

# How You Ask Matters! Adaptive RAG Robustness to Query Variations

Anonymous ACL submission

## Abstract

Adaptive Retrieval-Augmented Generation (RAG) promises accuracy and efficiency by dynamically triggering retrieval only when needed and is widely used in practice. However, real-world queries vary in surface form even with the same intent, and their impact on Adaptive RAG remains under-explored. We introduce the first large-scale benchmark of diverse yet semantically identical query variations, combining human-written and model-generated rewrites. Our benchmark enables systematic evaluation of Adaptive RAG robustness across answer, computational cost, and retrieval decisions. We discover a critical robustness gap, where small surface-level changes in queries dramatically alter retrieval behavior and accuracy. Although larger models show better performance, robustness does not improve accordingly. These findings reveal that Adaptive RAG methods are highly vulnerable to query variations that preserve identical semantics, exposing a critical robustness challenge.

## 1 Introduction

Retrieval-Augmented Generation (RAG) reduces hallucinations in large language models (LLMs) by grounding outputs in retrieved evidence, thereby improving factual accuracy in downstream tasks (Khandelwal et al., 2019; Lewis et al., 2020). While effective in reducing factual errors, conventional RAG often retrieves irrelevant or noisy documents, which can mislead reasoning and degrade answer quality (Li et al., 2023; Ding et al., 2025).

To address these limitations, Adaptive RAG conditions retrieval decisions on the input query and its evolving context (Jeong et al., 2024). Unlike conventional RAG with a fixed retrieval pattern, it adaptively decides *when* and *what* to retrieve as generation proceeds. For simple factoid questions, a small number of retrievals is sufficient, whereas multi-hop questions typically benefit from iterative subqueries that aggregate evidence across hops (Ho

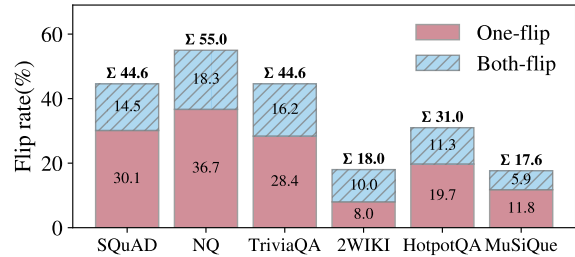


Figure 1: **Retrieval-decision flip rates** of the Qwen-32B model, measuring how often the model changes its retrieval judgment under meaning-preserving human query rewrites. Higher rates signal instability. Flips are categorized as *One-flip* (only one rewrite flips) and *Both-flip* (both rewrites flip).

et al., 2020; Trivedi et al., 2022, 2023). As a result, prior adaptive approaches can improve both accuracy and efficiency by allocating retrieval and LLM calls based on query difficulty and intermediate states (Su et al., 2024; Jiang et al., 2023; Ding et al., 2025). Due to their strong performance, these frameworks are increasingly deployed in real-world applications (Singh et al., 2025).

In real-world settings, users express the same intent in diverse forms, including paraphrases, stylistic variations, typos or other noise (Fu et al., 2025; Zhang et al., 2025; Cao et al., 2025). To examine the model’s decision stability under query variations, we measure how often its retrieval decision (retrieve or not) flips between the original query and a human rewrite. As shown in Figure 1, human rewrites trigger decision flips in up to 55% of cases, revealing a severe lack of robustness in the system’s retrieval decisions under human queries. Despite these risks, robustness to real-world query variation has not been systematically evaluated, which is particularly concerning in high-stakes settings where factual correctness and user trust are critical (Sharma et al., 2024; Oche et al., 2025).

This paper presents the first empirical study of Adaptive RAG robustness under realistic query

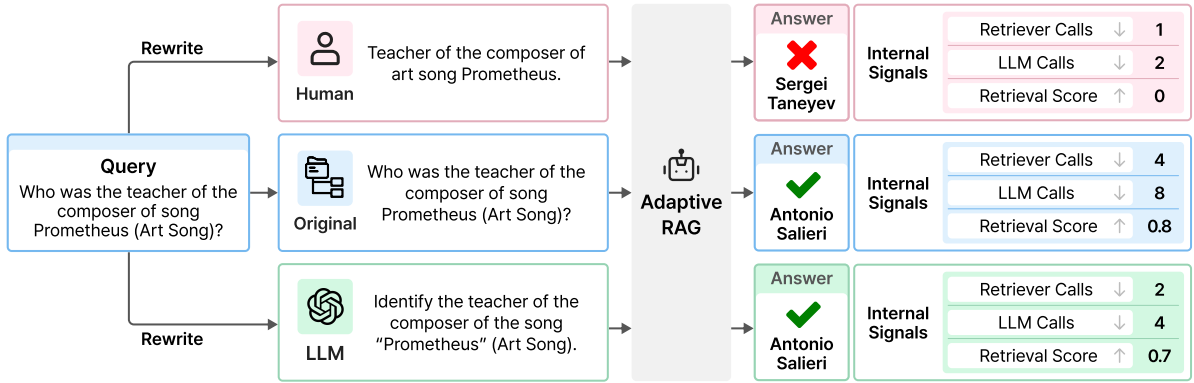


Figure 2: **Example of Adaptive RAG responses** under human, original, and LLM-generated query rewrites, highlighting differences in answer correctness, computation overhead, and retrieval score.

069 variations. To facilitate this, we introduce a bench- 102  
 070 mark spanning controlled linguistic paraphrases 103  
 071 and naturally written human queries. We evalu- 104  
 072 ate on six widely used question-answering (QA) 105  
 073 datasets covering both single-hop and multi-hop 106  
 074 questions and analyze how surface-level differ- 107  
 075 ences affect both intermediate retrieval trajec- 108  
 076 tories and final answers (Figure 2). For each 109  
 077 question, we include six model-generated vari- 110  
 078 ants and two human rewrites, yielding 27K QA 111  
 079 pairs spanning diverse levels of complexity. 112  
 080 Across these query variations, we comprehen- 113  
 081 sively assess model behavior along three dimen- 114  
 082 sions: answer quality, computa-  
 083 tional cost, and retrieval decisions.

This work presents the following key findings:

- 084 • Spelling errors cause the largest accuracy and 102  
 085 robustness drop at the cost of semantic distor- 103  
 086 tion, while human rewrites cause the second- 104  
 087 largest drop but preserve the original meaning. 105
- 088 • Improved performance does not imply greater 106  
 089 robustness; even strong baselines (e.g., the 107  
 090 generation-based method, QwQ-32B) show 108  
 091 substantially greater inconsistency in answer 109  
 092 semantics and computational cost. 110
- 093 • Even minor changes to the initial query can 111  
 094 alter the intermediate reasoning path and the 112  
 095 subqueries generation, leading to retrieval fail- 113  
 096 ure and ultimately lower-quality final answers. 114

## 097 2 Related Works

### 098 2.1 Adaptive RAG

099 Adaptive RAG implements dynamic, model-driven 132  
 100 decisions about when to retrieve and what to re- 133  
 101 trieve. We broadly categorize prior methods into 134  
 135

two main streams: logit-based and generation- 102  
 based methods. Logit-based methods monitor 103  
 model confidence and trigger retrieval when the 104  
 model appears uncertain, for example, using token- 105  
 level probabilities as in Jiang et al. (2023) or uncer- 106  
 tainty estimates based on token probabilities and 107  
 attention weights as in Su et al. (2024). Generation- 108  
 based methods, on the other hand, let the model 109  
 decide on the fly when it needs retrieval by generat- 110  
 ing search queries during the generation process (Li 111  
 et al., 2025). Our work focuses on approaches that 112  
 use a single model with a retriever and do not rely 113  
 on external classifiers, such as Jeong et al. (2024). 114

### 115 2.2 Query Robustness

A growing line of work studies robustness of RAG 116  
 systems. Several studies show that linguistic per- 117  
 turbations such as query-entry errors, grammatical 118  
 errors, typos, and stylistic variations affect QA per- 119  
 formance (Cao et al., 2025; Perçin et al., 2025; 120  
 Zeng et al., 2025). These findings demonstrate that 121  
 even minor modifications to a query can substan- 122  
 tially alter retrieved documents and downstream 123  
 answers (Zhang et al., 2025; Fu et al., 2025). How- 124  
 ever, prior work has two key limitations: (1) it fo- 125  
 cuses primarily on conventional RAG, leaving the 126  
 additional complexity introduced by adaptive RAG 127  
 largely unexplored; (2) it does not examine the ef- 128  
 fects of real human query reformulations across 129  
 different RAG frameworks. 130

## 131 3 Benchmark Construction

### 132 3.1 Dataset

We evaluate single-hop and multi-hop QA, using 133  
 Natural Questions (NQ) (Kwiatkowski et al., 2019), 134  
 TriviaQA (Joshi et al., 2017), and SQuAD v1.1 135

Original Query: Who is the child of the director of film Mukhyamantri (1996 Film)?		
Category	Variation	Sample Rewrites
Style	Formality (Form.) ↓	Yo, quick one — who’s the kid of the director behind the 1996 flick Mukhyamantri?
	Readability (Read.) ↓	By whom is the progeny of the director of Mukhyamantri (1996 Film) constituted?
Type	Declarative (Decl.)	I want to know who is the child of the director of the film Mukhyamantri (1996 Film).
	Imperative (Impr.)	Name the child of the director of film Mukhyamantri (1996 Film).
Error	Spelling (Spell.) ↓	Who is ght childer ar the director of filme Mukhyamantri (1996 Film) ?
	Grammar (Gram.) ↓	Who are the childs of the director from film Mukhyamantri 1996?

Table 1: **Examples of query rewrites** grouped by category, illustrating stylistic, structural, and error-based variations. ↓ denotes reduction in corresponding writing variation.

(Rajpurkar et al., 2016) for single-hop, and 2Wiki-MultiHopQA (2WIKI) (Ho et al., 2020), MuSiQue (Trivedi et al., 2022), and HotpotQA (Yang et al., 2018) for multi-hop evaluation.

### 3.2 Query Variation Generation

Building upon prior works (Moskvoretskii et al., 2025; Jeong et al., 2024; Trivedi et al., 2023), we evaluate the systems on the same 500 samples per dataset. To generate query variations, we employ GPT-5-mini and produce six distinct linguistic variants grouped into three categories: sentence style, sentence type, and error, as shown in Table 1. While earlier studies examined style and error (Zhang et al., 2025; Perçin et al., 2025; Cao et al., 2025), we extend this paradigm by incorporating sentence type variations, covering 24k queries.

**Sentence Style.** Human queries exhibit diverse linguistic styles, often characterised by less formality, simplified language, or omission of polite expressions. Following previous work (Cao et al., 2025), we introduce a set of natural queries with reduced formality and readability. Each variation is generated by LLM with prompts detailed in the Appendix B. For validation, we follow the process outlined in Cao et al. (2025) (see Appendix C.1). We formally define these variations as follows:

- **Formality:** We generate informal versions using simpler vocabulary, contractions, and more casual expressions.
- **Readability:** We reduce the readability of the queries, making them more difficult to understand. Although this may sacrifice clarity, it more accurately reflects how users often communicate in less structured contexts.

**Sentence Type.** Sentence type variations capture the grammatical function of the sentence, reflecting how the query is structured. Most current datasets cover interrogative questions, particularly in the form of wh-questions (e.g., who, what, when) (Rogers et al., 2023; Kwong and Yorke-Smith, 2012). Therefore, we challenge the model by generating two additional sentence types: declarative and imperative. We instruct GPT-5-mini to generate queries that fit each of these sentence types, and the used prompts are given in Appendix B.

- **Declarative sentence:** These queries make statements or assertions. They are often used to express knowledge, desire, or intention.
- **Imperative sentence:** These queries are commands or requests. They can suggest actions or instructions, typically used when the user wants the system to perform a task.

**Error.** We focus on the most frequent error types observed in human-generated queries: grammatical and spelling errors (Hagiwara and Mita, 2019). These kinds of errors commonly occur in user input due to hurried typing, typographical mistakes, and limited mastery of grammatical rules.

- **Spelling errors:** These errors occur when a user mistakenly types a word with incorrect letters or letter combinations, resulting in a non-standard word form. We use the python library *nlpaug* tool<sup>1</sup> to inject the spelling errors to the original query (Zhang et al., 2025).
- **Grammatical errors:** We implement these errors by prompting GPT-5-mini to inject grammatical mistakes. The prompts are provided in

<sup>1</sup>[pypi.org/project/nlpaug/0.0.5](https://pypi.org/project/nlpaug/0.0.5)

the Appendix B. After that, we use Language-Tool to check whether the query variations include at least one grammar error.

### 3.3 Human Rewrite Annotation

Human queries exhibit substantial variation in form and do not follow a single, consistent pattern; instead, they are a combination of multiple categories (Penha et al., 2021; Li et al., 2024). However, existing robustness evaluations for RAG systems largely overlook this dimension, as they do not incorporate human written queries (Cao et al., 2025; Perçin et al., 2025; Zhang et al., 2025). In this work, we collected two human rewrites for 600 queries (100 per dataset across six QA benchmarks), yielding 1,200 rewrites, as well as additional query reformulations from 12 annotators. Annotators were given the original query and its corresponding answer and were instructed to generate alternative formulations as if posing the question themselves. Detailed guidelines are provided in the Appendix C.2.

## 4 Experimental Setup

### 4.1 Methods

We consider two Adaptive RAG methods that opened new paradigms: logit-based control (LB) and generation-based control (GB). For the LB approach, we use DRAGIN (Su et al., 2024), which determines when and what to retrieve using token-level signals (e.g., probabilities and attention). For the GB approach, we use Search-o1 (Li et al., 2025), where the model uses its internal reasoning to construct a retrieval plan, including whether to initiate search, how to formulate queries, and how to derive the final answer. Unlike the threshold-style control of LB, Search-o1 performs retrieval control end-to-end without external triggers such as hand-tuned thresholds or auxiliary classifiers.

### 4.2 Metrics

**Answer Robustness.** We evaluate answer correctness using InAccuracy, following Baek et al. (2023), to check whether the gold answer is included in the model answer. In addition, we measure the semantic consistency of generated answers across variations using similarity and diversity.

1. **InAccuracy:** For each instance  $i$  and perturbation type  $v$ , let  $\hat{a}_{i,v}$  be the generated answer and  $g_i$  the gold answer. InAccuracy is reported as the dataset-level mean for each

perturbation type  $v$ .

$$\text{InAcc}_v = \frac{1}{N} \sum_{i=1}^N \mathbf{1}[g_i \subseteq \hat{a}_{i,v}]$$

2. **Similarity:** Let  $\phi(\cdot)$  denote an embedding function (SBERT). We denote the answers to the original and reformulated queries as  $\hat{a}_{i,\text{orig}}$  and  $\hat{a}_{i,v}$ , respectively, and define  $e_{i,v} = \phi(\hat{a}_{i,v})$ . We compute the cosine similarity between  $\hat{a}_{i,v}$  and  $\hat{a}_{i,\text{orig}}$ .

$$\text{Sim}_v = \frac{1}{N} \sum_{i=1}^N \cos(e_{i,v}, e_{i,\text{orig}})$$

3. **Diversity (Div):** Let  $\mathcal{V}$  be the set of variants (including the original) for instance  $i$ , and let  $K = |\mathcal{V}|$ . We quantify answer spread as:

$$\text{Div}_i = \frac{2}{K(K-1)} \sum_{v < w} (1 - \cos(e_{i,v}, e_{i,w}))$$

$$\text{Div} = \frac{1}{N} \sum_{i=1}^N \text{Div}_i$$

**Computation Robustness.** We measure call robustness for the LLM and the retriever using the average number of calls made for a single question. To assess robustness, we introduce an evaluation protocol that measures (i) the change in computational cost relative to the original queries and (ii) the consistency of cost across query variations. Let  $c_{i,v}$  denote the call count for instance  $i$  under variant  $v$ , with  $v = 0$  denoting the original query.

1. **Relative Error (RE):** For each instance  $i$  and variant  $v$ , we compute the deviation from the original-query call count:

$$\text{RE}_{i,v} = \frac{|c_{i,v} - c_{i,0}|}{\max(1, c_{i,0})}$$

We then report the dataset-level mean for each perturbation type  $v$ :

$$\text{RE}_v = \frac{1}{N} \sum_{i=1}^N \text{RE}_{i,v}$$

2. **Coefficient-of-Variation Robustness (CVR):** Let  $\mathcal{V}$  be the set of variants (including the original) and  $r_i = (c_{i,v})_{v \in \mathcal{V}}$ . We compute instance-wise variability across variants:

$$\text{CV}_i = \frac{\text{std}(r_i)}{\text{mean}(r_i)}, \quad \text{CVR}_i = \frac{1}{1 + \text{CV}_i}$$

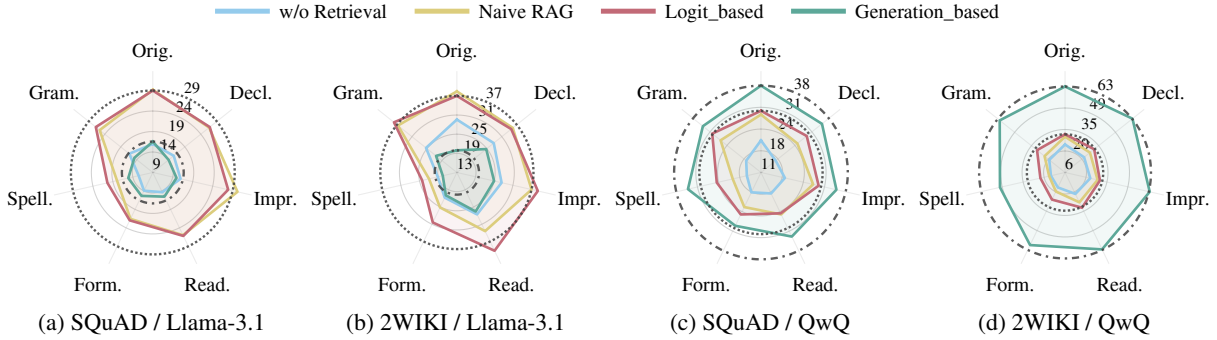


Figure 3: **InAccuracy results** across six perturbation types for a given dataset / model pair.

The final score is averaged over instances:

$$\text{CVR} = \frac{1}{N} \sum_{i=1}^N \text{CVR}_i$$

We report RE for each perturbation type and CVR aggregated across all perturbations, for both **retriever calls** and **LLM calls**.

**Retrieval Decision Robustness.** We evaluate whether the decision to retrieve is correct. The variable *act* indicates whether retrieval is triggered (*act* = 1 if triggered; *act* = 0 otherwise). The variable *need* is defined based on a no-retrieval baseline: if the model answers the question correctly without retrieval, then *need* = 0; otherwise, *need* = 1.

1. **Overconfidence:** The rate of cases where the model fails to retrieve despite retrieval being needed (i.e., *need* = 1, *act* = 0), following [Moskvoretskii et al. \(2025\)](#).
2. **Underconfidence:** The rate of cases where the model retrieves despite retrieval not being needed (i.e., *need* = 0, *act* = 1), following [Moskvoretskii et al. \(2025\)](#).

### 4.3 Implementation Details

We use two LLMs, QwQ-32B ([Team, 2025](#)) and Llama 3.1-8B-Instruct ([Grattafiori et al., 2024](#)). For retrieval, we employ the sparse retriever BM25, applied to a document pool constructed from Wikipedia documents ([Moskvoretskii et al., 2025](#)). See Appendix A for further implementation details.

## 5 Results

We present results on two representative datasets – SQuAD (single-hop) and 2WIKI (multi-hop) – and report others in the Appendix.

Dataset	Mtd.	Similarity ↑						Div. ↓
		Form.	Read.	Decl.	Impr.	Spell.	Gram.	
<b>Llama-3.1</b>								
SQuAD	LB	0.551	0.539	0.593	<b>0.634</b>	0.552	0.599	0.469
	GB	0.426	0.408	<b>0.433</b>	0.430	0.411	0.414	0.585
2WIKI	LB	0.560	0.628	0.649	<b>0.655</b>	0.537	0.627	0.414
	GB	0.423	<b>0.458</b>	0.300	0.403	0.313	0.385	0.667
<b>QwQ</b>								
SQuAD	LB	0.573	0.579	0.615	<b>0.631</b>	0.544	0.584	0.440
	GB	0.552	0.558	0.579	<b>0.589</b>	0.552	0.579	0.459
2WIKI	LB	0.500	0.545	0.551	<b>0.565</b>	0.462	0.519	0.511
	GB	0.520	0.557	0.554	<b>0.559</b>	0.499	0.548	0.465

Table 2: **Answer similarity and diversity** across query variations. Darker shades indicate more desirable outcomes: higher **similarity** and lower **diversity**.

### 5.1 Answer Robustness

**Model performance drastically drops on query variations.** We first analyze how query variations affect the answer accuracy. Figure 3 shows spelling error perturbation yields the largest accuracy drop across dataset/model pairs, indicating strong sensitivity to surface-form noise. Generation-based methods benefit from larger models and harder datasets. For example, on 2WIKI with QwQ, accuracy still degrades with changes in query phrasing (63.4% for the original query vs. 50.6% for spelling errors), with drops of up to 12.8% across variations. Results across all 6 datasets are shown in the Appendix D.1.

**Generation-based method shows inconsistent outputs.** Table 2 confirms the same trend: imperative rewrites maintain the highest answer similarity to the original, whereas spelling error consistently produces the lowest similarity. Moreover, generation-based methods exhibit higher diversity (lower is better), implying less stable outputs across perturbation.

Target	Mtd.	SQuAD							2WIKI							
		RE ↓						CVR ↑	RE ↓						CVR ↑	
		Form.	Read.	Decl.	Impr.	Spell.	Gram.		Form.	Read.	Decl.	Impr.	Spell.	Gram.		
Llama-3.1	RTR	LB	0.472	0.592	0.416	<b>0.338</b>	0.431	0.439	0.734	0.853	0.859	<b>0.455</b>	0.846	0.928	0.941	0.710
		GB	0.092	0.094	<b>0.086</b>	0.090	0.100	0.092	0.937	0.468	0.416	<b>0.402</b>	0.542	0.496	0.435	0.742
	LLM	LB	0.441	0.592	0.396	<b>0.313</b>	0.416	0.411	0.767	1.001	1.001	<b>0.427</b>	1.007	1.093	1.130	0.742
		GB	0.056	0.059	<b>0.048</b>	0.055	0.061	0.051	0.963	0.447	0.404	<b>0.401</b>	0.520	0.477	0.418	0.745
QwQ	RTR	LB	0.494	0.450	<b>0.405</b>	0.415	0.495	0.441	0.740	0.560	0.502	0.446	<b>0.443</b>	0.565	0.511	0.719
		GB	0.360	0.322	<b>0.298</b>	0.315	0.361	0.330	0.762	0.788	0.808	<b>0.780</b>	0.825	0.787	0.884	0.689
	LLM	LB	0.477	0.466	<b>0.389</b>	0.399	0.491	0.424	0.760	0.547	0.481	<b>0.426</b>	0.430	0.565	0.493	0.747
		GB	0.265	0.230	<b>0.225</b>	0.238	0.250	0.243	0.858	1.078	1.142	1.092	1.203	<b>1.049</b>	1.219	0.716

Table 3: **Computation robustness** across different model-generated query variations. Darker shading indicates more desirable outcomes: lower RE and higher CVR. RTR = Retriever. Bold denotes maximum across the variations.

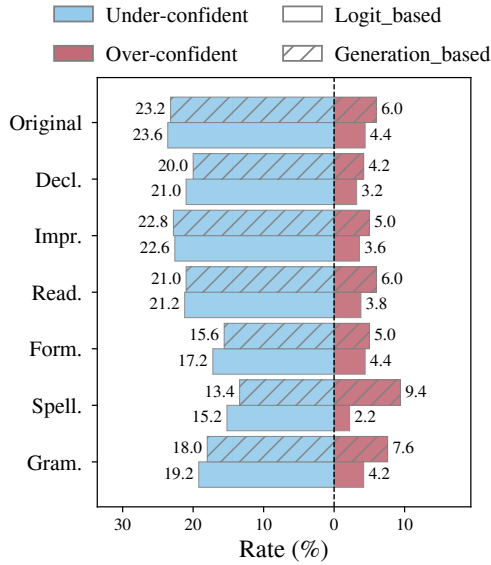


Figure 4: **Under- and over-confident rates** on model query variations.

## 5.2 Computation Robustness

**Computation robustness destabilizes on harder questions.** We evaluate computational robustness using RE and CVR to assess whether the model triggers the retriever and the LLM consistently across query variations. Table 3 shows that retrieval-call counts are most sensitive on 2WIKI, exhibiting larger RE across most perturbations. Spelling and grammatical errors produce the highest RE, indicating that surface-level noise most strongly shifts the retrieval-call budget.

LLM-call robustness follows a similar trend. On SQuAD, RE remains moderate and CVR stays high for Llama (up to 0.963 with the generation-based method), suggesting stable call counts under rewrites. In contrast, on 2WIKI, RE increases and CVR decreases, with the largest degradation under spelling errors.

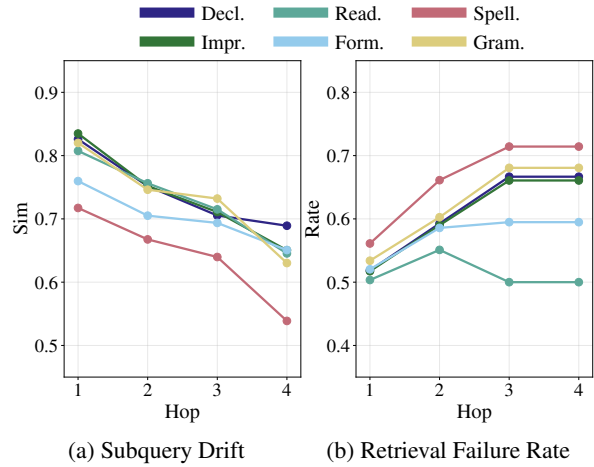


Figure 5: **Subquery drift effect on retrieval quality.** 2WIKI Dragin results (a) Subquery similarity to the original decreases across turns. (b) Retrieval failure rate (no gold document in top- $k$ ) across turns.

**Computation robustness does not imply higher accuracy.** In Figure 3, generation-based method achieves strong accuracy on harder datasets when paired with larger models. However, this trend reverses for computational robustness in Table 3. For 2WIKI with QwQ model, computational robustness degrades substantially. This contrast indicates that high accuracy does not necessarily correspond to strong computational robustness.

## 5.3 Retrieval Decision Robustness

**Query errors invoke over-confidence.** We investigate the retrieval decision correctness using under- and over-confidence metrics. Under-confident instances incur additional costs, while over-confidence harms accuracy by skipping necessary retrievals. Figure 4 shows spelling-error appear on high over-confident cases, reaching up to 9.4%. This tendency persists particularly in larger models.

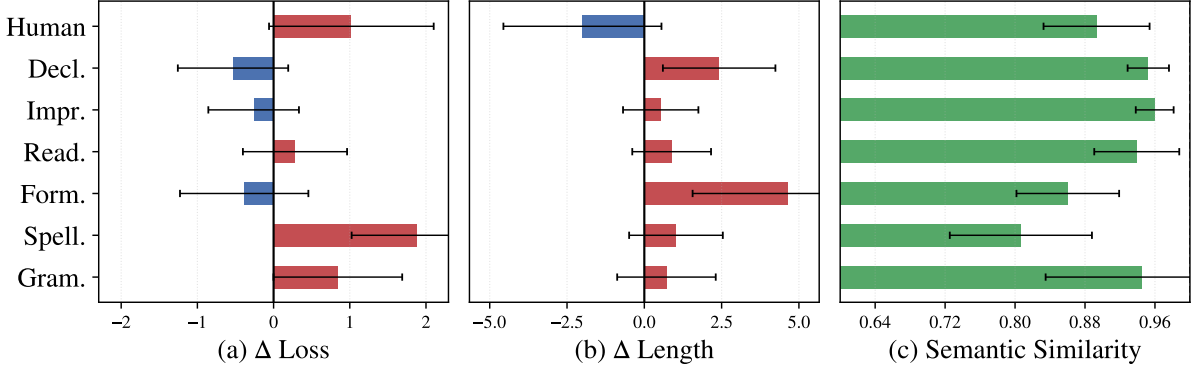


Figure 6: **Query analysis.** From left to right: (a) change in LM loss relative to the original query ( $\Delta$ Loss; nats/token), (b) change in query length ( $\Delta$ Length; words), and (c) semantic similarity to the original query (cosine similarity).

Mtd.	Similarity $\uparrow$							Div. $\downarrow$
	Human	Form.	Read.	Decl.	Impr.	Spell.	Gram.	
LB	0.518	0.528	0.570	0.573	0.612	0.522	0.572	0.477
GB	0.529	0.542	0.534	0.568	0.598	0.520	0.558	0.460

Table 4: **Answer Robustness with human queries.** Similarity to the original query and diversity on 2WIKI (QwQ-32B) including human rewrites.

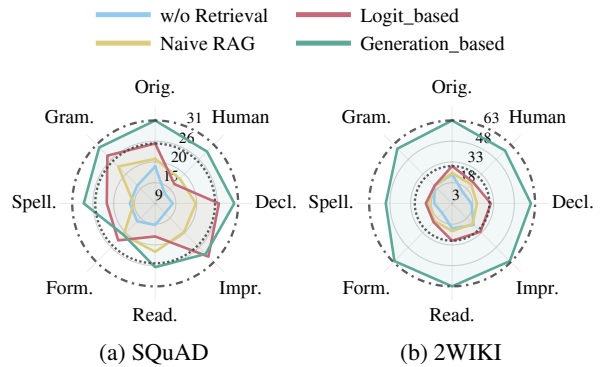


Figure 7: **InAccuracy results** with human rewrites

**Drifted subqueries miss gold documents.** We investigate retrieval-decision robustness in a multi-hop Adaptive RAG setting. At hop  $t$ , the model generates a subquery  $q_t$  and retrieves a top- $k$  document set  $D_t$ . We analyze how query variations perturb the multi-hop retrieval process by tracking (i) semantic drift in intermediate subqueries and (ii) hop-level retrieval failures rates. To quantify drift, we compute the mean subquery similarity (MSS), defined as the average semantic similarity between each variant subquery and its corresponding original subquery across hops. Retrieval failures are quantified using the retrieval failure rate (RFR), defined as the proportion of hops in which the top- $k$  retrieved documents contain no gold document, i.e.,  $\text{recall}@k(D_t^{\text{var}}, D_t^{\text{Gold}}) = 0$ .

As subqueries deviate from the original trajectory (lower MSS), the retrieved top- $k$  sets increasingly fail to cover gold evidence, resulting in a higher RFR. Figure 5 illustrates this trend: subquery similarity decreases while RFR increases across hops, indicating that subquery drift directly undermines retrieval quality.

## 6 Human Query Analysis

**Human rewrites are semantically faithful but harder for the model.** Figures 6(a) and (c) show

that human rewrites preserve semantic intent of the original queries, yet they induce substantially higher token-level loss. They are also shorter, frequently compressing or omitting surface cues as in Figure 6(b).

In Figure 7, human-authored rewrites lead to accuracy degradation compared to original query, in some cases comparable to that caused by spelling-errors. Table 4 further shows reduced answer stability relative to the original queries. This degradation is similar to that observed under simple surface-level perturbations, indicating that robustness failures are not limited to synthetic noise. To identify the source of errors, we trace a causal chain from properties of human rewrites to subquery drift, retrieval failures, and ultimately reduced end-to-end answer accuracy.

**Human query length and uncertainty are associated with subquery drift.** Shorter and more uncertain human rewrites are reflected in subquery drift. For Search-o1 on 2WIKI, query loss is negatively correlated with mean subquery similarity (MSS; Spearman’s  $\rho=-0.15$ ), while query length is

t	Original Subqueries	Human Subqueries	Sim	Hit(O)	Hit(H)
	Which film has the director born earlier, The Adventures Of Priscilla, Queen Of The Desert or Harvest: 3,000 Years?	Who was born first director of The Adventures Of Priscilla, Queen Of The Desert or Harvest: 3,000 Yearst			
1	director of the adventures of priscilla, queen of the desert	director of the adventures of priscilla queen of the desert	0.99	✓	✓
2	director of harvest 3000 years	director of harvest 3000 years	1.00	✓	✓
3	“harvest: 3,000 years” director birth year	harvest 3000 years movie director	0.85	✓	×
4	giuliano carnimeo birth year	harvest 3000 years korean movie director	0.13	✓	×

Table 5: **Per-hop comparison of original and human subqueries.** Sim is the semantic similarity between the two subqueries at hop  $t$ . Hit(O) indicates whether the retrieved top- $k$  set contains any ground-truth supporting evidence for the original(O) and human(H) query.

425 positively correlated with MSS ( $\rho=0.29$ ). These re-  
426 sults suggest that higher-loss, more compressed hu-  
427 man phrasing destabilizes intermediate subqueries,  
428 leading to lower MSS.

429 **Semantic drift leads to retrieval trajectory col-  
430 lapse and evidence failures.** Once intermediate  
431 subqueries drift, retrieval becomes unstable. For  
432 Search-01 on human rewrites, MSS is strongly  
433 negatively correlated with trajectory collapse ( $\rho=-$   
434 0.75), indicating that retrieved documents increas-  
435 ingly diverge from those for the original queries.  
436 MSS is also negatively correlated with RFR ( $\rho=-$   
437 0.35), suggesting that gold documents are more  
438 often missing from the retrieved set. These correla-  
439 tions imply that human rewrites leads not only to  
440 different retrieval trajectories but also to failures in  
441 retrieving necessary evidence.

442 Figure 8 visualizes this pattern: the gold-  
443 document inclusion rate decreases as multi-hop  
444 reasoning progresses, ultimately dropping well be-  
445 low the rates for the original and model queries.  
446 Table 5 illustrates the underlying mechanism with  
447 an example from QwQ outputs. In the initial hops,  
448 the model-generated subqueries remain highly sim-  
449 ilar to the original query (similarity  $\approx 0.99-1.00$ ),  
450 but they begin to drift starting at hop 3. This drift  
451 leads to retrieval failures, as reflected by the hit  
452 metric.

453 **Retrieval errors propagate.** Failures in interme-  
454 diate retrieval steps translate into end-to-end accu-  
455 racy degradation. Retrieval failure is negatively cor-  
456 related with the performance gap  $Y$  (Spearman’s  
457  $\rho = -0.20$  for Search-01). When the system fails  
458 to retrieve gold documents, it lacks the evidence  
459 required to answer the query, and downstream gen-  
460 eration rarely recovers.

461 Overall, our findings indicate that human  
462 rewrites preserve semantic meaning yet increase

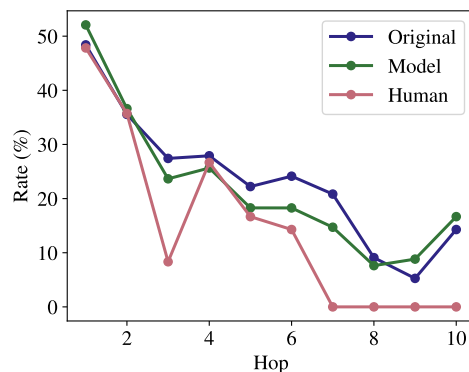


Figure 8: **Gold document inclusion rate across hops.** For each hop, we measure how often the top- $k$  retrieved results contain at least one gold document.

463 the difficulty for Adaptive RAG systems. Higher  
464 query loss and more compressed phrasing amplify  
465 subquery drift, which in turn destabilizes the re-  
466 trieval trajectory and increases retrieval failure, ul-  
467 timately degrading final answer accuracy.

## 468 7 Conclusion

469 In this work, we introduced a large-scale bench-  
470 mark of query variations to evaluate the robustness  
471 of Adaptive RAG systems. We comprehensively as-  
472 sessed three dimensions—answer robustness, com-  
473 putation robustness, and retrieval-decision robust-  
474 ness—covering both final answers and intermediate  
475 reasoning steps. We showed that both logit-based  
476 and generation-based adaptive methods were sen-  
477 sitive to query perturbations, which degraded ac-  
478 curacy, efficiency, and retrieval-decision stability.  
479 Human rewrites were particularly detrimental: they  
480 induced drift in intermediate subqueries and trig-  
481 gered retrieval-trajectory collapse, leading to sub-  
482 stantial end-to-end performance drops. Finally, we  
483 found that human rewrites differed systematically  
484 from model-generated variations, highlighting a  
485 realism gap in current robustness evaluations.

## 486 Limitations

487 Although our study offers a systematic view of  
488 Adaptive RAG responses under realistic query vari-  
489 ation, several aspects remain outside our current  
490 scope. First, our human rewrites are designed to be  
491 meaning-preserving paraphrases of existing ques-  
492 tions. This setting captures a common and well-  
493 controlled form of user rephrasing, but it does not  
494 include broader interactive behaviors such as multi-  
495 turn follow-up questions, conversational reformu-  
496 lations, or intentionally adversarial inputs. Second,  
497 we evaluate a specific set of retrievers and Adaptive  
498 RAG controllers. Other retrieval architectures or  
499 decision mechanisms may exhibit different robust-  
500 ness–efficiency profiles, and extending coverage  
501 to a wider range of systems is a straightforward  
502 next step. Finally, although our human rewrites  
503 are carefully collected and anonymized, they rep-  
504 resent only a small portion of the full diversity of  
505 real-world query phrasing. Scaling human-written  
506 evaluations across languages, domains, and user  
507 populations would further improve coverage.

## 508 References

509 Jinheon Baek, Soyeong Jeong, Minki Kang, Jong C.  
510 Park, and Sung Ju Hwang. 2023. [Knowledge-  
511 augmented language model verification](#). *Preprint*,  
512 arXiv:2310.12836.

513 Tianyu Cao, Neel Bhandari, Akhila Yerukola, Akari  
514 Asai, and Maarten Sap. 2025. [Out of style:  
515 Rag’s fragility to linguistic variation](#). *Preprint*,  
516 arXiv:2504.08231.

517 Hanxing Ding, Liang Pang, Zihao Wei, Huawei Shen,  
518 and Xueqi Cheng. 2025. [Rowen: Adaptive retrieval-  
519 augmented generation for hallucination mitigation in  
520 llms](#). *Preprint*, arXiv:2402.10612.

521 Daocheng Fu, Jianbiao Mei, Licheng Wen, Xuemeng  
522 Yang, Cheng Yang, Rong Wu, Tao Hu, Siqi Li, Yufan  
523 Shen, Xinyu Cai, Pinlong Cai, Botian Shi, Yong Liu,  
524 and Yu Qiao. 2025. [Re-searcher: Robust agentic  
525 search with goal-oriented planning and self-reflection](#).  
526 *Preprint*, arXiv:2509.26048.

527 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,  
528 Abhinav Pandey, Abhishek Kadian, Ahmad Al-  
529 Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,  
530 Alex Vaughan, and 1 others. 2024. The llama 3 herd  
531 of models. *arXiv preprint arXiv:2407.21783*.

532 Masato Hagiwara and Masato Mita. 2019. [Github  
533 typo corpus: A large-scale multilingual dataset  
534 of misspellings and grammatical errors](#). *CoRR*,  
535 abs/1911.12893.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, 536  
and Akiko Aizawa. 2020. [Constructing a multi- 537  
hop QA dataset for comprehensive evaluation of 538  
reasoning steps](#). In *Proceedings of the 28th Inter- 539  
national Conference on Computational Linguistics*, 540  
pages 6609–6625, Barcelona, Spain (Online). Inter- 541  
national Committee on Computational Linguistics. 542

Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju 543  
Hwang, and Jong C. Park. 2024. [Adaptive-rag: 544  
Learning to adapt retrieval-augmented large lan- 545  
guage models through question complexity](#). *Preprint*, 546  
arXiv:2403.14403. 547

Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, 548  
Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie 549  
Callan, and Graham Neubig. 2023. [Active retrieval 550  
augmented generation](#). *Preprint*, arXiv:2305.06983. 551

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke 552  
Zettlemoyer. 2017. [TriviaQA: A large scale distantly 553  
supervised challenge dataset for reading comprehen- 554  
sion](#). In *Proceedings of the 55th Annual Meeting of 555  
the Association for Computational Linguistics (Vol- 556  
ume 1: Long Papers)*, pages 1601–1611, Vancouver, 557  
Canada. Association for Computational Linguistics. 558

Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke 559  
Zettlemoyer, and Mike Lewis. 2019. [Generalization 560  
through memorization: Nearest neighbor language 561  
models](#). *CoRR*, abs/1911.00172. 562

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red- 563  
field, Michael Collins, Ankur Parikh, Chris Alberti, 564  
Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken- 565  
ton Lee, Kristina Toutanova, Llion Jones, Matthew 566  
Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob 567  
Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natu- 568  
ral questions: A benchmark for question answering 569  
research](#). *Transactions of the Association for Compu- 570  
tational Linguistics*, 7:453–466. 571

Helen Kwong and Neil Yorke-Smith. 2012. Detection of 572  
imperative and declarative question–answer pairs in 573  
email conversations. *AI Communications*, 25(4):271– 574  
283. 575

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio 576  
Petroni, Vladimir Karpukhin, Naman Goyal, Hein- 577  
rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock- 578  
täschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge- 579  
intensive NLP tasks](#). *CoRR*, abs/2005.11401. 580  
581

Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin 582  
Wang, Michal Lukasik, Andreas Veit, Felix Yu, and 583  
Sanjiv Kumar. 2023. [Large language models with 584  
controllable working memory](#). In *Findings of the As- 585  
sociation for Computational Linguistics: ACL 2023*, 586  
pages 1774–1793, Toronto, Canada. Association for 587  
Computational Linguistics. 588

Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, 589  
Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng 590  
Dou. 2025. [Search-o1: Agentic search-enhanced 591  
large reasoning models](#). *Preprint*, arXiv:2501.05366. 592

593	Zhicong Li, Jiahao Wang, Zhishu Jiang, Hangyu Mao, Zhongxia Chen, Jiazhen Du, Yuanxing Zhang, Fuzheng Zhang, Di Zhang, and Yong Liu. 2024. <a href="#">Dmqr-rag: Diverse multi-query rewriting for rag</a> . <i>Preprint</i> , arXiv:2411.13154.	648
594		649
595		650
596		651
597		652
598	Viktor Moskvoretskii, Maria Lysyuk, Mikhail Salnikov, Nikolay Ivanov, Sergey Pletenev, Daria Galimzianova, Nikita Krayko, Vasily Konovalov, Irina Nikishina, and Alexander Panchenko. 2025. <a href="#">Adaptive retrieval without self-knowledge? bringing uncertainty back home</a> . <i>Preprint</i> , arXiv:2501.12835.	653
599		654
600		655
601		656
602		657
603		
604	Agada Joseph Oche, Ademola Glory Folashade, Tirthankar Ghosal, and Arpan Biswas. 2025. <a href="#">A systematic review of key retrieval-augmented generation (rag) systems: Progress, gaps, and future directions</a> . <i>Preprint</i> , arXiv:2507.18910.	658
605		659
606		660
607		661
608		662
609	Gustavo Penha, Arthur Câmara, and Claudia Hauff. 2021. <a href="#">Evaluating the robustness of retrieval pipelines with query variation generators</a> . <i>CoRR</i> , abs/2111.13057.	663
610		664
611		665
612		666
613	Sezen Perçin, Xin Su, Qutub Sha Syed, Phillip Howard, Aleksei Kuvshinov, Leo Schwinn, and Kay-Ulrich Scholl. 2025. <a href="#">Investigating the robustness of retrieval-augmented generation at the query level</a> . In <i>Proceedings of the Fourth Workshop on Generation, Evaluation and Metrics (GEM<sup>2</sup>)</i> , pages 439–457, Vienna, Austria and virtual meeting. Association for Computational Linguistics.	667
614		
615		
616		
617		
618		
619		
620		
621	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <a href="#">Squad: 100,000+ questions for machine comprehension of text</a> . <i>Preprint</i> , arXiv:1606.05250.	
622		
623		
624		
625	Anna Rogers, Matt Gardner, and Isabelle Augenstein. 2023. <a href="#">Qa dataset explosion: A taxonomy of nlp resources for question answering and reading comprehension</a> . <i>ACM Computing Surveys</i> , 55(10):1–45.	
626		
627		
628		
629	Sanat Sharma, David Seunghyun Yoon, Franck Dernoncourt, Dewang Sultania, Karishma Bagga, Mengjiao Zhang, Trung Bui, and Varun Kotte. 2024. <a href="#">Retrieval augmented generation for domain-specific question answering</a> . <i>Preprint</i> , arXiv:2404.14760.	
630		
631		
632		
633		
634	Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talaie Khoei. 2025. <a href="#">Agentic retrieval-augmented generation: A survey on agentic rag</a> . <i>Preprint</i> , arXiv:2501.09136.	
635		
636		
637		
638	Weihang Su, Yichen Tang, Qingyao Ai, Zhijing Wu, and Yiqun Liu. 2024. <a href="#">Dragin: Dynamic retrieval augmented generation based on the information needs of large language models</a> . <i>Preprint</i> , arXiv:2403.10081.	
639		
640		
641		
642	Qwen Team. 2025. <a href="#">Qwq-32b: Embracing the power of reinforcement learning</a> .	
643		
644	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. <a href="#">Musique: Multi-hop questions via single-hop question composition</a> . <i>Preprint</i> , arXiv:2108.00573.	
645		
646		
647		

## A Implementation Details

We conduct our experiments with two LLMs: Llama 3.1-8B-Instruct and QwQ-32B. Llama-3.1-8B-Instruct is a lightweight model from the Llama-3.1 family, optimized for instruction-following tasks. QwQ-32B is part of the Qwen model family and is specifically optimized for reasoning, supporting long-context and multi-step inference with strong performance on mathematical, programming, and logic-intensive tasks. We generate responses using nucleus sampling with temperature  $\mathcal{T} = 0.7$  and  $\text{top-}p = 0.95$ , and report all results from a single run. We run inference locally using  $2 \times$  NVIDIA A100 GPUs with Intel(R) Xeon(R) Silver 4210R CPU @ 2.40GHz.

For the DRAGIN method, following prior work (Su et al., 2024), we tune the retrieval threshold by evaluating several high-performing candidate values on the Natural Questions (NQ) dataset and select 0.6. We then fix this setting and apply the same configuration to all remaining datasets to ensure a consistent evaluation of computational cost across benchmarks. For the Search-o1 method, maximum search limit hyperparameter is set 15 for the multi-hop QA and 5 for the single-hop QA. Rest of the hyperparameter follows the (Li et al., 2025) setting. All of the embedding similarity is calculated by the sentence-transformers/all-mpnet-base-v2, which is suitable for clustering and semantic search task.

## B Prompts

We report the prompts for model-generated query variations.

**Formality Variation.** We adopt the lower-formality rewriting prompt proposed in prior work (Cao et al., 2025) to rewrite each query in a more casual style while preserving its original meaning.

### Formality Variation Prompt Template

You are an AI assistant skilled at transforming formal queries into casual, everyday language. Rewrite the following query so that it sounds very informal. Experiment with different colloquial openings, varied sentence constructions, and a mix of slang, idioms, and casual expressions throughout the sentence. Avoid using the same phrase repeatedly (e.g., 'hey, so like') and ensure the meaning remains unchanged.

**Readability Variation.** We utilize the prompt

from Cao et al. (2025) shown in Table 6 to reduce the readability of the original query.

**Sentence Type Variation.** This prompt rewrites the given query into a specified sentence type (declarative or imperative) while preserving meaning. We implement this prompt for both declarative and imperative variations.

### Sentence Type Variation Prompt Template

Task: Rewrite Query as {*sentence\_type*} Sentence  
Take the given query and rewrite it as a {*sentence\_type*} sentence, changing only the sentence type while keeping the original meaning intact.  
Do not add extra details or change the intent.

**Grammar Error Variation.** This prompt injects at least one grammar error to the original query.

### Grammar Error Variation Prompt Template

You are an AI assistant skilled at generating realistic grammatical errors in queries.  
Rewrite the given query by introducing grammatical errors.  
Constraints:  
- Preserve the original meaning  
- Introduce at least three grammatical errors per query  
- Make it realistic and natural by varying the types of errors

## C Benchmark Construction

### C.1 Model-generated Queries

Style and sentence-type variations are generated using GPT-5-mini<sup>2</sup> with the prompts provided in Appendix B. For the error categories, we use external tools to inject or verify errors. Spelling errors are injected using the SpellingAug function in nlpaug. For the grammar error category, we use LanguageTool<sup>3</sup> to verify that each query contains at least one grammatical error. If a query is classified as having no grammar errors, we regenerate it using GPT-5-mini to inject a grammatical error.

For all query variations, the co-authors manually inspected sampled queries. The annotation process was conducted by two co-authors, both fluent in English. Annotators were instructed to focus on the quality of the rewritten queries, whether the intended perturbation was correctly applied while preserving the original meaning. We randomly sampled 100 examples across query variations, and the two annotators marked 96% and 98% as pass-

<sup>2</sup>[gpt-5-mini-2025-08-07](https://openai.com/gpt-5-mini)

<sup>3</sup><https://languagetool.org>

## Readability Variation Prompt Template

### 1. Task Definition:

You are rewriting a query to make it significantly less readable while preserving the original semantic meaning as closely as possible.

### 2. Constraints & Goals:

- Flesch Reading Ease Score: The rewritten text must have a Flesch score below 60 (preferably below 50).
- Semantic Similarity: The rewritten text must have SBERT similarity  $> 0.7$  compared with the original query.
- Length: The rewritten text must remain approximately the same length as the original query ( $\pm 10\%$ ).
- Preserve Domain Terminology: Do not remove or drastically change domain-specific words, abbreviations, or technical terms (e.g., "IRS," "distance," etc.).
- Abbreviation: Do not expand abbreviations unless the original query already used the expanded form.
- No New Information: You must not add additional details beyond what the original query states.
- Question Format: Retain the form of a question if the original is posed as a question.

### 3. How to Increase Complexity:

- Lexical Changes: Use advanced or academic synonyms only for common words. For domain or key terms (e.g., "distance," "IRS," "tax"), keep the original term or use a very close synonym if necessary to maintain meaning.
- Syntactic Complexity: Introduce passive voice, nominalizations, embedded clauses, and parenthetical or subordinate phrases. Ensure the sentence flow is more formal and convoluted without changing the core meaning.
- Redundancy & Formality: Employ circumlocution and excessively formal expressions (e.g., "due to the fact that" instead of "because") while avoiding any semantic drift.
- Dense, Indirect Construction: Favor longer phrases, indirect references, and wordiness. Avoid direct or simple phrasing.

Table 6: Prompt for lower-readability query generation.

ing, indicating that the rewrites are largely meaning-preserving and adhere well to the intended query categories.

## C.2 Human Annotation Details

We recruited 12 volunteer annotators to produce human query rewrites. Annotation guidelines are provided in Table 7. All volunteers were fluent in English and between 20 and 40 years old; with 5 female and 7 male participants. All annotators were graduate students or postdoctoral researchers with research experience in NLP or closely related fields. They reported high familiarity with AI tools (e.g., LLM-based assistants) and used them frequently in everyday tasks. Participation was unpaid and voluntary. Prior to the annotation process, all annotators provided informed consent for their rewrites to be collected and used for research purposes.

## D Robustness Results

In this section, we report results for all six datasets and model variants in our benchmark.

### D.1 Answer Robustness

**InAccuracy.** Figure 9 and 10 show consistent degradation under spelling errors. In contrast, sentence-type variations (imperative and declarative) achieve performance comparable to the original, suggesting that these rewrites have little effect

on accuracy. For human rewrites, Figures 11 and 12 show consistently lower performance across most datasets.

**Answer Consistency.** Tables 8 and 9 report model-generated answer consistency across all six datasets. Full results including human query variations are shown in Tables 10 and 11. Overall, human rewrites tend to yield answers that are less similar to those produced for the original queries.

### D.2 Computation Robustness

We measure computation robustness in terms of the computational cost incurred by the LLM and the retriever. We report the actual number of calls for each method–dataset pair. Figures 13 and 14 visualize the number of LLM calls across query variations. Figures 15 and 16 report the corresponding retriever calls for each model. As dataset difficulty increases, the generation\_based method adaptively increases both LLM and retriever calls.

Robustness across model-generated query variations is reported in Tables 12 and 13. Results including human queries are shown in Tables 14 and 15. Overall, increased query complexity and higher call counts are associated with lower computational robustness.

## 793 **E Self-Knowledge Results**

794 We visualize the under- and over-confidence results  
795 for model-generated queries in Figures 17 and 18.  
796 As dataset difficulty increases, the system exhibits  
797 more overconfident cases. Figures 19 and 20 report  
798 the results when including human rewrites as well.

## 799 **F Human Query Analysis**

800 In this section, we visualize human query features  
801 in terms of loss, length, and semantic similarity. As  
802 shown in Figure 22, human queries are consistently  
803 shorter than other query variations. Consistent with  
804 their shorter length, Figure 21 shows that human,  
805 readability, spelling, and grammar-error variations  
806 tend to have higher loss. Lastly, semantic similarity  
807 to the original query shows a similar trend across  
808 datasets, indicating that human queries preserve  
809 the original query’s meaning adequately.

## 810 **G Usage of GenAI**

811 We used AI-assisted coding tools to support data  
812 analysis and visualization (e.g., drafting and debug-  
813 ging scripts for metric aggregation and figure/table  
814 generation). All code and reported results were  
815 reviewed and validated by the authors.

---

## Human Annotation Instruction

---

### Task Overview

1. You will be given a query (question) and its corresponding answer.
2. Your task is to rewrite the given query — imagine how you would ask the same question if you were speaking to a large language model (like ChatGPT).

### What You Should Do

1. Understand the original query.
  - Identify what the user is trying to find out.
  - Grasp the intent and main topic clearly.
2. Rewrite the query in a new way while keeping the same meaning.
  - You may change the tone, phrasing, or structure.
  - Information loss, addition, or simplification is acceptable.
  - The goal is to make the rewritten query similar in meaning but different in expression.
  - Make sure that your rewritten query can still be answered appropriately with the same answer as the original query.
3. Incorporate your personal style.
  - You can make the question sound more casual, detailed, concise, or clear.
  - Different writing styles and tones are encouraged.
4. Avoid copying the original query.
  - Do not simply rephrase it with only minor surface changes.
  - The rewritten query should feel genuinely reworded and natural.
5. Use the answer only for context.
  - The answer helps you understand what the original query intended to ask.
  - You should not change your rewrite based on the answer's content.
  - Your task is to rewrite the question, not the answer.

---

Table 7: Human Annotation Instruction: Query Rewriting Task.

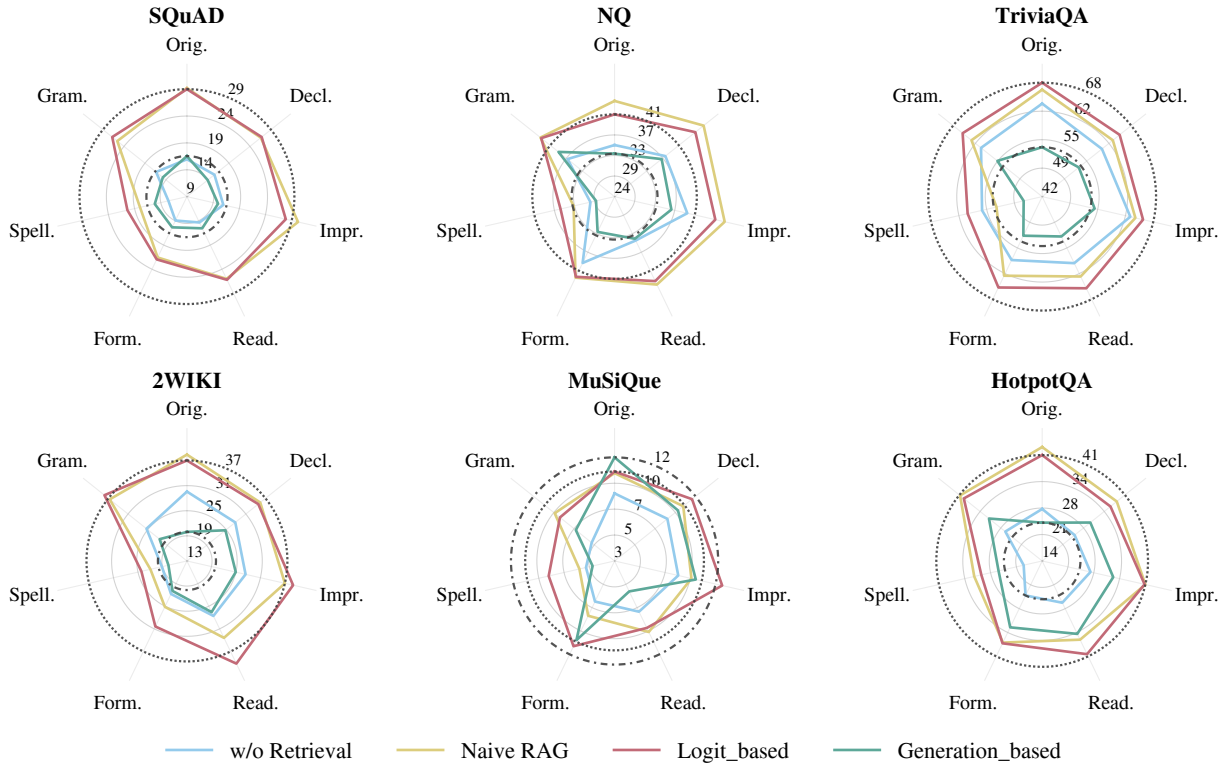


Figure 9: **InAccuracy (Llama-3.1-8B)**. Performance across datasets for four Adaptive RAG methods. Each panel reports InAccuracy on seven query variations: the original query and six model-generated perturbations.

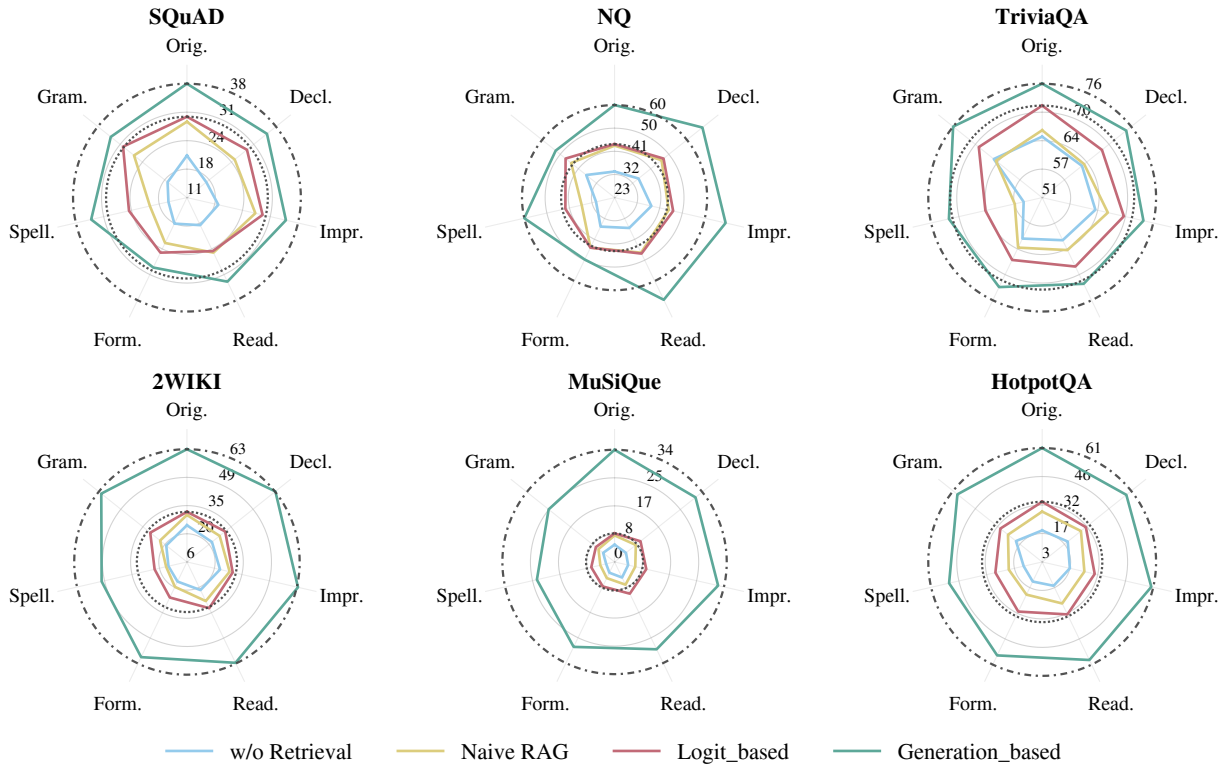


Figure 10: **InAccuracy (QwQ-32B)**. Performance across datasets for four Adaptive RAG methods. Each panel reports InAccuracy on seven query variations: the original query and six model-generated perturbations).

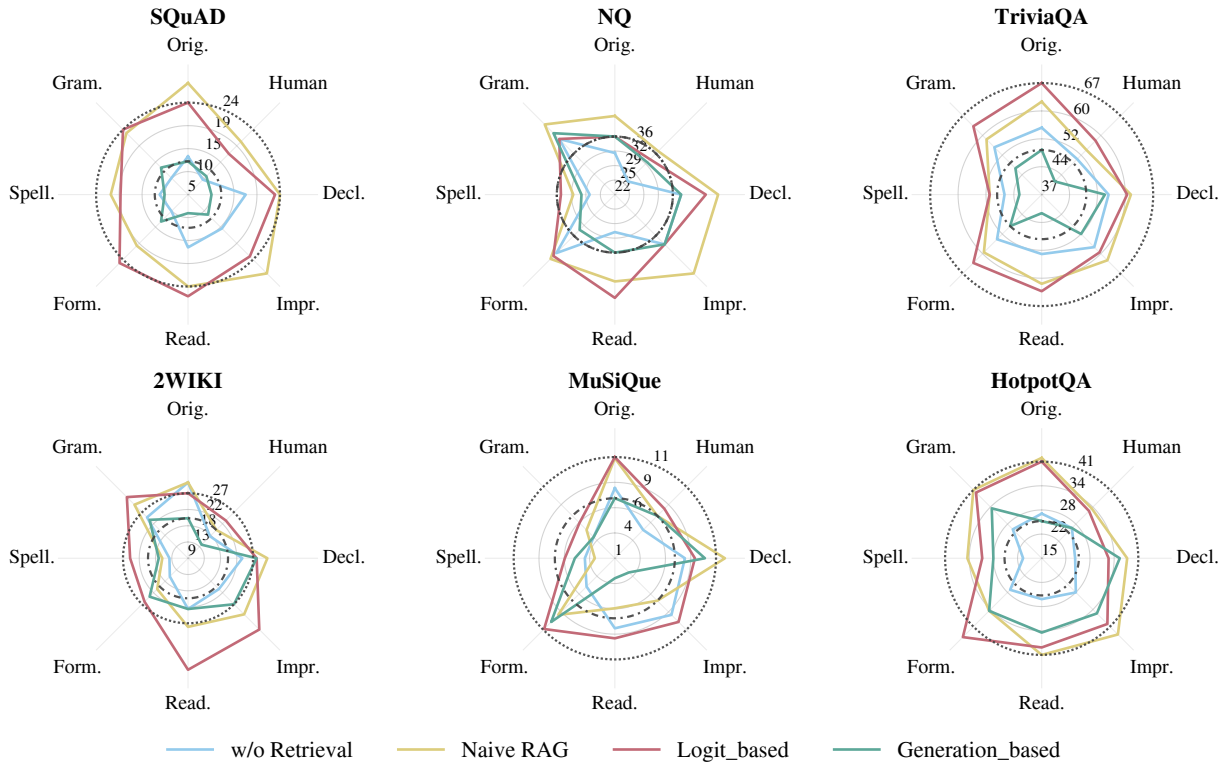


Figure 11: **InAccuracy (Llama-3.1-8B) including human rewrites.** Performance across datasets for four Adaptive RAG methods, evaluated on a matched subset of instances that have human rewrites. For each panel, InAccuracy is reported for eight query variations (original, human rewrite, and six model-generated perturbations), all computed on the same human-rewrite subset.

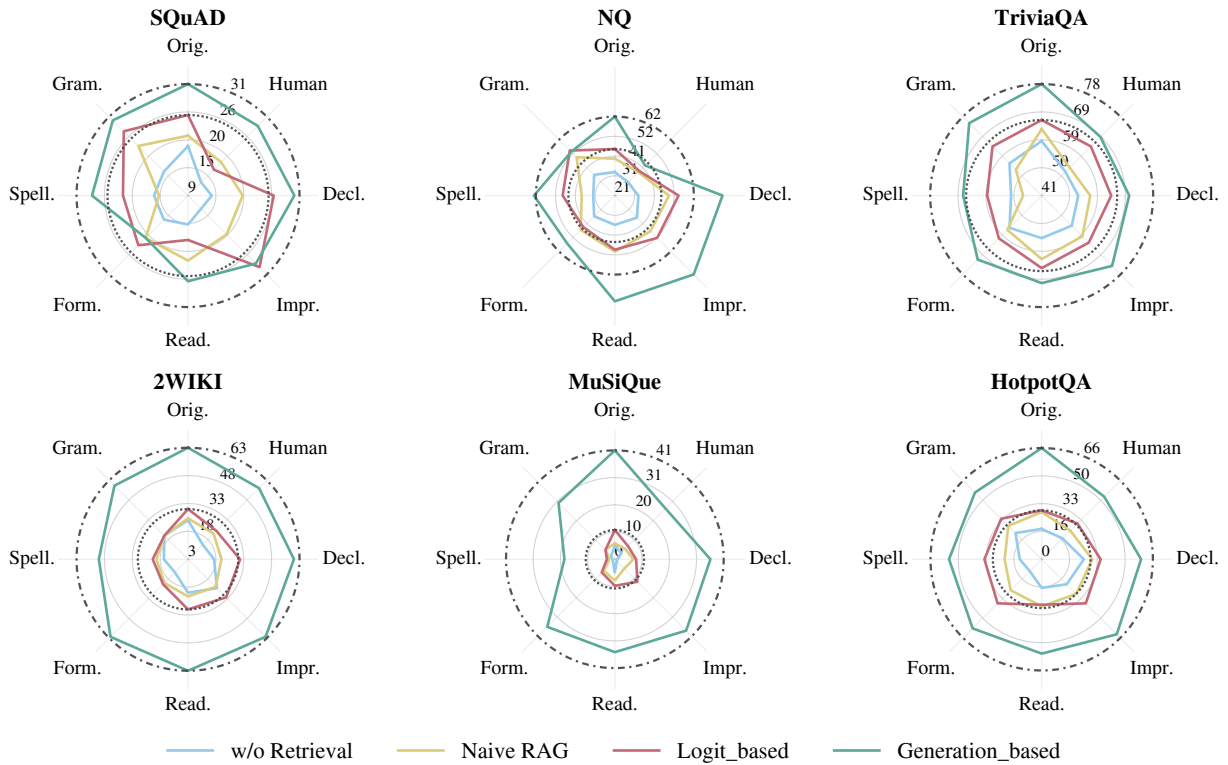


Figure 12: **InAccuracy (QwQ-32B) including human rewrites.** Performance across datasets for four Adaptive RAG methods, evaluated on a matched subset of instances that have human rewrites. For each panel, InAccuracy is reported for eight query variations (original, human rewrite, and six model-generated perturbations), all computed on the same human-rewrite subset.

Dataset	Method	Similarity $\uparrow$						Diversity $\downarrow$
		Form.	Read.	Decl.	Impr.	Spell.	Gram.	
<b>Llama-3.1-8B</b>								
SQuAD	WR	0.568	0.580	0.598	<b>0.635</b>	0.538	0.584	<b>0.459</b>
	NR	0.509	0.552	0.606	<b>0.631</b>	0.466	0.578	<b>0.499</b>
	LB	0.551	0.539	0.593	<b>0.634</b>	0.552	0.599	<b>0.469</b>
	GB	0.426	0.408	<b>0.433</b>	0.430	0.411	0.414	<b>0.585</b>
NQ	WR	0.672	0.670	0.652	<b>0.712</b>	0.610	0.675	<b>0.373</b>
	NR	0.311	0.483	<b>0.499</b>	0.494	0.425	0.321	<b>0.621</b>
	LB	0.706	0.686	0.708	<b>0.729</b>	0.626	0.724	<b>0.346</b>
	GB	0.382	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	0.435	<b>0.347</b>
TriviaQA	WR	0.767	0.773	0.762	<b>0.787</b>	0.734	0.763	<b>0.268</b>
	NR	0.700	0.716	0.725	<b>0.758</b>	0.648	0.738	<b>0.327</b>
	LB	0.781	<b>0.799</b>	0.769	0.798	0.727	0.785	<b>0.259</b>
	GB	0.528	0.537	0.543	0.547	0.504	<b>0.549</b>	<b>0.467</b>
2WIKI	WR	0.561	0.595	0.630	<b>0.660</b>	0.520	0.593	<b>0.442</b>
	NR	0.568	0.628	0.644	<b>0.651</b>	0.530	0.641	<b>0.432</b>
	LB	0.560	0.628	0.649	<b>0.655</b>	0.537	0.627	<b>0.414</b>
	GB	0.423	<b>0.458</b>	0.300	0.403	0.313	0.385	<b>0.667</b>
HotpotQA	WR	0.536	0.566	0.547	<b>0.580</b>	0.509	0.557	<b>0.484</b>
	NR	0.532	0.542	0.546	0.543	0.500	<b>0.564</b>	<b>0.500</b>
	LB	0.594	<b>0.628</b>	0.591	0.625	0.557	0.617	<b>0.415</b>
	GB	0.479	<b>0.502</b>	0.481	0.483	0.456	0.489	<b>0.523</b>
MuSiQue	WR	0.469	0.489	0.491	<b>0.512</b>	0.423	0.480	<b>0.560</b>
	NR	0.427	0.454	<b>0.471</b>	0.459	0.385	0.452	<b>0.588</b>
	LB	0.513	0.480	0.502	<b>0.528</b>	0.428	0.481	<b>0.537</b>
	GB	0.424	0.426	0.419	<b>0.438</b>	0.389	0.392	<b>0.576</b>

Table 8: Answer similarity to the original query and diversity across perturbation types (excluding human rewrites). We compare w/o retrieval (WR), naive RAG (NR), logit-based (LB), and generation-based (GB). Darker shading indicates more desirable outcomes: higher Similarity and lower Diversity.

Dataset	Method	Similarity $\uparrow$						Diversity $\downarrow$
		Form.	Read.	Decl.	Impr.	Spell.	Gram.	
<b>QwQ-32B</b>								
<b>SQuAD</b>	WR	0.566	0.576	0.602	<b>0.634</b>	0.549	0.593	<b>0.438</b>
	NR	0.511	0.570	0.625	<b>0.649</b>	0.488	0.599	<b>0.475</b>
	LB	0.573	0.579	0.615	<b>0.631</b>	0.544	0.584	<b>0.440</b>
	GB	0.552	0.558	0.579	<b>0.589</b>	0.552	0.579	<b>0.459</b>
<b>NQ</b>	WR	0.660	0.634	0.671	<b>0.678</b>	0.583	0.673	<b>0.373</b>
	NR	0.367	0.518	0.375	<b>0.520</b>	0.483	0.389	<b>0.566</b>
	LB	0.344	0.304	0.321	0.309	0.305	<b>0.353</b>	<b>0.662</b>
	GB	0.438	0.769	0.783	<b>0.797</b>	0.739	0.442	<b>0.428</b>
<b>TriviaQA</b>	WR	0.743	0.735	<b>0.759</b>	0.759	0.699	0.752	<b>0.288</b>
	NR	0.698	0.738	0.735	<b>0.759</b>	0.671	0.747	<b>0.309</b>
	LB	0.770	0.775	0.785	0.779	0.730	<b>0.788</b>	<b>0.255</b>
	GB	0.608	0.588	0.603	<b>0.626</b>	0.580	0.604	<b>0.395</b>
<b>2WIKI</b>	WR	0.528	0.554	0.589	<b>0.591</b>	0.481	0.483	<b>0.531</b>
	NR	0.443	0.500	0.520	<b>0.531</b>	0.422	0.514	<b>0.547</b>
	LB	0.500	0.545	0.551	<b>0.565</b>	0.462	0.519	<b>0.511</b>
	GB	0.520	0.557	0.554	<b>0.559</b>	0.499	0.548	<b>0.465</b>
<b>HotpotQA</b>	WR	0.359	<b>0.402</b>	0.385	0.401	0.333	0.394	<b>0.639</b>
	NR	0.377	<b>0.414</b>	0.408	0.401	0.357	0.397	<b>0.624</b>
	LB	0.444	0.452	<b>0.468</b>	0.458	0.406	0.444	<b>0.568</b>
	GB	0.582	0.585	0.581	<b>0.608</b>	0.563	0.580	<b>0.429</b>
<b>MuSiQue</b>	WR	0.394	0.390	<b>0.408</b>	0.405	0.371	0.408	<b>0.629</b>
	NR	0.387	0.377	0.396	<b>0.405</b>	0.374	0.383	<b>0.634</b>
	LB	0.395	0.400	0.414	<b>0.422</b>	0.380	0.404	<b>0.610</b>
	GB	0.593	<b>0.607</b>	0.606	0.607	0.589	0.582	<b>0.422</b>

Table 9: Answer similarity to the original query and diversity across perturbation types (excluding human rewrites). We compare w/o retrieval (WR), naive RAG (NR), logit-based (LB), and generation-based (GB). Darker shading indicates more desirable outcomes: higher Similarity and lower Diversity.

Dataset	Method	Similarity $\uparrow$							Diversity $\downarrow$
		Human	Form.	Read.	Decl.	Impr.	Spell.	Gram.	
<b>Llama-3.1-8B</b>									
SQuAD	WR	0.542	0.579	0.610	0.644	<b>0.653</b>	0.542	0.600	<b>0.446</b>
	NR	0.505	0.580	0.544	0.613	<b>0.622</b>	0.479	0.607	<b>0.503</b>
	LB	0.527	0.557	0.557	0.599	<b>0.626</b>	0.566	0.584	<b>0.469</b>
	GB	0.447	0.416	0.421	0.423	<b>0.449</b>	0.427	0.449	<b>0.583</b>
NQ	WR	0.624	0.698	0.646	0.683	<b>0.748</b>	0.586	0.692	<b>0.385</b>
	NR	0.292	0.300	0.470	<b>0.507</b>	0.500	0.385	0.319	<b>0.616</b>
	LB	0.642	0.674	0.662	0.661	0.686	0.577	<b>0.691</b>	<b>0.370</b>
	GB	0.356	0.382	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	0.435	<b>0.416</b>
TriviaQA	WR	0.660	0.734	0.753	0.741	<b>0.755</b>	0.672	0.700	<b>0.331</b>
	NR	0.634	0.673	0.691	0.690	<b>0.743</b>	0.635	0.725	<b>0.370</b>
	LB	0.697	0.735	<b>0.776</b>	0.718	0.729	0.708	0.735	<b>0.314</b>
	GB	0.475	0.499	0.458	0.541	<b>0.555</b>	0.460	0.538	<b>0.517</b>
2WIKI	WR	0.540	0.516	0.567	0.628	<b>0.628</b>	0.511	0.607	<b>0.480</b>
	NR	0.488	0.516	0.585	<b>0.641</b>	0.592	0.453	0.615	<b>0.473</b>
	LB	0.550	0.541	<b>0.617</b>	0.616	0.610	0.502	0.590	<b>0.451</b>
	GB	0.279	0.423	<b>0.458</b>	0.300	0.403	0.313	0.385	<b>0.647</b>
HotpotQA	WR	0.514	0.544	0.545	0.555	<b>0.592</b>	0.511	0.565	<b>0.482</b>
	NR	0.524	0.568	0.578	0.574	0.583	0.520	<b>0.616</b>	<b>0.465</b>
	LB	0.543	0.630	0.617	0.586	0.634	0.588	<b>0.652</b>	<b>0.401</b>
	GB	0.462	0.477	0.484	0.452	0.465	0.489	<b>0.515</b>	<b>0.507</b>
MuSiQue	WR	0.459	0.417	0.511	0.471	<b>0.525</b>	0.403	0.454	<b>0.564</b>
	NR	0.407	0.415	<b>0.476</b>	0.433	0.452	0.382	0.400	<b>0.600</b>
	LB	0.504	0.528	<b>0.578</b>	0.486	0.576	0.436	0.487	<b>0.551</b>
	GB	0.360	0.423	<b>0.442</b>	0.380	0.415	0.372	0.370	<b>0.590</b>

Table 10: Answer similarity to the original query and diversity across perturbation types (including human rewrites). We compare w/o retrieval (WR), naive RAG (NR), logit-based (LB), and generation-based (GB). Darker shading indicates more desirable outcomes: higher Similarity and lower Diversity.

Dataset	Method	Similarity $\uparrow$							Diversity $\downarrow$
		Human	Form.	Read.	Decl.	Impr.	Spell.	Gram.	
<b>QwQ-32B</b>									
SQuAD	WR	0.585	0.567	0.608	0.630	<b>0.642</b>	0.583	0.605	<b>0.448</b>
	NR	0.556	0.514	0.619	0.655	<b>0.664</b>	0.518	0.584	<b>0.459</b>
	LB	0.565	0.608	0.590	0.608	<b>0.635</b>	0.577	0.599	<b>0.435</b>
	GB	0.544	0.549	0.560	0.596	<b>0.601</b>	0.575	0.588	<b>0.461</b>
NQ	WR	0.594	0.631	0.612	0.649	0.635	0.573	<b>0.660</b>	<b>0.374</b>
	NR	0.370	0.358	<b>0.517</b>	0.374	0.510	0.461	0.374	<b>0.551</b>
	LB	0.330	0.348	0.294	0.319	0.346	0.300	<b>0.387</b>	<b>0.628</b>
	GB	0.438	0.438	0.769	0.783	<b>0.797</b>	0.739	0.442	<b>0.454</b>
TriviaQA	WR	0.684	<b>0.718</b>	0.678	0.696	0.705	0.661	0.687	<b>0.336</b>
	NR	0.646	0.649	0.683	0.687	<b>0.720</b>	0.597	0.664	<b>0.363</b>
	LB	0.697	0.730	0.712	<b>0.745</b>	0.693	0.665	0.717	<b>0.307</b>
	GB	0.548	0.551	0.531	0.539	<b>0.582</b>	0.546	0.558	<b>0.440</b>
2WIKI	WR	0.560	0.542	0.556	0.623	<b>0.634</b>	0.519	0.570	<b>0.496</b>
	NR	0.511	0.492	0.545	<b>0.579</b>	0.569	0.482	0.551	<b>0.516</b>
	LB	0.518	0.528	0.570	0.573	<b>0.612</b>	0.522	0.572	<b>0.486</b>
	GB	0.529	0.542	0.534	0.568	<b>0.598</b>	0.520	0.558	<b>0.463</b>
HotpotQA	WR	0.356	0.394	<b>0.455</b>	0.412	0.416	0.368	0.421	<b>0.636</b>
	NR	0.379	0.366	0.405	0.398	<b>0.412</b>	0.364	0.381	<b>0.621</b>
	LB	0.420	0.467	0.417	<b>0.486</b>	0.476	0.428	0.468	<b>0.560</b>
	GB	0.616	0.595	0.596	0.598	<b>0.629</b>	0.573	0.612	<b>0.426</b>
MuSiQue	WR	0.395	0.384	0.388	<b>0.411</b>	0.402	0.372	0.395	<b>0.633</b>
	NR	0.363	<b>0.427</b>	0.396	0.425	0.404	0.354	0.382	<b>0.627</b>
	LB	0.397	<b>0.416</b>	0.371	0.403	0.408	0.389	0.390	<b>0.619</b>
	GB	0.584	0.582	0.604	<b>0.627</b>	0.603	0.595	0.573	<b>0.424</b>

Table 11: Answer similarity to the original query and diversity across perturbation types (including human rewrites). We compare w/o retrieval (WR), naive RAG (NR), logit-based (LB), and generation-based (GB). Darker shading indicates more desirable outcomes: higher Similarity and lower Diversity.

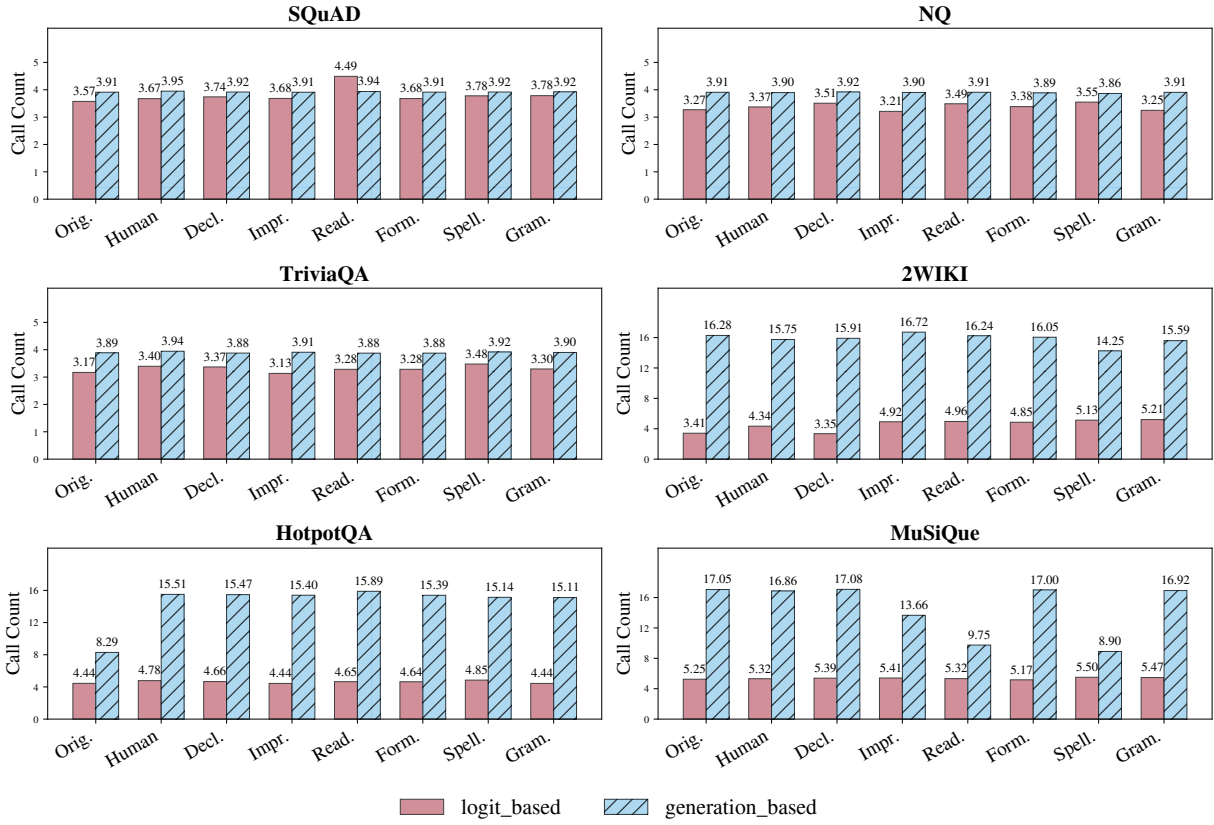


Figure 13: LLM Call on Llama-3.1-8B on both model-generated and human queries.

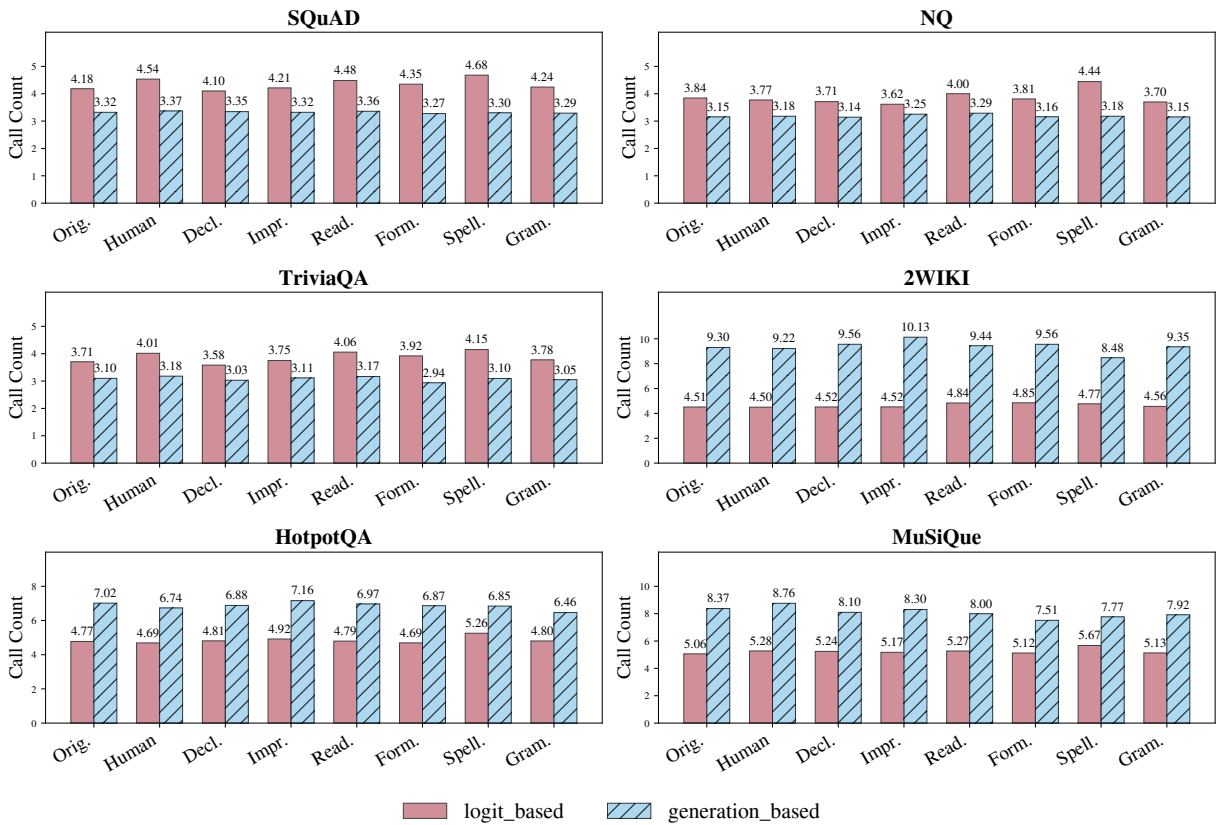


Figure 14: LLM Call on QwQ-32B on both model-generated and human queries.

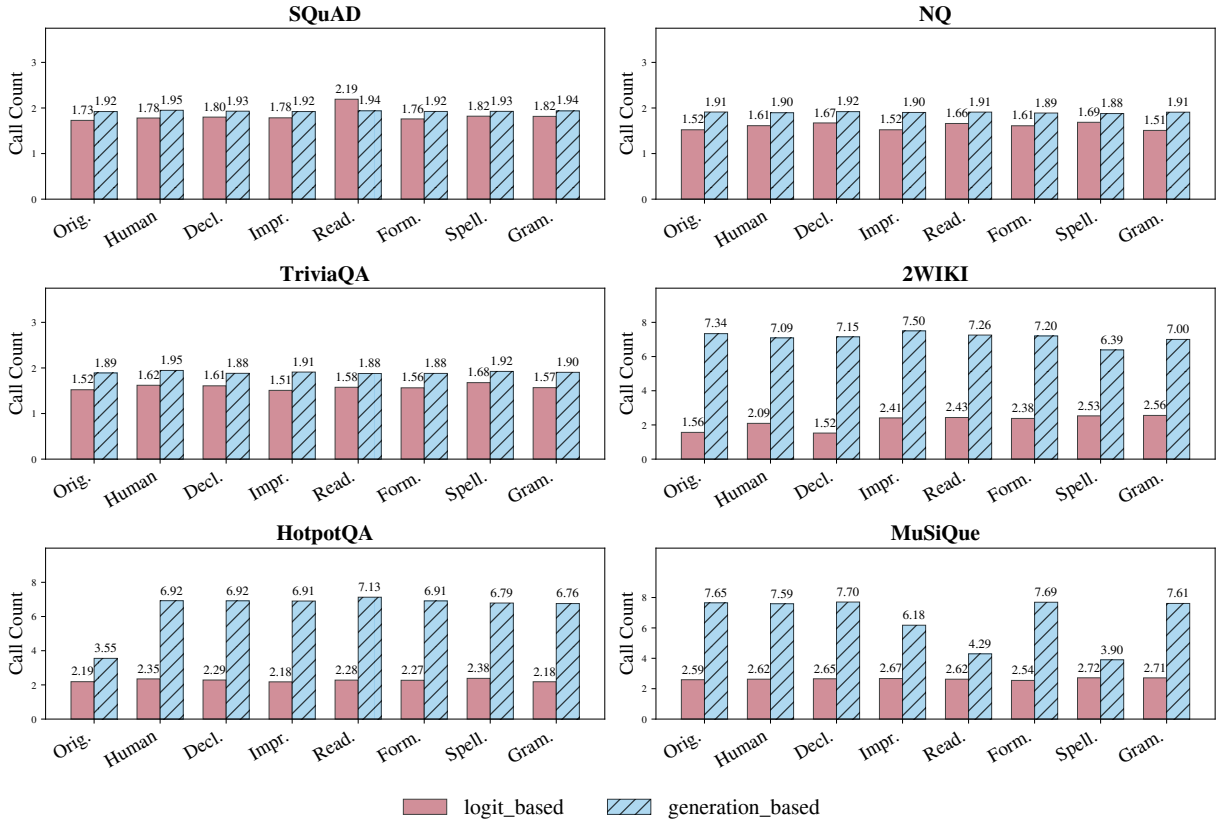


Figure 15: Retriever Call on Llama-3.1-8B on both model-generated and human queries.

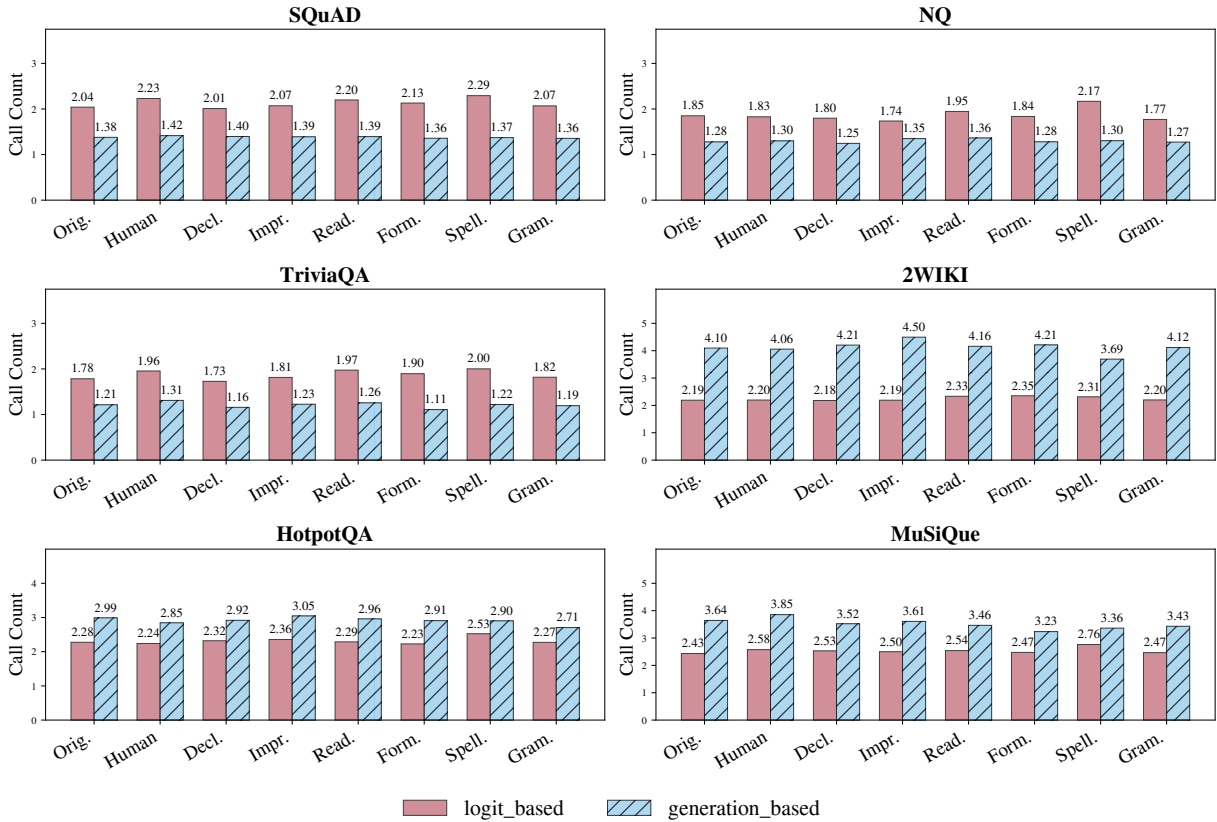


Figure 16: Retriever Call on QwQ-32B on both model-generated and human queries.

Dataset	Target	Method	RE ↓						CVR ↑
			Form.	Read.	Decl.	Impr.	Spell.	Gram.	
<b>Llama-3.1-8B</b>									
SQuAD	Retriever	LB	0.472	0.592	0.416	0.338	0.431	0.439	0.734
		GB	<b>0.092</b>	<b>0.094</b>	<b>0.086</b>	<b>0.090</b>	<b>0.100</b>	<b>0.092</b>	0.937
	LLM	LB	0.441	0.592	0.396	0.313	0.416	0.411	0.767
		GB	<b>0.056</b>	<b>0.059</b>	<b>0.048</b>	<b>0.055</b>	<b>0.061</b>	<b>0.051</b>	0.963
NQ	Retriever	LB	0.520	0.485	0.576	0.465	0.578	0.456	0.702
		GB	<b>0.096</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.096</b>	0.949
	LLM	LB	0.499	0.462	0.541	0.428	0.545	0.440	0.750
		GB	<b>0.038</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.042</b>	0.973
TriviaQA	Retriever	LB	0.452	0.442	0.493	0.436	0.536	0.471	0.712
		GB	<b>0.137</b>	<b>0.120</b>	<b>0.121</b>	<b>0.122</b>	<b>0.124</b>	<b>0.122</b>	0.917
	LLM	LB	0.463	0.420	0.472	0.410	0.527	0.466	0.747
		GB	<b>0.054</b>	<b>0.047</b>	<b>0.048</b>	<b>0.046</b>	<b>0.047</b>	<b>0.047</b>	0.956
2WIKI	Retriever	LB	0.853	0.859	0.455	0.846	0.928	0.941	0.710
		GB	<b>0.468</b>	<b>0.416</b>	<b>0.402</b>	<b>0.542</b>	<b>0.496</b>	<b>0.435</b>	0.742
	LLM	LB	1.001	1.001	0.427	1.007	1.093	1.130	0.742
		GB	<b>0.447</b>	<b>0.404</b>	<b>0.401</b>	<b>0.520</b>	<b>0.477</b>	<b>0.418</b>	0.745
HotpotQA	Retriever	LB	<b>0.486</b>	<b>0.461</b>	<b>0.428</b>	<b>0.440</b>	<b>0.487</b>	<b>0.428</b>	0.747
		GB	1.914	1.992	1.850	1.848	1.924	1.825	0.687
	LLM	LB	<b>0.490</b>	<b>0.439</b>	<b>0.421</b>	<b>0.427</b>	<b>0.502</b>	<b>0.418</b>	0.761
		GB	1.560	1.626	1.525	1.517	1.565	1.495	0.694
MuSiQue	Retriever	LB	<b>0.404</b>	<b>0.417</b>	<b>0.382</b>	<b>0.345</b>	<b>0.399</b>	<b>0.340</b>	0.753
		GB	0.521	0.525	0.533	0.446	0.574	0.510	0.691
	LLM	LB	<b>0.365</b>	<b>0.398</b>	<b>0.353</b>	<b>0.323</b>	<b>0.388</b>	<b>0.318</b>	0.763
		GB	0.502	0.502	0.519	0.429	0.550	0.484	0.698

Table 12: **Computation robustness across query variations (Llama-3.1-8B; no human rewrites)**. Darker shading indicates more desirable outcomes: lower RE and higher CVR.

Dataset	Target	Method	RE ↓						CVR ↑
			Form.	Read.	Decl.	Impr.	Spell.	Gram.	
<b>QwQ-32B</b>									
SQuAD	Retriever	LB	0.494	0.450	0.405	0.415	0.495	0.441	0.740
		GB	<b>0.360</b>	<b>0.322</b>	<b>0.298</b>	<b>0.315</b>	<b>0.361</b>	<b>0.330</b>	0.762
	LLM	LB	0.477	0.466	0.389	0.399	0.491	0.424	0.760
		GB	<b>0.265</b>	<b>0.230</b>	<b>0.225</b>	<b>0.238</b>	<b>0.250</b>	<b>0.243</b>	0.858
NQ	Retriever	LB	0.551	0.586	0.525	0.513	0.673	0.520	0.689
		GB	<b>0.411</b>	<b>0.402</b>	<b>0.374</b>	<b>0.416</b>	<b>0.439</b>	<b>0.452</b>	0.721
	LLM	LB	0.546	0.598	0.511	0.504	0.679	0.516	0.724
		GB	<b>0.379</b>	<b>0.375</b>	<b>0.358</b>	<b>0.372</b>	<b>0.411</b>	<b>0.424</b>	0.820
TriviaQA	Retriever	LB	0.581	0.563	0.495	0.504	0.587	0.546	0.697
		GB	<b>0.362</b>	<b>0.377</b>	<b>0.347</b>	<b>0.374</b>	<b>0.349</b>	<b>0.387</b>	0.721
	LLM	LB	0.564	0.565	0.470	0.509	0.601	0.530	0.727
		GB	<b>0.329</b>	<b>0.373</b>	<b>0.331</b>	<b>0.348</b>	<b>0.324</b>	<b>0.348</b>	0.818
2WIKI	Retriever	LB	<b>0.560</b>	<b>0.502</b>	<b>0.446</b>	<b>0.443</b>	<b>0.565</b>	<b>0.511</b>	0.719
		GB	0.788	0.808	0.780	0.825	0.787	0.884	0.689
	LLM	LB	<b>0.547</b>	<b>0.481</b>	<b>0.426</b>	<b>0.430</b>	<b>0.565</b>	<b>0.493</b>	0.747
		GB	1.078	1.142	1.092	1.203	1.049	1.219	0.716
HotpotQA	Retriever	LB	<b>0.589</b>	<b>0.620</b>	<b>0.580</b>	<b>0.590</b>	<b>0.671</b>	<b>0.613</b>	0.707
		GB	0.806	0.774	0.788	0.818	0.894	0.663	0.654
	LLM	LB	<b>0.516</b>	<b>0.520</b>	<b>0.487</b>	<b>0.504</b>	<b>0.605</b>	<b>0.505</b>	0.743
		GB	0.968	0.920	0.902	1.000	1.078	0.778	0.695
MuSiQue	Retriever	LB	<b>0.632</b>	<b>0.660</b>	<b>0.598</b>	<b>0.607</b>	<b>0.702</b>	<b>0.603</b>	0.711
		GB	0.811	0.890	0.767	0.883	0.845	0.881	0.644
	LLM	LB	<b>0.538</b>	<b>0.572</b>	<b>0.537</b>	<b>0.515</b>	<b>0.607</b>	<b>0.523</b>	0.740
		GB	1.032	1.168	0.984	1.203	1.100	1.117	0.679

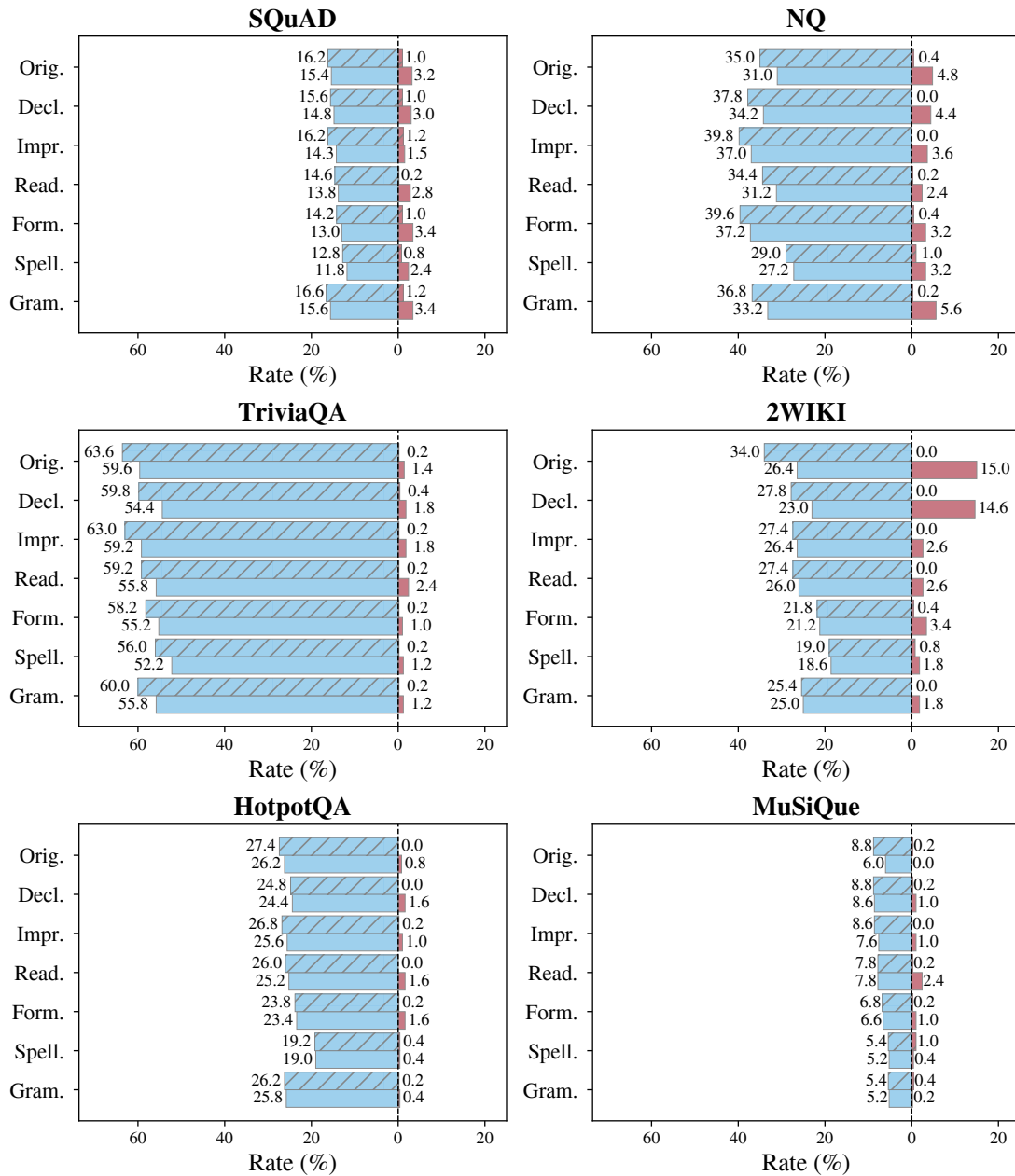
Table 13: **Computation robustness across query variations (QwQ-32B; no human rewrites)**. Darker shading indicates more desirable outcomes: lower RE and higher CVR.

Dataset	Target	Method	RE ↓						CVR ↑	
			Human	Form.	Read.	Decl.	Impr.	Spell.		Gram.
<b>Llama-3.1-8B</b>										
SQuAD	Retriever	LB	0.415	0.452	0.577	0.423	0.344	0.460	0.450	0.748
		GB	<b>0.105</b>	<b>0.090</b>	<b>0.100</b>	<b>0.105</b>	<b>0.090</b>	<b>0.120</b>	<b>0.120</b>	0.946
	LLM	LB	0.407	0.418	0.587	0.398	0.306	0.414	0.414	0.778
		GB	<b>0.061</b>	<b>0.055</b>	<b>0.060</b>	<b>0.065</b>	<b>0.048</b>	<b>0.068</b>	<b>0.067</b>	0.969
NQ	Retriever	LB	0.549	0.495	0.487	0.562	0.481	0.551	0.425	0.696
		GB	<b>0.077</b>	<b>0.096</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.096</b>	0.942
	LLM	LB	0.542	0.485	0.475	0.546	0.410	0.498	0.427	0.739
		GB	<b>0.029</b>	<b>0.038</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.042</b>	0.969
TriviaQA	Retriever	LB	0.492	0.508	0.457	0.500	0.397	0.535	0.472	0.709
		GB	<b>0.065</b>	<b>0.095</b>	<b>0.105</b>	<b>0.090</b>	<b>0.080</b>	<b>0.090</b>	<b>0.095</b>	0.924
	LLM	LB	0.496	0.563	0.466	0.537	0.378	0.532	0.468	0.747
		GB	<b>0.024</b>	<b>0.039</b>	<b>0.044</b>	<b>0.041</b>	<b>0.030</b>	<b>0.037</b>	<b>0.041</b>	0.959
2WIKI	Retriever	LB	0.745	0.862	0.870	0.510	1.000	0.890	0.960	0.717
		GB	<b>0.506</b>	<b>0.468</b>	<b>0.416</b>	<b>0.402</b>	<b>0.542</b>	<b>0.496</b>	<b>0.435</b>	0.740
	LLM	LB	0.827	1.065	1.041	0.526	1.259	1.126	1.195	0.747
		GB	<b>0.480</b>	<b>0.447</b>	<b>0.404</b>	<b>0.401</b>	<b>0.520</b>	<b>0.477</b>	<b>0.418</b>	0.743
HotpotQA	Retriever	LB	<b>0.426</b>	<b>0.495</b>	<b>0.429</b>	<b>0.350</b>	<b>0.413</b>	<b>0.496</b>	<b>0.415</b>	0.747
		GB	1.962	1.895	1.928	1.627	1.697	1.805	1.832	0.684
	LLM	LB	<b>0.426</b>	<b>0.479</b>	<b>0.396</b>	<b>0.328</b>	<b>0.381</b>	<b>0.498</b>	<b>0.401</b>	0.757
		GB	1.668	1.524	1.556	1.326	1.389	1.481	1.527	0.693
MuSiQue	Retriever	LB	<b>0.284</b>	<b>0.275</b>	<b>0.377</b>	<b>0.448</b>	<b>0.343</b>	<b>0.296</b>	<b>0.219</b>	0.751
		GB	0.597	0.568	0.544	0.554	0.438	0.586	0.585	0.704
	LLM	LB	<b>0.268</b>	<b>0.262</b>	<b>0.331</b>	<b>0.410</b>	<b>0.307</b>	<b>0.274</b>	<b>0.221</b>	0.761
		GB	0.528	0.510	0.504	0.491	0.415	0.542	0.505	0.710

Table 14: **Computation robustness across query variations (Llama-3.1-8B; including human rewrites).** Darker shading indicates more desirable outcomes: lower RE and higher CVR.

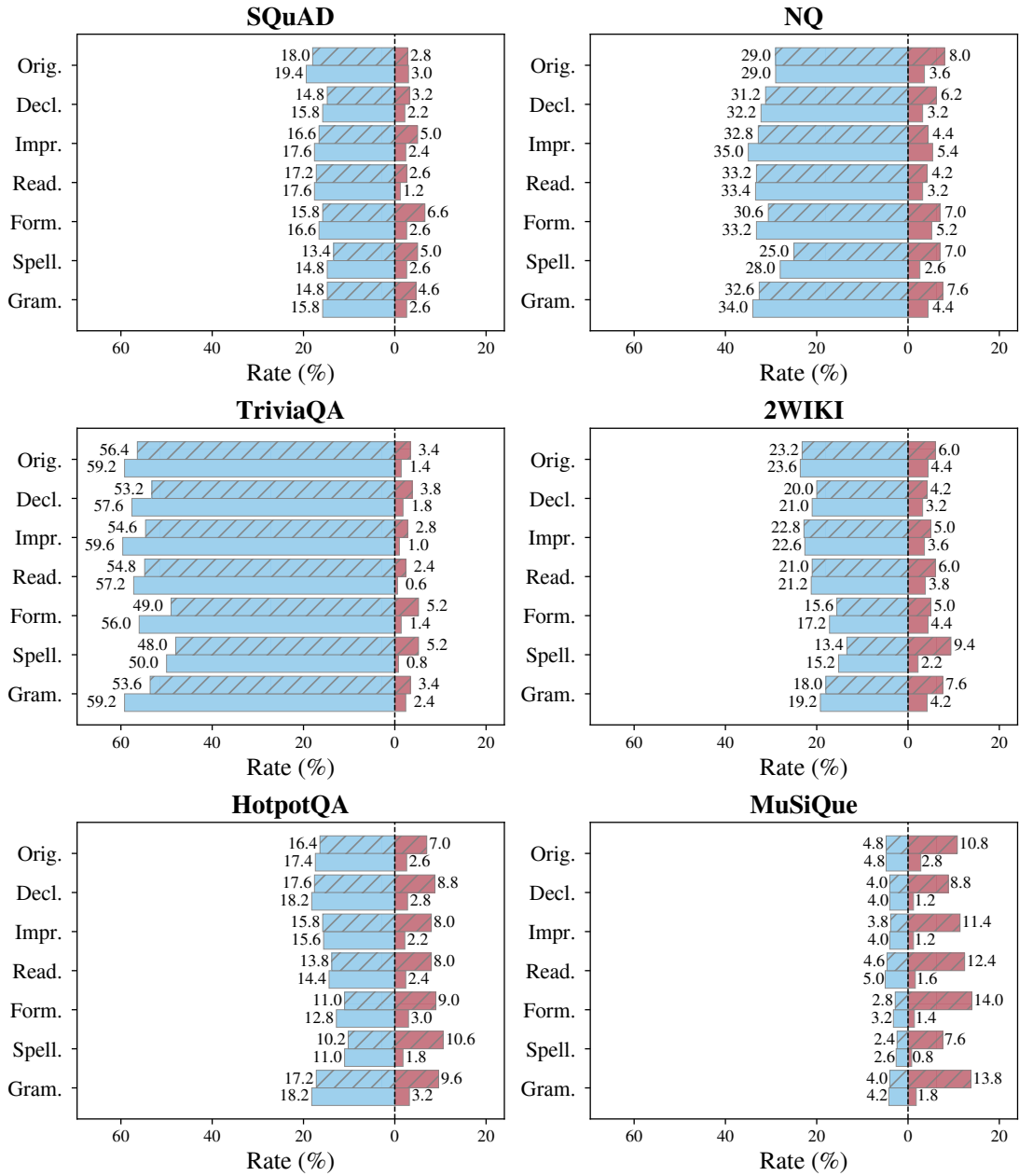
Dataset	Target	Method	RE ↓						CVR ↑	
			Human	Form.	Read.	Decl.	Impr.	Spell.		Gram.
<b>QwQ-32B</b>										
SQuAD	Retriever	LB	0.439	0.500	0.568	0.446	0.404	0.451	0.547	0.741
		GB	<b>0.370</b>	<b>0.295</b>	<b>0.280</b>	<b>0.310</b>	<b>0.240</b>	<b>0.335</b>	<b>0.290</b>	0.766
	LLM	LB	0.461	0.468	0.597	0.426	0.389	0.417	0.494	0.757
		GB	<b>0.285</b>	<b>0.252</b>	<b>0.233</b>	<b>0.242</b>	<b>0.210</b>	<b>0.273</b>	<b>0.242</b>	0.859
NQ	Retriever	LB	0.565	0.485	0.462	0.533	0.558	0.662	0.537	0.692
		GB	<b>0.404</b>	<b>0.411</b>	<b>0.402</b>	<b>0.374</b>	<b>0.416</b>	<b>0.439</b>	<b>0.452</b>	0.715
	LLM	LB	0.554	0.456	0.492	0.510	0.555	0.669	0.533	0.725
		GB	<b>0.370</b>	<b>0.379</b>	<b>0.375</b>	<b>0.358</b>	<b>0.372</b>	<b>0.411</b>	<b>0.424</b>	0.816
TriviaQA	Retriever	LB	0.488	0.510	0.438	0.538	0.514	0.600	0.557	0.700
		GB	<b>0.455</b>	<b>0.405</b>	<b>0.420</b>	<b>0.430</b>	<b>0.395</b>	<b>0.440</b>	<b>0.440</b>	0.716
	LLM	LB	0.473	0.459	<b>0.438</b>	0.483	0.530	0.590	0.486	0.727
		GB	<b>0.472</b>	<b>0.417</b>	0.499	<b>0.481</b>	<b>0.407</b>	<b>0.448</b>	<b>0.448</b>	0.811
2WIKI	Retriever	LB	<b>0.407</b>	<b>0.531</b>	<b>0.486</b>	<b>0.461</b>	<b>0.437</b>	<b>0.546</b>	<b>0.510</b>	0.720
		GB	0.746	0.686	0.643	0.760	0.757	0.751	0.794	0.679
	LLM	LB	<b>0.437</b>	<b>0.506</b>	<b>0.480</b>	<b>0.459</b>	<b>0.425</b>	<b>0.548</b>	<b>0.502</b>	0.749
		GB	0.911	0.939	0.942	1.083	1.130	0.996	1.135	0.708
HotpotQA	Retriever	LB	<b>0.500</b>	<b>0.497</b>	<b>0.552</b>	<b>0.508</b>	<b>0.451</b>	<b>0.570</b>	<b>0.465</b>	0.706
		GB	1.023	0.928	0.712	0.910	0.901	0.877	0.716	0.632
	LLM	LB	<b>0.446</b>	<b>0.462</b>	<b>0.448</b>	<b>0.440</b>	<b>0.411</b>	<b>0.575</b>	<b>0.392</b>	0.746
		GB	1.146	1.092	0.776	0.976	1.105	1.083	0.874	0.678
MuSiQue	Retriever	LB	<b>0.652</b>	<b>0.629</b>	<b>0.657</b>	0.602	0.689	<b>0.655</b>	<b>0.602</b>	0.719
		GB	0.770	0.797	0.767	<b>0.548</b>	<b>0.642</b>	0.713	0.709	0.645
	LLM	LB	<b>0.557</b>	<b>0.501</b>	<b>0.537</b>	<b>0.512</b>	<b>0.579</b>	<b>0.570</b>	<b>0.510</b>	0.745
		GB	0.888	0.908	0.916	0.584	0.718	0.767	0.719	0.679

Table 15: **Computation robustness across query variations (QwQ-32B; including human rewrites).** Darker shading indicates more desirable outcomes: lower RE and higher CVR.



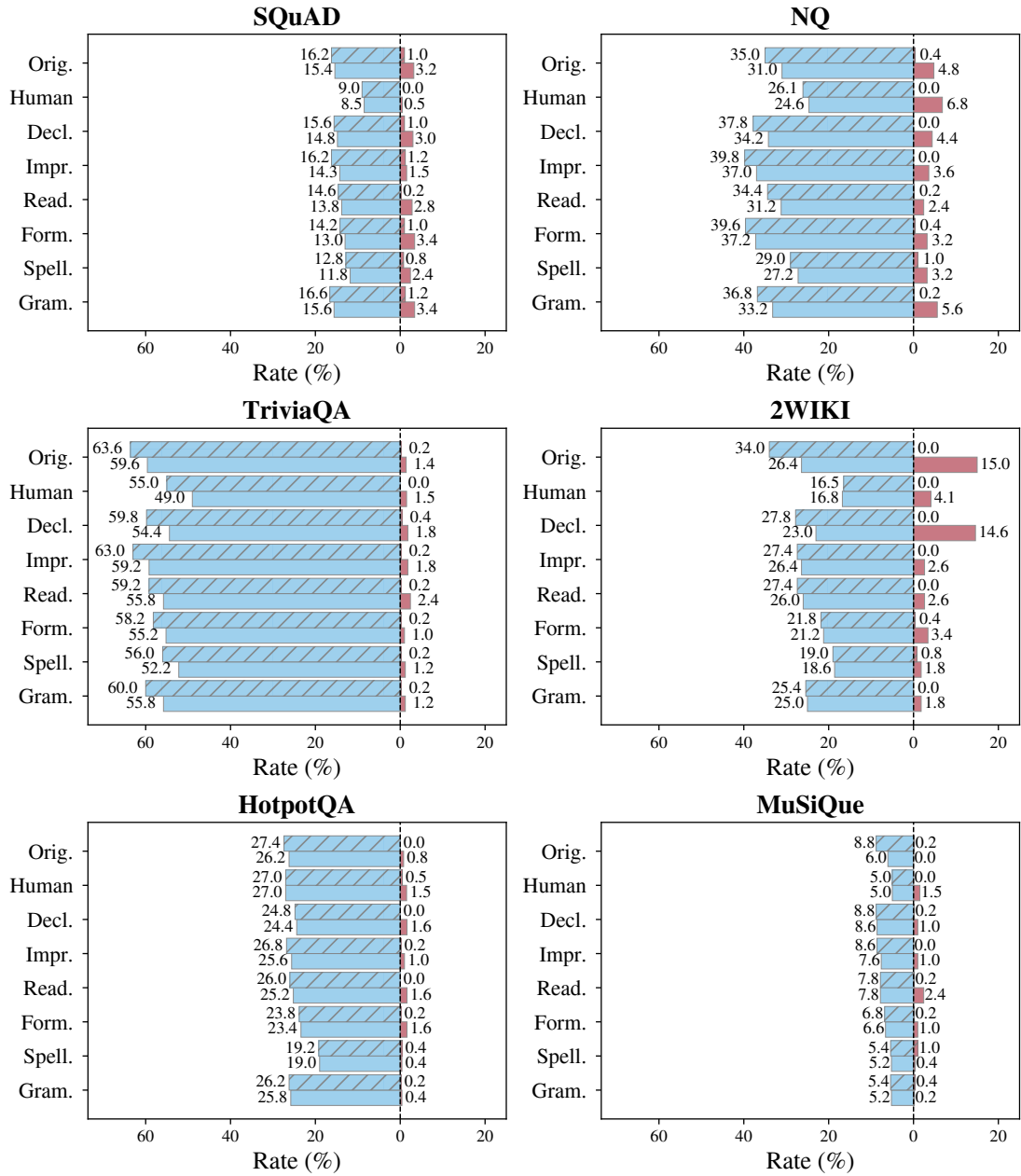
Under-confident
  Over-confident
  Logit\_based
  Generation\_based

Figure 17: Under- and Over-confidence results on Llama-3.1-8B on model-generated queries.



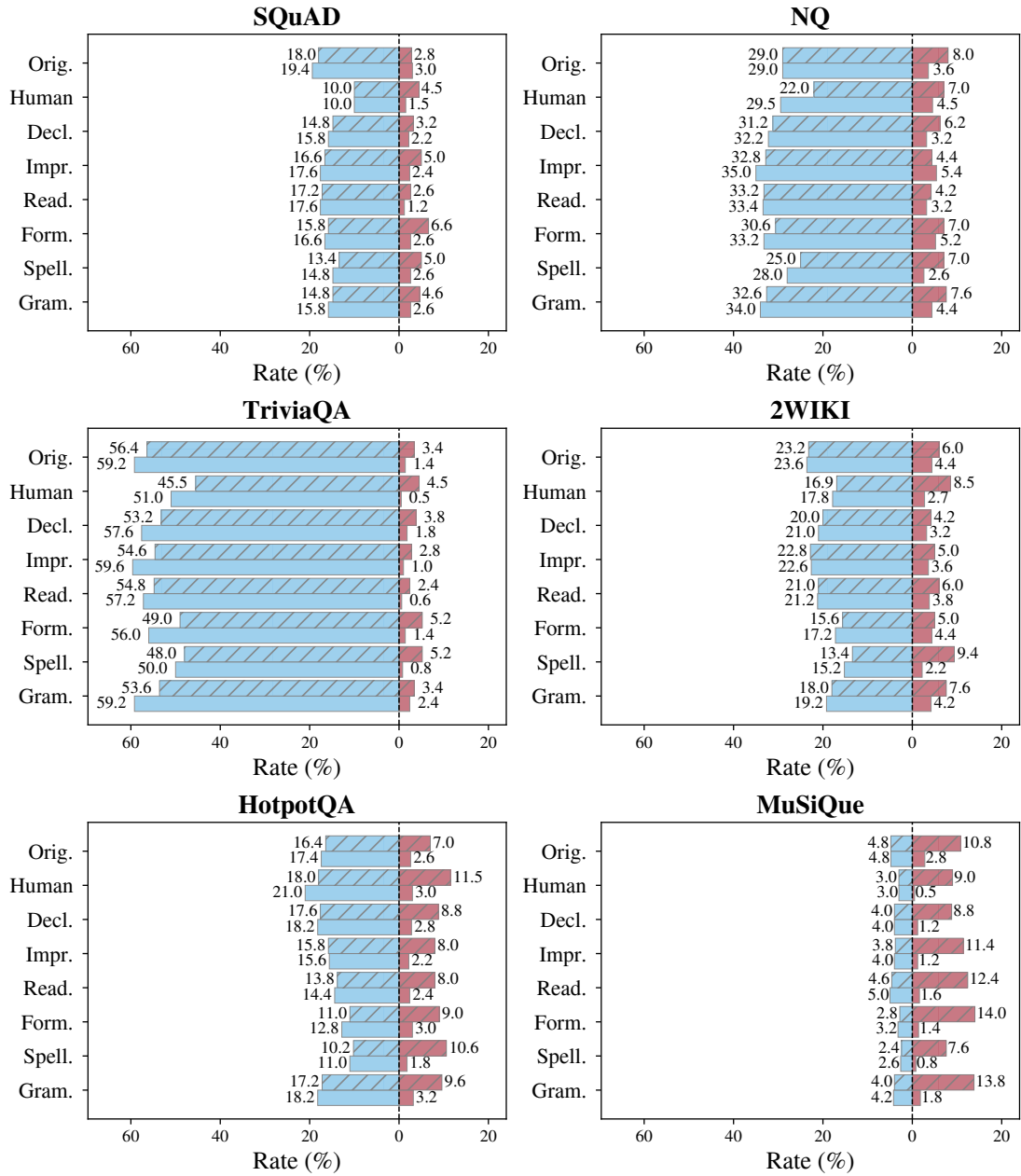
Under-confident    Over-confident    Logit\_based    Generation\_based

Figure 18: Under- and Over-confidence results on QwQ-32B on model-generated queries.



Under-confident    Over-confident    Logit\_based    Generation\_based

Figure 19: Under- and Over-confidence results on Llama-3.1-8B on including human-written queries.



Under-confident    Over-confident    Logit\_based    Generation\_based

Figure 20: Under- and Over-confidence results on Llama-3.1-8B on including human-written queries.

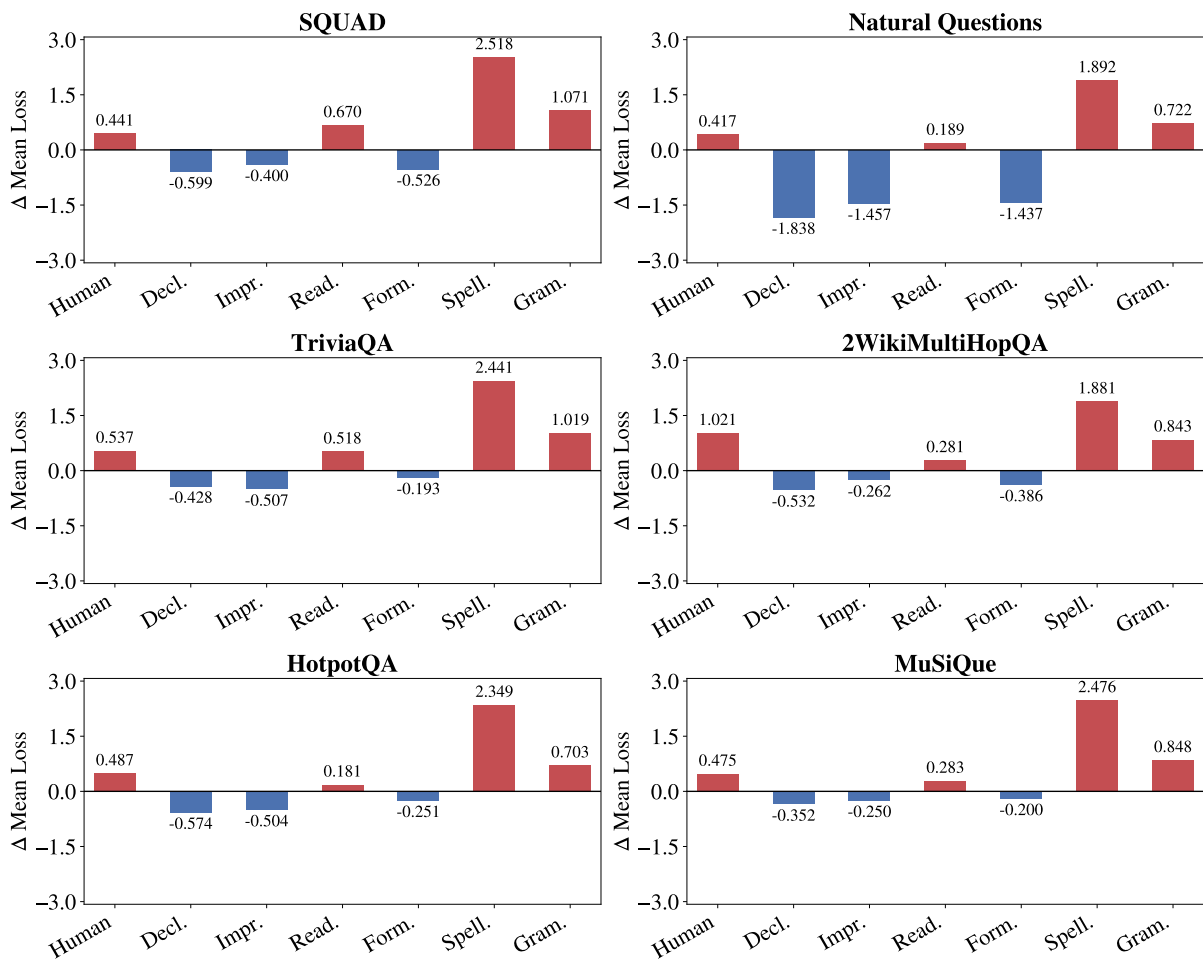


Figure 21: Loss score across all datasets

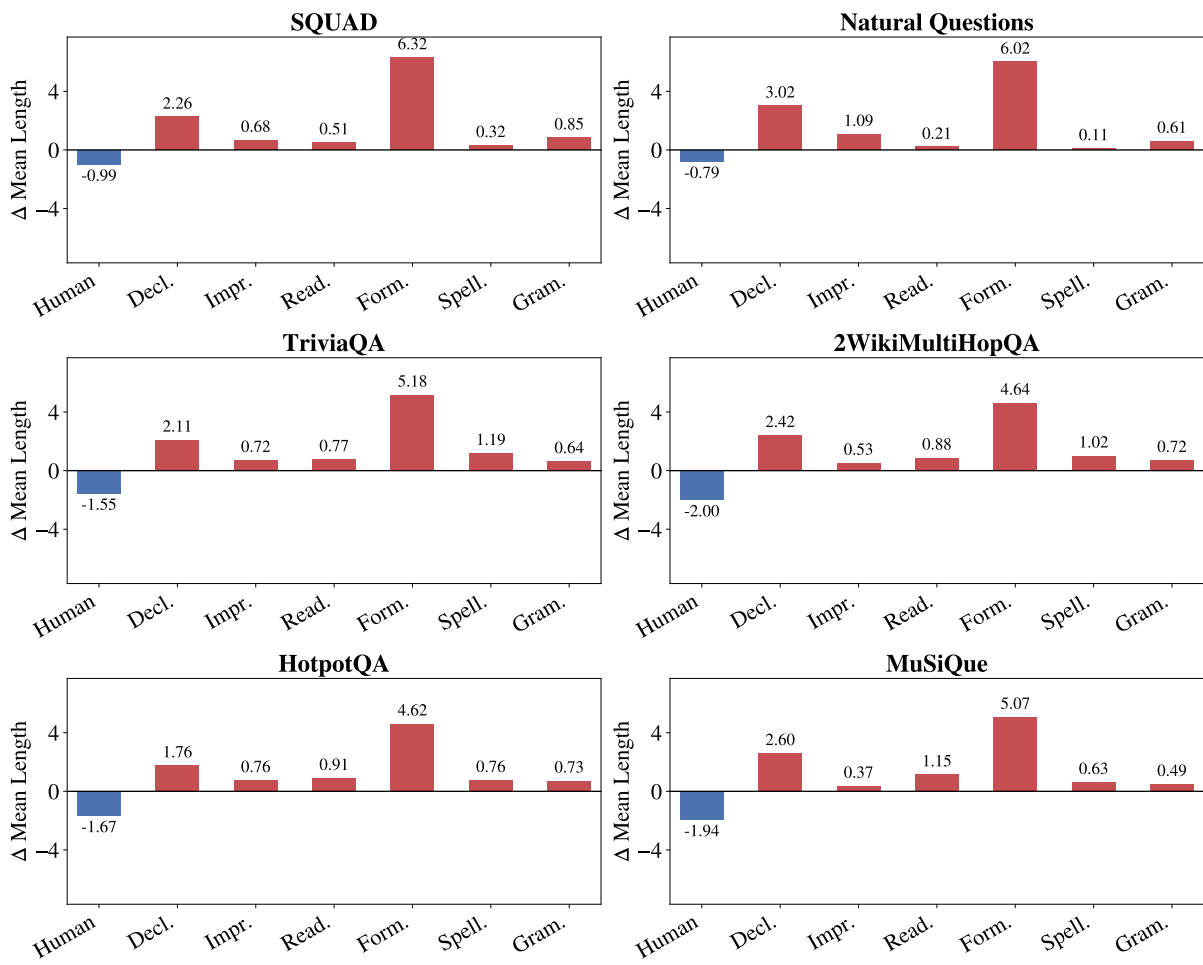


Figure 22: Query length across all datasets

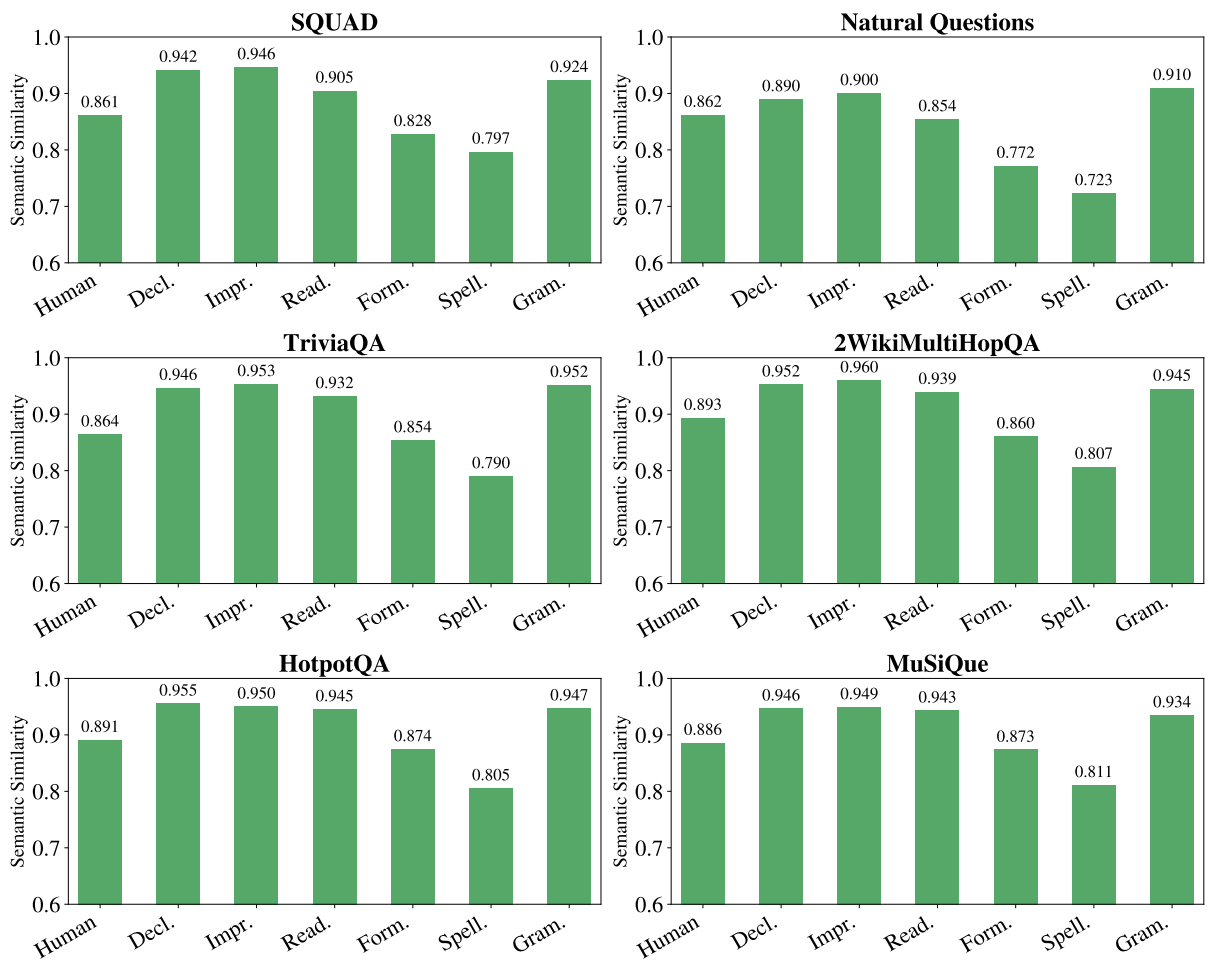


Figure 23: Semantic similarity to original query across all datasets