# The Great Nugget Recall: Automating Fact Extraction and RAG Evaluation with Large Language Models

Ronak Pradeep
University of Waterloo
Ontario, Canada

Nandan Thakur
University of Waterloo
Ontario, Canada

Shivani Upadhyay
University of Waterloo
Ontario, Canada

Daniel Campos
Snowflake
San Mateo, USA

Nick Craswell
Microsoft
Seattle, USA

Jimmy Lin
University of Waterloo
Ontario, Canada

## Abstract

Large Language Models (LLMs) have significantly enhanced the capabilities of information access systems, especially with retrieval-augmented generation (RAG). Nevertheless, the *evaluation* of RAG systems remains a barrier to continued progress, a challenge we tackle in this work by proposing an automatic evaluation framework that is validated against human annotations. We believe that the nugget evaluation methodology provides a solid foundation for evaluating RAG systems. This approach, originally developed for the TREC Question Answering (QA) Track in 2003, evaluates systems based on atomic facts that should be present in good answers. Our efforts focus on "refactoring" this methodology, where we describe the AutoNuggetizer framework that specifically applies LLMs to both automatically *create* nuggets and automatically *assign* nuggets to system answers. In the context of the TREC 2024 RAG Track, we calibrate a fully automatic approach against strategies where nuggets are created manually or semi-manually by human assessors and then assigned manually to system answers. Based on results from a community-wide evaluation, we observe strong agreement at the run level between scores derived from fully automatic nugget evaluation and human-based variants. The agreement is stronger when individual framework components such as nugget assignment are automated independently. This suggests that our evaluation framework provides tradeoffs between effort and quality that can be used to guide the development of future RAG systems. However, further research is necessary to refine our approach, particularly in establishing robust per-topic agreement to diagnose system failures effectively.

## 1 Introduction

This paper tackles the challenge of evaluating long-form responses from retrieval-augmented generation (RAG) systems built on top of large language models (LLMs), as part of the TREC 2024 RAG Track. There is, obviously, tremendous excitement and interest in RAG, but we feel that the *evaluation* of RAG output remains deficient from many perspectives. Furthermore, the lack of standardized evaluations presents a barrier to continued progress in information access, and more broadly, NLP as well as AI.

Our central hypothesis is that the nugget evaluation methodology [45, 46] provides a solid foundation for evaluating RAG systems. This methodology was originally articulated more than two decades ago in the context of the TREC Question Answering (QA) Track for evaluating answers to free-form "definition questions". This matches our RAG setting, where for a given query, there are atomic facts (i.e., nuggets) from different documents that a system must synthesize into a fluent and cohesive natural language response.

Given this starting point, our efforts have focused on "refactoring" the original nugget evaluation methodology in light of LLMs. Specifically, we leverage LLMs to both automatically *create* nuggets (nugget creation) and automatically *assign* nuggets (nugget assignment) to system-generated answers. We implement our nugget evaluation methodology in the AutoNuggetizer framework. While this is not the first attempt to automate nugget evaluations [21, 22], the advent of LLMs provides opportunities that did not exist before.

We explored the following research questions:

**RQ1** For RAG systems, how well does our fully automatic nugget evaluation framework correlate with different manual nugget evaluation strategies (that vary in human involvement)?

**RQ2** Does automating only nugget assignment result in stronger agreement with manual evaluations compared to fully automating the entire evaluation (including nugget creation)?

**RQ3** Are there any noticeable differences between nugget assignments by humans versus LLMs?

In this work, we analyzed diverse runs from the TREC 2024 RAG Track, a community-wide evaluation that attracted participation from dozens of teams around the world. Our findings can be summarized as follows:

- For **RQ1**, we find that scores from our fully automatic (end-to-end) nugget evaluation framework show strong correlations with manual nugget evaluations at the run level. This suggests that our approach can potentially serve as a good surrogate for manual evaluations in assessing RAG systems.

- For **RQ2**, we find that automating only nugget assignment leads to stronger agreement with manual evaluations, compared to a fully automated evaluation where nuggets are constructed automatically.

- For **RQ3**, our analyses suggest that LLM assessors appear to be more strict than NIST assessors in nugget assignment. Additionally, the use of LLMs to provide "draft" nuggets in the nugget creation step does not appear to noticeably increase alignment with human nugget assignment.

The biggest advantage of our approach is minimal reinvention of the wheel, in that we can leverage the existing body of work that has gone into exploring the nugget evaluation methodology, e.g., [7, 23, 24]. For aspects that are not directly impacted by our use of LLMs, we can continue to assert findings from the literature without needing to carefully validate those claims again. Furthermore, and unique to the TREC setup, we calibrate our fully automatic

evaluation results against alternative strategies that involve NIST assessors to varying extents, representing different design choices in an overall evaluation framework (e.g., fully manual nugget creation vs. post-editing automatically generated nuggets).

The contribution of this work is, to our knowledge, the first large-scale study of automating the nugget evaluation methodology using LLMs for RAG systems. We demonstrate that system rankings derived from our fully automatic nugget evaluation process strongly correlate with those obtained from manual or semi-manual evaluations by human assessors. This validation of automatic nugget creation and assignment represents an advance in the practical evaluation of RAG systems, offering a scalable and efficient alternative to labor-intensive manual assessments. Preliminary findings from our nugget evaluation were shared with the community in Pradeep et al. [32], but here we provide a much more thorough analysis of evaluation results.

It is important to emphasize that this work focuses exclusively on evaluating the recall of information nuggets within RAG answers. There are, of course, other important aspects of RAG evaluation, such as the assessment of citation support and the fluency of answers. These are important to the overall evaluation picture, but beyond the scope of this paper, as it is impossible to tackle everything all at once.

## 2 Background and Related Work

*RAG Systems and Benchmarks.* RAG in the context of LLMs has been used to generate entirely new, non-extractive answers to user queries [3, 12, 16, 20], where integrating external knowledge improves factual accuracy and mitigates issues such as hallucination [11, 18]. To better understand the nuances of these systems, RAG benchmarks have been developed to evaluate their effectiveness on different tasks, ranging from factoid QA [6, 50] to long-form answer generation [13, 49]. Additionally, RAG benchmarks have been developed to explore various sources: knowledge graphs [10, 30], multilingual documents [39, 40], coding [41, 48], structured tabular data [17], and images [15].

*LLM-as-a-Judge Evaluation.* Using LLM judges [51] in place of human judges has gained popularity for evaluating different aspects of RAG, such as factuality [19], retrieval quality [36], and robustness [38] for specific aspects such as contextual faithfulness [28]. Researchers have also taken advantage of LLMs to directly perform pairwise comparisons between RAG responses [13, 31, 33] and have suggested incorporating grading rubrics with strong LLMs [47].

Existing RAG frameworks assess responses along multiple dimensions. For example, RAGAs [8] introduced several metrics to evaluate faithfulness, answer relevance, and context relevance in RAG, leveraging LLMs (such as ChatGPT) for automatic evaluation based on custom prompts. ARES [35] demonstrated that fine-tuning an open-source LLM with a small amount of training data suffices to produce an effective automatic judge. Related, eRAG [36] utilizes the LLM in the RAG system itself to evaluate retrieval effectiveness.

*Nugget Evaluation Methodology.* Introduced in the TREC QA Track in 2003 [46], the nugget evaluation methodology emphasized the identification of *nuggets*, or essential facts relevant to an answer. Voorhees [45] demonstrated the importance of constructing nuggets to evaluate definition questions, thereby aligning complex question answering with retrieval-augmented generation. In both cases, queries require coherent answers drawn from information present in multiple relevant documents.

Further studies have refined and extended the basic nugget evaluation methodology. For example, nuggets served as the basis for Summarization Content Units (SCUs) [27] to evaluate summaries in the Document Understanding Conferences. Further refinements for question answering evaluation were introduced in later work [7, 23, 24]. Lin and Demner-Fushman [21, 22] explored automated approaches for nugget evaluation. Pavlu et al. [29] introduced a nugget-based methodology for evaluating retrieval, addressing scalability and reusability issues inherent in traditional document-level relevance judgments.

Recently, the advent of LLMs has enabled researchers and practitioners to perform more reliable nugget evaluation automatically, an observation noted in concurrent papers [1, 25, 41]. In this work, we incorporate these ideas into the TREC 2024 RAG Track.

*Related Nugget Evaluation Frameworks.* We can identify several attempts to build RAG evaluation frameworks around the notion of nuggets. For example, FactScore [26] proposed to evaluate system-generated biographies from Wikipedia in terms of the percentage of atomic facts that are supported. Farzi and Dietz [9] introduced the LLM-based RUBRIC framework, which evaluates the quality of system answers based on a rubric. In a similar vein, Arabzadeh and Clarke [2] used *subtopics*, which are corpus-independent predicates generated by an LLM, to evaluate relevance, and Mayfield et al. [25] used question and answer pairs as nuggets, manually identified by assessors. Concurrently, Thakur et al. [41] used an LLM to generate nuggets that decompose reference answers, which are then used to assess document relevance. As most of these methods were primarily validated against traditional *ad hoc* retrieval tasks, it is unclear if they can be directly applied to RAG answer evaluation in our context.

In contrast, our nugget evaluation framework evaluates how many nuggets—gathered from a pool of relevant documents—are present in a system answer. Our experiments with AutoNuggetizer are grounded in the MS MARCO V2.1 segment collection, reflecting real-world search more accurately than (for example) Wikipedia alone (the target corpus of many other studies). Crucially, our study includes correlations against high-quality human annotations from NIST assessors, similar to the established TREC QA Tracks, and features a "fresh" set of queries (reducing the risk of data pollution). The human annotations painstakingly gathered in the TREC 2024 RAG Track provide valuable data for validating and calibrating automated LLM-based nugget creation and assignment.

## 3 The Setup of the TREC 2024 RAG Track

The TREC 2024 RAG Track comprises three distinct but interconnected tasks: Retrieval (R), Augmented Generation (AG), and full Retrieval-Augmented Generation (RAG). Participants were given 301 queries (topics); the system task is to return, for the AG and RAG tasks, well-formed answers for each individual query (up to a maximum of 400 words). The Retrieval (R) task can be viewed as an intermediate product in a full RAG pipeline.

| |
|---|
| **Query**: how did african rulers contribute to the triangle trade |
| **Answer**: African rulers played a significant role in the triangular trade by capturing and supplying slaves to European traders. They actively participated in the trade by capturing individuals from inland areas and transporting them to the coast, where they were sold to European traders. [ … ] In summary, African rulers contributed to the triangular trade by capturing and supplying slaves to European traders, driven by the economic benefits they received in exchange. Their involvement was crucial for the success and expansion of the transatlantic slave trade. |

**Table 1: Sample answer from the TREC 2024 RAG Track for topic 2024-35227. Note that here we have purposely omitted citations in the answer since this paper does not cover support evaluation. The answer has been elided for length considerations.**

Throughout this paper, we use topic 2024-35227 "how did african rulers contribute to the triangle trade" as a running example. A system-generated answer is provided in Table 1; in this case, the answer was generated using GPT-4o [31]. We have purposely omitted citations from this answer: actual submissions took the form of structured JSON data wherein systems explicitly attempted to link each answer sentence to passages in the corpus that (purportedly) provide evidence for the assertions made in the sentence. This paper does *not* consider the evaluation of support, i.e., attempts to assess the veracity of the citations.

The Retrieval (R) task adopts a standard *ad hoc* retrieval setup, where systems are tasked with returning ranked lists of relevant passages from the MS MARCO V2.1 segment collection, comprised of 113 million passages [31, 32, 42]. The top-ranked passages retrieved for each query can then serve as grounding in the prompt input for LLM-based answer-generation systems. In the Augmented Generation (AG) task, the organizers provided the participants with a fixed list of 100 retrieved results from which to generate their answers. In the retrieval-augmented generation (RAG) task, participants were free to do whatever they wished in terms of retrieval from the corpus. This paper focuses on analyzing the results of system-generated answers from only the AG and RAG tasks.

Our evaluation used queries (topics) sourced from Bing search logs, specifically selecting non-factoid queries that require comprehensive, multi-faceted, and potentially subjective responses [34]. The set of test queries was constructed near the submission deadline (late July 2024), rather than reusing queries that had already been publicly released in Rosset et al. [34]. This strategy primarily addresses concerns about relevance judgment contamination, though we note that the corpus is likely to be already contained in LLM pretraining data given its web-based nature and open-source availability. After additional curation by NIST, we arrived at the final test set of 301 queries.

## 4 The Nugget Evaluation Methodology

The nugget evaluation methodology comprises two main steps, which we paraphrase from Voorhees [46], developed for the TREC 2003 QA Track:

(1) The assessor first creates a list of "information nuggets", where an information nugget is defined as a fact for which the assessor can make a binary decision as to whether a response contains the nugget. At the end of this step, the assessor decides which nuggets are vital—nuggets that must appear in a response for that response to be good.

(2) The assessor then proceeds to the second step once the nugget list is created. In this step, the assessor goes through each of the system answers and marks (whether or not) each nugget appears in the response.

Our AutoNuggetizer framework represents a "refactoring" of this decades-old nugget evaluation methodology, incorporating the latest advances in LLMs. While there have been previous attempts to automate nugget evaluations [21, 22], this effort represents the first since the advent of modern LLMs, which provide capabilities that did not previously exist. In more detail, the nugget evaluation methodology comprises two main steps:

**Nugget Creation.** This corresponds to the first step in the nugget evaluation methodology described by Voorhees [46], often dubbed "nuggetization". These nuggets capture elements of what should be in a good answer, divided into "vital" and "okay" categories. "Vital" nuggets are those, as Voorhees [46] articulates, must be present in a good response, while "okay" nuggets are "good to have", but are not absolutely necessary.

Nuggets must be created from an input set of documents under consideration; the input set is a design choice that we detail below. In the original implementation of the evaluation methodology from 2003, NIST assessors manually formulated these nuggets based on documents in the pool that they assessed to be relevant. That is, nugget creation followed relevance assessment (via pooling). In our AutoNuggetizer framework, we have attempted to automate this entire process, using UMBRELA [43] to generate relevance assessments for the pool and LLMs to generate what we call Auto-Nuggets (see Section 4.1.1).

It is important to note that while nuggets manifest as short natural language phrases or sentences, they are formulated at the semantic or conceptual level, and thus may or may not correspond to phrases or other lexical realizations in the sources.

**Nugget Assignment.** This corresponds to the second step in the nugget evaluation methodology described by Voorhees [46]. After the nuggets have been created, the list can be treated like an answer key. The assessor then reads the answer of each system to perform nugget assignment, which is a determination of whether each nugget appears in the response. That is: Does a system answer contain this particular nugget? In our AutoNuggetizer framework, nugget assignment is performed automatically by an LLM (see Section 4.2.1).

It is important to note that the nugget assignment process is performed at the semantic or conceptual level, not merely based on lexical matching (and hence requires understanding, inference, and

**Automatic nugget creation (AutoNuggets) using UMBRELA qrels**

African rulers captured and sold slaves to Europeans (vital)
African rulers waged wars to capture more slaves (vital)
African rulers exchanged slaves for firearms (vital)
African rulers' involvement was crucial for the trade's scale (vital)
African rulers' trade led to increased internal slavery (okay)

**Automatic nugget creation (AutoNuggets) using NIST qrels**

African rulers captured and sold slaves to European traders (vital)
African rulers exchanged slaves for firearms and goods (vital)
African rulers' involvement was crucial for the transatlantic slave trade (vital)
African rulers' cooperation enabled large-scale slave trade (vital)
African rulers encouraged European traders to come to their ports (okay)

**Post-edited nuggets (AutoNuggets+Edits) using NIST qrels**

African rulers sold slaves to European traders (vital)
African rulers transported captives to coastal slave forts (vital)
African rulers formed alliances with European traders (vital)
African rulers' cooperation enabled large-scale slave trade (vital)
African rulers' actions had a lasting negative impact on Africa (okay)

**Table 2: Comparison of five sample nuggets for the query "how did african rulers contribute to the triangle trade" from three conditions: AutoNuggets using UMBRELA and NIST qrels, and AutoNuggets+Edits (i.e., after human post-editing).**

| Nuggets | Assignment |
|---|---|
| **Automatic nugget creation (AutoNuggets) using UMBRELA qrels** | **AutoAssign** |
| African rulers captured and sold slaves to Europeans (vital) | Support |
| African rulers waged wars to capture more slaves (vital) | Not Support |
| African rulers exchanged slaves for firearms (vital) | Partial Support |
| African rulers' involvement was crucial for the trade's scale (vital) | Support |
| African rulers' trade led to increased internal slavery (okay) | Partial Support |
| **Post-edited nuggets (AutoNuggets+Edits) using NIST qrels** | **ManualAssign** |
| African rulers sold slaves to European traders (vital) | Support |
| African rulers sold war captives to European traders (okay) | Not Support |
| African rulers sold criminals to European traders (okay) | Not Support |
| African rulers sold debtors to European traders (okay) | Not Support |
| African rulers' actions had a lasting negative impact on Africa (okay) | Not Support |

**Table 3: Comparison of nugget assignment across different nugget creation approaches for the answer in Table 1.**

reasoning). In particular, a nugget can be assigned to an answer (i.e., appears in the answer) even if there is no lexical overlap between the nugget itself and the system answer.

As a concrete example, the nuggets for topic 2024-35227 "how did african rulers contribute to the triangle trade" are shown in Table 2 (top). In this case, the nuggets were automatically created from documents considered by UMBRELA to be at least related to the topic (more details in Section 4.1.1).

After the nugget assignment process, we arrive at a record of which nuggets are found in which systems' answers. For the sample answer in Table 1, the outcome of this is shown in Table 3 (top). At this point, it is straightforward to compute various metrics to quantify the quality of a response and the overall score of a run. We refer interested readers to Voorhees [46] for details on how final scores were computed previously; we adopt a different approach to quantifying answer quality in our evaluation (see Section 4.4).

## 4.1 Nugget Creation

*4.1.1 Automatic Nugget Creation.* The AutoNuggetizer framework begins by extracting a list of nuggets, or atomic information units, from a set of input documents. This nugget creation process, dubbed

---

**Figure 1: Prompt for the iterative nuggetization at turn $i$.**

"nuggetization", characterizes the information that should be contained in a high-quality answer to the user query.

To perform nuggetization, we employ GPT-4o through the Azure endpoint. This process is run over all documents that are judged to be at least "related" to the query (grade $\geq 1$). Note that, depending on the actual evaluation condition (see below), these relevance judgments are either provided by NIST assessors (manual) or by UMBRELA [43] (automatic). Details about relevance assessments are discussed in Upadhyay et al. [42]. Examples of these differences are shown in Table 2, comparing the top vs. the middle portions of the table. There are minor differences, but overall the resulting nuggets are quite similar.

Using the prompt in Figure 1, the LLM iteratively updates a list of nuggets that collectively represent the key facts required to fully answer the query, conditioned on the provided contexts (passages). The first iteration for each query begins with an empty list. Our prompt design encourages the model to produce comprehensive and diverse nuggets, ensuring broad coverage of different aspects of the user's information need. This iterative approach allows for the generation of a rich set of nuggets, capturing both explicit and implicit information requirements derived from the query. Auto-Nuggetizer aims to generate nuggets that are neither too broad nor too specific. Informed by previous implementations of nugget evaluations, we limited generation to at most 30 nuggets.

Once we have generated a set of nuggets for a query, the next step is to assess the importance of each nugget. We again use GPT-4o: the LLM is asked to assign an importance label of either "vital" or "okay" (see Section 4) using the prompt shown in Figure 2.

At this point, we sort the nuggets in descending order of importance and select the first 20 nuggets. This approach usually trims a few "okay" nuggets and, less frequently, some "vital" nuggets (when

**Figure 2: Prompt for determining the importance of nuggets. At each turn, at most 10 nuggets are passed to the LLM.**

there are over 20 of them). Note that these nuggets are ordered by decreasing importance, as imposed by the prompt in Figure 1. The resulting nugget set, we dub AutoNuggets.

*4.1.2 Semi-Manual Nugget Creation.* In the condition that we denote as AutoNuggets+Edits, NIST assessors post-edit nuggets that have been proposed by AutoNuggetizer as a "rough draft". Here, we start with the set of documents that have been manually judged by NIST assessors to be at least "related" to the query (grade ≥ 1), which serves as the input to AutoNuggetizer to create nuggets automatically (per above). Note that in this case, the input set of documents has already been judged by humans to be at least related to the user's query, unlike with the UMBRELA labels, which are generated automatically.

The automatically generated nuggets that we provide to NIST assessors are prepared in a slightly different way: for each query, we created a set of 30 nuggets without any indication of importance. This over-generation is deliberate. These nuggets are then edited by NIST assessors, who may add, eliminate, or combine nuggets. The NIST assessors perform this task by concurrently considering the list of relevant documents for that topic. Based on our logs, this process takes roughly an hour per topic on average, which suggests that the NIST assessors are not merely "rubber stamping" automatically generated nuggets, but are actually carefully deliberating over how they are formulated.

An example of this process is shown in Table 2 for our running example. On the top, we have the automatically generated nuggets using automatically generated UMBRELA relevance judgments; in the middle, we have automatically generated nuggets using manually generated relevance labels; and on the bottom, we have post-edited nuggets generated using manually generated relevance labels. Despite some differences, resulting nuggets from the different approaches are quite similar semantically.

*4.1.3 Manual Nugget Creation.* In the condition that we denote as ManualNuggets, NIST assessors examine the documents that have been manually judged to be at least "related" to the query (grade ≥ 1) to create a set of at most 20 nuggets with their importance. This followed the same procedure as the TREC 2003 QA Track; manual nugget creation is known to be a time-consuming and

**Figure 3: Prompt for nugget assignment. At each turn, at most 10 nuggets are passed to the LLM.**

labor-intensive task. Each topic required substantial assessor effort to thoroughly analyze the relevant passages and to synthesize the contents of those passages into a concise set of nuggets. This took roughly 2.5 hours per topic, based on our measurements.

## 4.2 Nugget Assignment

*4.2.1 Automatic Nugget Assignment.* The final component of our AutoNuggetizer framework, AutoAssign, automatically assigns nuggets to systems' answers. We adopt a listwise approach to nugget assignment, where the LLM is used to analyze an answer and determine if each nugget is fully supported (support), partially supported (partial_support), or not supported (not_support) by the answer. Again, we employ GPT-4o through the Azure endpoint, with the prompt shown in Figure 3. We iteratively prompt the model with at most 10 nuggets to evaluate assignment at each step.

*4.2.2 Manual Nugget Assignment.* In the manual nugget assignment process, we leveraged the expertise of the original NIST assessor involved in creating the list of nuggets (i.e., the same assessor performed nugget creation *and* nugget assignment). The assessor examines each answer, assigning each nugget one of three labels based on the extent of its support: support, partial_support, and not_support (same as with automatic assignment). It is worth clarifying that assessors in this step are not shown any LLM output, so this is a fully manual process.

By involving the same assessor responsible for nugget creation, we ensure consistency in nugget interpretation. This continuity helps maintain reliability across the nugget assignment process, as the assessor applies the original context and intent behind each nugget to the assessment process.

An example is shown in Table 3: on the top, we show fully automatic assignment (AutoAssign) of automatically created nuggets (AutoNuggets) using UMBRELA qrels; on the bottom, we show manual assignment (ManualAssign) of post-edited nuggets (AutoNuggets+Edits) using NIST qrels. The answer being evaluated is the one shown in Table 1.

## 4.3 Evaluation Strategies

Within our general AutoNuggetizer framework, we can instantiate a number of variants combining different design choices for nugget creation (Section 4.1) and nugget assignment (Section 4.2). In this work, we considered the following combinations:

(1) Automatic nugget creation with manual post-editing and manual nugget assignment (AutoNuggets+Edits / ManualAssign). Nuggets are initially generated automatically from passages judged relevant by NIST assessors and then refined by human assessors. Nugget assignment is performed manually by NIST assessors. This represents semi-manual nugget creation with manual assignment.

(2) Manual nugget creation and assignment (ManualNuggets / ManualAssign). NIST assessors manually create nuggets from relevant passages and manually assign them to system answers. This is a fully manual evaluation baseline, mirroring the original TREC QA approach.

(3) Fully automatic nugget creation and assignment (AutoNuggets / AutoAssign). Nuggets are automatically generated from documents assessed relevant by UMBRELA, and nugget assignment is performed automatically using AutoAssign. This represents fully automated, end-to-end evaluation.

(4) Semi-manual nugget creation with automatic nugget assignment. In this condition, we use the nuggets from conditions (1) and (2), which involve human assessors in nugget creation (AutoNuggets+Edits and ManualNuggets, respectively). However, we replace manual nugget assignment with automatic nugget assignment using AutoAssign, i.e., AutoNuggets+Edits / AutoAssign and ManualNuggets / AutoAssign, respectively. This condition isolates the impact of automating only the nugget assignment step.

## 4.4 Scoring

At this point in the evaluation, we have already completed nugget creation and nugget assignment. For each query, we have a list of nuggets, and for each system answer, we have a record of which nuggets it contains, in terms of a three-way judgment: support, partial_support, and not_support (either manual or automatic nugget assignment).

The final step is to compute the score for a system answer to a query $q$. The score of a run is simply the mean of the score for each query. We compute the following scores per query:

*All Strict ($A_{strict}$).* We calculate a score based on all nuggets in an answer, but with strict nugget matching. For a given nugget $i$:

$$p_i = \begin{cases} 1 & \text{if assignment = support} \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

The "All Strict" score is then calculated as:

$$A_{\text{strict}} = \frac{\sum_i p_i}{N_{\text{nuggets}}}$$

*Vital Strict ($V_{strict}$).* This score, used as the primary evaluation metric, applies the same strict matching criterion as $p_i$ above, but only over the vital nuggets, $p_i^{\text{v}}$. The score is calculated as:

$$V_{\text{strict}} = \frac{\sum_i p_i^{\text{v}}}{N_{\text{vital nuggets}}}$$

We emphasize that the scores presented here focus exclusively on evaluating the recall of nuggets in RAG answers. This analysis does not encompass other crucial aspects of RAG evaluation, such as the assessment of citation support or fluency, which are important but are outside the scope of this paper.

## 5 Results

### 5.1 Descriptive Statistics

For the TREC 2024 RAG Track, NIST received 93 runs from 20 groups for the RAG task and 53 runs from 11 groups for the AG task. Given resource constraints, the NIST annotators were able to evaluate only the two highest priority submissions from each group across the RAG and AG tasks, which translates into 31 runs from 18 groups for RAG and 14 runs from 9 groups for AG.

To be clear, our analyses are over 56 topics that have been fully judged across the 31 RAG runs and 14 AG runs discussed above. Of these, 36 topics were under the AutoNuggets+Edits / ManualAssign condition and 20 topics were under the ManualNuggets / ManualAssign condition. Descriptive statistics for nugget creation and assignment are presented in Table 4, broken down by the various experimental conditions. Naturally, since condition (3) in Section 4.3, corresponding to row (4) in the table, is fully automated, we were able to evaluate all topics, whereas for the other conditions we were constrained by available resources.

Examining the nugget characteristics, we observe that conditions employing AutoNuggets, rows (2a), (2b), and (4), tend to generate a higher average number of nuggets per topic (around 18) compared to conditions using AutoNuggets+Edits or ManualNuggets (around 14). Interestingly, the average nugget length remains relatively consistent across all conditions, hovering around 7–8 tokens, suggesting a stable granularity in nugget extraction regardless of the automation level in nugget creation. Furthermore, humans tend to be more conservative in assigning the label of "vital" importance, likely due to a stricter evaluation of what constitutes essential information within the context of a topic.

Focusing on nugget assignment, a clear distinction emerges between manual and automatic assignment. Conditions employing manual assignment, (1a) and (1b), exhibit a substantially lower percentage of nuggets categorized as "Partial Support". In contrast, the automatic assignment conditions, especially the directly comparable ones, (3a) and (3b), see higher percentages. This shift suggests that automatic assignment tends to distribute some human-deemed "Support" assignments to "Partial Support", potentially indicating a different interpretation or application of support criteria compared to manual assignment. However, this behavior could simply be the result of the current prompt that we are using.

### 5.2 Fully Automatic Nugget Evaluation

To address **RQ1**, we examine the correlation between our fully automatic nugget evaluation methodology and variants that require manual intervention to various degrees, as detailed in conditions

| | Condition | # topics | avg nuggets | avg nugget length | % nuggets | | % assigned | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | vital | okay | NS | PS | S |
| (1a) | AutoNuggets+Edits / ManualAssign | 36 | 14.1 | 7.0 | 61.1 | 38.9 | 57.9 | 6.6 | 35.6 |
| (1b) | ManualNuggets / ManualAssign | 20 | 13.7 | 8.5 | 59.0 | 41.0 | 62.9 | 3.5 | 33.6 |
| (2a) | AutoNuggets / AutoAssign (1a subset) | 36 | 18.2 | 6.8 | 66.2 | 33.8 | 50.9 | 22.0 | 27.0 |
| (2b) | AutoNuggets / AutoAssign (1b subset) | 20 | 18.7 | 7.2 | 70.3 | 29.7 | 47.1 | 24.9 | 28.0 |
| (3a) | AutoNuggets+Edits / AutoAssign | 36 | 14.1 | 7.0 | 61.1 | 38.9 | 54.1 | 18.7 | 27.2 |
| (3b) | ManualNuggets / AutoAssign | 20 | 13.7 | 8.5 | 59.0 | 41.0 | 57.8 | 17.6 | 24.6 |
| (4) | AutoNuggets / AutoAssign (Overall) | 301 | 18.7 | 6.8 | 72.5 | 27.5 | 50.9 | 23.6 | 25.5 |

**Table 4: Descriptive statistics for nugget creation and assignment under different conditions.**
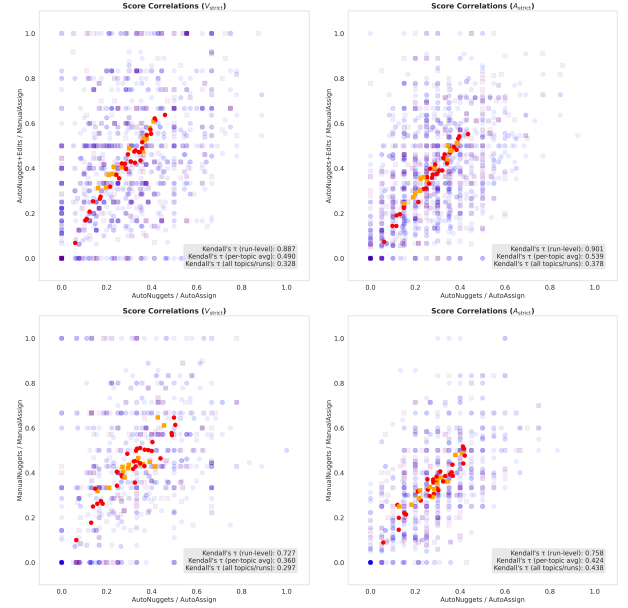
(1) and (2) in Section 4.3. We reiterate that the AutoNuggets / Auto-Assign condition is fully automatic, requiring no manual intervention (specifically, nuggetization is based on automatically generated relevance labels).

Here, we adopt the meta-evaluation methodology common in the information retrieval literature dating back several decades [4, 5, 14, 37, 44, 52]. We evaluate runs both with our fully automatic nugget evaluation methodology and with the other variants involving humans. Rank correlations (in terms of Kendall's $\tau$) between the system orderings induced by the different approaches capture the quality of our automatic metric. Our goal is to obtain high rank correlations against human-based metrics.

Figure 4 presents scatter plots of scores derived from the various experimental conditions. In the left column, we show $V_{strict}$ and in the right column, $A_{strict}$. In all these plots, the $x$-axes show the Auto-Nuggets / AutoAssign score. The top row shows the comparison against the AutoNuggets+Edits / ManualAssign condition, where nuggets are automatically generated and post-edited by humans, while assignment is manual. The bottom row shows the comparison against the ManualNuggets / ManualAssign condition, where both nugget creation and assignment are fully manual. Red circles and orange squares represent run-level scores for RAG and AG tasks, respectively. Blue circles (RAG runs) and purple squares (AG runs) represent all topic/run combinations for the same.

We compute Kendall's $\tau$ in three different ways: (1) run-level correlations (i.e., we first compute run scores by averaging topic scores, and then compute correlations); (2) average of per-topic correlations (per-topic avg), where we compute correlations per topic, and then average those; (3) correlation across all topic/run combinations (all topics/runs), where each topic/run combination is considered an independent observation, and we compute Kendall's $\tau$ across all these observations.
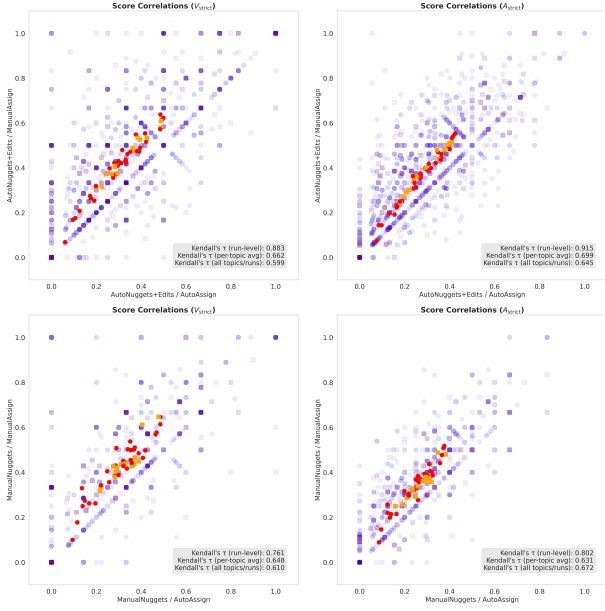
For the AutoNuggets+Edits / ManualAssign condition (top row), we observe a strong run-level Kendall's $\tau$ correlation of 0.887 for $V_{strict}$ and 0.901 for $A_{strict}$. In IR meta-evaluations, these would be sufficient to "validate" a metric. For the ManualNuggets / Manual-Assign condition (bottom row), the run-level Kendall's $\tau$ correlations drop to 0.727 for $V_{strict}$ and 0.758 for $A_{strict}$. The correlations are slightly lower, but still good. This can be partly explained by the much smaller topic count of 20: experimenting with a 20 topic subset of AutoNuggets+Edits results in correlations of 0.826 for $V_{strict}$ and 0.838 for $A_{strict}$. This drop could perhaps be explained by fully manual nuggets "looking different" due to the assessors



**Figure 4: Scatter plots between manual vs. automatic $V_{strict}$ and $A_{strict}$ scores for AG and RAG runs. The $x$ axes show scores from AutoNuggets / AutoAssign and the $y$ axes show scores from different manual conditions. Red circles (RAG runs) and orange squares (AG runs) represent run-level scores. Blue circles (RAG runs) and purple squares (AG runs) show all topic/run combinations. The bottom-right box reports Kendall's $\tau$ correlations at the run level (red circles/orange squares), over all topic/run combinations (blue circles/purple squares), average of per-topic correlations.**

being anchored by the automatically generated nuggets (in the post-editing condition). Nevertheless, results suggest good agreement between automatic and manual metrics at the run level.

However, examining the scatter of blue circles and purple squares representing individual topic/run combinations in all the scatter plots, we observe a greater dispersion compared to the run-level scores. The Kendall's $\tau$ correlation across all topic/run combinations is much lower, ranging from 0.297 to 0.438. The average per-topic Kendall's $\tau$ correlations fall to between 0.360 and 0.539. This indicates that while run-level rankings are well-preserved by our automatic evaluation, there is more variability at the topic level.

Figure 5: Scatter plots showing correlations designed to answer RQ2, isolating the effects of nugget assignment. Overall organization is identical to Figure 4.

In summary, to answer **RQ1**:

**RQ1** We find that scores from our fully automatic (end-to-end) nugget evaluation framework show strong correlations with manual nugget evaluations at the run level. This suggests that our approach can potentially serve as a good surrogate for manual evaluations in assessing RAG systems.
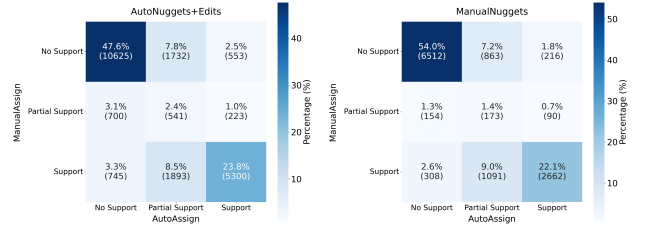
We observe stronger correlations with AutoNuggets+Edits than ManualNuggets, but Kendall's $\tau$ remains quite high even with fully manual nuggets. We emphasize that while we achieve good *run-level* correlations, *per-topic* correlations are quite low, suggesting that our evaluation framework is inadequate for fine-grained debugging of individual answers.

## 5.3 Automatic Nugget Assignment

To investigate **RQ2**, we analyze whether automating nugget assignment in isolation, an expensive annotation step, leads to stronger agreement with manual evaluations compared to fully automating the entire framework. That is, we hold the nuggetization approach constant, but vary the assignment method.

Figure 5 compares scores with automatic vs. manual nugget assignment. The top row shows the AutoNuggets+Edits / AutoAssign condition against AutoNuggets+Edits / ManualAssign, isolating the effect of automating assignment when nuggets are post-edited by humans. The bottom row compares ManualNuggets / AutoAssign against ManualNuggets / ManualAssign, isolating the effect of automating assignment when nuggets are manually created.

Comparing the top row of Figure 5 to the top row of Figure 4, we observe that automating only nugget assignment (AutoNuggets-+Edits / AutoAssign vs. AutoNuggets+Edits / ManualAssign) continues to show strong agreement at the run level. Additionally, it



Figure 6: Confusion matrices showing agreement between automatic and manual nugget assignment strategies.

results in a higher per-topic average and all topics/runs Kendall's $\tau$ correlations for both $V_{\text{strict}}$ and $A_{\text{strict}}$. Specifically, for $V_{\text{strict}}$, the per-topic average Kendall's $\tau$ increases from 0.490 in the fully automatic (AutoNuggets / AutoAssign in Figure 4) to 0.662 when only nugget assignment is automated (AutoNuggets+Edits / AutoAssign in Figure 5). Similarly, the all topics/runs Kendall's $\tau$ improves from 0.328 to 0.599. For $A_{\text{strict}}$, we observe comparable improvements. These metrics are arguably more sensitive to the nuances of topic-level evaluation, as they capture agreement across individual queries rather than just at the aggregate run level. The substantial increase in correlations suggests that automating nugget assignment while maintaining human-curated nuggets is perhaps a way forward to evaluate a large number of system responses, given that the assignment step occupies a lot of the assessment budget.

When we examine the bottom row of Figure 5, where nuggets are manually created, we see a similar pattern of improvement. Comparing ManualNuggets / AutoAssign to ManualNuggets / ManualAssign, we observe increases in per-topic average and all topics/runs Kendall's $\tau$ correlations for both $V_{\text{strict}}$ and $A_{\text{strict}}$. This further reinforces the finding that automating nugget assignment, regardless of whether nuggets are automatically generated and post-edited or manually created from scratch, leads to a better alignment with manual evaluations, particularly at the topic level.

To answer **RQ2**:

**RQ2** We find that automating only nugget assignment leads to stronger agreement with manual evaluations, compared to a fully automated evaluation where nuggets are constructed automatically.

Our framework is not an "all or nothing" proposition. Instead, there is a vast design space that supports different tradeoffs between effort and quality at the run- and topic-level.

## 5.4 Nugget Assignment Confusion Matrices

To address **RQ3**, we investigate quantitative differences between human and LLM nugget assignment. Figure 6 presents confusion matrices that quantify agreements (and disagreements) in nugget assignment labels between automatic (AutoAssign) and manual (ManualAssign) approaches. The left matrix compares AutoAssign against ManualAssign when using AutoNuggets+Edits nuggets, and the right matrix compares them when using ManualNuggets. Each matrix shows the percentages and counts of nugget assignments falling into different categories.

Examining the diagonal elements of both confusion matrices, we observe substantial agreement between automatic and manual

nugget assignment, particularly for the "Support" and "No Support" labels. For AutoNuggets+Edits nuggets (left matrix), 47.6% of the (response, nugget) pairs resulted in both AutoAssign and Manual-Assign labeling a "No Support" and 23.8% a "Support". For Manual-Nuggets (right matrix), we see 54.0% of the distribution with "No Support" and 22.1% with "Support". Partial support appears to be a source of disagreement in both conditions, as evidenced by only a small mass of the distribution contained in the central element of the matrices.

Off-diagonal elements reveal disagreements. In both matrices, we observe a tendency for AutoAssign to assign "Partial Support" more frequently than ManualAssign. In general, we see a slightly larger mass of the distribution in the lower triangle part of both matrices compared to the upper triangle, indicating that LLM assessors are a bit stricter compared to NIST assessors.

The use of LLMs to provide "draft" nuggets in the prior nuggetization step (AutoNuggets+Edits vs. ManualNuggets) does not appear to drastically alter the alignment across labeling in nugget assignment, as the confusion matrices for both nugget types exhibit similar patterns of agreement and disagreement.

To answer **RQ3**:

**RQ3** Our analyses suggest that LLM assessors appear to be more strict than NIST assessors in nugget assignment. Additionally, the use of LLMs to provide "draft" nuggets in the nugget creation step does not appear to noticeably increase alignment with human nugget assignment.

Two points of caution here, though: These findings are specific to the current implementation. For example, it is entirely possible that a different LLM or a different prompt might elicit different behavior. Another gap in our current understanding is the relationship between inter-annotator agreement and LLM–human differences. This analysis only compares *one* set of nugget assignments between an assessor and a particular LLM implementation. We do not know the variations exhibited by different human assessors and different LLMs, and how they compare, which would be an interesting question to tackle in the future.

## 6  Conclusions

While we are certainly not the first to propose a RAG evaluation methodology, we view our efforts as having two main distinguishing characteristics: First, by building on the nugget evaluation methodology dating back over two decades, we minimize reinvention of the wheel. The information retrieval literature has a long tradition of deliberate and careful meta-evaluations that validate evaluation methodologies, and much work has examined different aspects of nugget evaluations. For aspects of the evaluation that are not dependent on LLMs, we can simply build on existing findings. Second, we have demonstrated that our AutoNuggetizer framework, which leverages LLMs for automatic nugget creation and assignment, achieves strong run-level agreement with manual and semi-manual evaluations, suggesting its viability as a scalable alternative for assessing RAG systems.

Under the AutoNuggets / AutoAssign condition in our Auto-Nuggetizer framework, we are able to provide evaluation scores for *all* runs submitted to the TREC 2024 RAG Track, across all 301 topics—at only the cost of LLM inference, while doing the same

with human assessors remains prohibitively expensive. Our insights reveal that automating nugget assignment in isolation yields stronger agreement, especially at the per-topic level, with manual evaluations than with full framework automation, highlighting the potential of a hybrid approach. Furthermore, we find that LLMs apply stricter assignment labels than human assessors, indicating the need for calibration. Our framework has the potential to enable rapid iteration on improving RAG systems with varying degrees of human involvement, while providing some confidence that the automatically generated metrics have some degree of correlation to answer quality as determined by human assessors. For RAG, we can now potentially climb hills quickly and in a meaningful way!

## References

[1] Marwah Alaofi, Negar Arabzadeh, Charles LA Clarke, and Mark Sanderson. 2024. Generative information retrieval evaluation. In *Information Access in the Era of Generative AI*. Springer, 135–159.

[2] Negar Arabzadeh and Charles LA Clarke. 2024. A Comparison of Methods for Evaluating Generative IR. *arXiv:2404.04044* (2024).

[3] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, et al. 2022. Improving Language Models by Retrieving from Trillions of Tokens. In *Proceedings of the 39th International Conference on Machine Learning (ICML 2022)*. Baltimore, Maryland, 2206–2240.

[4] Chris Buckley and Ellen M. Voorhees. 2000. Evaluating Evaluation Measure Stability. In *Proceedings of the 23rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2000)*. Athens, Greece, 33–40.

[5] Chris Buckley and Ellen M. Voorhees. 2004. Retrieval Evaluation with Incomplete Information. In *Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2004)*. Sheffield, United Kingdom, 25–32.

[6] Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking Large Language Models in Retrieval-Augmented Generation. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI-24)*. Vancouver, Canada, 17754–17762.

[7] Hoa Trang Dang and Jimmy Lin. 2007. Different Structures for Evaluating Answers to Complex Questions: Pyramids Won't Topple, and Neither Will Human Assessors. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics (ACL 2007)*. Prague, Czech Republic, 768–775.

[8] Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. 2024. RAGAs: Automated Evaluation of Retrieval Augmented Generation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*. St. Julians, Malta, 150–158.

[9] Naghmeh Farzi and Laura Dietz. 2024. Pencils Down! Automatic Rubric-based Evaluation of Retrieve/Generate Systems. In *Proceedings of the 2024 ACM SIGIR International Conference on Theory of Information Retrieval (ICTIR '24)*. Washington, D.C., 175–184.

[10] Yanlin Feng, Simone Papicchio, and Sajjadur Rahman. 2024. CypherBench: Towards Precise Retrieval over Full-scale Modern Knowledge Graphs in the LLM Era. *arXiv:2412.18702* (2024).

[11] Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling Large Language Models to Generate Text with Citations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Singapore, 6465–6488.

[12] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. REALM: Retrieval-Augmented Language Model Pre-training. In *Proceedings*

*of the 37th International Conference on Machine Learning (ICML 2020)*. 3929–3938.

[13] Rujun Han, Yuhao Zhang, Peng Qi, Yumo Xu, Jenyuan Wang, Lan Liu, William Yang Wang, Bonan Min, and Vittorio Castelli. 2024. RAG-QA Arena: Evaluating Domain Robustness for Long-form Retrieval Augmented Question Answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Miami, Florida, 4354–4374.

[14] Donna Harman. 2011. *Information Retrieval Evaluation*. Morgan & Claypool Publishers.

[15] Wenbo Hu, Jia-Chen Gu, Zi-Yi Dou, Mohsen Fayyaz, Pan Lu, Kai-Wei Chang, and Nanyun Peng. 2024. MRAG-Bench: Vision-Centric Evaluation for Retrieval-Augmented Multimodal Models. *arXiv:2410.08182* (2024).

[16] Gautier Izacard and Edouard Grave. 2021. Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. 874–880.

[17] Xingyu Ji, Aditya Parameswaran, and Madelon Hulsebos. 2024. TARGET: Benchmarking Table Retrieval for Generative Tasks. In *Proceedings of the NeurIPS 2024 Third Table Representation Learning Workshop*.

[18] Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through Memorization: Nearest Neighbor Language Models. In *Proceedings of the 8th International Conference on Learning Representations (ICLR 2020)*. Addis Ababa, Ethiopia.

[19] Satyapriya Krishna, Kalpesh Krishna, Anhad Mohananey, Steven Schwarcz, Adam Stambler, Shyam Upadhyay, and Manaal Faruqui. 2024. Fact, Fetch, and Reason: A Unified Evaluation of Retrieval-Augmented Generation. *arXiv:2409.12941* (2024).

[20] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems 33*. 9459–9474.

[21] Jimmy Lin and Dina Demner-Fushman. 2005. Automatically Evaluating Answers to Definition Questions. In *Proceedings of the 2005 Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP 2005)*. Vancouver, Canada, 931–938.

[22] Jimmy Lin and Dina Demner-Fushman. 2006. Methods for Automatically Evaluating Answers to Complex Questions. *Information Retrieval* 9, 5 (2006), 565–587.

[23] Jimmy Lin and Dina Demner-Fushman. 2006. Will Pyramids Built of Nuggets Topple Over?. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*. New York, New York, 383–390.

[24] Jimmy Lin and Pengyi Zhang. 2007. Deconstructing Nuggets: The Stability and Reliability of Complex Question Answering Evaluation. In *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2007)*. Amsterdam, the Netherlands, 327–334.

[25] James Mayfield, Eugene Yang, Dawn Lawrie, Sean MacAvaney, Paul McNamee, Douglas W. Oard, Luca Soldaini, Ian Soboroff, Orion Weller, Efsun Kayi, Kate Sanders, Marc Mason, and Noah Hibbler. 2024. On the Evaluation of Machine-Generated Reports. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2024)*. Washington, D.C., 1904–1915.

[26] Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Singapore, 12076–12100.

[27] Ani Nenkova and Rebecca Passonneau. 2004. Evaluating Content Selection in Summarization: The Pyramid Method. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*. Boston, Massachusetts, 145–152.

[28] Xuan-Phi Nguyen, Shrey Pandit, Senthil Purushwalkam, Austin Xu, Hailin Chen, Yifei Ming, Zixuan Ke, Silvio Savarese, Caiming Xong, and Shafiq Joty. 2024. SFR-RAG: Towards Contextually Faithful LLMs. *arXiv:2409.09916* (2024).

[29] Virgil Pavlu, Shahzad Rajput, Peter B. Golbus, and Javed A. Aslam. 2012. IR System Evaluation using Nugget-based Test Collections. In *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining (WSDM 2012)*. Seattle, Washington, 393–402.

[30] Ronak Pradeep, Daniel Lee, Ali Mousavi, Jeffrey Pound, Yisi Sang, Jimmy Lin, Ihab Ilyas, Saloni Potdar, Mostafa Arefiyan, and Yunyao Li. 2024. ConvKG-Yarn: Spinning Configurable and Scalable Conversational Knowledge Graph QA Datasets with Large Language Models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*. Miami, Florida, 1176–1206.

[31] Ronak Pradeep, Nandan Thakur, Sahel Sharifymoghaddam, Eric Zhang, Ryan Nguyen, Daniel Campos, Nick Craswell, and Jimmy Lin. 2025. Ragnarök: A Reusable RAG Framework and Baselines for TREC 2024 Retrieval-Augmented Generation Track. In *Proceedings of the 47th European Conference on Information Retrieval (ECIR 2025), Part I*. Lucca, Italy, 132–148.

[32] Ronak Pradeep, Nandan Thakur, Shivani Upadhyay, Daniel Campos, Nick Craswell, and Jimmy Lin. 2024. Initial Nugget Evaluation Results for the TREC 2024 RAG Track with the AutoNuggetizer Framework. *arXiv:2411.09607* (2024).

[33] Zackary Rackauckas, Arthur Câmara, and Jakub Zavrel. 2024. Evaluating RAG-Fusion with RAGElo: an Automated Elo-based Framework. *arXiv:2406.14783* (2024).

[34] Corby Rosset, Ho-Lam Chung, Guanghui Qin, Ethan C. Chau, Zhuo Feng, Ahmed Awadallah, Jennifer Neville, and Nikhil Rao. 2024. Researchy Questions: A Dataset of Multi-Perspective, Decompositional Questions for LLM Web Agents. *arXiv:2402.17896* (2024).

[35] Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. 2024. ARES: An Automated Evaluation Framework for Retrieval-Augmented Generation Systems. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. Mexico City, Mexico, 338–354.

[36] Alireza Salemi and Hamed Zamani. 2024. Evaluating Retrieval Quality in Retrieval-Augmented Generation. *arXiv:2404.13781* (2024).

[37] Mark Sanderson and Justin Zobel. 2005. Information Retrieval System Evaluation: Effort, Sensitivity, and Reliability. In *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2005)*. Salvador, Brazil, 162–169.

[38] Shangeetha Sivasothy, Scott Barnett, Stefanus Kurniawan, Zafaryab Rasool, and Rajesh Vasa. 2024. RAGProbe: An Automated Approach for Evaluating RAG Applications. *arXiv:2409.19019* (2024).

[39] Nandan Thakur, Luiz Bonifacio, Crystina Zhang, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Boxing Chen, Mehdi Rezagholizadeh, and Jimmy Lin. 2024. "Knowing When You Don't Know": A Multilingual Relevance Assessment Dataset for Robust Retrieval-Augmented Generation. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. Miami, Florida, 12508–12526.

[40] Nandan Thakur, Suleman Kazi, Ge Luo, Jimmy Lin, and Amin Ahmad. 2025. MIRAGE-Bench: Automatic Multilingual Benchmark Arena for Retrieval-Augmented Generation Systems. *arXiv:2410.13716* (2025).

[41] Nandan Thakur, Jimmy Lin, Sam Havens, Michael Carbin, Omar Khattab, and Andrew Drozdov. 2025. FreshStack: Building Realistic Benchmarks for Evaluating Retrieval on Technical Documents. *arXiv:2504.13128* (2025).

[42] Shivani Upadhyay, Ronak Pradeep, Nandan Thakur, Daniel Campos, Nick Craswell, Ian Soboroff, Hoa Trang Dang, and Jimmy Lin. 2024. A Large-Scale Study of Relevance Assessments with Large Language Models: An Initial Look. *arXiv:2411.08275* (2024).

[43] Shivani Upadhyay, Ronak Pradeep, Nandan Thakur, Nick Craswell, and Jimmy Lin. 2024. UMBRELA: UMbrela is the (Open-Source Reproduction of the) Bing RELevance Assessor. *arXiv:2406.06519* (2024).

[44] Ellen M. Voorhees. 1998. Variations in Relevance Judgments and the Measurement of Retrieval Effectiveness. In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 1998)*. Melbourne, Australia, 315–323.

[45] Ellen M. Voorhees. 2003. Evaluating Answers to Definition Questions. In *Companion Volume of the Proceedings of HLT-NAACL 2003 — Short Papers*. Edmonton, Canada, 109–111.

[46] Ellen M. Voorhees. 2003. Overview of the TREC 2003 Question Answering Track. In *Proceedings of the Twelfth Text REtrieval Conference (TREC 2003)*. Gaithersburg, Maryland.

[47] Yang Wang, Alberto Garcia Hernandez, Roman Kyslyi, and Nicholas Kersting. 2024. Evaluating Quality of Answers for Retrieval-Augmented Generation: A Strong LLM Is All You Need. *arXiv:2406.18064* (2024).

[48] Zora Zhiruo Wang, Akari Asai, Xinyan Velocity Yu, Frank F. Xu, Yiqing Xie, Graham Neubig, and Daniel Fried. 2025. CodeRAG-Bench: Can Retrieval Augment Code Generation? *arXiv:2406.14497* (2025).

[49] Fangyuan Xu, Yixiao Song, Mohit Iyyer, and Eunsol Choi. 2023. A Critical Evaluation of Evaluations for Long-form Question Answering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Toronto, Canada, 3225–3245.

[50] Xiao Yang, Kai Sun, Hao Xin, Yushi Sun, Nikita Bhalla, Xiangsen Chen, Sajal Choudhary, Rongze Daniel Gui, Ziran Will Jiang, Ziyu Jiang, Lingkun Kong, Brian Moran, Jiaqi Wang, Yifan Ethan Xu, An Yan, Chenyu Yang, Eting Yuan, Hanwen Zha, Nan Tang, Lei Chen, Nicolas Scheffer, Yue Liu, Nirav Shah, Rakesh Wanga, Anuj Kumar, Wen tau Yih, and Xin Luna Dong. 2024. CRAG – Comprehensive RAG Benchmark. *arXiv:2406.04744* (2024).

[51] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. In *Advances in Neural Information Processing Systems 36 (NeurIPS 2023) Datasets and Benchmarks Track*. New Orleans, Louisiana, 46595–46623.

[52] Justin Zobel. 1998. How Reliable Are the Results of Large-Scale Information Retrieval Experiments?. In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 1998)*. Melbourne, Australia, 307–314.