

RAC: Retrieval-Augmented Clarification for Faithful Conversational Search

Ahmed Rayane Kebir^{1,2}[0009–0009–2512–832X], Vincent Guigue³[0000–0002–1450–5566], Lynda Said Lhadj²[0009–0005–3850–9229], and Laure Soulier¹[0000–0001–9827–7400]

¹ Sorbonne Université, CNRS, ISIR, F-75005 Paris, France

² Ecole nationale Supérieure d’Informatique (ESI), Algeria

³ AgroParisTech, UMR MIA-PS, Palaiseau, France

Abstract. Clarification questions help conversational search systems resolve ambiguous or underspecified user queries. While prior work has focused on fluency and alignment with user intent, especially through facet extraction, much less attention has been paid to grounding clarifications in the underlying corpus. Without such grounding, systems risk asking questions that cannot be answered from the available documents. We introduce RAC (**R**etrieval-**A**ugmented **C**larification), a framework for generating corpus-faithful clarification questions. After comparing several indexing strategies for retrieval, we fine-tune a large language model to make optimal use of research context and to encourage the generation of evidence-based question. We then apply contrastive preference optimization to favor questions supported by retrieved passages over ungrounded alternatives. Evaluated on four benchmarks, RAC demonstrate significant improvements over baselines. In addition to LLM-as-Judge assessments, we introduce novel metrics derived from NLI and data-to-text to assess how well questions are anchored in the context, and we demonstrate that our approach consistently enhances faithfulness.

Keywords: Conversational Search · Clarifying Questions · RAG.

1 Introduction

In open-domain information-seeking tasks, user queries are often short, ambiguous, or under-specified. Such characteristics make it difficult for traditional search systems to accurately capture user intent, as they typically provide only a ranked list of documents or passages without engaging in clarifying interactions [22]. Recent work has explored generating clarifying questions that are relevant, diverse, and human-plausible [8,28,30]. However, little attention has been given to whether these questions are grounded in the document corpus, even though unsupported clarifications may mislead users and harm retrieval effectiveness [10,18].

This is the author’s version of the work. It is posted here for your personal use. The definitive version is published in: *Proc of the 48th European Conference on Information Retrieval (ECIR ’26), 29 March–2 April, 2026, Delft, Netherlands*

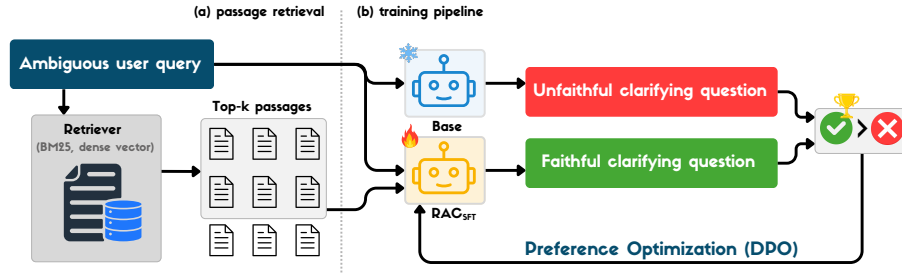


Fig. 1: Overview of RAC. Given an ambiguous user query, the system first retrieves the top- k passages ((a) passage retrieval). A mixture of the fine-tuned model and the base model is then used to generate unfaithful clarifying questions. Both faithful and unfaithful clarifying questions are subsequently leveraged for preference optimization via the DPO algorithm ((b) training pipeline). During inference, the trained model directly generates faithful clarifying questions.

Early approaches to clarifying question generation in conversational search largely relied on facet-based methods. These methods extracted candidate facets from the document collection to produce clarifying questions via templates or sequence-to-sequence models [1,2]. While this offered a basic form of corpus grounding, the reliance on coarse-grained facets proved reductive.

The advent of large language models (LLMs) enabled more fluent generation, with systems either conditioning on extracted facets to produce natural clarifications or directly deriving facets from queries before turning them into questions. Yet the task remains split into two stages—facet identification and question generation—creating bottlenecks in facet extraction and risks of hallucination when clarifications introduce content unsupported by the corpus [27,28].

In this work, we build on the retrieval-augmented generation (RAG) paradigm [12] to ground clarifications directly in the corpus, focusing on answers supported by the documents. Facet extraction is performed implicitly by supplying the top- k retrieved passages to the LLM, which then generates the clarifying question. The first contribution of this article is to propose a fine-tuning of conditional clarification generation, which greatly improves the quality of the questions. To further mitigate entity-level hallucinations, we also introduce a faithfulness reinforcement mechanism that steers the model to rely on the retrieved inputs rather than its internal knowledge, following the approach of [6].

Thus, we aim to address the following research questions:

RQ1. How can relevant passages be selected from the corpus, and how many should be used to optimally guide clarification?

RQ2. How does conditioning on these relevant passages affect the generation of clarifying questions?

RQ3. How can the faithfulness of clarifying question generation be improved when conditioned on relevant passages?

We introduce RAC, a framework for generating clarifying questions grounded in relevant retrieved passages, and train a large language model to prioritize faithful questions using preference tuning and contrastive learning, as illustrated in Fig. 1. We validate our approach on conversational search and open-domain Question Answering datasets through automatic metrics and LLM-as-Judge evaluations. Results show that RAC consistently enhances both the quality and faithfulness of clarifying questions, outperforming existing baselines.

2 Related Work

Query Clarification in Conversational Search. Asking clarifying questions enables users to actively participate in query disambiguation, with the goal of better capturing their information intent [1,2,7,11]. Prior work in this area has primarily focused on two tasks: predicting the need for clarification and generating clarifying questions. In this paper, we focus on the latter. Recent studies have increasingly explored large language model based approaches. For instance, Sekulić et al. [27] conditioned an LLM on specific facets; however, such facets are not always readily available and often require external extraction tools. Siro et al. [28] leveraged temperature control and facet information to generate diverse clarifications, while Wang et al. [30] introduced a zero-shot clarifying-question generator using fixed templates and query facets. More recently, Tang et al. [29] proposed a prompting strategy grounded in an ambiguity taxonomy to improve handling of ambiguous queries. Although these methods produce plausible and diverse clarifications, they remain prone to hallucination, frequently generating questions about aspects unsupported by the underlying corpus. Additionally, the reliance on explicit facets limits applicability when facets are difficult to extract or unavailable.

Retrieval Augmented Generation. Since the original article [12], several variants have been proposed, first for question answering [9] and later for clarification. Early studies primarily examined the role of the retriever in selecting corpus-grounded clarifications among candidate suggestions [18], whereas more recent work has shifted the focus toward generation [10], with particular attention to maximizing faithfulness during inference. In addition, [26] demonstrates that the RAG paradigm can be combined with knowledge bases to enhance disambiguation in domain-specific applications. However, these approaches rely on a zero-shot paradigm, whereas we demonstrate the benefit of fine-tuning the generator to better exploit the retrieved passages.

Preference Tuning. Reinforcement learning from human feedback was introduced to align LLMs with human preferences [19], but reward-model methods were costly and often unstable. More recent techniques such as direct preference optimization (DPO) [23] and extensions [32] have improved efficiency by learning directly from pairwise comparisons. Beyond general alignment, generating both faithful and unfaithful baseline sentences allows contrastive learning algorithms

to be effectively applied for improving text generation. Such approaches have demonstrated strong performance in tasks such as automatic summarization [4] and data-to-text generation [6]. To be useful, text variants must be generated carefully, and previous work has relied on mixture-of-logits decoding. Such techniques are directly relevant to conversational search, where clarifying questions must remain faithful to the corpus.

Faithfulness Evaluation. Faithfulness measures whether generated text remains consistent with its input. In summarization, state-of-the-art approaches employ entailment-based metrics that leverage NLI models to score the consistency of summaries with source documents (e.g., RoBERTa-based entailment [16]). These methods provide fine-grained judgments of factual alignment on a continuous scale. In data-to-text generation, metrics such as PAR-ENT [5] evaluate whether candidate outputs faithfully express entities and relations from structured inputs. By contrast, clarifying question generation has not been systematically assessed for faithfulness. Existing evaluations rely mainly on reference-based metrics (e.g., BLEU, METEOR) or indirect retrieval-based proxies [2,24], which do not directly measure factual consistency with the input context. In this work, we adapt entailment-based and data-grounding approaches from summarization and data-to-text to develop faithfulness evaluations tailored to clarifying question generation.

3 Methodology

Our RAC framework follows a two-stage training pipeline as illustrated in Fig. 2. In the first stage, a large language model p_{LM} is fine-tuned on existing clarification datasets along two axes to generate: (1) factual questions conditioned by user queries and retrieved passages p_{θ_0} and (2) less factual questions unconditioned by passages p_{uncond} . The two levels of question quality are then passed to a preference learning algorithm (contrastive) that encourages the model to rank faithful, evidence-grounded clarifications higher than unsupported or hallucinated alternatives.

We formulate the task of generating clarifying questions Cq as a retrieval-augmented generation task. The initial user query U_q enables the retrieval of a set of relevant passages $\mathcal{D} = \{d_1, \dots, d_N\}$, which will be used as context for the generation. We assume all queries to be ambiguous, focusing on clarifying

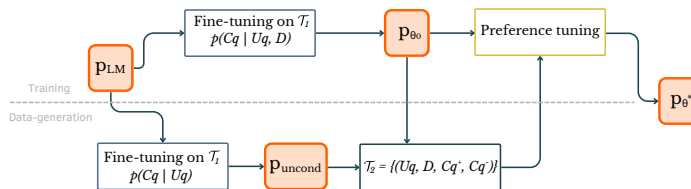


Fig. 2: Overview of our proposed training pipeline.

question generation rather than clarification need prediction [17]. Each passage may capture different semantic facets of the query, but we restrict to a single-turn setup, generating one clarifying question targeting the most useful facet.

3.1 Supervised Clarifying Question Generation

Retrieval-augmented generation (RAG) has shown that conditioning large language models on retrieved passages improves factual grounding and reduces reliance on parametric memory [9,12]. However, previous work has focused on generating direct zero-shot answers. Our contribution is to propose a fine-tuned model (twice) to better exploit the retrieved passages for the clarification task. To this end, we employ supervised fine-tuning (SFT) as the first stage of training: a large language model is trained to generate clarifying questions C_q conditioned on both the user query U_q and the corresponding retrieved passages \mathcal{D} (leading to p_{θ_0}). Given a dataset \mathcal{T}_1 of query–passage–ground-truth-question tuples $(U_q, \mathcal{D}, C_q^+)$, the model is optimized with the negative log-likelihood objective:

$$\mathcal{L}_{\text{SFT}}(\theta) = -\mathbb{E}_{\sim \mathcal{T}_1} \left[\sum_{t=1}^{|C_q^+|} \log p_{\theta}(C_{q,t} \mid U_q, \mathcal{D}, C_{q,<t}) \right] \quad (1)$$

Here, each token of the clarifying question $C_{q,t}$ is predicted sequentially, conditioned on the user query, the retrieved passages, and the previously generated tokens (denoted $C_{q,<t}$).

SFT establishes a strong baseline for clarification. By learning to ask questions supported by retrieved passages, the model reduces ambiguity in user intent and provides an evidence-aligned starting point for the subsequent preference based alignment stage. This further improves faithfulness and mitigates hallucinations.

3.2 Faithfulness Alignment

Although the p_{θ_0} model is already fine-tuned to generate clarifying questions that are much more relevant than the initial p_{LM} model, one of its main limitations is its tendency to hallucinate: it may generate details that are absent from the retrieved passages \mathcal{D} .

Preference tuning. To mitigate this, we introduce a second training stage focused on faithfulness. We augment the training data with pairs of faithful (C_q^+) and unfaithful (C_q^-) clarifying questions and apply a contrastive learning approach. In particular, we employ DPO [23] over a dataset $\mathcal{T}_2 = \{(U_q, \mathcal{D}, C_q^+, C_q^-)\}$, where the model is explicitly trained to prefer faithful clarifying questions over unfaithful ones. In DPO, the learning objective (Eq. 2) aligns a policy model p_{θ} with a preference signal, favoring C_q^+ over C_q^- , given the same input (U_q, \mathcal{D}) , as defined below:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{\sim \mathcal{T}_2} \left[\log \sigma \left(\beta \log \frac{p_{\theta}(C_q^+ \mid U_q, \mathcal{D})}{p_{\theta_0}(C_q^+ \mid U_q, \mathcal{D})} - \beta \log \frac{p_{\theta}(C_q^- \mid U_q, \mathcal{D})}{p_{\theta_0}(C_q^- \mid U_q, \mathcal{D})} \right) \right] \quad (2)$$

Unfaithful clarifying questions generation. Preference-based alignment requires faithful–unfaithful question pairs, but manual creation is costly and automatic detection remains difficult. We propose an unsupervised method that simulates unfaithful questions by injecting controlled noise during decoding. Our method adapts the noisy decoding strategy of Duong et al. [6] to the clarification setting. The approach relies on two complementary models:

Grounded model p_{θ_0} : obtained by fine-tuning a pretrained base model p_{LM} on half of the dataset \mathcal{T}_1 . Given a query and retrieved passages (U_q, \mathcal{D}) , it outputs generally faithful clarifying questions $C_q \sim p_{\theta_0}(\cdot \mid U_q, \mathcal{D})$, though minor inaccuracies remain.

Ungrounded model p_{uncond} : obtained by fine-tuning the same base model but conditioned only on the user query U_q , i.e., $C_q \sim p_{\text{uncond}}(\cdot \mid U_q)$. It produces fluent and relevant clarifying questions, yet these are not guaranteed to be grounded in the retrieved passages \mathcal{D} .

While p_{uncond} produces overly unconstrained questions and p_{θ_0} tends to remain faithful, their combination yields plausible but unfaithful clarifying questions (the balance is critical, as highlighted by Duong et al. [6]). Specifically, we decode token-by-token from a mixture distribution (Eq. 3), using stochastic decoding (temperature and top- k sampling) to promote diversity and encourage hallucinated tokens.

$$C_{q,t} \sim (1 - \alpha_t) p_{\theta_0}(\cdot \mid C_{q,<t}, U_q, \mathcal{D}) + \alpha_t p_{\text{uncond}}(\cdot \mid C_{q,<t}, U_q), \quad (3)$$

where $\alpha_t \sim \text{Bernoulli}(\alpha)$ controls the injection of ungrounded content. The noise parameter $\alpha \in [0, 1]$ determines the faithfulness–fluency trade-off: $\alpha = 0$ recovers clarifying questions from p_{θ_0} , whereas $\alpha = 1$ generates ungrounded ones from p_{uncond} . The resulting questions remain fluent but contain ungrounded spans, yielding both intrinsic errors (contradictions with retrieved passages) and extrinsic hallucinations (additions not inferable from \mathcal{D}). These are used as unfaithful clarifying questions C_q^- in the augmented dataset $\mathcal{T}_2 = (U_q, \mathcal{D}, C_q^+, C_q^-)$, enabling preference optimization for faithfulness alignment.

3.3 Joint Training Objective

Supervised fine-tuning and preference optimization address complementary objectives: supervised fine-tuning operates at the token level, teaching the model to produce clarifying questions, while preference optimization encourages it to prefer faithful outputs over unfaithful ones. To leverage both, we propose a combined training objective: $\mathcal{L}_{\text{RAC}}(\theta) = \gamma \cdot \mathcal{L}_{\text{DPO}}(\theta) + (1 - \gamma) \cdot \mathcal{L}_{\text{SFT}}(\theta)$.

4 Experimental Setup

4.1 Datasets and Evaluation

Datasets. We evaluate RAC on four datasets across conversational search and open-retrieval QA. For search, we use Qulac (derived from TREC Web Track

2009–2012) [2] and the filtered version of ClariQ proposed by Sekulic et al. [27], which maps clarifying questions to facets. For QA, we use PaQa (AmbigNQ with GPT-3 clarifications) [8] and CambigNQ (AmbigNQ queries augmented with human-validated clarifications) [11].

Adapting Datasets for Retrieval-Augmented Clarification. Existing clarification datasets (Qulac, ClariQ) lack passage-level grounding, as their relevance labels are assigned at the document level and not explicitly tied to the clarifying question. To bridge this gap, we derive passage-level supervision through a three-stage pipeline: (i) *Passage Indexing*: we segment Clueweb09-12⁴ into 250-token passages, following TREC CAsT [20], and index them with Pyserini [15]; (ii) *Query Rewriting*: for each ambiguous query–clarification pair (U_q, C_q) , we generate a facet-specific reformulation U_q^r by incorporating C_q using an LLM, yielding sharper retrieval intents than U_q alone; (iii) *Pseudo-Relevance Retrieval*: we employ BM25 [25] over the passage index to retrieve the top- k passages \mathcal{D} for U_q^r , treating them as pseudo-relevant evidence. This produces training tuples (U_q, \mathcal{D}, C_q) that support retrieval-conditioned clarification generation.

Metrics. We employ both reference-based and reference-free metrics to evaluate the quality of generated clarifying questions. Reference-based metrics measure similarity to gold questions, while reference-free metrics assess faithfulness to the input query and associated passages. In addition, we use GPT-4 to assess faithfulness, serving as a model-based proxy for human judgment.

Reference-based evaluation. We report BLEU [21], ROUGE-L [14], METEOR [3], and BERTScore [33]. BLEU and ROUGE-L capture n-gram and longest common subsequence overlap, respectively, while METEOR accounts for synonym and stem matches. BERTScore computes semantic similarity via contextualized token embeddings, providing a finer-grained assessment of meaning preservation. These metrics are consistent with prior work in clarification question generation and facilitate direct comparison.

Faithfulness evaluation. We evaluate faithfulness using PARENT Recall (PAR) [5] and AlignScore (AL) [13]. PAR, originally proposed for data-to-text generation, computes n-gram recall against both the input and the reference, serving as a proxy for input-groundedness. To apply it to unstructured passages, we adapt the metric by extracting named entities, multi-word noun phrases, and subject–verb–object triples with SpaCy⁵, allowing content-level overlap measurement without reliance on structured data. AL is an entailment-based metric built on RoBERTa [16] and trained on multiple NLI datasets. Because clarifying questions are often interrogative and not well-suited for direct entailment evaluation, we convert them into declarative statements by removing question templates, retaining only content-bearing tokens, and filtering query overlaps.

⁴ <https://lemurproject.org/clueweb09/>

⁵ <https://spacy.io/>

This yields hypotheses compatible with AL’s premise–hypothesis structure while preserving the semantic content of the questions.

4.2 Baselines

We evaluate RAC against several baselines. First, we include (AT-CoT), the ambiguity taxonomy chain-of-thought prompting baseline of Tang et al. [29], which applies few-shot prompting conditioned only on the query. Following Sekulic et al. [27], we use the widely adopted (Q-Cond) fine-tuned model, which generates clarifications from the query alone. To assess the impact of supervision, we compare RAC to a (QP-Zero_{shot}) variant conditioned on both query and passages in a zero-shot setting. Finally, on ClariQ, where facet annotations are available, we also report results for the template-based (TB) and facet-based (QF-Cond) baselines of Sekulic et al. [27]. For LLM-based methods, we use the same underlying model to ensure a fair comparison.

4.3 Implementation Details and Hyperparameters

We build on the pre-trained LLaMA3.1-8B-base checkpoint from the Hugging-Face Hub, using the Transformers and TRL libraries [31]. For supervised fine-tuning (SFT), we train for 2 epochs with a learning rate of 1×10^{-5} , batch size 32, and a linear learning rate schedule. For direct preference optimization (DPO), we use 2 epochs with a learning rate of 2×10^{-6} , batch size 32, and $\beta = 0.1$. In our joint loss, we set $\gamma = 0.5$, based on ablation results. Zero-shot baselines rely on the Instruct variant of the base model. All experiments are run on NVIDIA A100 GPUs (80GB). Source code is available at: <https://github.com/RayaneA7/RAC-Retrieval-augmented-clarification>.

5 Results

5.1 Main Results

The main evaluation results are reported in Table 1. We find that *RAC* significantly outperforms the baselines across all metrics and datasets, confirming that passage conditioning substantially improves clarifying question generation, answering **RQ2**.

Moreover, results show that reference-based measures fail to capture the gains from preference tuning, consistent with prior findings [4,6,23]. In contrast, reference-free evaluation –reported only for models conditioned with passages– reveals that *RAC*_{DPO} achieves better performance over *RAC*_{SFT}. This demonstrates that preference-based optimization enhances corpus faithfulness beyond supervised fine-tuning, directly addressing **RQ3**.

The fact that QP-Zero performs significantly worse than Q-cond highlights the importance of learning the form of a clarification question, independently of its content.

Table 1: Evaluation scores of RAC variants against different baselines, with $\beta = 0.1$ and for mixture $\alpha = 0.7$. Bold values indicate best performance, and \dagger indicates a statistically significant improvement (Welch’s t-test, $p < 0.001$).

Dataset	Model	ROUGE-L \uparrow	BLEU \uparrow	METEOR \uparrow	BERTScore (F1) \uparrow	ALScore \uparrow	Par-R \uparrow
Conversational Search Datasets							
Qulac	AT-CoT	17.97	2.77	20.81	84.72	–	–
	Q-Cond	29.44	10.51	25.92	88.24	–	–
	QP-Zero _{shot}	27.39	5.68	33.33	87.20	–	–
	RAC_{SFT}(ours)	33.14\dagger	12.59\dagger	31.30\dagger	89.34\dagger	79.14	42.53
	+ RAC_{DPO}(ours)	32.42\dagger	11.52\dagger	31.48\dagger	88.92\dagger	81.73	44.83
ClariQ	AT-CoT	18.63	3.49	21.19	84.74	–	–
	Q-Cond	28.68	11.19	25.47	88.16	–	–
	TB	35.50	0.28	24.26	87.65	–	–
	QF-Cond	33.70	2.20	37.56	89.08	–	–
	QP-Zero _{shot}	26.03	4.99	31.81	86.59	–	–
	RAC_{SFT}(ours)	36.25\dagger	14.88\dagger	34.01\dagger	89.52\dagger	51.32	53.15
	+ RAC_{DPO}(ours)	35.52\dagger	14.86\dagger	33.84\dagger	89.39\dagger	52.41	55.77
Question Answering Datasets							
PaQa	AT-CoT	23.59	7.07	22.93	85.97	–	–
	Q-Cond	42.46	16.62	41.58	90.12	–	–
	QP-Zero _{shot}	33.79	10.42	35.84	88.66	–	–
	RAC_{SFT}(ours)	46.83\dagger	20.17\dagger	47.97\dagger	90.85\dagger	43.36	27.62
	+ RAC_{DPO}(ours)	45.26\dagger	18.32\dagger	46.40\dagger	90.41\dagger	45.75	28.54
CAmbigNQ	AT-CoT	10.33	2.10	8.53	84.02	–	–
	Q-Cond	28.41	8.90	33.06	87.17	–	–
	QP-Zero _{shot}	18.20	4.27	19.48	85.15	–	–
	RAC_{SFT}(ours)	36.66\dagger	14.81\dagger	43.37\dagger	88.93\dagger	47.62	87.99
	+ RAC_{DPO}(ours)	35.47\dagger	14.40\dagger	41.99\dagger	88.89\dagger	49.95	88.05

These findings highlight both the benefit of passage conditioning and the added value of preference-based optimization. We further validate these results through qualitative analysis and LLM-based judgments in subsequent experiments.

5.2 LLM-based Evaluation

To further address **RQ2**, we evaluate the faithfulness of our approach using GPT-4 as a evaluator, comparing RAC_{DPO} against RAC_{SFT} . Results are shown in Table 2. Across all datasets, RAC_{DPO} achieves higher win rates compared to RAC_{SFT} , in some cases by more than a factor of two, whereas a large fraction of outputs are judged as ties. These results suggest that supervised fine-tuning already provides a strong baseline, preference optimization yield further gains on harder cases, reinforcing **RQ3** by enhancing faithfulness beyond supervised training.

5.3 Impact of the Number of Input Passages

We next examine the impact of the number and quality of retrieved passages on RAC. Because RAC relies on retrieval to expose potential ambiguities, both the quantity and relevance of the input passages directly affect its ability to generate effective clarifications.

Table 2: GPT-4 preference results comparing RAC_{DPO} and RAC_{SFT} . Results with * are statistically significantly different based on the one-sided McNemar’s test with $p < 0.05$.

Dataset	$RAC_{DPO}\%$	Tie%	$RAC_{SFT}\%$
Qulac	28.88*	50.56	20.56
ClariQ	28.36*	48.79	22.85
PaQa	36.24	30.36	33.40
CAMbigNQ	16.27*	72.23	11.50

As shown in Fig. 3, performance improves as the number of passages increases, but the effect saturates after approximately four passages, suggesting that the most salient query-related ambiguities are typically captured within the top-ranked results. Passage quality is equally important: using random passages results in performance close to the “*Q-Cond*” baseline, whereas BM25 and dense retrievers achieve substantially higher scores. BM25’s advantage is likely due to a domain mismatch, since the dense retriever is trained on MS MARCO, whose passage structure and content differ from the chunked ClueWeb passages used in our setting. These findings indicate that RAC benefits from informative retrieval signals and can extract relevant facets from high-quality passages rather than relying on arbitrary content, thereby addressing **RQ1**.

5.4 Impact of the Quality of Noisy Generated Elements

We study the effect of our noisy generation method by comparing p_{uncond} , finetuned with only the query as input, with p_{LM} , the initial language model. We then measure their impact on preference tuning (positive samples always being generated by p_{θ_0}). Table 3 shows that p_{uncond} provides more effective negative samples than p_{LM} . Unlike the approach of Duong et al. [6], which relies on generic

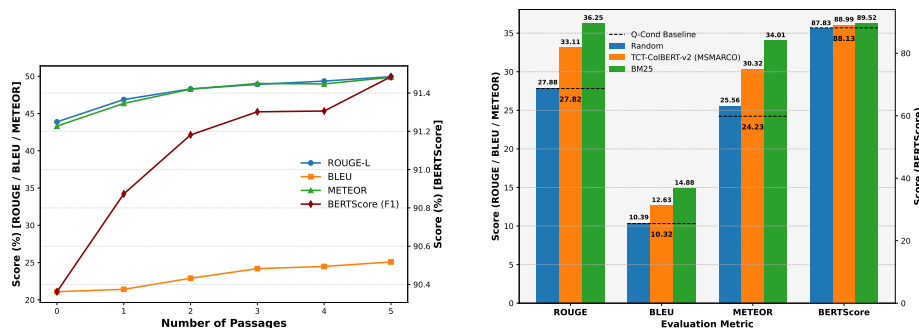


Fig. 3: NLG metrics on ClariQ: impact of varying the number of passages (left) and comparison of retrieval strategies (BM25, TCT, random) using the top 5 retrieved passages (right).

Table 3: ClariQ validation results using different negative generation methods.

Method	ROUGE-L	BLEU	METEOR	BERTScore (F1)	ALScore \uparrow	Par-R \uparrow
$RAC_{DPO}, C_q^- \sim p_{LM}$	33.84	12.93	30.79	89.25	50.81	50.73
$RAC_{DPO}, C_q^- \sim p_{uncond}$	35.52	14.86	33.84	89.39	52.41	55.77

noise injection, p_{uncond} generates clarifications that are structurally well-formed but factually misaligned. This contrast makes them harder negatives and better training signals for preference optimization. By comparison, samples from p_{LM} often fail to resemble clarifications at all, limiting their usefulness. These results highlight the importance of tailoring noise generation to the clarification format rather than reusing generic base-model outputs.

5.5 Qualitative Analysis

Noisy clarifying questions. We qualitatively assess the effect of mixture between the conditioned & unconditioned models p_{θ_0} and p_{uncond} , controlled by α (Eq.3). At $\alpha = 0$, outputs come solely from p_{θ_0} ; at $\alpha = 1$, from p_{uncond} . Table 4 shows an example from ClariQ, where noise increases with α . Irrelevant spans (highlighted in red) illustrate how higher α degrades faithfulness. For preference learning, selecting intermediate α values yields negative examples that are challenging yet informative, avoiding both trivial and overly noisy supervision.

Generated Clarifying questions. We compare clarifications from RAC_{SFT} and RAC_{DPO} on ClariQ validation data (Table 5). Faithful content is highlighted in yellow, hallucinations in red. RAC_{SFT} exhibits occasional grounding failures, such as introducing unsupported entities (e.g., hallucinating “season 17”) or omitting relevant evidence (e.g., overlooking “season 16”). In contrast, RAC_{DPO} produces questions more tightly grounded in passages. This aligns

Table 4: Noisy generation at different α interpolation weights of p_{uncond} and p_{θ_0}

α	Noisy Generation
0.0	are you interested in the causes of angular cheilitis
0.1	do you want to read an article on that, do you want to be diagnosed
0.2	do you have any recommendations for treating this condition
0.3	do you want to see a picture of angular cheilitis
0.4	would you like to know about symptoms, risk factors, or a way to treat or manage your condition
0.5	would recommend a topical or internal medicine to cure your specific symptom for you based off its symptoms or how do I treat your specific angular cheilitis?
0.6	do any methods have been tested?
0.7	What works to treat this
0.8	Use warm soapy water, soft cotton balls, an old razor blade. If you do not have any of the aforementioned items, you can use a regular nail clipper.
0.9	How many times a day should I apply the treatment?

Table 5: Qualitative comparison of *RAC* against baseline models.

Input	<i>RAC_{SFT}</i>	<i>RAC_{DPO}</i>
Query: When does the new family guy season come out? Passages: ['... Ed O Neill Liam Payne Louis Tomlinson and Neil deGrasse Tyson During this season the guys head to South Korea after', ' Family Guy (season 15) Family Guy season 15 Family Guy fifteenth season premiered on Fox in the United States on September 25 2016 and ended on May 21 2017 The season contained 20 episodes The series follows the dysfunctional Griffin family consisting of father Peter mother Lois daughter Meg son Chris baby ...re Appel and Callaghan Guest stars for the season include Kyle Chandler Stephen Curry Flea Rob Gronkowski']	Are you looking for the release date of Family Guy season 15, 16 or 17 ?	Are you looking for the release dates of Family Guy seasons 14, 15 or 16 ?

with quantitative gains reflecting in more relevant and grounded clarifying questions, demonstrating that DPO improves reliability in ambiguous cases where SFT fails to capture the core ambiguity.

6 Conclusion

In this work, we introduced clarification question generation as a retrieval conditioned generation task, where questions are generated based on both the user query and retrieved passages. This formulation ensures that clarifications are grounded in information the system can realistically access. Our RAC framework combines retrieval context with preference tuning to improve both the relevance and corpus-faithfulness of generated questions. Experiments on four benchmarks demonstrate that both *RAC_{SFT}* and *RAC_{DPO}* significantly outperform existing baselines, Q-Cond and QP-Zero_{shot}, across all reference-based metrics (ROUGE-L, BLEU, METEOR, and BERTScore). We further employ LLM-as-Judge evaluations and novel metrics derived from NLI and data-to-text to quantify the gains in faithfulness to retrieved content of *RAC_{DPO}* over *RAC_{SFT}*, which is critical for conversational search, where the objective is to disambiguate and answer user queries based on retrieved evidence rather than knowledge internal to the language model. As future work, we plan to extend this task to multi-turn clarification and evaluate its impact on downstream retrieval performance.

Acknowledgments. The authors acknowledge the ANR – FRANCE (French National Research Agency) for its financial support of the GUIDANCE project n°ANR-23-IAS1-0003 as well as the Chaire Multi-Modal/LLM ANR Cluster IA ANR-23-IACL-0007. This work was granted access to the HPC resources of IDRIS under the allocation AD011016470 made by GENCI.

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article.

References

1. Aliannejadi, M., Kiseleva, J., Chuklin, A., Dalton, J., Burtsev, M.: Building and evaluating open-domain dialogue corpora with clarifying questions. arXiv preprint arXiv:2109.05794 (2021)
2. Aliannejadi, M., Zamani, H., Crestani, F., Croft, W.B.: Asking clarifying questions in open-domain information-seeking conversations. In: Proceedings of the 42nd international acm sigir conference on research and development in information retrieval. pp. 475–484 (2019)
3. Banerjee, S., Lavie, A.: Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In: Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization. pp. 65–72 (2005)
4. Choi, J., Chae, K., Song, J., Jo, Y., Kim, T.: Model-based preference optimization in abstractive summarization without human feedback. In: Al-Onaizan, Y., Bansal, M., Chen, Y.N. (eds.) Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing. pp. 18837–18851. Association for Computational Linguistics, Miami, Florida, USA (Nov 2024). <https://doi.org/10.18653/v1/2024.emnlp-main.1048>, <https://aclanthology.org/2024.emnlp-main.1048/>
5. Dhingra, B., Faruqui, M., Parikh, A., Chang, M.W., Das, D., Cohen, W.W.: Handling divergent reference texts when evaluating table-to-text generation. arXiv preprint arXiv:1906.01081 (2019)
6. Duong, S., Bronnec, F.L., Allauzen, A., Guigue, V., Lumbreras, A., Soulier, L., Gallinari, P.: SCOPE: A self-supervised framework for improving faithfulness in conditional text generation. In: The Thirteenth International Conference on Learning Representations (2025), <https://openreview.net/forum?id=dTkqaCKLPp>
7. Erbacher, P., Nie, J.Y., Preux, P., Soulier, L.: Augmenting ad-hoc ir dataset for interactive conversational search. Transactions on Machine Learning Research (2024), <https://openreview.net/forum?id=z8d7nT1HWw>
8. Erbacher, P., Nie, J.Y., Preux, P., Soulier, L.: Paqa: toward proactive open-retrieval question answering. arXiv preprint arXiv:2402.16608 (2024)
9. Izacard, G., Grave, E.: Leveraging passage retrieval with generative models for open domain question answering. In: Merlo, P., Tiedemann, J., Tsarfaty, R. (eds.) Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume. pp. 874–880. Association for Computational Linguistics, Online (Apr 2021). <https://doi.org/10.18653/v1/2021.eacl-main.74>, <https://aclanthology.org/2021.eacl-main.74/>
10. Krasakis, A.M., Yates, A., Kanoulas, E.: Corpus-informed retrieval augmented generation of clarifying questions. arXiv preprint arXiv:2409.18575 (2024)
11. Lee, D., Kim, S., Lee, M., Lee, H., Park, J., Lee, S.W., Jung, K.: Asking clarification questions to handle ambiguity in open-domain qa. arXiv preprint arXiv:2305.13808 (2023)
12. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.t., Rocktäschel, T., et al.: Retrieval-augmented generation for knowledge-intensive nlp tasks. Advances in neural information processing systems **33**, 9459–9474 (2020)
13. Li, Y.Z.Y.Y.R., Hu, Z.: Alignscore: Evaluating factual consistency with a unified alignment function
14. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. In: Text summarization branches out. pp. 74–81 (2004)

15. Lin, J., Ma, X., Lin, S.C., Yang, J.H., Pradeep, R., Nogueira, R.: Pyserini: A python toolkit for reproducible information retrieval research with sparse and dense representations. In: Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval. pp. 2356–2362 (2021)
16. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 (2019)
17. Lu, L., Meng, C., Ravenda, F., Aliannejadi, M., Crestani, F.: Zero-shot and efficient clarification need prediction in conversational search. In: European Conference on Information Retrieval. pp. 389–404. Springer (2025)
18. Mass, Y., Cohen, D., Yehudai, A., Konopnicki, D.: Conversational search with mixed-initiative - asking good clarification questions backed-up by passage retrieval. In: Feng, S., Wan, H., Yuan, C., Yu, H. (eds.) Proceedings of the Second DialDoc Workshop on Document-grounded Dialogue and Conversational Question Answering. pp. 65–71. Association for Computational Linguistics, Dublin, Ireland (May 2022). <https://doi.org/10.18653/v1/2022.dialdoc-1.7>, <https://aclanthology.org/2022.dialdoc-1.7/>
19. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al.: Training language models to follow instructions with human feedback. *Advances in neural information processing systems* **35**, 27730–27744 (2022)
20. Owoicho, P., Dalton, J., Aliannejadi, M., Azzopardi, L., Trippas, J.R., Vakulenko, S.: Trec cast 2022: Going beyond user ask and system retrieve with initiative and response generation. In: TREC (2022)
21. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting of the Association for Computational Linguistics. pp. 311–318 (2002)
22. Radlinski, F., Craswell, N.: A theoretical framework for conversational search. In: Proceedings of the 2017 conference on conference human information interaction and retrieval. pp. 117–126 (2017)
23. Rafailov, R., Sharma, A., Mitchell, E., Manning, C.D., Ermon, S., Finn, C.: Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems* **36**, 53728–53741 (2023)
24. Rao, S., Daumé, H.: Learning to ask good questions: Ranking clarification questions using neural expected value of perfect information. In: Annual Meeting of the Association for Computational Linguistics (2018)
25. Robertson, S., Zaragoza, H., et al.: The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval* **3**(4), 333–389 (2009)
26. Sahay, R., Tekumalla, L.S., Aggarwal, P., Jain, A., Saladi, A.: Ask: Aspects and retrieval based hybrid clarification in task oriented dialogue systems. In: Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 6: Industry Track). pp. 881–895 (2025)
27. Sekulić, I., Aliannejadi, M., Crestani, F.: Towards facet-driven generation of clarifying questions for conversational search. In: Proceedings of the 2021 ACM SIGIR international conference on theory of information retrieval. pp. 167–175 (2021)
28. Siro, C., Yuan, Y., Aliannejadi, M., de Rijke, M.: AGENT-CQ: Automatic Generation and Evaluation of Clarifying Questions for Conversational Search with LLMs. arXiv preprint arXiv:2410.19692 (2024), <http://arxiv.org/abs/2410.19692>

29. Tang, A., Soulier, L., Guigue, V.: Clarifying ambiguities: on the role of ambiguity types in prompting methods for clarification generation. In: Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 20–30 (2025)
30. Wang, Z., Tu, Y., Rosset, C., Craswell, N., Wu, M., Ai, Q.: Zero-shot clarifying question generation for conversational search. In: Proceedings of the ACM web conference 2023. pp. 3288–3298 (2023)
31. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., et al.: Transformers: State-of-the-art natural language processing. In: Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations. pp. 38–45 (2020)
32. Yang, X., Tan, Z., Li, H.: Ipo: Iterative preference optimization for text-to-video generation. arXiv preprint arXiv:2502.02088 (2025)
33. Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y.: Bertscore: Evaluating text generation with bert. arXiv preprint arXiv:1904.09675 (2019)