

CCRS: A Zero-Shot LLM-as-a-Judge Framework for Comprehensive RAG Evaluation

Aashiq Muhamed
 Carnegie Mellon University
 Pittsburgh, Pennsylvania, USA
 amuhamed@andrew.cmu.edu

Abstract

Retrieval-Augmented Generation (RAG) systems enhance Large Language Models (LLMs) by incorporating external knowledge, which is crucial for domains that demand factual accuracy and up-to-date information. However, evaluating the multifaceted quality of RAG outputs – spanning aspects such as contextual coherence, query relevance, factual correctness, and informational completeness—poses significant challenges. Existing evaluation methods often rely on simple lexical overlap metrics, which are inadequate for capturing these nuances, or involve complex multi-stage pipelines with intermediate steps like claim extraction or require finetuning specialized judge models, hindering practical efficiency. To address these limitations, we propose CCRS (Contextual Coherence and Relevance Score), a novel suite of five metrics that utilizes a single, powerful, pretrained LLM (Llama 70B) as a zero-shot, end-to-end judge. CCRS evaluates: Contextual Coherence (CC), Question Relevance (QR), Information Density (ID), Answer Correctness (AC), and Information Recall (IR). We apply CCRS to evaluate six diverse RAG system configurations (varying retrievers and readers) on the challenging BioASQ biomedical question-answering dataset. Our analysis demonstrates that CCRS effectively discriminates between system performances, confirming, for instance, that the Mistral-7B reader outperforms Llama variants and that the E5 neural retriever enhances QR and IR for Llama models in this task. We provide a detailed analysis of CCRS metric properties, including score distributions, convergent/discriminant validity, tie rates, population statistics, and discriminative power, finding QR highly discriminative overall, while observing a strong AC-IR correlation. Compared to the complex RAGChecker framework, CCRS offers comparable or superior discriminative power for key aspects like recall and faithfulness, while being significantly more computationally efficient. CCRS thus provides a practical, comprehensive, and efficient framework for evaluating and iteratively improving RAG systems.

CCS Concepts

- **Information systems → Evaluation of retrieval systems; Question answering; Language models; Retrieval effectiveness.**

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

LLM4Eval@SIGIR '25, Padua, Italy

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
 ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

Keywords

Retrieval-Augmented Generation, RAG Evaluation, LLM-as-a-judge, Evaluation Metrics, Contextual Coherence, Answer Correctness, Information Retrieval

ACM Reference Format:

Aashiq Muhamed. 2025. CCRS: A Zero-Shot LLM-as-a-Judge Framework for Comprehensive RAG Evaluation. In *Proceedings of the Third Workshop on Large Language Models for Evaluation in Information Retrieval (LLM4Eval@SIGIR '25), July 13–17, 2025, Padua, Italy*. ACM, New York, NY, USA, 25 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 Introduction

RAG systems [1, 6] represent a significant advancement in Large Language Model (LLM) technology by integrating external knowledge retrieval with generative capabilities. This synergy aims to produce responses that are more factual, current, and contextually appropriate than those generated by LLMs relying solely on their internal parametric knowledge. The ability to ground generated content in external evidence is particularly vital in specialized domains such as medicine, finance, and law, where precision, reliability, and grounding in specific, often dynamic, knowledge bases are paramount [24]. However, the effectiveness of RAG systems hinges on a complex interplay between the quality of the retrieved information (relevance, accuracy, completeness) and the LLM’s ability to comprehend, synthesize, and faithfully generate responses based on that information. This inherent complexity poses significant evaluation challenges.

Failures in RAG systems can manifest in various ways, hindering their trustworthiness and utility. These failures include the retrieval of irrelevant or contradictory documents, the LLM’s inability to discern correct information amidst noisy or conflicting context, unresolved conflicts between retrieved information and the LLM’s internal parametric knowledge, or the generation of fluent text that nonetheless misrepresents, omits crucial details from, or hallucinates information beyond the provided context [7, 9, 10]. Consequently, evaluating RAG systems requires a multifaceted approach that moves beyond assessing simple output quality. It must consider multiple critical dimensions, including faithfulness to the provided context, relevance to the user’s query, factual accuracy against ground truth, completeness of information, logical coherence, and conciseness [4, 18].

Traditional NLP metrics, such as BLEU [15], ROUGE [11], and even more advanced semantic similarity measures like BERTScore [22], primarily rely on comparing the generated output to one or more static reference texts. These metrics are often insufficient for RAG evaluation because they fail to adequately capture critical aspects like whether the generated text is factually supported by the

retrieved context—a core requirement for trustworthy RAG systems [4, 18]. They cannot reliably penalize plausible-sounding but unfaithful or factually inaccurate generation, nor can they effectively assess the quality of the retrieval process itself.

The advent of increasingly powerful LLMs has spurred interest in using them as evaluators, often referred to as the *LLM-as-a-judge* paradigm [2, 23]. This approach offers a potential solution for assessing the nuanced quality dimensions required for RAG evaluation. Several RAG-specific evaluation frameworks leveraging this paradigm have been proposed. RAGAS [4] introduces metrics like faithfulness and answer relevance, but its faithfulness calculation involves a multi-step pipeline: first decomposing the generated answer into individual statements and then performing Natural Language Inference (NLI) checks for each statement against the retrieved context. ARES [18] utilizes finetuned smaller LLMs as judges to predict scores for context relevance, answer faithfulness, and answer relevance. While potentially accurate due to supervised training, this approach demands considerable effort in generating large synthetic training datasets and obtaining human validation data for calibration, often using complex statistical techniques like Prediction-Powered Inference (PPI). CRUD-RAG [13] focuses on completeness and accuracy by using an LLM to first generate questions from the ground truth and then assessing the RAG system’s response based on its ability to answer these generated questions. RAGChecker [17] offers highly fine-grained analysis by extracting claims from the response, ground truth, and context, followed by pairwise entailment checking. This allows for detailed metrics like precision, recall, faithfulness, hallucination rate, and noise sensitivity, providing comprehensive insights but at the cost of significant computational complexity and sensitivity to potential errors in the intermediate claim extraction and entailment steps.

While these frameworks represent significant progress, their reliance on complex multi-stage pipelines (RAGAS, RAGChecker, CRUD-RAG) or extensive setup requirements involving data generation and fine-tuning (ARES) can limit their practical applicability, especially for rapid iteration during development or for large-scale evaluations across many system variants. There remains a need for a robust, comprehensive, yet efficient RAG evaluation method.

To address this gap, we introduce CCRS (Contextual Coherence and Relevance Score), a novel suite of five evaluation metrics. CCRS employs a single, powerful, pretrained LLM (specifically, Llama 70B-Instruct [19]) as a zero-shot judge. It directly evaluates the end-to-end RAG output across five key dimensions without requiring intermediate processing steps (like claim extraction or question generation) or specialized fine-tuning. The CCRS metrics are designed to capture critical aspects of RAG quality:

- (1) **Contextual Coherence (CC)**: Assesses the logical consistency and flow of the response with respect to the provided context.
- (2) **Question Relevance (QR)**: Evaluates the appropriateness and directness of the response in addressing the user’s query.
- (3) **Information Density (ID)**: Measures the balance between conciseness and informativeness in the response.
- (4) **Answer Correctness (AC)**: Determines the factual accuracy of the response compared to a ground truth answer, considering the context.

- (5) **Information Recall (IR)**: Assesses the completeness of the response in capturing essential information present in the ground truth answer.

This framework aims to harness the nuanced understanding capabilities of large LLMs in an efficient, end-to-end manner, providing a practical tool for RAG evaluation. Our main contributions are:

- (1) We propose and define CCRS, a zero-shot, end-to-end LLM-as-a-judge framework with five distinct metrics designed for practical and comprehensive RAG evaluation (subsection 2.3).
- (2) We conduct a comprehensive empirical evaluation using the challenging BioASQ biomedical question-answering dataset, comparing six diverse RAG configurations (varying retrievers and readers) and rigorously testing hypotheses about the impact of system components (subsection 4.2).
- (3) We perform a detailed analysis of the properties of the CCRS metrics, including their distributions, validity (convergent and discriminant), tie rates, population statistics, and discriminative power. We also compare their effectiveness and computational efficiency against the complex, state-of-the-art RAGChecker framework (subsection 4.1).

This work contributes a practical evaluation methodology and valuable empirical insights into RAG system performance.

2 Related Work and Framework

2.1 Limitations of Traditional Metrics

Evaluating the output quality of generative systems like RAG requires moving beyond metrics designed for tasks like machine translation or summarization, which primarily assess surface-level similarity to reference texts. Traditional metrics such as BLEU [15], which measures n-gram precision against references (often calculated as $\text{BLEU} = \text{BP} \cdot \exp(\sum_{n=1}^N w_n \log p_n)$, where BP is a brevity penalty), and ROUGE [11], particularly ROUGE-L which uses the Longest Common Subsequence (LCS) ($F_L = \frac{(1+\beta^2)R_L P_L}{\beta^2 R_L + P_L}$), focus on lexical overlap. While useful for assessing fluency and lexical similarity, they fail to capture deeper semantic meaning or factual correctness.

Metrics like BERTScore [22] represent an improvement by leveraging contextual embeddings ($\sum_{x_i \in \text{gen}} \max_{y_j \in \text{ref}} \text{sim}(x_i, y_j)$) to compute semantic similarity between generated and reference texts. However, even these advanced semantic metrics have significant limitations in the context of RAG evaluation [4, 12]. They cannot directly evaluate the faithfulness of the generated response to the *retrieved context*, nor can they assess factual accuracy independently of lexical or semantic similarity to a potentially incomplete or imperfect reference answer. They also do not penalize plausible but unfaithful hallucinations or factually incorrect statements that might still share some semantic similarity with the reference.

2.2 LLM-based Evaluation Frameworks

Recognizing the limitations of traditional metrics, recent research has increasingly focused on leveraging the sophisticated language understanding capabilities of LLMs themselves to perform evaluation, often termed *LLM-as-a-judge* [2, 23]. Several frameworks

specifically designed for RAG evaluation have emerged from this line of work:

- **TruLens** [5]: Proposes a conceptual *RAG Triad* consisting of groundedness (faithfulness to context), answer relevance (to the query), and context relevance (of retrieved documents to the query), which are assessed using LLM prompts. While conceptually aligned with key RAG quality dimensions, the framework remains somewhat underspecified in its published form and lacks extensive validation studies demonstrating its reliability and effectiveness compared to other methods.
- **RAGAS** [4]: Introduces several metrics, including *faithfulness*, *answer relevancy*, and *context relevancy*. Its faithfulness calculation requires a multi-step pipeline: an LLM first extracts discrete claims from the generated response (r), and then another LLM (or the same one) performs Natural Language Inference (NLI) checks for each claim against the retrieved context (C). The final score is the proportion of verified claims ($\text{Faithfulness} = \frac{|\{s_i \in \text{claims}(r) | \text{verify}(s_i, C)\}|}{|\text{claims}(r)|}$). This reliance on intermediate claim extraction and verification steps can make the process computationally demanding and potentially brittle if the intermediate steps fail.
- **ARES** [18]: Employs finetuned smaller LLMs (e.g., DeBERTa variants) as specialized judges to predict scores for context relevance, answer faithfulness, and answer relevance. Its potential strength lies in achieving high accuracy through supervised training. However, this comes at the cost of significant upfront investment: creating large-scale synthetic datasets (often using powerful generator models like FLAN-T5 XXL) to train the judges, and obtaining human validation data for calibration using sophisticated statistical techniques like Prediction-Powered Inference (PPI). This setup complexity can be a barrier to adoption.
- **CRUD-RAG** [13]: Focuses specifically on evaluating the completeness and accuracy of the RAG response (r) relative to a ground truth answer (g). It uses an LLM to first generate a set of questions $Q(g)$ based on the ground truth. It then assesses the RAG response by checking if it can answer these generated questions ($\text{Recall} = |\{q \in Q(g) | A_r(q) \neq \emptyset\}| / |Q(g)|$) and measures the accuracy of the answers extracted from r compared to those derived from g ($\text{Precision} \approx \text{avg}(\text{Sim}(A_r(q), A_g(q)))$). This approach introduces dependencies on the quality and coverage of the intermediate question generation and answer extraction steps.
- **RAGChecker** [17]: Provides arguably the most fine-grained analysis through claim-level checking. It involves extracting claims from the model’s response (m), the ground truth (gt), and the retrieved context chunks ($\{chunk_j\}$). It then performs pairwise entailment checks between these sets of claims. This enables the computation of detailed metrics such as Precision ($\frac{|\{c_i^{(m)} | c_i^{(m)} \in gt\}|}{|\{c_i^{(m)}\}|}$), Recall ($\frac{|\{c_i^{(gt)} | c_i^{(gt)} \in m\}|}{|\{c_i^{(gt)}\}|}$), Faithfulness ($\frac{|\{c_i^{(m)} | c_i^{(m)} \in \{chunk_j\}\}|}{|\{c_i^{(m)}\}|}$), Hallucination Rate, Noise Sensitivity, and Context Utilization. While extremely comprehensive, the multi-step process involving claim extraction and numerous entailment checks makes RAGChecker computationally intensive and potentially sensitive to errors propagating from the intermediate steps.

2.3 The CCRS Framework

CCRS (Contextual Coherence and Relevance Score) is designed to provide a comprehensive yet efficient evaluation of RAG systems, balancing the depth of multi-dimensional assessment with the practicality of a streamlined workflow. Unlike the related works discussed, CCRS seeks an effective balance between the simplicity of direct LLM prompting (akin to aspects of TruLens or general reference-free LLM evaluation) and the multi-dimensional assessment capability offered by more complex frameworks like RAGAS, ARES, CRUD-RAG, and RAGChecker. The core hypothesis underlying CCRS is that a single, sufficiently powerful, pretrained LLM (in our case, Llama 70B-Instruct) can perform nuanced judgments across multiple quality dimensions in a zero-shot, end-to-end manner, directly evaluating the final RAG output without needing intermediate processing steps or specialized training.

By avoiding the pipeline complexity inherent in claim extraction/verification (RAGAS, RAGChecker) or question generation/answering (CRUD-RAG), and bypassing the extensive training data generation and calibration overhead required by supervised judge models (ARES), CCRS aims to offer a practical and computationally efficient solution suitable for rapid development cycles and large-scale comparative evaluations.

2.3.1 Task Definition & Constructs. We focus on evaluating RAG systems designed for Question Answering (QA). The input to the RAG system is $x = (q, D)$, where q is the user’s question and D is the document collection. The output is $y = (C, r)$, where $C \subseteq D$ is the set of retrieved context passages and r is the generated response. The evaluation task is to assess the quality of the generated response r , given the input question q , the retrieved context C , and a ground truth answer g .

Our primary goal is to develop metrics that serve as effective proxies for **user satisfaction** with the RAG system’s output. We operationalize this goal by measuring five key quality constructs, defined in Table 1. These constructs were chosen to cover essential aspects determining whether a RAG response is helpful, reliable, and addresses the user’s information need effectively.

2.3.2 CCRS Components and Metric Calculation. CCRS utilizes Llama-70B-Instruct as the LLM_{Judge} to perform zero-shot evaluations for each of the five constructs. Specific prompts, detailed in Appendix B, are used to guide the LLM judge. The judge is instructed to output a score between 0 and 100 for each dimension, which we then normalize to a $[0, 1]$ range for consistency.

Contextual Coherence (CC): Calculated as $CC(r, C) = LLM_{\text{Judge}}(r, C, \text{prompt}_{CC}) / 100$. This metric directly assesses whether the generated text r logically follows from and avoids contradicting the provided context C . A high CC score indicates that the response is well-grounded in the evidence and integrates information coherently. The corresponding prompt (Appendix B) explicitly asks the judge to evaluate logical consistency and coherence.

Question Relevance (QR): Calculated as $QR(r, q) = LLM_{\text{Judge}}(r, q, \text{prompt}_{QR}) / 100$. This metric directly evaluates if the response r effectively answers the specific question q posed by the user. High relevance is fundamental to the perceived usefulness of the response. The prompt (Appendix B) asks the judge *how well the response addresses the user’s query*.

Table 1: CCRS Construct Definitions and Motivations

Construct	Definition & Motivation
Contextual Coherence (CC)	Logical consistency and flow between the response (r) and the retrieved context (C). (<i>Motivation</i> : Ensures the response is well-supported by the evidence and internally consistent with it, crucial for groundedness.)
Question Relevance (QR)	Appropriateness and directness of the response (r) to the user's query (q). (<i>Motivation</i> : Ensures the response actually addresses the user's specific information need, fundamental to usefulness.)
Information Density (ID)	Optimal balance between informativeness and conciseness in the response (r). (<i>Motivation</i> : Avoids overly verbose answers that overwhelm the user or overly brief answers that lack necessary detail, impacting user experience.)
Answer Correctness (AC)	Factual accuracy of claims in the response (r) compared to the ground truth answer (g). (<i>Motivation</i> : Critical for trust and reliability, especially in high-stakes domains like medicine or finance.)
Information Recall (IR)	Coverage of essential information from the ground truth answer (g) within the response (r). (<i>Motivation</i> : Ensures the answer is sufficiently complete and provides necessary details, complementing correctness.)

Information Density (ID): Calculated as $ID(r, C, q) = LLM_{\text{Judge}}(r, C, q, \text{prompt}_{\text{ID}})/100$. This assesses the communication efficiency of the response r . It's important for user experience, aiming to avoid information overload from excessive verbosity or frustration from underspecification. The prompt (Appendix B) asks the judge to assess the *balance of conciseness and informativeness*.

Answer Correctness (AC): Calculated as $AC(r, g, C) = (\lambda \cdot EM(r, g) + (1 - \lambda) \cdot LLM_{\text{Judge}}(r, g, C, \text{prompt}_{\text{AC}})/100)$. Measures the factual alignment of the response r with the ground truth answer g , which is crucial for reliability. We combine a strict exact match (EM) check ($EM(r, g) = 1$ if $r = g$, else 0) with the LLM's judgment of semantic correctness, allowing for paraphrasing while considering the retrieved context C . We use a weight $\lambda = 0.7$ to emphasize exact matches, reflecting the high precision often required in domains like biomedical QA, while still incorporating semantic evaluation. The prompt (Appendix B) asks for *factual accuracy... compared to the ground truth answer, considering the context*.

Information Recall (IR): Calculated as $IR(r, g, C) = LLM_{\text{Judge}}(r, g, C, \text{prompt}_{\text{IR}})/100$. Measures the completeness of the response r in covering the essential information present in the ground truth answer g . This complements AC by ensuring that key details required for a full answer are not omitted. The prompt (Appendix B) asks *how much of the essential information from the ground truth is captured*.

By combining these five metrics, CCRS provides a multi-faceted evaluation of RAG system output quality in an efficient, zero-shot manner.

3 Experimental Setup

We designed our experiments to rigorously evaluate the CCRS framework and compare the performance of different RAG system configurations on a challenging task.

3.1 Dataset: BioASQ

We utilized the BioASQ dataset [20], a widely recognized benchmark for biomedical question answering. We used a publicly available version containing N=4,719 expert-curated question-answer pairs along with associated PubMed passages¹. BioASQ questions are formulated by domain experts and cover a range of types (factoid, list, summary, yes/no) and complexities, requiring both accurate retrieval from biomedical literature and precise generation. Its domain specificity and complexity make it a suitable and challenging testbed for evaluating RAG systems and the metrics designed to assess them. Example data points can be found in Appendix A.

3.2 Systems Compared

To assess CCRS's ability to differentiate between systems and to investigate the impact of RAG components, we evaluated six distinct RAG system configurations. These configurations were created by combining two different retrieval methods with three different reader LLMs:

Retrievers:

- **BM25:** A traditional sparse retrieval algorithm based on lexical matching [16], implemented using OpenSearch.
- **E5-Mistral:** A state-of-the-art neural dense retriever [21], implemented using OpenSearch's approximate k-NN search capabilities.

Readers (Generators):

- **Mistral-7B:** Mistral-7B-Instruct-v0.3 [8].
- **Llama3-8B:** Meta-Llama-3-8B-Instruct [14].
- **Llama3.2-3B:** Meta-Llama-3.2-3B-Instruct [3].

This resulted in the following six system configurations, labeled A through F:

- **System A:** BM25 Retriever + Mistral-7B Reader
- **System B:** E5-Mistral Retriever + Mistral-7B Reader
- **System C:** BM25 Retriever + Llama3-8B Reader
- **System D:** E5-Mistral Retriever + Llama3-8B Reader
- **System E:** BM25 Retriever + Llama3.2-3B Reader
- **System F:** E5-Mistral Retriever + Llama3.2-3B Reader

CCRS Judge Model: For all CCRS metric calculations, we used Meta-Llama-3-70B-Instruct as the LLM_{Judge} .

3.3 Implementation Details

All experiments were executed on institutional compute clusters, utilizing nodes equipped with 8x A100 80GB GPUs for parallel model inference.

Document Processing and Retrieval: Documents from the BioASQ corpus were chunked into segments of 300 tokens with an overlap of 20% (60 tokens). We consistently used the tokenizer associated with the E5-Mistral model for chunking across both BM25 and E5 experiments to ensure comparability. For each query, the top-k=20 chunks as ranked by the respective retriever (BM25

¹<https://huggingface.co/datasets/rag-datasets/rag-mini-bioasq>

or E5-Mistral) were retrieved and concatenated to form the context provided to the reader LLM.

Response Generation: The reader LLMs generated responses based on the input query and the retrieved context ($k=20$ chunks). To ensure deterministic and comparable outputs, we used a temperature setting of 0.0 for generation. The maximum length for generated responses was set to 2048 tokens. The prompt template used for instructing the reader models is shown in Figure 1.

```
Please answer the given question based on the context. <context>
<content>chunk_1</content>    <content>chunk_2</content> ...
<content>chunk_k</content> </context>
Question: {question}
Please answer the question and tag your answer with <answer></answer>.
```

Figure 1: Default prompt template used for RAG response generation by the reader models (Mistral-7B, Llama3-8B, Llama3.2-3B).

Evaluation: CCRS metrics were computed for each generated response using the Llama-70B-Instruct judge model. For the comparison with RAGChecker (subsubsection 4.1.4), we also used Llama-70B-Instruct for the claim extraction step to maintain consistency in the LLM used, although the original RAGChecker paper used GPT-4o. Statistical analyses were performed using custom Python scripts leveraging libraries like NumPy and SciPy, with the specific implementation detailed in Appendix F.

4 Results and Analysis

In this section, we present the results of our experiments. We first analyze the properties of the CCRS metrics themselves, including their distributions, validity, tie rates, and discriminative power, along with a comparison to the RAGChecker framework. We then evaluate the performance of the six RAG systems based on the CCRS metrics and test our predefined hypotheses.

4.1 CCRS Metric Properties

Understanding the behavior and characteristics of the CCRS metrics is crucial for interpreting the RAG system evaluation results.

4.1.1 Score Distributions and Boundedness. We analyzed the distributions of scores for each CCRS metric across all 4,719 queries and 6 systems. The histograms of the mean scores per query (averaged over the 6 systems) are shown in Figure 2. Detailed per-system histograms are available in Appendix Figures 6-10.

Most metrics exhibit a negative skew (left tail), with scores tending towards the upper end of the scale (1.0). This suggests that, on average, the evaluated RAG systems perform reasonably well across these dimensions on the BioASQ dataset.

Question Relevance (QR) shows a particularly pronounced negative skew and a strong ceiling effect. As detailed in Table 5, between 65% and 85% of responses received a perfect score of 1.0 for QR, depending on the system. This indicates that QR is effective at identifying poorly relevant answers (assigning low scores) but offers limited granularity for differentiating among highly relevant responses.

Answer Correctness (AC) scores are notably low, with a mean around 0.2 across systems (Table 4). Crucially, AC never reached the

maximum score of 1.0 in our observations (Table 6). This is partly due to the $\lambda = 0.7$ weight given to the strict Exact Match component (meaning perfect semantic match without exact match yields max 0.3) and partly reflects the inherent difficulty of achieving perfect factual accuracy in the complex biomedical domain.

Information Recall (IR) displays the widest and most symmetric distribution among the metrics (Figure 2(e)). This suggests that IR effectively captures a broad range of completeness levels in the generated answers, making it potentially useful for identifying systems that excel or struggle with providing comprehensive information.

Contextual Coherence (CC) shows distributions whose shapes are highly dependent on the reader model (Figure 6 in Appendix). Systems using Llama readers (C, D, E, F) exhibit a higher proportion of scores at the lower bound (0.0) compared to Mistral-based systems (A, B), suggesting more frequent coherence failures (hallucinations or contradictions relative to context) in the Llama models tested. This highlights CC's ability to capture differences in generation quality related to grounding.

The observed bounds (Table 5, Table 6 in Appendix) confirm these trends, showing high frequencies of perfect scores for QR and CC (especially for Mistral), zero perfect scores for AC, and a wider spread for IR.

4.1.2 Validity (Convergent & Discriminant). To assess whether the CCRS metrics measure distinct constructs (discriminant validity), and whether related metrics correlate as expected (convergent validity), we computed Pearson, Spearman, and Kendall correlations between all pairs of metrics. The averaged Pearson correlations (using Fisher's Z-transform) across all six systems are shown in Figure 3(a). Detailed per-system correlation matrices are in subsection C.2.

We observe strong **convergent validity** between Answer Correctness (AC) and Information Recall (IR), with an average Pearson correlation of $r=0.756$. This strong positive relationship is expected, as responses that are more complete (higher IR) are often also more likely to contain the correct facts (higher AC), and vice-versa.

We find evidence for **discriminant validity**, particularly for Contextual Coherence (CC). CC shows only weak to moderate correlations with the other metrics (average Pearson r ranging from 0.201 with AC to 0.452 with IR). This suggests that CC captures a distinct aspect of quality—logical consistency with the context—that is not strongly captured by metrics focusing on relevance (QR), density (ID), accuracy (AC), or completeness (IR).

Question Relevance (QR) and Information Density (ID) show a moderate positive correlation ($r=0.494$), suggesting some relationship: highly relevant answers might also tend to be appropriately dense, or judges might implicitly link these qualities. However, the correlation is not strong enough to suggest they measure the same construct entirely.

Per-system analyses (subsection C.2) reveal some architectural influences. For instance, the correlation between CC and IR is stronger for Llama-based systems (e.g., $r=0.528$ for System C) than for Mistral-based systems (e.g., $r=0.377$ for System B), suggesting coherence and recall might be more tightly coupled in Llama architectures.

4.1.3 Discriminative Power (DP) and Ties. An important property of an evaluation metric is its ability to reliably distinguish between

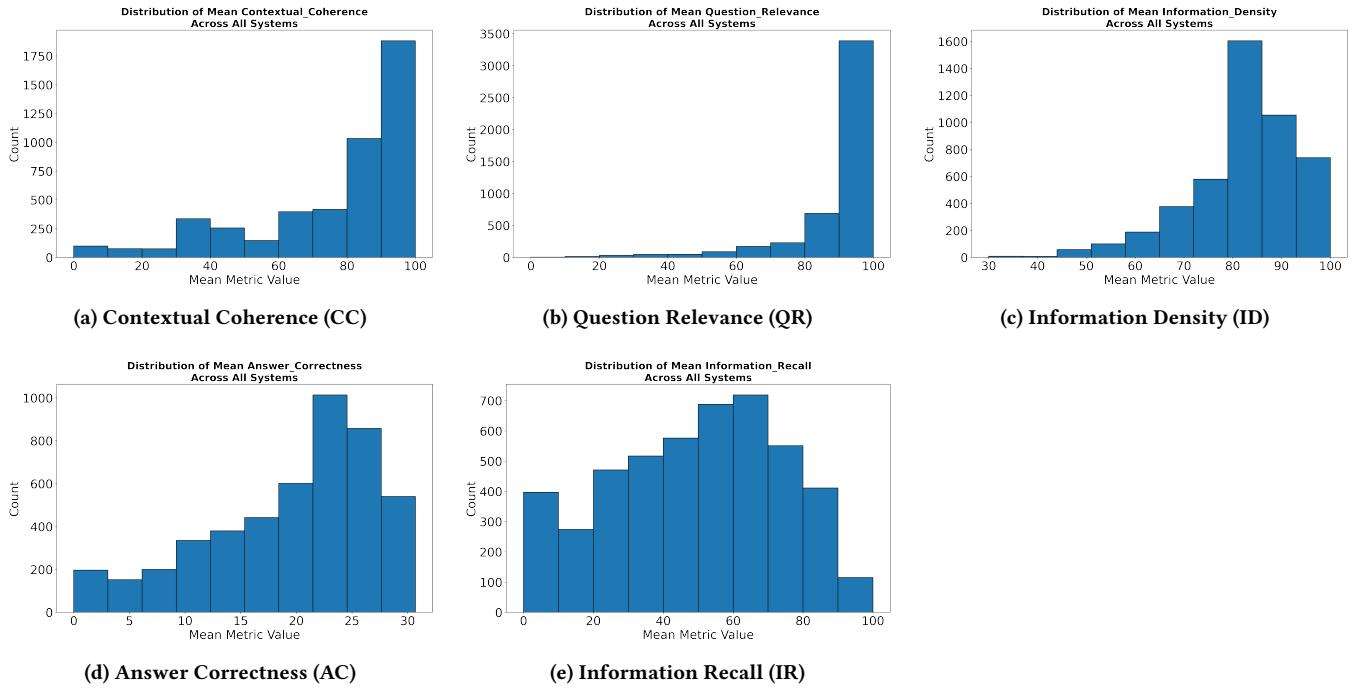


Figure 2: Distribution of Mean CCRS Metric Values Across All Systems (N=4719 queries). These histograms show the frequency (Y-axis) of mean scores (X-axis, range 0-1, averaged across the 6 systems for each query) for each of the five CCRS metrics (a-e), illustrating overall scoring trends.

systems that perform differently (discriminative power). We assessed DP using pairwise statistical tests (Tukey's HSD with randomization, B=10,000, $\alpha = 0.05$) between all $\binom{6}{2} = 15$ pairs of systems for each metric. DP is calculated as the fraction of pairs for which a statistically significant difference was found. Results are summarized in Table 2, and visualized using Achieved Significance Level (ASL) curves in Figure 3(b).

Question Relevance (QR) exhibits the highest DP (0.9333), detecting significant differences in 14 out of 15 system pairs. This indicates high sensitivity to performance variations in relevance. **Information Recall (IR)** follows closely with DP=0.8667 (13/15 pairs), and **Answer Correctness (AC)** also shows good DP=0.8000 (12/15 pairs). **Contextual Coherence (CC)** and **Information Density (ID)** have lower DP (both 0.7333, 11/15 pairs), suggesting they are slightly less sensitive in distinguishing between the systems evaluated here.

However, DP must be considered alongside the rate of ties (Table 13). QR, despite its high DP, suffers from a very high empirical tie rate (ranging from 47% to 74% depending on the system) due to the ceiling effect discussed earlier. This means that while QR effectively separates systems with different average relevance, it offers poor granularity for distinguishing between individual responses that are all deemed highly relevant.

Conversely, AC and IR have the lowest tie rates (AC: 15-18%, IR: 15-18%). Their combination of good DP and low ties suggests they offer better potential for fine-grained differentiation between RAG outputs, capturing variations in accuracy and completeness more

granularly than QR captures variations in relevance at the top end of the scale.

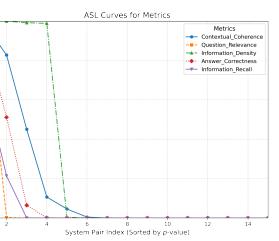
Averaged Pearson Correlations Across Systems

Using Fisher's Z-transformation for averaging (N=6 systems)

	CC	QR	ID	AC	IR
CC	1.000	0.305	0.429	0.217	0.464
QR	0.305	1.000	0.582	0.439	0.475
ID	0.429	0.582	1.000	0.335	0.370
AC	0.217	0.439	0.335	1.000	0.761
IR	0.464	0.475	0.370	0.761	1.000

Method: $z = 0.5 * \ln((1+r)/(1-r))$, averaged, then back-transformed

(a) Avg. Pearson Correlation



(b) ASL Curves (DP)

Figure 3: CCRS Metric Validity (left) assessed via averaged Pearson correlations across 6 systems, and Discriminative Power (right) visualized with Achieved Significance Level (ASL) curves.

4.1.4 Comparison with RAGChecker. To benchmark CCRS against a state-of-the-art but complex RAG evaluation framework, we compared it to RAGChecker [17]. We focused on RAGChecker's core end-to-end metrics: Precision (P), Recall (R), and Faithfulness (F), which rely on claim extraction and entailment checking. We used

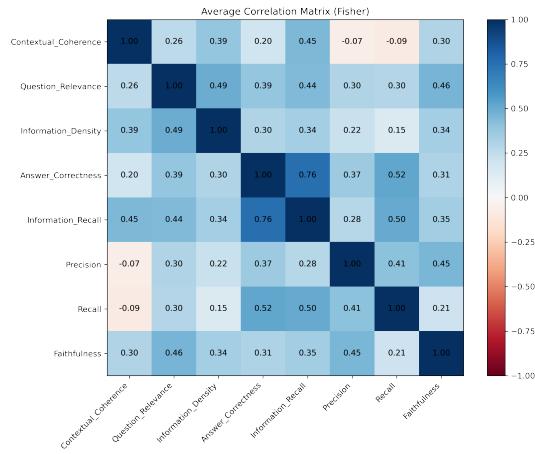
Table 2: Discriminative Power (DP) of CCRS Metrics based on pairwise Tukey's HSD tests ($\alpha = 0.05$) across 6 systems (15 pairs).

Metric	Significant Pairs	DP
Question Relevance (QR)	14 / 15	0.9333
Information Recall (IR)	13 / 15	0.8667
Answer Correctness (AC)	12 / 15	0.8000
Contextual Coherence (CC)	11 / 15	0.7333
Information Density (ID)	11 / 15	0.7333

Llama-70B-Instruct for claim extraction in RAGChecker to maintain consistency with the CCRS judge model. Full details, including distributions, bounds, ties, and population stats for RAGChecker metrics, are in Appendix D.

Distributions/Bounds/Ties: RAGChecker Faithfulness (F) exhibited a strong ceiling effect, similar to CCRS QR, with high mean scores, 64-74% of responses achieving a perfect score of 1.0 (Table 17), and high tie rates (44-54%, Table 18). In contrast, RAGChecker Precision (P) and Recall (R) showed more balanced distributions (Table 19) and low tie rates (P: 13-20%, R: 9-10%), behaving similarly to CCRS AC and IR.

Validity: We examined the correlation between CCRS and RAGChecker metrics (Figure 4). RAGChecker Recall (R) showed moderate-to-strong correlations with CCRS AC ($r=0.518$) and CCRS IR ($r=0.504$), confirming they measure related aspects of correctness and completeness. RAGChecker Faithfulness (F) correlated most strongly with CCRS QR ($r=0.463$) and moderately with CCRS CC ($r=0.304$), suggesting overlap in assessing relevance and coherence or grounding. RAGChecker Precision (P) showed moderate correlation with CCRS AC ($r=0.366$) but weaker links to other CCRS metrics.

**Figure 4: Averaged Pearson Correlation Matrix comparing CCRS metrics (CC, QR, ID, AC, IR) and RAGChecker metrics (P, R, F) across 6 systems.**

Discriminative Power (DP): We compared the DP of key metrics (Table 3). CCRS QR (DP=0.933) and IR (DP=0.867) significantly outperformed RAGChecker’s Recall (R, DP=0.800) and Faithfulness (F, DP=0.800). CCRS AC (DP=0.800) matched the DP of RAGChecker

R and F. RAGChecker Precision (P) showed very poor discriminative power (DP=0.200) in our experiments.

Table 3: Discriminative Power (DP) Comparison: CCRS vs. RAGChecker (Core Metrics).

CCRS Metric	DP	RAGChecker Metric	DP
QR	0.9333	Faithfulness (F)	0.8000
IR	0.8667	Recall (R)	0.8000
AC	0.8000	Precision (P)	0.2000
CC	0.7333		
ID	0.7333		

Efficiency: A significant advantage of CCRS is its computational efficiency. In our setup, running the five zero-shot CCRS evaluations using Llama-70B was approximately **5 times faster** than executing the RAGChecker pipeline (claim extraction using Llama-70B + entailment checks) for P, R, and F. This difference stems from CCRS avoiding the computationally intensive claim extraction and multiple pairwise entailment steps required by RAGChecker.

Conclusion vs. RAGChecker: CCRS offers comparable or superior discriminative power for key aspects of RAG evaluation (especially relevance via QR, completeness via IR, and correctness via AC) compared to RAGChecker’s core metrics (F, R, P), while being significantly more computationally efficient and simpler to implement due to its zero-shot, end-to-end nature.

4.2 RAG System Performance and Hypothesis Testing

We used the CCRS framework to evaluate the six RAG system configurations (A-F). Mean performance scores are presented in Table 4, and performance distributions are visualized using box plots in Figure 5. We tested three pre-defined hypotheses using paired Tukey’s HSD tests (B=10,000 permutations) with Holm-Bonferroni correction across hypotheses ($\alpha = 0.05$). Detailed statistical test results are available in Appendix E. All three hypotheses were supported by the data.

Table 4: Mean Performance Scores (%) by System and CCRS Metric

Sys.	Config.	CC	QR	ID	AC	IR
A	Mistral+BM25	87.95	92.83	86.81	21.10	54.96
B	Mistral+E5	88.61	95.40	86.90	21.81	58.28
C	Llama8B+BM25	67.89	87.07	79.09	18.54	43.00
D	Llama8B+E5	69.29	90.55	80.41	19.24	45.44
E	Llama3.2+BM25	68.88	88.28	80.23	18.41	42.04
F	Llama3.2+E5	70.88	90.83	80.43	18.76	44.79

H1: Mistral-7B Reader Superiority (Adjusted $p < 0.0001$):

The results overwhelmingly support the hypothesis that Mistral-7B significantly outperforms both Llama3-8B and Llama3.2-3B readers across all five CCRS metrics, regardless of the retriever used. Comparing systems B (Mistral+E5) and D (Llama8B+E5), Mistral showed large advantages in CC (+19.3 points), QR (+4.9 points), ID (+6.5 points), AC (+2.6 points), and IR (+12.8 points). Similar significant advantages were observed when comparing A vs C, A vs E, and B vs F (see Table 20 to Table 24 for all pairwise p-values).

The box plots in Figure 5 visually confirm these findings, showing consistently higher median scores and often tighter interquartile ranges (IQRs) for Mistral-based systems (A, B) compared to Llama-based systems (C, D, E, F).

H2: E5 Retriever Advantage for Llama Models (Adjusted p < 0.0001)

The hypothesis that the E5 neural retriever provides benefits over BM25 specifically for QR and IR when used with Llama readers was strongly supported. Comparing System D (Llama8B+E5) vs. System C (Llama8B+BM25), E5 yielded significant improvements in QR (+3.48 points, $p < 0.0001$) and IR (+2.44 points, $p < 0.0001$).

Similarly, comparing System F (Llama3.2+E5) vs. System E (Llama3.2+BM25), E5 led to significant gains in QR (+2.55 points, $p < 0.0001$) and IR (+2.74 points, $p < 0.0001$). These results highlight the targeted benefits of neural retrieval for enhancing relevance and completeness, particularly with Llama models in this task.

H3: Performance Differences Between Llama Models (Final adj. p = 0.0070): The hypothesis that there exists a significant performance difference between the Llama3-8B and Llama3.2-3B readers was supported. The differences varied across metrics.

Llama3.2-3B showed significant advantages in QR (System E > C, $p=0.0037$), ID (System E > C, $p=0.0032$), and CC (System F > D, $p=0.0417$). Conversely, Llama3-8B demonstrated a significant advantage in AC (System D > F, $p=0.0014$). No significant difference was found for IR ($p=0.2171$ for C vs E, $p=0.6328$ for D vs F). This indicates a nuanced trade-off, with the smaller Llama3.2 potentially benefiting from higher-quality training data for coherence and relevance, while the larger Llama3-8B maintains an edge in factual accuracy.

5 Discussion

Our evaluation using the CCRS framework on the BioASQ dataset provides several important insights into RAG system performance and evaluation methodologies.

Interpreting RAG System Performance: The results strongly underscore the critical role of the generator (reader) LLM in determining overall RAG performance. Mistral-7B's consistent and statistically significant superiority across all five CCRS dimensions (supporting H1) suggests that its architecture or training provides substantial advantages in synthesizing information, maintaining coherence, and extracting relevant facts from retrieved context compared to the Llama 3 and 3.2 models tested. The performance gaps were particularly striking for Contextual Coherence (approx. +19-20 points) and Information Recall (approx. +10-13 points), indicating that the choice of reader model can have a profound impact on the quality and completeness of the generated output, potentially outweighing differences stemming from the retrieval component.

The targeted benefits observed for the E5 neural retriever (supporting H2) highlight the value of semantic retrieval. E5 significantly improved Question Relevance and Information Recall specifically for the Llama-based systems. This suggests E5 successfully retrieved context that was more topically aligned with the query and contained more of the necessary factual elements present in the ground truth. However, a crucial finding is that these improvements in retrieval quality (QR, IR) did not translate into statistically significant gains in Answer Correctness (AC) for the same Llama systems (D vs C, F vs E). This points towards a potential bottleneck at the

generation stage: providing more relevant or complete context does not automatically guarantee a more factually accurate generated response. The LLM's ability to faithfully extract, synthesize without distortion, and avoid introducing its own errors or hallucinations remains a distinct challenge, emphasizing the need for evaluation frameworks like CCRS that assess both faithfulness/coherence (CC) and factual accuracy (AC) independently.

The comparison between Llama3-8B and Llama3.2-3B (supporting H3) revealed a nuanced relationship between model scale, training data, and performance dimensions. The smaller Llama3.2-3B, potentially benefiting from higher-quality (possibly multimodal) training data, excelled in QR, ID, and CC. This suggests specialized training can enhance aspects related to query understanding, structured text generation, and conciseness, even in smaller models. Conversely, the larger Llama3-8B maintained an edge in AC, possibly due to better factual recall or more precise generation capabilities associated with its scale. This complex outcome suggests that optimizing RAG systems might involve trade-offs, requiring component selection based on which quality dimensions (e.g., coherence vs. factual precision) are most critical for the specific application.

Evaluating the Task and Dataset: The BioASQ dataset served as an effective and challenging benchmark, capable of discriminating clearly between the performance levels of different RAG systems. The overall performance profile observed—generally high scores for QR and CC, but significantly lower scores for AC and IR—underscores the inherent difficulty of high-fidelity question answering in specialized domains like biomedicine. Generating fluent and relevant text appears considerably less challenging for current models than ensuring absolute factual correctness and comprehensive coverage of information according to an expert standard. This highlights the importance of including metrics like AC and IR alongside QR and CC for a holistic evaluation, particularly in domains where reliability and completeness are paramount.

Reflecting on CCRS Metrics and Framework: CCRS proved to be a viable, efficient, and effective evaluation framework for RAG systems. Its multi-dimensional nature successfully captured distinct aspects of response quality, as validated by the correlation analysis (e.g., the strong AC-IR link indicating convergent validity, and CC's weak correlations with others suggesting discriminant validity). The discriminative power analysis identified QR and IR as highly sensitive metrics in our experimental setup, capable of detecting subtle system differences, while AC also performed well.

The comparison with the RAGChecker framework (subsubsection 4.1.4) showed that CCRS achieved comparable or superior discriminative power for core aspects of RAG quality (approximating recall via IR, correctness via AC, and aspects of faithfulness via CC/QR) compared to RAGChecker's key metrics (R, F, P). Crucially, CCRS achieved this with significantly reduced computational overhead (5x faster in our setup) and implementation complexity, as it avoids the need for intermediate claim extraction and entailment checking. This strongly supports the potential and practicality of using zero-shot, end-to-end LLM judgments for comprehensive RAG evaluation.

However, CCRS is not without limitations and challenges. The pronounced ceiling effect observed for QR, while indicating good performance on relevance, limits its granularity for differentiating between top-performing systems. The consistently low absolute

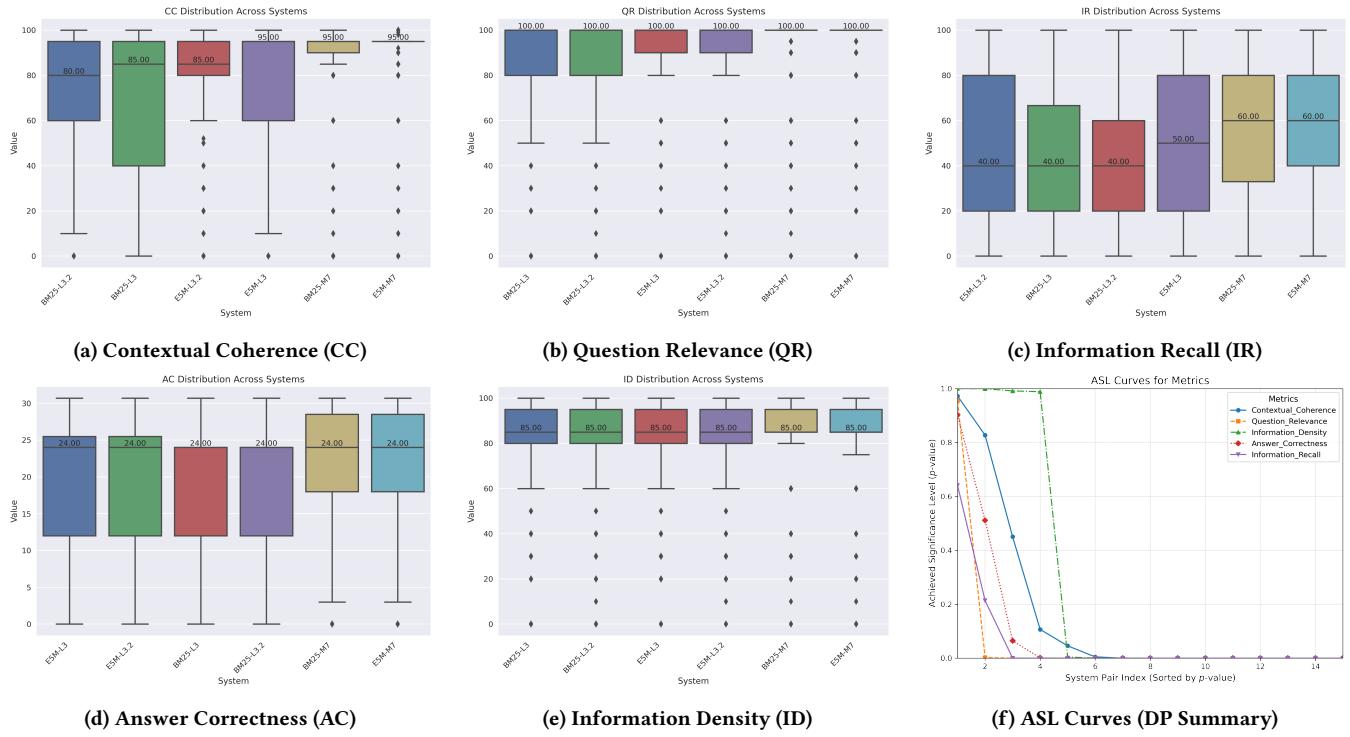


Figure 5: Box plots illustrating performance distributions for all CCRS metrics across the six RAG systems (A-F), plus ASL curves summarizing discriminative power. Systems within each plot are sorted by median performance. System Key: A(M+BM25), B(M+E5), C(L8+BM25), D(L8+E5), E(L3.2+BM25), F(L3.2+E5).

scores and absence of perfect scores for AC warrant further investigation into its calibration, the influence of the EM component, and the fundamental difficulty of exact factual matching in complex domains. The lower DP observed for CC and ID might reflect limitations in the judge LLM’s ability to consistently discern subtle differences in these more subjective qualities, or it might simply indicate smaller actual performance gaps between the systems on these dimensions. A significant limitation of the current study is the lack of direct human correlation data; the alignment between CCRS scores and human perception of quality remains to be formally established. Furthermore, interpretability remains a general challenge for all LLM-based metrics, including CCRS. Understanding *why* the judge assigned a particular score can be difficult.

Future Directions: Based on our findings and limitations, priority future directions include: (1) Validating CCRS across diverse domains and datasets to assess its generalizability. (2) Conducting rigorous human correlation studies to establish the alignment between CCRS metrics and human judgments of RAG quality. (3) Analyzing the sensitivity of CCRS scores to the choice of judge LLM and prompt variations. (4) Developing methods to enhance the interpretability of CCRS scores, potentially by prompting the judge for justifications alongside scores. (5) Exploring principled approaches to combine the five CCRS metrics into a single composite score, possibly weighted based on user studies or task requirements.

6 Conclusion

Evaluating the multifaceted quality of RAG systems remains a significant challenge, requiring methods that go beyond traditional metrics. To address the need for comprehensive yet efficient evaluation, we introduced CCRS, a novel suite of five metrics (Contextual Coherence, Question Relevance, Information Density, Answer Correctness, Information Recall). CCRS leverages a single, powerful, pretrained LLM (Llama 70B) as a zero-shot, end-to-end judge, thereby avoiding the complexities of multi-stage pipelines or extensive fine-tuning required by previous frameworks.

References

- [1] Akari Asai, Zexuan Zhong, Danqi Chen, Pang Wei Koh, Luke Zettlemoyer, Hanneh Hajishirzi, and Wen-tau Yih. 2024. Reliable, adaptable, and attributable language models with retrieval. *arXiv preprint arXiv:2403.03187* (2024).
- [2] Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937* (2023).
- [3] Abhishek Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
- [4] Shahul Es, Jithin James, Luis Espinosa-Anke, and Steven Schockaert. 2023. Ra-gas: Automated evaluation of retrieval augmented generation. *arXiv preprint arXiv:2309.15217* (2023).
- [5] Joe Ferrara, Ethan-Tonic, and Olguzhan Mete Ozturk. 2024. *The RAG Triad*. https://www.trulens.org/trulens_eval/core_concepts_rag_triad/.
- [6] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinlin Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* (2023).
- [7] Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2023. Large Language Models Cannot

- Self-Correct Reasoning Yet. arXiv:2310.01798 [cs.CL]
- [8] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lampe, Lucile Saulnier, et al. 2023. Mistral 7B. *arXiv preprint arXiv:2310.06825* (2023).
 - [9] Zhioran Jin, Pengfei Cao, Yubo Chen, Kang Liu, Xiaojian Jiang, Jiebin Xu, Qiuxia Li, and Jun Zhao. 2024. Tug-of-war between knowledge: Exploring and resolving knowledge conflicts in retrieval-augmented language models. *arXiv preprint arXiv:2402.14409* (2024).
 - [10] Evgenii Kortukov, Alexander Rubinstein, Elisa Nguyen, and Seong Joon Oh. 2024. Studying Large Language Model Behaviors Under Realistic Knowledge Conflicts. *arXiv preprint arXiv:2404.16032* (2024).
 - [11] Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.
 - [12] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment. *arXiv preprint arXiv:2303.16634* (2023).
 - [13] Yuanjie Lyu, Zhiyu Li, Simin Niu, Feiyu Xiong, Bo Tang, Wenjin Wang, Hao Wu, Huanyong Liu, Tong Xu, and Enhong Chen. 2024. CRUD-RAG: A comprehensive chinese benchmark for retrieval-augmented generation of large language models. *arXiv preprint arXiv:2401.17043* (2024).
 - [14] Meta. 2024. Introducing Meta Llama 3: The most capable openly available LLM to date. (2024). <https://ai.meta.com/blog/meta-llama-3/>
 - [15] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.
 - [16] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.
 - [17] Dongyu Ru, Lin Qiu, Xiangkun Hu, Tianhang Zhang, Peng Shi, Shuaichen Chang, Jiayang Cheng, Cuxiang Wang, Shichao Sun, Huanyu Li, Zizhao Zhang, Binjie Wang, Jiayong Jiang, Tong He, Zhiguo Wang, Pengfei Liu, Yue Zhang, and Zheng Zhang. 2024. RAGCHECKER: A Fine-grained Framework for Diagnosing Retrieval-Augmented Generation. *arXiv preprint arXiv:2408.08067* (2024).
 - [18] Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. 2023. ARES: An Automated Evaluation Framework for Retrieval-Augmented Generation Systems. *arXiv:2311.09476* [cs.CL]
 - [19] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
 - [20] George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R Alvers, Dirk Weissenborn, Anastasia Krithara, Sergios Petridis, Dimitris Polychronopoulos, et al. 2015. An overview of the BIOASQ large-scale biomedical semantic indexing and question answering competition. *BMC bioinformatics* 16 (2015), 1–28.
 - [21] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368* (2023).
 - [22] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTScore: Evaluating text generation with BERT. In *International Conference on Learning Representations*.
 - [23] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghai Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems* 36 (2023), 46595–46623.
 - [24] Caleb Ziems, William Held, Omar Shaikh, Jiao Chen, Zhehao Zhang, and Diyi Yang. 2024. Can Large Language Models Transform Computational Social Science? *arXiv:2305.03514* [cs.CL]

A BioASQ Dataset Examples

Examples illustrating the structure of the BioASQ dataset used in our experiments.

Listing 1: Example BioASQ Query-Answer Pairs

- ```

1 [
2 {
3 "query_id": 0,
4 "text": "Is Hirschsprung disease a mendelian
5 or a multifactorial disorder?",
```

```

6 "metadata": {},
7 "gt_answer": "Coding sequence mutations in RET
8 , GDNF, EDNRB, EDN3, and SOX10 are
9 involved in the development of
10 Hirschsprung disease. The majority of
11 these genes was shown to be related to
12 Mendelian syndromic forms of Hirschsprung'
13 s disease, whereas the non-Mendelian
14 inheritance of sporadic non-syndromic
15 Hirschsprung disease proved to be complex;
16
17 },
18 {
19 "query_id": 1,
20 "text": "List signaling molecules (ligands)
21 that interact with the receptor EGFR?",
```

```

22 "metadata": {},
23 "gt_answer": "The 7 known EGFR ligands are:
24 epidermal growth factor (EGF),
25 betacellulin (BTC), epiregulin (EPR),
26 heparin-binding EGF (HB-EGF), transforming
27 growth factor-\alpha [TGF-\alpha],
28 amphiregulin (AREG) and epigen (EPG)."
```

**Listing 2: Example BioASQ Supporting Document Snippets**

```

1 [
2 {
3 "title": "",
4 "text": "INTRODUCTION: The majority of
5 patients with type 1 diabetes mellitus (T1
6 DM) do not achieve glycemic targets. In
7 addition, treatment with insulin is
8 associated with increased risk for
9 hypoglycemia and weight gain. Accordingly,
10 there is an unmet need for new safe and
11 effective glucose-lowering agents in this
12 population. Sotagliflozin, a dual
13 inhibitor of sodium-glucose co-
14 transporters 1 and 2, has been recently
15 approved for use in patients with T1DM."
16 "metadata": {},
17 "doc_id": "33108240"
18 },
19 {
20 "title": "",
21 "text": "The World Health Organization is
22 still revising the epidemiology of multi-
23 system inflammatory syndrome in children ((
24 MIS-C) and the preliminary case definition
25 , although there is a dearth of robust
26 evidence regarding the clinical
27 presentations, severity, and outcomes.
28 Researchers, epidemiologists, and
29 clinicians are struggling to characterize
30 and describe the disease phenomenon while
31 taking care of the diseased persons at the
32 forefronts."
33 "metadata": {},
34 "doc_id": "33110725"
35 }
36]
37
```

## B CCRS Evaluation Prompts

The specific prompts used to instruct the Llama-70B-Instruct judge model for each CCRS metric calculation.

### Contextual Coherence (CC) Evaluation Prompt

Evaluate the Contextual Coherence between the generated response and the retrieved context.

**Retrieved Context (C):** {context}

**Generated Response (r):** {response}

**Task:** Assess the logical consistency and coherence of the response with respect to the provided context. Ensure the response logically follows from and does not contradict the context.

**Scoring:** Score from 0 to 100, where 0 is completely incoherent or contradictory, and 100 is perfectly coherent and consistent. If no response is generated, the score is 0.

**Output:** Only print the score and nothing else.

### Question Relevance (QR) Evaluation Prompt

Evaluate the Question Relevance of the generated response to the user query.

**User Query (q):** {question}

**Generated Response (r):** {response}

**Task:** Assess how well the response directly addresses the user's query. Consider if the response answers the question asked.

**Scoring:** Score from 0 to 100, where 0 is completely irrelevant and 100 is perfectly relevant and directly answers the query. If no response is generated, the score is 0.

**Output:** Only print the score and nothing else.

### Information Density (ID) Evaluation Prompt

Evaluate the Information Density of the generated response.

**User Query (q):** {question}

**Retrieved Context (C):** {context}

**Generated Response (r):** {response}

**Task:** Assess the balance of conciseness and informativeness in the response, considering both the context and the query. The response should provide necessary information without being overly verbose or unnecessarily brief.

**Scoring:** Score from 0 to 100, where 0 is either too verbose (contains excessive irrelevant detail) or uninformative (lacks necessary information), and 100 is optimally concise and informative for the query. If no response is generated, the score is 0.

**Output:** Only print the score and nothing else.

### Answer Correctness (AC) Evaluation Prompt (LLM part)

Evaluate the Answer Correctness of the generated response.

**Retrieved Context (C):** {context}

**Generated Response (r):** {response}

**Ground Truth Answer (g):** {ground\_truth}

**Task:** Assess the factual accuracy of the information presented in the response compared to the ground truth answer, considering the provided context. Do not penalize differences in phrasing if the core factual meaning is preserved and accurate according to the ground truth.

**Scoring:** Score from 0 to 100, where 0 is completely incorrect or contains significant factual errors, and 100 represents perfect factual accuracy (semantically equivalent to the ground truth). If no response is generated, the score is 0.

**Output:** Only print the score and nothing else.

Note: The final AC score combines this LLM score (weighted 0.3) with an Exact Match check (weighted 0.7).

### Information Recall (IR) Evaluation Prompt

Evaluate the Information Recall of the generated response.

**Retrieved Context (C):** {context}

**Generated Response (r):** {response}

**Ground Truth Answer (g):** {ground\_truth}

**Task:** Assess how much of the essential information present in the ground truth answer is captured in the generated response, considering the provided context. Focus on whether key facts or points from the ground truth are included.

**Scoring:** Score from 0 to 100, where 0 means no essential information from the ground truth is recalled, and 100 means all essential information is fully captured. If no response is generated, the score is 0.

**Output:** Only print the score and nothing else.

## C Detailed Metric Property Data

### C.1 Observed Metric Bounds

Tables showing the frequency and percentage of scores reaching the minimum (0) and maximum (1) bounds for each CCRS metric across the 6 systems (N=4,719).

**Table 5: Observed Bounds (Count/Percentage) of CCRS Components - Part 1 (N=4719)**

| System       | Bound | CC           | QR           | ID         |
|--------------|-------|--------------|--------------|------------|
| A: M+BM25    | 0     | 186 / 3.9%   | 61 / 1.3%    | 14 / 0.3%  |
|              | 1     | 779 / 16.5%  | 3698 / 78.4% | 188 / 4.0% |
| C: L8+BM25   | 0     | 1070 / 22.7% | 229 / 4.8%   | 177 / 3.7% |
|              | 1     | 630 / 13.3%  | 3064 / 64.9% | 136 / 2.9% |
| E: L3.2+BM25 | 0     | 892 / 18.9%  | 151 / 3.2%   | 29 / 0.6%  |
|              | 1     | 368 / 7.8%   | 3104 / 65.8% | 182 / 3.9% |
| B: M+E5      | 0     | 186 / 3.9%   | 43 / 0.9%    | 16 / 0.3%  |
|              | 1     | 724 / 15.3%  | 4031 / 85.4% | 132 / 2.8% |
| D: L8+E5     | 0     | 995 / 21.1%  | 139 / 2.9%   | 82 / 1.7%  |
|              | 1     | 618 / 13.1%  | 3361 / 71.2% | 132 / 2.8% |
| F: L3.2+E5   | 0     | 819 / 17.4%  | 100 / 2.1%   | 17 / 0.4%  |
|              | 1     | 419 / 8.9%   | 3381 / 71.6% | 147 / 3.1% |

**Table 6: Observed Bounds (Count/Percentage) of CCRS Components - Part 2 (N=4719)**

| System       | Bound | AC          | IR           |
|--------------|-------|-------------|--------------|
| A: M+BM25    | 0     | 377 / 8.0%  | 610 / 12.9%  |
|              | 1     | 0 / 0.0%    | 254 / 5.4%   |
| C: L8+BM25   | 0     | 614 / 13.0% | 1015 / 21.5% |
|              | 1     | 0 / 0.0%    | 146 / 3.1%   |
| E: L3.2+BM25 | 0     | 592 / 12.5% | 1007 / 21.3% |
|              | 1     | 0 / 0.0%    | 102 / 2.2%   |
| B: M+E5      | 0     | 321 / 6.8%  | 483 / 10.2%  |
|              | 1     | 0 / 0.0%    | 228 / 4.8%   |
| D: L8+E5     | 0     | 530 / 11.2% | 887 / 18.8%  |
|              | 1     | 0 / 0.0%    | 138 / 2.9%   |
| F: L3.2+E5   | 0     | 562 / 11.9% | 895 / 19.0%  |
|              | 1     | 0 / 0.0%    | 112 / 2.4%   |

**Table 8: Correlations for System B (Mistral+E5)**

|                 | CC    | QR    | ID    | AC    | IR    |
|-----------------|-------|-------|-------|-------|-------|
| <b>Pearson</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.205 | 0.321 | 0.226 | 0.377 |
| QR              | 0.205 | 1.000 | 0.437 | 0.413 | 0.411 |
| ID              | 0.321 | 0.437 | 1.000 | 0.253 | 0.259 |
| AC              | 0.226 | 0.413 | 0.253 | 1.000 | 0.814 |
| IR              | 0.377 | 0.411 | 0.259 | 0.814 | 1.000 |
| <b>Spearman</b> |       |       |       |       |       |
| CC              | 1.000 | 0.294 | 0.324 | 0.324 | 0.402 |
| QR              | 0.294 | 1.000 | 0.391 | 0.354 | 0.403 |
| ID              | 0.324 | 0.391 | 1.000 | 0.281 | 0.292 |
| AC              | 0.324 | 0.354 | 0.281 | 1.000 | 0.785 |
| IR              | 0.402 | 0.403 | 0.292 | 0.785 | 1.000 |
| <b>Kendall</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.269 | 0.299 | 0.278 | 0.341 |
| QR              | 0.269 | 1.000 | 0.359 | 0.305 | 0.350 |
| ID              | 0.299 | 0.359 | 1.000 | 0.238 | 0.247 |
| AC              | 0.278 | 0.305 | 0.238 | 1.000 | 0.708 |
| IR              | 0.341 | 0.350 | 0.247 | 0.708 | 1.000 |

## C.2 Detailed CCRS Validity Analysis (Per System Correlations)

Correlation matrices (Pearson, Spearman, Kendall) for each of the six RAG systems, illustrating system-specific relationships between CCRS metrics.

**Table 7: Correlations for System A (Mistral+BM25)**

|                 | CC    | QR    | ID    | AC    | IR    |
|-----------------|-------|-------|-------|-------|-------|
| <b>Pearson</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.211 | 0.280 | 0.218 | 0.343 |
| QR              | 0.211 | 1.000 | 0.451 | 0.469 | 0.496 |
| ID              | 0.280 | 0.451 | 1.000 | 0.299 | 0.299 |
| AC              | 0.218 | 0.469 | 0.299 | 1.000 | 0.822 |
| IR              | 0.343 | 0.496 | 0.299 | 0.822 | 1.000 |
| <b>Spearman</b> |       |       |       |       |       |
| CC              | 1.000 | 0.284 | 0.324 | 0.300 | 0.362 |
| QR              | 0.284 | 1.000 | 0.458 | 0.421 | 0.495 |
| ID              | 0.324 | 0.458 | 1.000 | 0.370 | 0.368 |
| AC              | 0.300 | 0.421 | 0.370 | 1.000 | 0.809 |
| IR              | 0.362 | 0.495 | 0.368 | 0.809 | 1.000 |
| <b>Kendall</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.255 | 0.297 | 0.255 | 0.303 |
| QR              | 0.255 | 1.000 | 0.416 | 0.362 | 0.425 |
| ID              | 0.297 | 0.416 | 1.000 | 0.312 | 0.310 |
| AC              | 0.255 | 0.362 | 0.312 | 1.000 | 0.736 |
| IR              | 0.303 | 0.425 | 0.310 | 0.736 | 1.000 |

**Table 9: Correlations for System C (Llama8B+BM25)**

|                 | CC    | QR    | ID    | AC    | IR    |
|-----------------|-------|-------|-------|-------|-------|
| <b>Pearson</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.365 | 0.514 | 0.237 | 0.528 |
| QR              | 0.365 | 1.000 | 0.687 | 0.479 | 0.492 |
| ID              | 0.514 | 0.687 | 1.000 | 0.420 | 0.444 |
| AC              | 0.237 | 0.479 | 0.420 | 1.000 | 0.734 |
| IR              | 0.528 | 0.492 | 0.444 | 0.734 | 1.000 |
| <b>Spearman</b> |       |       |       |       |       |
| CC              | 1.000 | 0.354 | 0.402 | 0.234 | 0.514 |
| QR              | 0.354 | 1.000 | 0.611 | 0.362 | 0.527 |
| ID              | 0.402 | 0.611 | 1.000 | 0.339 | 0.466 |
| AC              | 0.234 | 0.362 | 0.339 | 1.000 | 0.716 |
| IR              | 0.514 | 0.527 | 0.466 | 0.716 | 1.000 |
| <b>Kendall</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.301 | 0.328 | 0.203 | 0.425 |
| QR              | 0.301 | 1.000 | 0.538 | 0.304 | 0.445 |
| ID              | 0.328 | 0.538 | 1.000 | 0.275 | 0.374 |
| AC              | 0.203 | 0.304 | 0.275 | 1.000 | 0.642 |
| IR              | 0.425 | 0.445 | 0.374 | 0.642 | 1.000 |

**Table 10: Correlations for System D (Llama8B+E5)**

|                 | CC    | QR    | ID    | AC    | IR    |
|-----------------|-------|-------|-------|-------|-------|
| <b>Pearson</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.334 | 0.478 | 0.195 | 0.511 |
| QR              | 0.334 | 1.000 | 0.585 | 0.363 | 0.436 |
| ID              | 0.478 | 0.585 | 1.000 | 0.281 | 0.362 |
| AC              | 0.195 | 0.363 | 0.281 | 1.000 | 0.713 |
| IR              | 0.511 | 0.436 | 0.362 | 0.713 | 1.000 |
| <b>Spearman</b> |       |       |       |       |       |
| CC              | 1.000 | 0.387 | 0.367 | 0.215 | 0.511 |
| QR              | 0.387 | 1.000 | 0.552 | 0.223 | 0.459 |
| ID              | 0.367 | 0.552 | 1.000 | 0.207 | 0.363 |
| AC              | 0.215 | 0.223 | 0.207 | 1.000 | 0.684 |
| IR              | 0.511 | 0.459 | 0.363 | 0.684 | 1.000 |
| <b>Kendall</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.335 | 0.298 | 0.189 | 0.423 |
| QR              | 0.335 | 1.000 | 0.488 | 0.189 | 0.390 |
| ID              | 0.298 | 0.488 | 1.000 | 0.167 | 0.290 |
| AC              | 0.189 | 0.189 | 0.167 | 1.000 | 0.617 |
| IR              | 0.423 | 0.390 | 0.290 | 0.617 | 1.000 |

**Table 12: Correlations for System F (Llama3.2+E5)**

|                 | CC    | QR    | ID    | AC    | IR    |
|-----------------|-------|-------|-------|-------|-------|
| <b>Pearson</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.220 | 0.354 | 0.170 | 0.477 |
| QR              | 0.220 | 1.000 | 0.382 | 0.275 | 0.367 |
| ID              | 0.354 | 0.382 | 1.000 | 0.245 | 0.315 |
| AC              | 0.170 | 0.275 | 0.245 | 1.000 | 0.726 |
| IR              | 0.477 | 0.367 | 0.315 | 0.726 | 1.000 |
| <b>Spearman</b> |       |       |       |       |       |
| CC              | 1.000 | 0.340 | 0.353 | 0.252 | 0.526 |
| QR              | 0.340 | 1.000 | 0.502 | 0.227 | 0.420 |
| ID              | 0.353 | 0.502 | 1.000 | 0.226 | 0.344 |
| AC              | 0.252 | 0.227 | 0.226 | 1.000 | 0.701 |
| IR              | 0.526 | 0.420 | 0.344 | 0.701 | 1.000 |
| <b>Kendall</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.296 | 0.289 | 0.221 | 0.433 |
| QR              | 0.296 | 1.000 | 0.438 | 0.193 | 0.357 |
| ID              | 0.289 | 0.438 | 1.000 | 0.179 | 0.270 |
| AC              | 0.221 | 0.193 | 0.179 | 1.000 | 0.630 |
| IR              | 0.433 | 0.357 | 0.270 | 0.630 | 1.000 |

**Table 11: Correlations for System E (Llama3.2+BM25)**

|                 | CC    | QR    | ID    | AC    | IR    |
|-----------------|-------|-------|-------|-------|-------|
| <b>Pearson</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.225 | 0.356 | 0.159 | 0.462 |
| QR              | 0.225 | 1.000 | 0.361 | 0.315 | 0.421 |
| ID              | 0.356 | 0.361 | 1.000 | 0.315 | 0.337 |
| AC              | 0.159 | 0.315 | 0.315 | 1.000 | 0.701 |
| IR              | 0.462 | 0.421 | 0.337 | 0.701 | 1.000 |
| <b>Spearman</b> |       |       |       |       |       |
| CC              | 1.000 | 0.279 | 0.366 | 0.221 | 0.490 |
| QR              | 0.279 | 1.000 | 0.474 | 0.267 | 0.464 |
| ID              | 0.366 | 0.474 | 1.000 | 0.307 | 0.387 |
| AC              | 0.221 | 0.267 | 0.307 | 1.000 | 0.672 |
| IR              | 0.490 | 0.464 | 0.387 | 0.672 | 1.000 |
| <b>Kendall</b>  |       |       |       |       |       |
| CC              | 1.000 | 0.238 | 0.300 | 0.195 | 0.406 |
| QR              | 0.238 | 1.000 | 0.412 | 0.227 | 0.392 |
| ID              | 0.300 | 0.412 | 1.000 | 0.247 | 0.307 |
| AC              | 0.195 | 0.227 | 0.247 | 1.000 | 0.612 |
| IR              | 0.406 | 0.392 | 0.307 | 0.612 | 1.000 |

**Table 13: Empirical Tie Probabilities (%) for CCRS Metrics Across Systems**

| Metric | A    | B    | C    | D    | E    | F    |
|--------|------|------|------|------|------|------|
| CC     | 33.9 | 40.7 | 22.6 | 24.0 | 23.0 | 24.1 |
| QR     | 62.5 | 73.5 | 45.9 | 53.7 | 46.7 | 54.2 |
| ID     | 31.1 | 33.1 | 21.1 | 22.1 | 21.4 | 21.5 |
| AC     | 16.3 | 16.2 | 16.6 | 16.5 | 16.0 | 15.6 |
| IR     | 17.0 | 18.2 | 15.4 | 14.9 | 15.4 | 14.8 |

#### C.4 Population Level Statistics

Detailed descriptive statistics (central tendency, variability, shape) for each CCRS metric across the six systems.

**Table 14: Central Tendency Measures Across Systems and Metrics (%)**

| System       | Metric | Mean  | GeoMean* | Median | Midhinge |
|--------------|--------|-------|----------|--------|----------|
| A: M+BM25    | AC     | 21.10 | 0.00     | 24.00  | 23.25    |
|              | CC     | 87.95 | 0.00     | 95.00  | 92.50    |
|              | ID     | 86.81 | 0.00     | 85.00  | 90.00    |
|              | IR     | 54.96 | 0.00     | 60.00  | 56.50    |
|              | QR     | 92.83 | 0.00     | 100.00 | 100.00   |
| B: M+E5      | AC     | 21.81 | 0.00     | 24.00  | 23.25    |
|              | CC     | 88.61 | 0.00     | 95.00  | 95.00    |
|              | ID     | 86.90 | 0.00     | 85.00  | 90.00    |
|              | IR     | 58.28 | 0.00     | 60.00  | 60.00    |
|              | QR     | 95.40 | 0.00     | 100.00 | 100.00   |
| C: L8+BM25   | AC     | 18.54 | 0.00     | 24.00  | 18.00    |
|              | CC     | 67.89 | 0.00     | 85.00  | 67.50    |
|              | ID     | 79.09 | 0.00     | 85.00  | 87.50    |
|              | IR     | 43.00 | 0.00     | 40.00  | 43.33    |
|              | QR     | 87.07 | 0.00     | 100.00 | 90.00    |
| D: L8+E5     | AC     | 19.24 | 0.00     | 24.00  | 18.75    |
|              | CC     | 69.29 | 0.00     | 95.00  | 77.50    |
|              | ID     | 80.41 | 0.00     | 85.00  | 87.50    |
|              | IR     | 45.44 | 0.00     | 50.00  | 50.00    |
|              | QR     | 90.55 | 0.00     | 100.00 | 95.00    |
| E: L3.2+BM25 | AC     | 18.41 | 0.00     | 24.00  | 18.00    |
|              | CC     | 68.88 | 0.00     | 80.00  | 77.50    |
|              | ID     | 80.23 | 0.00     | 85.00  | 87.50    |
|              | IR     | 42.04 | 0.00     | 40.00  | 40.00    |
|              | QR     | 88.28 | 0.00     | 100.00 | 90.00    |
| F: L3.2+E5   | AC     | 18.76 | 0.00     | 24.00  | 18.75    |
|              | CC     | 70.88 | 0.00     | 85.00  | 87.50    |
|              | ID     | 80.43 | 0.00     | 85.00  | 87.50    |
|              | IR     | 44.79 | 0.00     | 40.00  | 50.00    |
|              | QR     | 90.83 | 0.00     | 100.00 | 95.00    |

\*Geometric Mean is 0 due to presence of 0 scores.

**Table 15: Variability Measures Across Systems and Metrics**

| System       | Metric | Var     | Min | Max   | Range | IQR  |
|--------------|--------|---------|-----|-------|-------|------|
| A: M+BM25    | AC     | 75.72   | 0.0 | 30.7  | 30.7  | 10.5 |
|              | CC     | 424.25  | 0.0 | 100.0 | 100.0 | 5.0  |
|              | ID     | 123.58  | 0.0 | 100.0 | 100.0 | 10.0 |
|              | IR     | 918.67  | 0.0 | 100.0 | 100.0 | 47.0 |
|              | QR     | 340.18  | 0.0 | 100.0 | 100.0 | 0.0  |
| B: M+E5      | AC     | 69.38   | 0.0 | 30.7  | 30.7  | 10.5 |
|              | CC     | 423.80  | 0.0 | 100.0 | 100.0 | 0.0  |
|              | ID     | 123.40  | 0.0 | 100.0 | 100.0 | 10.0 |
|              | IR     | 834.55  | 0.0 | 100.0 | 100.0 | 40.0 |
|              | QR     | 227.31  | 0.0 | 100.0 | 100.0 | 0.0  |
| C: L8+BM25   | AC     | 93.68   | 0.0 | 30.7  | 30.7  | 12.0 |
|              | CC     | 1527.75 | 0.0 | 100.0 | 100.0 | 55.0 |
|              | ID     | 502.50  | 0.0 | 100.0 | 100.0 | 15.0 |
|              | IR     | 973.87  | 0.0 | 100.0 | 100.0 | 46.7 |
|              | QR     | 633.71  | 0.0 | 100.0 | 100.0 | 20.0 |
| D: L8+E5     | AC     | 89.46   | 0.0 | 30.7  | 30.7  | 13.5 |
|              | CC     | 1494.39 | 0.0 | 100.0 | 100.0 | 35.0 |
|              | ID     | 392.57  | 0.0 | 100.0 | 100.0 | 15.0 |
|              | IR     | 970.01  | 0.0 | 100.0 | 100.0 | 60.0 |
|              | QR     | 440.64  | 0.0 | 100.0 | 100.0 | 10.0 |
| E: L3.2+BM25 | AC     | 93.49   | 0.0 | 30.7  | 30.7  | 12.0 |
|              | CC     | 1330.38 | 0.0 | 100.0 | 100.0 | 35.0 |
|              | ID     | 335.04  | 0.0 | 100.0 | 100.0 | 15.0 |
|              | IR     | 944.44  | 0.0 | 100.0 | 100.0 | 40.0 |
|              | QR     | 537.50  | 0.0 | 100.0 | 100.0 | 20.0 |
| F: L3.2+E5   | AC     | 92.92   | 0.0 | 30.7  | 30.7  | 13.5 |
|              | CC     | 1298.15 | 0.0 | 100.0 | 100.0 | 15.0 |
|              | ID     | 321.37  | 0.0 | 100.0 | 100.0 | 15.0 |
|              | IR     | 952.22  | 0.0 | 100.0 | 100.0 | 60.0 |
|              | QR     | 408.58  | 0.0 | 100.0 | 100.0 | 10.0 |

**Table 16: Shape Measures Across Systems and Metrics**

| System       | Metric | Skewness | Kurtosis |
|--------------|--------|----------|----------|
| A: M+BM25    | AC     | -1.211   | 0.486    |
|              | CC     | -3.425   | 11.379   |
|              | ID     | -3.686   | 20.748   |
|              | IR     | -0.503   | -0.882   |
|              | QR     | -3.278   | 11.045   |
| B: M+E5      | AC     | -1.387   | 1.047    |
|              | CC     | -3.537   | 11.916   |
|              | ID     | -3.949   | 22.791   |
|              | IR     | -0.695   | -0.565   |
|              | QR     | -4.363   | 20.605   |
| C: L8+BM25   | AC     | -0.766   | -0.672   |
|              | CC     | -1.030   | -0.744   |
|              | ID     | -2.221   | 4.479    |
|              | IR     | 0.020    | -1.260   |
|              | QR     | -2.412   | 5.112    |
| D: L8+E5     | AC     | -0.872   | -0.460   |
|              | CC     | -1.092   | -0.615   |
|              | ID     | -2.262   | 5.052    |
|              | IR     | -0.078   | -1.271   |
|              | QR     | -3.009   | 9.181    |
| E: L3.2+BM25 | AC     | -0.719   | -0.738   |
|              | CC     | -1.165   | -0.358   |
|              | ID     | -2.071   | 4.284    |
|              | IR     | 0.050    | -1.262   |
|              | QR     | -2.527   | 5.965    |
| F: L3.2+E5   | AC     | -0.783   | -0.674   |
|              | CC     | -1.256   | -0.140   |
|              | ID     | -2.056   | 4.160    |
|              | IR     | -0.071   | -1.281   |
|              | QR     | -2.968   | 9.049    |

## D RAGChecker Comparison Details

Supporting data for the comparison between CCRS and RAGChecker (Precision, Recall, Faithfulness).

**Table 17: Observed Bounds (Count/Percentage) for RAGChecker Metrics (N=4719)**

| System       | Bound | Precision    | Recall       | Faithfulness |
|--------------|-------|--------------|--------------|--------------|
| A: M+BM25    | 0     | 568 / 12.0%  | 767 / 16.2%  | 209 / 4.4%   |
|              | 1     | 1471 / 31.2% | 1003 / 21.2% | 3158 / 66.9% |
| B: M+E5      | 0     | 500 / 10.6%  | 587 / 12.4%  | 128 / 2.7%   |
|              | 1     | 1475 / 31.3% | 1075 / 22.8% | 3469 / 73.5% |
| C: L8+BM25   | 0     | 971 / 20.6%  | 783 / 16.6%  | 564 / 11.9%  |
|              | 1     | 1710 / 36.2% | 1037 / 22.0% | 3175 / 67.3% |
| D: L8+E5     | 0     | 839 / 17.8%  | 640 / 13.6%  | 454 / 9.6%   |
|              | 1     | 1788 / 37.9% | 1062 / 22.5% | 3386 / 71.8% |
| E: L3.2+BM25 | 0     | 947 / 20.1%  | 860 / 18.2%  | 505 / 10.7%  |
|              | 1     | 1788 / 37.9% | 960 / 20.3%  | 3062 / 64.9% |
| F: L3.2+E5   | 0     | 857 / 18.2%  | 689 / 14.6%  | 391 / 8.3%   |
|              | 1     | 1773 / 37.6% | 1018 / 21.6% | 3353 / 71.1% |

**Table 18: Empirical Tie Probabilities (%) for RAGChecker Metrics**

| System       | Precision | Recall | Faithfulness |
|--------------|-----------|--------|--------------|
| A: M+BM25    | 13.0      | 9.3    | 45.5         |
| B: M+E5      | 12.6      | 8.9    | 54.4         |
| C: L8+BM25   | 18.5      | 9.6    | 47.0         |
| D: L8+E5     | 18.7      | 9.1    | 52.6         |
| E: L3.2+BM25 | 19.5      | 9.4    | 43.6         |
| F: L3.2+E5   | 18.5      | 8.9    | 51.4         |

**Table 19: Population Statistics for RAGChecker Metrics (F/P/R, Scores 0-100)**

| System       | M | Mean | Median | Var  | Skew  | Kurt  |
|--------------|---|------|--------|------|-------|-------|
| A: M+BM25    | F | 87.3 | 100.0  | 6.1  | -2.34 | 4.93  |
|              | P | 57.2 | 55.6   | 12.9 | -0.15 | -1.35 |
|              | R | 49.9 | 50.0   | 12.3 | 0.08  | -1.28 |
| B: M+E5      | F | 91.3 | 100.0  | 4.2  | -3.08 | 9.74  |
|              | P | 58.0 | 57.1   | 12.5 | -0.17 | -1.32 |
|              | R | 53.7 | 50.0   | 11.6 | -0.05 | -1.22 |
| C: L8+BM25   | F | 82.5 | 100.0  | 11.0 | -1.82 | 1.76  |
|              | P | 55.8 | 53.3   | 15.9 | -0.17 | -1.54 |
|              | R | 50.3 | 50.0   | 12.6 | 0.06  | -1.31 |
| D: L8+E5     | F | 85.8 | 100.0  | 9.2  | -2.20 | 3.32  |
|              | P | 58.2 | 60.0   | 15.3 | -0.26 | -1.48 |
|              | R | 53.2 | 50.0   | 11.9 | -0.05 | -1.25 |
| E: L3.2+BM25 | F | 81.9 | 100.0  | 10.5 | -1.76 | 1.64  |
|              | P | 57.2 | 60.0   | 15.8 | -0.22 | -1.52 |
|              | R | 48.0 | 46.7   | 12.7 | 0.14  | -1.30 |
| F: L3.2+E5   | F | 86.3 | 100.0  | 8.5  | -2.25 | 3.70  |
|              | P | 58.0 | 60.0   | 15.2 | -0.25 | -1.47 |
|              | R | 51.2 | 50.0   | 12.1 | 0.04  | -1.27 |

## E Detailed Statistical Test Results

Complete pairwise comparison results from Tukey's HSD tests ( $\alpha = 0.05$ , B=10,000 permutations) for each CCRS metric. Diff = (Mean1 - Mean2) in percentage points. Significant p-values (<0.05) are bolded.

Table 20: Contextual Coherence (CC) Pairwise Comparisons

| Comparison | Mean Diff | P-value           | Significant? |
|------------|-----------|-------------------|--------------|
| A vs B     | -0.66     | 0.8236            | No           |
| A vs C     | +20.06    | <b>&lt;0.0001</b> | Yes          |
| A vs D     | +18.66    | <b>&lt;0.0001</b> | Yes          |
| A vs E     | +19.07    | <b>&lt;0.0001</b> | Yes          |
| A vs F     | +17.07    | <b>&lt;0.0001</b> | Yes          |
| B vs C     | +20.72    | <b>&lt;0.0001</b> | Yes          |
| B vs D     | +19.32    | <b>&lt;0.0001</b> | Yes          |
| B vs E     | +19.73    | <b>&lt;0.0001</b> | Yes          |
| B vs F     | +17.73    | <b>&lt;0.0001</b> | Yes          |
| C vs D     | -1.40     | 0.1017            | No           |
| C vs E     | -0.99     | 0.4446            | No           |
| C vs F     | -2.99     | <b>&lt;0.0001</b> | Yes          |
| D vs E     | +0.41     | 0.9736            | No           |
| D vs F     | -1.59     | <b>0.0417</b>     | Yes          |
| E vs F     | -2.00     | <b>0.0036</b>     | Yes          |

**Table 23: Answer Correctness (AC) Pairwise Comparisons**

| Comparison | Mean Diff | P-value         | Significant? |
|------------|-----------|-----------------|--------------|
| A vs B     | -0.71     | < <b>0.0001</b> | Yes          |
| A vs C     | +2.56     | < <b>0.0001</b> | Yes          |
| A vs D     | +1.86     | < <b>0.0001</b> | Yes          |
| A vs E     | +2.69     | < <b>0.0001</b> | Yes          |
| A vs F     | +2.34     | < <b>0.0001</b> | Yes          |
| B vs C     | +3.27     | < <b>0.0001</b> | Yes          |
| B vs D     | +2.57     | < <b>0.0001</b> | Yes          |
| B vs E     | +3.40     | < <b>0.0001</b> | Yes          |
| B vs F     | +3.06     | < <b>0.0001</b> | Yes          |
| C vs D     | -0.70     | < <b>0.0001</b> | Yes          |
| C vs E     | +0.13     | 0.9025          | No           |
| C vs F     | -0.22     | 0.5122          | No           |
| D vs E     | +0.83     | < <b>0.0001</b> | Yes          |
| D vs F     | +0.48     | <b>0.0020</b>   | Yes          |
| E vs F     | -0.35     | 0.0620          | No           |

**Table 21: Question Relevance (QR) Pairwise Comparisons**

| Comparison | Mean Diff | P-value | Significant? |
|------------|-----------|---------|--------------|
| A vs B     | -2.57     | <0.0001 | Yes          |
| A vs C     | +5.76     | <0.0001 | Yes          |
| A vs D     | +2.28     | <0.0001 | Yes          |
| A vs E     | +4.55     | <0.0001 | Yes          |
| A vs F     | +2.00     | <0.0001 | Yes          |
| B vs C     | +8.33     | <0.0001 | Yes          |
| B vs D     | +4.85     | <0.0001 | Yes          |
| B vs E     | +7.12     | <0.0001 | Yes          |
| B vs F     | +4.57     | <0.0001 | Yes          |
| C vs D     | -3.48     | <0.0001 | Yes          |
| C vs E     | -1.21     | 0.0037  | Yes          |
| C vs F     | -3.76     | <0.0001 | Yes          |
| D vs E     | +2.27     | <0.0001 | Yes          |
| D vs F     | -0.28     | 0.9553  | No           |
| E vs F     | -2.55     | <0.0001 | Yes          |

Table 24: Information Recall (IR) Pairwise Comparisons

| Comparison | Mean Diff | P-value         | Significant? |
|------------|-----------|-----------------|--------------|
| A vs B     | -3.32     | < <b>0.0001</b> | Yes          |
| A vs C     | +11.96    | < <b>0.0001</b> | Yes          |
| A vs D     | +9.52     | < <b>0.0001</b> | Yes          |
| A vs E     | +12.92    | < <b>0.0001</b> | Yes          |
| A vs F     | +10.17    | < <b>0.0001</b> | Yes          |
| B vs C     | +15.28    | < <b>0.0001</b> | Yes          |
| B vs D     | +12.84    | < <b>0.0001</b> | Yes          |
| B vs E     | +16.24    | < <b>0.0001</b> | Yes          |
| B vs F     | +13.49    | < <b>0.0001</b> | Yes          |
| C vs D     | -2.44     | < <b>0.0001</b> | Yes          |
| C vs E     | +0.96     | 0.2123          | No           |
| C vs F     | -1.78     | <b>0.0006</b>   | Yes          |
| D vs E     | +3.39     | < <b>0.0001</b> | Yes          |
| D vs F     | +0.65     | 0.6418          | No           |
| E vs F     | -2.74     | < <b>0.0001</b> | Yes          |

**Table 22: Information Density (ID) Pairwise Comparisons**

| Comparison | Mean Diff | P-value           | Significant? |
|------------|-----------|-------------------|--------------|
| A vs B     | -0.09     | 0.9998            | No           |
| A vs C     | +7.72     | <b>&lt;0.0001</b> | Yes          |
| A vs D     | +6.40     | <b>&lt;0.0001</b> | Yes          |
| A vs E     | +6.58     | <b>&lt;0.0001</b> | Yes          |
| A vs F     | +6.38     | <b>&lt;0.0001</b> | Yes          |
| B vs C     | +7.81     | <b>&lt;0.0001</b> | Yes          |
| B vs D     | +6.49     | <b>&lt;0.0001</b> | Yes          |
| B vs E     | +6.67     | <b>&lt;0.0001</b> | Yes          |
| B vs F     | +6.47     | <b>&lt;0.0001</b> | Yes          |
| C vs D     | -1.32     | <b>0.0004</b>     | Yes          |
| C vs E     | -1.14     | <b>0.0032</b>     | Yes          |
| C vs F     | -1.34     | <b>0.0004</b>     | Yes          |
| D vs E     | +0.18     | 0.9932            | No           |
| D vs F     | -0.02     | 1.0000            | No           |
| E vs F     | -0.20     | 0.9891            | No           |

## F Statistical Testing Implementation Code

The Python code used for performing the Tukey's HSD with randomization and the Holm-Bonferroni correction across hypotheses.

### **Listing 3: Statistical Testing Implementation for RAG Hypotheses**

```
1 import numpy as np
2 from scipy import stats
3 import json
4 from typing import Dict, List, Tuple, Any
5 from tqdm import tqdm
6 import argparse # Added for file arguments
7
8 def load_system_data(files: Dict[str, str]) ->
9 Dict[str, np.ndarray]:
10 """Load data for all systems and extract
11 metrics."""
12 print("\nLoading system data...")
13 data_systems = {}
```

```

12 # Ensure metrics match the order used in
13 # analysis consistently
14 metrics = ['Contextual_Coherence', '
15 Question_Relevance', 'Information_Density'
16 ,
17 'Answer_Correctness', '
18 Information_Recall']
19
20 for sys_label, file_path in files.items():
21 print(f"Loading system {sys_label} from {
22 file_path}")
23 try:
24 with open(file_path, 'r') as f:
25 data = json.load(f)
26 # Check if 'ccrs_results' key
27 # exists
28 if 'ccrs_results' not in data:
29 print(f"Error: 'ccrs_results'
30 key not found in {
31 file_path}")
32 return None
33
34 results_list = data['ccrs_results'
35]
36 if not isinstance(results_list,
37 list):
38 print(f"Error: 'ccrs_results'
39 is not a list in {
40 file_path}")
41 return None
42
43 processed_data = []
44 for i, m in enumerate(results_list
45):
46 if not isinstance(m, dict) or
47 'metrics' not in m:
48 print(f"Warning: Invalid
49 item format at index {
50 i} in {file_path}.
51 Skipping.")
52 continue
53 metric_values = []
54 valid_item = True
55 for metric in metrics:
56 if metric not in m['
57 metrics']:
58 print(f"Warning:
59 Metric '{metric}'"
60 not found for item
61 {i} in {file_path}
62). Skipping item."
63)
64 valid_item = False
65 break
66 # Ensure score is numeric
67 # and handle potential
68 # None or non-numeric
69 # types
70 score = m['metrics'][
71 metric]
72 if isinstance(score, (int,
73 float)):
74
75 metric_values.append(
76 score / 100.0) #
77 Normalize
78 else:
79 print(f"Warning: Non-
80 numeric score '{
81 score}' for metric
82 '{metric}' item {
83 i} in {file_path}.
84 Skipping item.")
85 valid_item = False
86 break
87 if valid_item:
88 processed_data.append(
89 metric_values)
90
91 if not processed_data: # If no
92 valid data was processed
93 print(f"Error: No valid data
94 processed from {file_path
95 }")
96 return None
97
98 data_systems[sys_label] = np.array(
99 processed_data)
100
101 print(f"System {sys_label} loaded:
102 shape {data_systems[sys_label].
103 shape}")
104
105 except FileNotFoundError:
106 print(f"Error: File not found {
107 file_path}")
108 return None
109 except json.JSONDecodeError:
110 print(f"Error: Could not decode JSON
111 from {file_path}")
112 return None
113 except Exception as e:
114 print(f"An unexpected error occurred
115 loading {file_path}: {e}")
116 return None
117
118 # Final checks after loading all files
119 if len(data_systems) != len(files):
120 print("Error: Could not load all specified
121 system files.")
122 return None
123 if any(arr.size == 0 for arr in data_systems.
124 values()):
125 print("Error: At least one system resulted
126 in no valid data after processing.")
127 return None
128
129 # Check for shape consistency (number of
130 # samples and metrics)
131 first_shape = next(iter(data_systems.values()))
132 .shape
133 if len(first_shape) != 2:
134
135

```



```

Adjust for one-sided tests
if one_sided:
 if expected_directions is None: raise
 ValueError("Expected directions needed
 for one-sided test.")
print("Adjusting for one-sided test...")
for pair_key in list(p_values.keys()):
 if np.isnan(p_values[pair_key]):
 continue

 i_str, j_str = pair_key.split(',')
 i, j = int(i_str), int(j_str)
 lookup_key = f"{i},{j}"
 expected_dir = expected_directions.get(
 lookup_key)
 if expected_dir is None:
 # Try reverse pair if direction
 # not defined for i,j
 reverse_lookup_key = f"{j},{i}"
 expected_dir =
 expected_directions.get(
 reverse_lookup_key)
 if expected_dir is not None:
 expected_dir *= -1 # Invert
 direction
 else: raise KeyError(f"Direction
 for pair {lookup_key} or {
 reverse_lookup_key} not found
 .")

 observed_diff = observed_diffs[
 pair_key]
 if np.isnan(observed_diff): continue

 current_p = p_values[pair_key]
 # Check if observed difference matches
 # expected direction (use tolerance
)
 is_correct_direction = (expected_dir >
 0 and observed_diff > 1e-9) or \
 (expected_dir <
 0 and
 observed_diff
 < -1e-9)
 is_zero_diff = abs(observed_diff) < 1e
 -9

 if is_correct_direction:
 # If direction is correct, p-value
 # is halved
 p_values[pair_key] = current_p /
 2.0
 elif expected_dir != 0 and
 is_zero_diff:
 # If difference is effectively
 # zero for a directional test
 p_values[pair_key] = 0.5
 elif expected_dir == 0:
 # Keep two-sided p-value if no
 # direction expected (should
 # not happen here)
 pass

```

```

239 reject[name] = True
240 significant_found = True
241 else:
242 print(f" -> Fail to reject H0 for {
243 name} (and subsequent hypotheses)"
244)
245 # Once we fail to reject, stop
246 # rejecting for subsequent
247 # hypotheses
248 break # Optimization: no need to check
249 further hypotheses
250
251 # Calculate final adjusted p-values ensuring
252 # monotonicity
253 # Adjusted p-value for hypothesis i is max(p_j
254 # * (n - j + 1)) for j <= i (where p_j is
255 # sorted raw p)
256 last_p_adj = 0.0
257 for i, (p_sort, p_orig, idx, name) in
258 enumerate(sorted_pairs):
259 holm_p = min(1.0, p_sort * (n - i)) #
260 Adjusted p for this step
261 last_p_adj = max(last_p_adj, holm_p) #
262 Enforce monotonicity
263 adjusted_p_values[name] = last_p_adj if
264 not np.isnan(p_orig) else 1.0 #
265 Assign adjusted p, keep NaN as 1.0
266
267 if not significant_found:
268 print(" No hypotheses rejected.")
269
270 return reject, adjusted_p_values
271
272 def test_all_hypotheses(data_systems: Dict[str, np
273 .ndarray],
274 alpha: float = 0.05) ->
275 Dict[str, Any]:
276 """Run all hypothesis tests with correct
277 multiple testing correction."""
278 metric_names = ['CC', 'QR', 'ID', 'AC', 'IR']
279 h2_metrics = ['QR', 'IR'] # Metrics relevant
280 for H2
281 hypothesis_names = ['H1', 'H2', 'H3']
282 system_labels = ['A', 'B', 'C', 'D', 'E', 'F']
283 sys_idx = {label: i for i, label in enumerate(
284 system_labels)}
285
286 # Define comparison pairs indices based on
287 # system_labels A=0, B=1, ... F=5
288 h1_pairs = [(sys_idx['A'], sys_idx['C']), (
289 sys_idx['A'], sys_idx['E']),
290 (sys_idx['B'], sys_idx['D']), (
291 sys_idx['B'], sys_idx['F'])]
292 h2_pairs = [(sys_idx['D'], sys_idx['C']), (
293 sys_idx['F'], sys_idx['E'])]
294 h3_pairs = [(sys_idx['C'], sys_idx['E']), (
295 sys_idx['D'], sys_idx['F'])]
296
297 # Expected directions for one-sided tests (H1:
298 # Mistral > Llama, H2: E5 > BM25)
299
300 h1_directions = {f"{i},{j}": 1 for i,j in
301 h1_pairs} # Assuming mean(A/B) > mean(C/D/
302 E/F)
303 h2_directions = {f"{i},{j}": 1 for i,j in
304 h2_pairs} # Assuming mean(D/F) > mean(C/E)
305
306 # Storage for aggregated p-values per
307 # hypothesis
308 h1_pvals_all_metrics_pairs = []
309 h2_pvals_all_metrics_pairs = []
310 h3_min_p_by_metric = {}
311 results_by_metric = {}
312
313 # Combine data for easier slicing: shape (
314 n_samples, n_metrics, n_systems)
315 all_system_data_array = np.stack([data_systems
316 [lbl] for lbl in system_labels], axis=2)
317
318 for metric_idx, metric_name in enumerate(
319 metric_names):
320 print(f"\n===== Testing metric: {
321 metric_name} =====")
322 # Extract data for the current metric:
323 # shape (n_samples, n_systems)
324 metric_data = all_system_data_array[:, ,
325 metric_idx, :]
326
327 # --- H1: Mistral Superiority ---
328 print("\n--- H1 (Mistral > Llama) ---")
329 h1_pvals_pairs, h1_diffs, _ =
330 tukeys_hsd_randomization(metric_data,
331 h1_pairs, one_sided=True,
332 expected_directions=h1_directions)
333 # Collect valid p-values for H1
334 aggregation
335 h1_pvals_all_metrics_pairs.extend([p for p
336 in h1_pvals_pairs.values() if p is
337 not None and not np.isnan(p)])
338
339 # --- H2: E5 Advantage (Llama only) ---
340 h2_pvals_pairs_metric = {}
341 h2_diffs_metric = {}
342 if metric_name in h2_metrics:
343 print(f"\n--- H2 (E5 > BM25 for Llama)
344 ---")
345 h2_pvals_pairs_metric, h2_diffs_metric
346 , _ = tukeys_hsd_randomization(
347 metric_data, h2_pairs, one_sided=
348 True, expected_directions=
349 h2_directions)
350 # Collect valid p-values for H2
351 aggregation
352 h2_pvals_all_metrics_pairs.extend([p
353 for p in h2_pvals_pairs_metric.
354 values() if p is not None and not
355 np.isnan(p)])
356 else: # Initialize placeholders if metric
357 not relevant for H2
358 h2_pvals_pairs_metric = {f"{i},{j}": np.nan
359 for i,j in h2_pairs}
360 h2_diffs_metric = {f"{i},{j}": np.nan
361 for i,j in h2_pairs}

```

```

310 # --- H3: Llama Model Differences ---
311 print(f"\n--- H3 (Llama 8B vs 3.2B) ---")
312 # Two-sided test for H3
313 h3_pvals_pairs, h3_diffs, means =
314 tukeys_hsd_randomization(metric_data,
315 h3_pairs, one_sided=False)
316 # Find the minimum non-Nan p-value for
317 # this metric for H3 aggregation
318 valid_h3_pvals = [p for p in
319 h3_pvals_pairs.values() if p is not
320 None and not np.isnan(p)]
321 h3_min_p_by_metric[metric_name] = min(
322 valid_h3_pvals) if valid_h3_pvals else
323 1.0
324
325 # Store detailed results per metric
326 results_by_metric[metric_name] = {
327 'h1_pairs_results': {f"{system_labels[i]}_vs_{system_labels[j]}": {
328 'p_value': h1_pvals_pairs.get(f"{i},{j}", np.nan), 'difference':
329 h1_diffs.get(f"{i},{j}", np.nan)} for i,j in h1_pairs},
330 'h2_pairs_results': {f"{system_labels[i]}_vs_{system_labels[j]}": {
331 'p_value': h2_pvals_pairs_metric.get(f"{i},{j}", np.nan), 'difference':
332 h2_diffs_metric.get(f"{i},{j}", np.nan)} for i,j in h2_pairs} if metric_name in
333 h2_metrics else None,
334 'h3_pairs_results': {f"{system_labels[i]}_vs_{system_labels[j]}": {
335 'p_value': h3_pvals_pairs.get(f"{i},{j}", np.nan), 'difference':
336 h3_diffs.get(f"{i},{j}", np.nan)} for i,j in h3_pairs},
337 'means': (means * 100).round(3).tolist()
338 () if not np.isnan(means).all() else []
339 }
340 # print(f"\n{metric_name} Summary: H1 Max(p)={max(h1_pvals_pairs.values() if
341 h1_pvals_pairs else [1.0]):.6f}" +
342 f", H2 Max(p)={max(
343 h2_pvals_pairs_metric.values() if
344 h2_pvals_pairs_metric and any(p is not
345 None and not np.isnan(p) for p in
346 h2_pvals_pairs_metric.values()) else
347 [1.0]):.6f}" if metric_name in
348 h2_metrics else "") +
349 f", H3 Min(p)={h3_min_p_by_metric[metric_name]:.6f}")
350
351 # Aggregate p-values across metrics for each
352 # hypothesis based on logic
353 # H1/H2 (AND condition): max p-value across
354 # all relevant comparisons
355 h1_final_p = max(h1_pvals_all_metrics_pairs)
356 if h1_pvals_all_metrics_pairs else 1.0
357
358 h2_final_p = max(h2_pvals_all_metrics_pairs)
359 if h2_pvals_all_metrics_pairs else 1.0
360 # H3 (OR condition): min p-value across
361 # metrics, then Bonferroni correct
362 h3_min_p_across_metrics = min(
363 h3_min_p_by_metric.values()) if
364 h3_min_p_by_metric else 1.0
365 h3_final_p_bonferroni = min(1.0,
366 h3_min_p_across_metrics * len(metric_names))
367) # Bonferroni correction factor = number
368 of metrics
369
370 print("\n==== Aggregated Hypothesis p-values
371 (before Holm-Bonferroni) ====")
372 print(f"H1 (Max p across all comparisons): {h1_final_p:.6f}")
373 print(f"H2 (Max p across relevant comparisons):
374 : {h2_final_p:.6f}")
375 print(f"H3 (Bonf-corr Min p across metrics): {h3_final_p_bonferroni:.6f}")
376
377 # Apply Holm-Bonferroni correction across the
378 # three final hypothesis p-values
379 hypothesis_p_values = [h1_final_p, h2_final_p,
380 h3_final_p_bonferroni]
381 significant, final_adjusted_p_values =
382 apply_holm_bonferroni(hypothesis_p_values,
383 hypothesis_names, alpha)
384
385 # Compile final results object
386 results = {'hypothesis_results': {}, 'metric_results': results_by_metric}
387 for h_idx, h_name in enumerate(
388 hypothesis_names):
389 results['hypothesis_results'][h_name] = {
390 'significant': significant[h_name],
391 'raw_aggregated_p_value':
392 hypothesis_p_values[h_idx],
393 'final_adjusted_p_value':
394 final_adjusted_p_values[h_name],
395 'conclusion': 'Reject H0' if
396 significant[h_name] else 'Fail to
397 reject H0'
398 }
399 if h_name == 'H3': # Add extra info for
400 H3
401 results['hypothesis_results'][h_name][
402 'min_p_across_metrics_raw'] =
403 h3_min_p_across_metrics
404 results['hypothesis_results'][h_name][
405 'metric_min_p_values'] =
406 h3_min_p_by_metric
407
408 print("\n==== Final Conclusions after Holm-
409 Bonferroni Correction ====")
410 for h_name, h_result in results['hypothesis_results'].items():
411 p_val_desc = 'Max Raw p' if h_name in ['H1',
412 'H2'] else 'Bonf-corr Min Raw p'

```

```

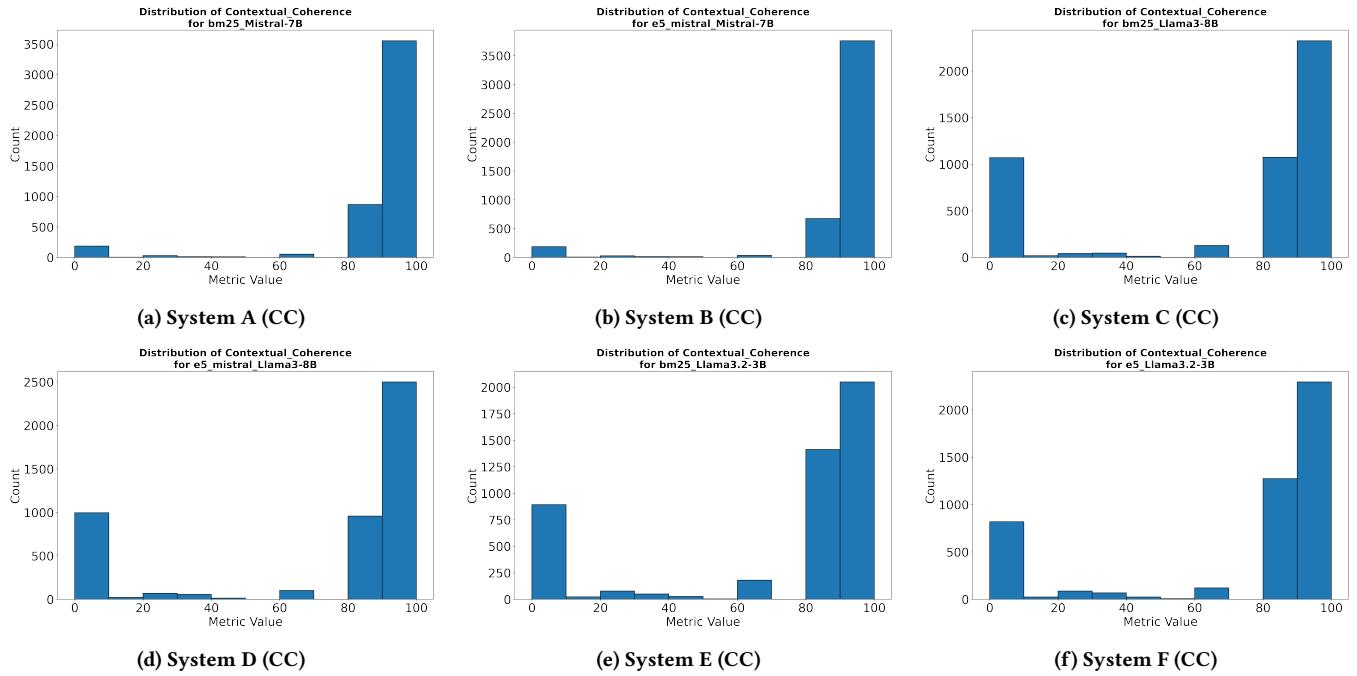
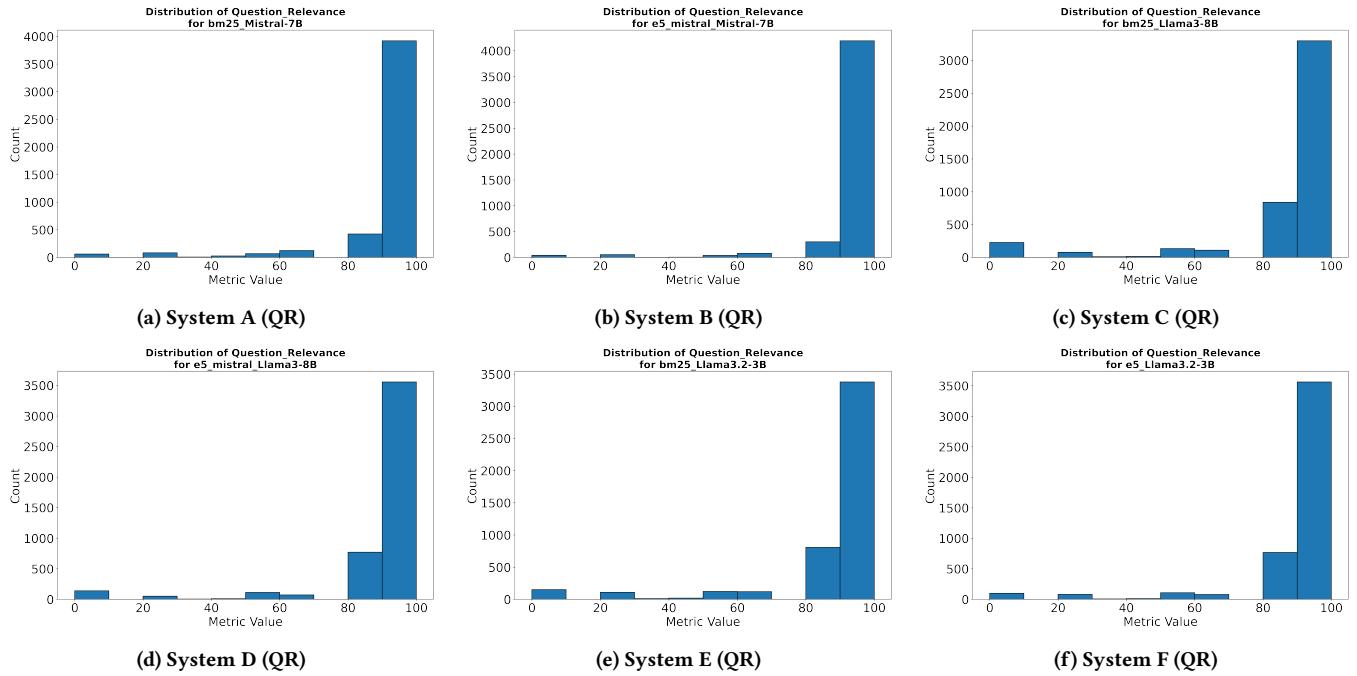
363 print(f"\n{h_name}: {p_val_desc}={h_result
364 ['raw_aggregated_p_value']:.6f}; Final
365 adj. p={h_result['
366 final_adjusted_p_value']:.6f} -> {
367 h_result['conclusion']}")"
368
369 return results
370
371 def main():
372 parser = argparse.ArgumentParser(description="
373 Run RAG system hypothesis tests.")
374 # Add arguments for file paths if needed,
375 # otherwise use defaults
376 parser.add_argument('--file_a', default='
377 modified_bioasq_bm25_100_0.2_k20_Mistral-7
378 B-results.json')
379 parser.add_argument('--file_b', default='
380 modified_bioasq_e5_mistral_100_0.2
381 _k20_Mistral-7B-results.json')
382 parser.add_argument('--file_c', default='
383 modified_bioasq_bm25_100_0.2_k20_Llama3-8B
384 -results.json')
385 parser.add_argument('--file_d', default='
386 modified_bioasq_e5_mistral_100_0.2
387 _k20_Llama3-8B-results.json')
388 parser.add_argument('--file_e', default='
389 modified_bioasq_bm25_100_0.2_k20_Llama3
390 .2-3B-results.json')
391 parser.add_argument('--file_f', default='
392 modified_bioasq_e5_mistral_100_0.2
393 _k20_Llama3.2-3B-results.json')
394 parser.add_argument('--alpha', type=float,
395 default=0.05, help='Significance level
396 alpha.')
397 parser.add_argument('--output', default='
398 hypothesis_test_results_full.json', help='
399 Output JSON file name.')
400
401 args = parser.parse_args()
402
403 files = {
404 'A': args.file_a, 'B': args.file_b, 'C':
405 args.file_c,
406 'D': args.file_d, 'E': args.file_e, 'F':
407 args.file_f
408 }
409
410 print("Starting hypothesis testing...")
411 data_systems = load_system_data(files)
412 if data_systems is None:
413 print("Failed to load data. Exiting.")
414 return
415
416 results = test_all_hypotheses(data_systems,
417 alpha=args.alpha)
418
419 # Helper function for JSON serialization of
420 # numpy types
421 def convert_numpy(obj):
422 if isinstance(obj, np.ndarray): return obj
423 .tolist()

```

```

398 if isinstance(obj, np.generic): return obj
399 .item()
400 if isinstance(obj, (np.bool_,)): return
401 bool(obj)
402 if isinstance(obj, dict): return {k:
403 convert_numpy(v) for k, v in obj.items()
404 }
405 if isinstance(obj, list): return [
406 convert_numpy(i) for i in obj]
407 # Handle NaNs gracefully for JSON output
408 if isinstance(obj, float) and np.isnan(obj):
409 return None # Represent NaN as null
410 return obj
411
412 try:
413 serializable_results = convert_numpy(
414 results)
415 with open(args.output, 'w') as f:
416 json.dump(serializable_results, f,
417 indent=2)
418 print(f"\nComplete results saved to {args.
419 output}")
420 except Exception as e:
421 print(f"Error saving results to JSON: {e}"
422)
423
424 if __name__ == "__main__":
425 main()

```

**Figure 6: Appendix: Distribution of Contextual Coherence (CC) scores for each RAG system.****Figure 7: Appendix: Distribution of Question Relevance (QR) scores for each RAG system.**

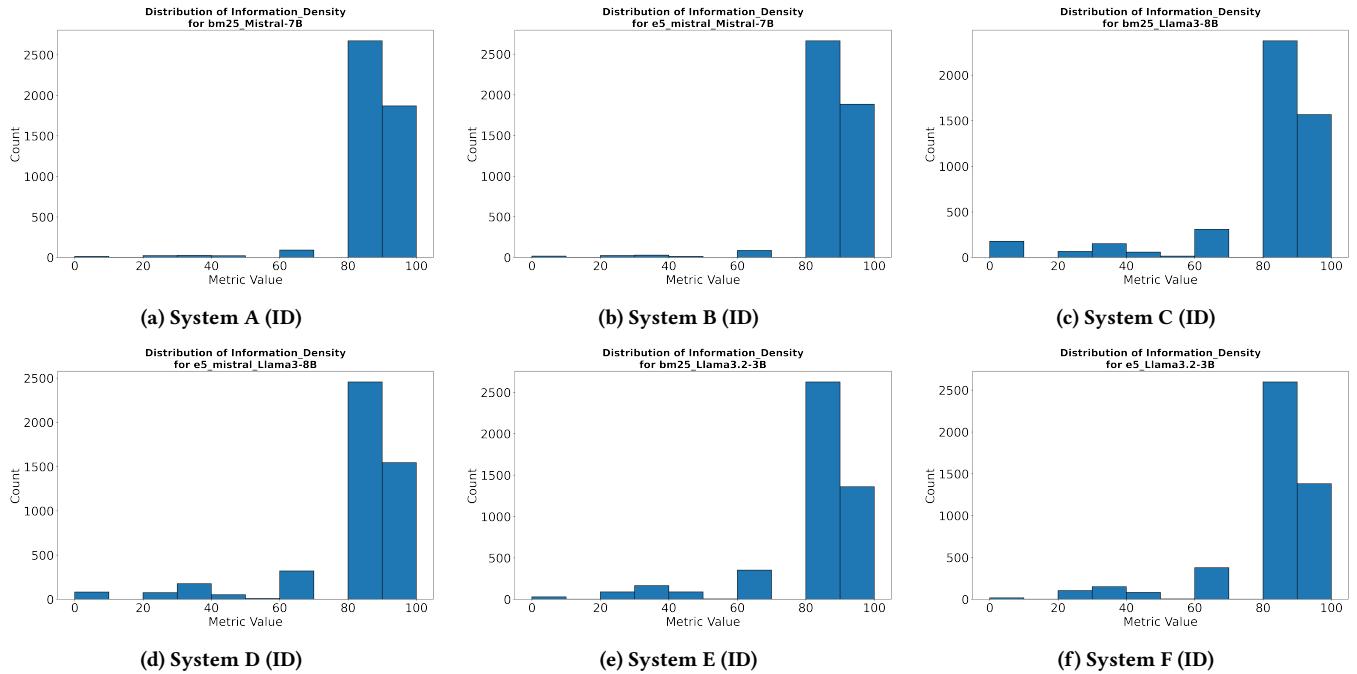


Figure 8: Appendix: Distribution of Information Density (ID) scores for each RAG system.

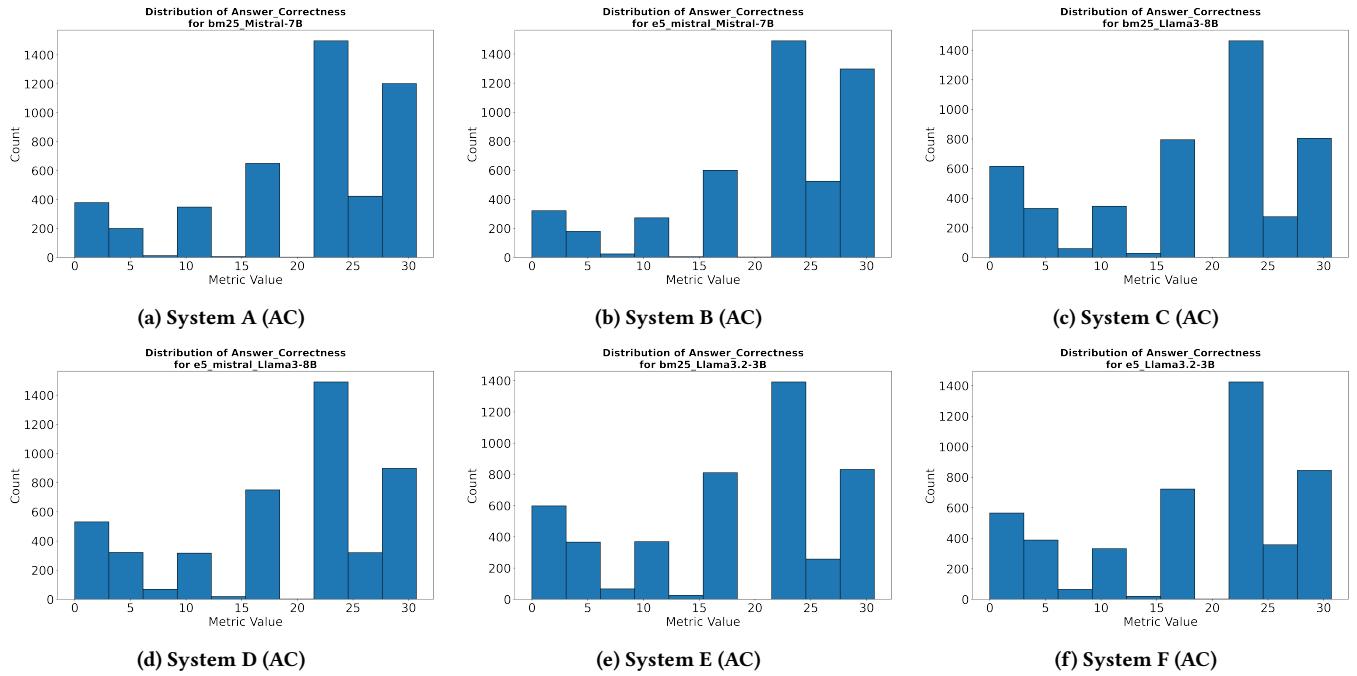
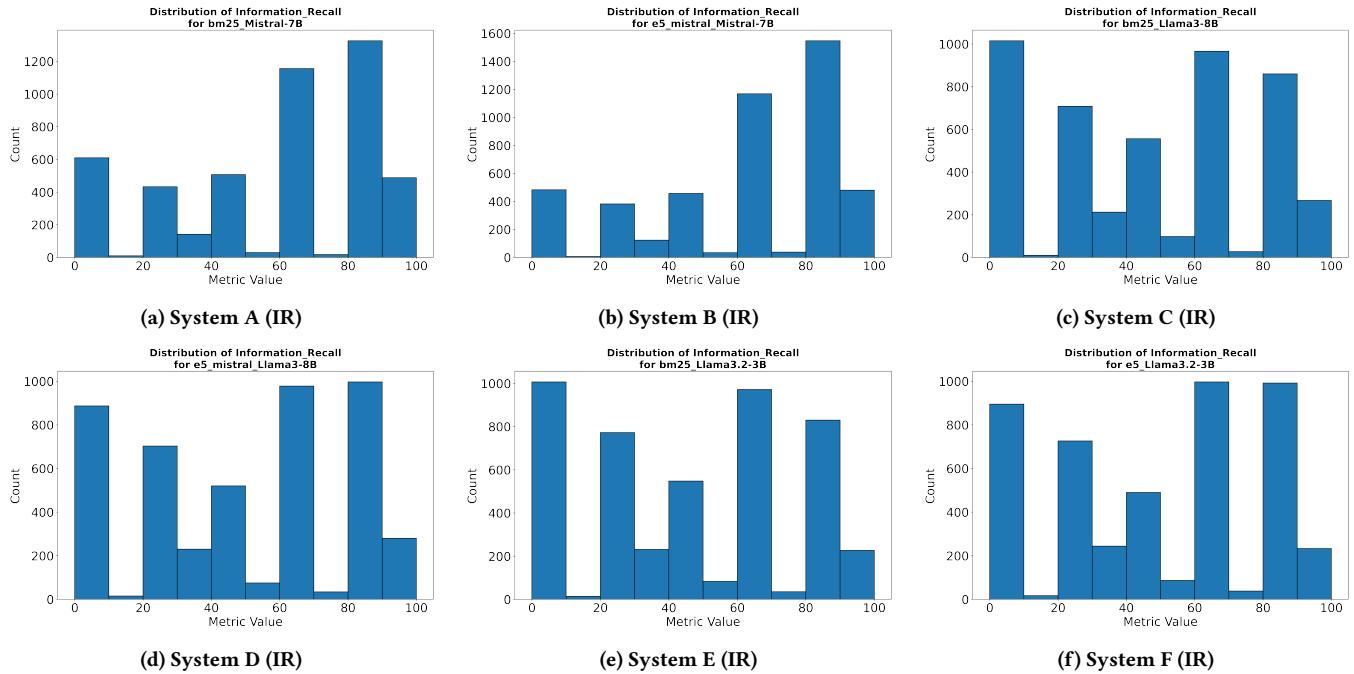


Figure 9: Appendix: Distribution of Answer Correctness (AC) scores for each RAG system.



**Figure 10: Appendix: Distribution of Information Recall (IR) scores for each RAG system.**