

mTRAG: A Multi-Turn Conversational Benchmark for Evaluating Retrieval-Augmented Generation Systems

Yannis Katsis*, Sara Rosenthal*, Kshitij Fadnis, Chulaka Gunasekara, Young-Suk Lee, Lucian Popa, Vraj Shah, Huaiyu Zhu, Danish Contractor, Marina Danilevsky

IBM Research, USA

yannis.katsis@ibm.com, sjrosenthal@us.ibm.com, kpfadnis@us.ibm.com, chulaka.gunasekara@ibm.com, ysuklee@us.ibm.com, lpopa@us.ibm.com, vraj@ibm.com, huaiyu@us.ibm.com, Danish.Contractor@ibm.com, mdanile@us.ibm.com

Abstract

Retrieval-augmented generation (RAG) has recently become a very popular task for Large Language Models (LLMs). Evaluating them on *multi-turn* RAG conversations, where the system is asked to generate a response to a question in the context of a preceding conversation, is an important and often overlooked task with several additional challenges. We present mTRAG, an end-to-end human-generated multi-turn RAG benchmark that reflects several real-world properties across diverse dimensions for evaluating the full RAG pipeline. mTRAG contains 110 conversations averaging 7.7 turns each across four domains for a total of 842 tasks. We also explore automation paths via synthetic data and LLM-as-a-Judge evaluation. Our human and automatic evaluations show that even state-of-the-art LLM RAG systems struggle on mTRAG. We demonstrate the need for strong retrieval and generation systems that can handle later turns, unanswerable questions, non-standalone questions, and multiple domains. mTRAG is available at <https://github.com/ibm/mt-rag-benchmark>.

1 Introduction

Large Language Models (LLMs) play an important role as chat-based assistants (Lin et al., 2023). Relying on knowledge-sources during the conversation is an important task that helps improve answer reliability and trust, and hence retrieval-augmented generation (RAG) has become an important and popular field in recent years (Lewis et al., 2020; Chen et al., 2024). The primary focus of RAG benchmarks has been on single turn (Friel et al., 2024; Niu et al., 2024; Yang et al., 2024), which LLMs have become proficient

at (Kuo et al., 2024); however, multi-turn RAG, where a turn is defined as a question-response pair, has been largely overlooked as a benchmark, and presents additional challenges not covered in single-turn RAG. A multi-turn conversation benchmark should sufficiently cover several challenging aspects to effect a holistic evaluation of the full RAG pipeline:

Retrieval The relevant passages should change during the conversation causing repeated retrieval.

Generation The generator should struggle to answer many of the questions correctly, particularly questions that refer to and rely on previous turns.

We present mTRAG, a diverse and representative multi-turn RAG benchmark of human-generated conversations across 4 different domains that vary in style, topic and source. Our conversations comprise turns that vary along the dimensions of *question type*, *multi-turn*, and *answerability*.

Our benchmark is constructed using a novel process where human annotators simulate a real-world conversation, by actually interacting with a live RAG system via a custom chat application and improving the output in real time. Annotators took care to diversify their questions across the aforementioned different dimensions, including referencing earlier turns, while ensuring a flowing and natural conversation. At every turn, after issuing their questions, annotators checked the passages retrieved by the RAG system and modified the passage set to improve relevance and diversity. Next, they reviewed and repaired the generated response to improve its quality. Figure 1 shows part of a conversation from the benchmark to illustrate the output of our data creation process (described in detail in Section 4). The resulting

*These authors contributed equally to this work.

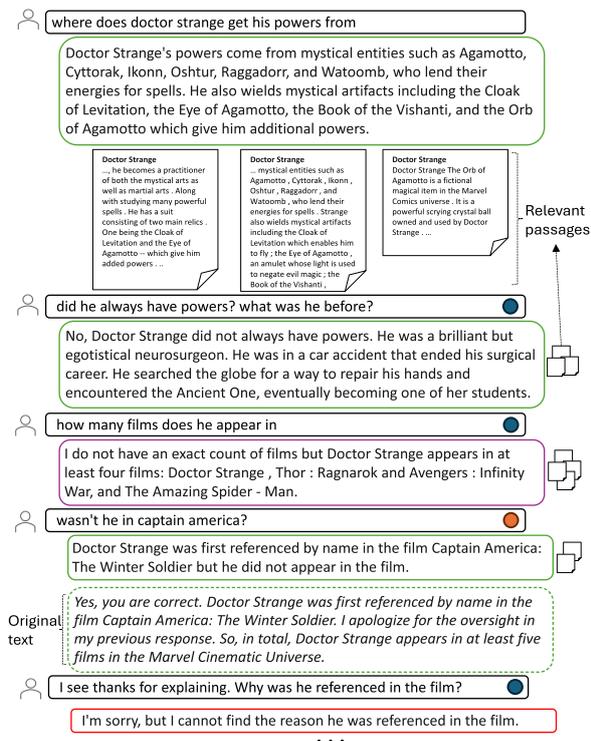


Figure 1: 5/8 turns of a conversation from the CLAP_{NO} domain. The conversation is enriched with question dimensions, passage diversity, and repair. The answerability is shown using the assistant response color: answerable, unanswerable, and partial. The multi-turn type is shown using the question circle: follow-up and clarification. The different relevant passages highlight diversity and the original text shows a repair from the model response.

conversations average 7–8 turns in length and 16.9 unique relevant passages per conversation.

We evaluate our mTRAG benchmark on the retrieval and generation components of RAG systems. We examine retrieval performance of lexical, sparse and dense retrieval under two settings (last turn and query rewrite) and analyze generative performance of 9 LLMs under three retrieval settings (reference, reference+RAG, and full RAG). Our rigorous human evaluation of model responses demonstrates that all models struggle on our tasks, especially on unanswerable questions, and in later turns.

As human data generation and evaluation does not scale well, we systematically explore automation paths. We identify several automated evaluation metrics that correlate with human scores (and many that do not), and demonstrate the need for more work in automatic evaluation. We also construct a companion benchmark, mTRAG-S, of synthetically generated conversa-

tions. By providing both human-generated and synthetic conversations over the same corpora, we aim to help the community analyze and understand the relative advantages of the two types of data.

Our contributions are as follows:

1. We present mTRAG, a comprehensive and diverse human-generated multi-turn RAG benchmark, accompanied by four document corpora. The benchmark is available at <https://github.com/ibm/mt-rag-benchmark>.
2. We evaluate lexical, dense, and sparse retrieval models and 9 large language models and conduct a rigorous human evaluation of model responses.
3. We systematically explore automation paths for data and metrics, demonstrating gaps in current automatic evaluation and constructing a companion benchmark, mTRAG-S, of synthetically generated conversations.

To the best of our knowledge, mTRAG is the first end-to-end human-generated multi-turn RAG benchmark that reflects real-world properties of multi-turn conversations.

2 Related Work

Question Answering (QA) and Information Retrieval (IR) have been popular tasks for many years. Prior focus has been on Extractive QA (Rajpurkar et al., 2016; Kwiatkowski et al., 2019), Long Form Question Answering (Fan et al., 2019; Stelmakh et al., 2022; Rosenthal et al., 2025), and Open Domain QA (Lee et al., 2019). There has also been work on multi-turn conversations, primarily through datasets like MT-Bench (Zheng et al., 2023), which focuses on tasks without retrieval, and Wizard of Wikipedia (Dinan et al., 2019), which is conversational rather than information seeking. Several surveys exist which summarize prior approaches and datasets available to the community (Bai and Wang, 2022; Cambazoglu et al., 2021; Wang, 2022; Rogers et al., 2023).

We compare our benchmark to existing multi-turn RAG datasets in Table 1. Except for FaithDial (Dziri et al., 2022), prior work keeps the retrieval component fixed. Retrieval is either performed once at the beginning, restricting the

Benchmark	Active Retrieval	Long Answer	Unanswerable	Multi-Domain
QuAC (Choi et al., 2018)	✗	✗	✓	✗
Or-QuAC (Qu et al., 2020)	✗	✗	✓	✗
CoQA (Reddy et al., 2019)	✗	✗	✓	✗
ShARC (Saeidi et al., 2018)	✗	✗	✗	✗
MD2Dial (Feng et al., 2021)	✗	✗	✗	✗
FaithDial (Dziri et al., 2022)	✓	✗	✓	✗
iKAT (Aliannejadi et al., 2024)	✗	✓	✗	✗
RADBench (Kuo et al., 2024)	✗	✓	✗	✓
mTRAG (Our work)	✓	✓	✓	✓

Table 1: mTRAG compared to prior multi-turn RAG benchmarks.

entire conversation to the initial passages; or retrieval is only deployed to find evidence for an existing conversation, which can result in strange mismatches between answers and verifiability. In contrast, we perform *active retrieval* (Jiang et al., 2023), where ongoing passage retrieval influences both follow-up questions and provided answers, more closely reflecting real life scenarios. Excluding RAD-Bench (Kuo et al., 2024) and iKAT (Aliannejadi et al., 2024), most prior datasets focus on extractive or short answers (1–2 sentences) limiting the kind of questions that can be asked. Further, many existing datasets ignore unanswerable questions - a ripe source of hallucinations in LLMs (Song et al., 2025). Finally, most datasets focus on a single domain/topic, while we explore several domains of different types. Our mTRAG benchmark reflects real-world properties of multi-turn conversations, including active retrieval, long-form answers, unanswerable questions, and multiple domains.

3 mTRAG Benchmark

We now describe mTRAG’s characteristics, before delving into how it was created in Section 4.

3.1 Dimensions

To ensure that mTRAG is representative of real RAG use cases, we designed it to be diverse across several important dimensions. For the detailed definitions, see Appendix A.

Question Types: Similar to contemporaneous work on datasets for RAG (Yang et al., 2024; Lee et al., 2024) conversations in mTRAG contain questions of diverse types, including *factoid*, *comparison*, *explanation*, *keyword* questions, and others. Each question has one or more question type labels.

Multi-turn: In addition to question types that apply to individual questions regardless of the surrounding conversation, there are different types of multi-turn questions. A question can be a *follow-up* or *clarification*. All questions in mTRAG beyond the first turn are labeled with one multi-turn type.

Answerability: Models often struggle when confronted with questions or problems that cannot be answered (Rajpurkar et al., 2018; Sun et al., 2024). Based on this observation, we include *answerable*, *partially answerable*, *unanswerable* based on the corpora, and *conversational* statements (e.g., “Hi”, “That’s interesting”, “Thank you”).

Domain: To test RAG systems over different types of documents, mTRAG is created over four domains and corresponding document corpora (discussed in more detail in Section 4.2).

3.2 Conversation Properties

To ensure that mTRAG is a challenging benchmark, the following properties were incorporated during conversation creation.

High Quality Responses: Following and extending prior work (Rosenthal et al., 2025; Es et al., 2024; Aliannejadi et al., 2024), each reference answer is written to satisfy the following properties: (a) *Faithfulness*: The answer is faithful to the passages or earlier turns, (b) *Appropriateness*: It is appropriate/relevant to the question, (c) *Naturalness*: The answer sounds natural, and (d) *Completeness*: Includes all information in the passages relevant to the question. We refer to these properties as *FANC*.

Passage Diversity: Conversations contain questions that are diverse enough that the relevant passages do not remain static through the conversation. Our conversations have on average 16.9 unique relevant passages and 20.9 relevant passages in total.

Answer Repair: During conversation creation we employ an LLM to create the initial response, which is then repaired by the annotator as needed. A response that does not require repair can be considered an indication that the question is not challenging for the LLM. Our conversations contain repairs on 92% of the turns.

Corpus	Domain	Provenance	Documents (D)	Passages (P)	Avg P/D
CLAP _{NQ}	Wikipedia	(Rosenthal et al., 2025)	4,293	183,408	42.7
FiQA	Finance	(Maia et al., 2018)	57,638*	61,022	1.1
Govt	Government	New corpus	7,661	49,607	6.5
Cloud	Technical documentation	New corpus	8,578	72,442	8.4

Table 2: Statistics of document corpora included in the mTRAG benchmark. (*) For FiQA we report the number of individual forum posts, as the dataset does not have the notion of a document.

Non-standalone: The questions occasionally rely on prior turns in the conversation (e.g., by using co-references to prior questions or answers or by clarifying previous turns). On average 1.3 questions per conversation include co-references.

4 Benchmark Creation

4.1 Annotators

The annotators who contributed to this work are highly skilled individuals hired solely to perform language annotation tasks and paid well above minimum wage. Unless otherwise noted, the annotators were used for all annotation tasks in this paper. Great care was taken to ensure random assignment.

The core task of creating conversations was carried out by 13 annotators in total. Each conversation typically took between 30 and 60 minutes to create, highlighting the complexity of the task and the need for highly skilled individuals.

4.2 Document Corpora

The first step in creating the benchmark was to assemble the document corpora over which the conversations would be built. mTRAG consists of four document corpora covering a diverse set of domains:¹

CLAP_{NQ} (Rosenthal et al., 2025): a subset of Wikipedia pages,

FiQA (Maia et al., 2018): a set of StackExchange posts discussing financial advice,

Govt: the crawled contents of select web-pages under the .gov and .mil domains, and

Cloud: the crawled contents of select technical documentation pages of a major cloud provider.

CLAP_{NQ} and FiQA are existing corpora from QA/IR datasets, while Govt and Cloud are new corpora assembled specifically for this benchmark.

The decision to create new corpora was motivated by the need to also have corpora specifically designed for creating diverse multi-turn conversations (i.e., conversations that touch several different passages). This was accomplished by designing the new Govt and Cloud corpora to contain sets of inter-connected pages (see Appendix B.1 for details).

Each corpus was indexed using Elser from ElasticSearch.² During ingestion, documents were split into passages of 512 tokens with an overlap stride of 100 tokens. The indexes are used to perform retrieval during conversation creation and all RAG experiments. Table 2 summarizes the corpora.

Finally, for each corpus we also assembled a set of seed questions to help human annotators bootstrap the conversation generation process, described in Section 4.3. For CLAP_{NQ} and FiQA we leveraged questions from the corresponding QA datasets, while for Govt and Cloud we selected a set of seed documents from each corpus and asked annotators to write seed questions based on them.

4.3 Human-generated Conversations

Annotators were asked to create multi-turn conversations over the four corpora. To aid in this process, we developed a custom chat application allowing them to interact with a live RAG system consisting of an ELSERV1 (ElasticSearch 8.10)¹ retriever and Mixtral 8×7b Instruct (Jiang et al., 2024) generator as the assistant and correcting the retriever and generator outputs as needed. Using the application, annotators created multi-turn conversations by performing the following actions at every turn: Write a question, adjust the set of retrieved passages (by discarding irrelevant passages and/or bringing in additional relevant passages using a custom passage search functionality), edit assistant response, and enrich

¹All corpora are publicly released as part of mTRAG.

²<https://www.elastic.co/guide/en/machine-learning/current/ml-nlp-elser.html>.

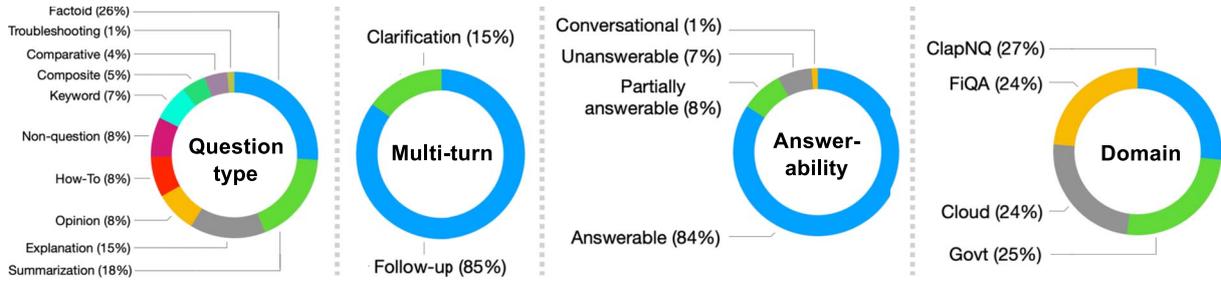


Figure 2: Distribution of tasks in mTRAG based on each of the benchmark’s dimensions.

with question dimensions (i.e., question type, multi-turn type, and answerability). Annotators were directed to write 6–9 (or more) turns per conversation.

(i) **Training:** Prior to the main creation task we provided annotators with interactive training sessions, extensive documentation with examples, and pilot tasks with feedback. Once they were comfortable with the task and producing conversations that met our criteria, we began the full annotation task.

(ii) **Creation:** Annotators started a conversation using a seed question. For subsequent turns, they were encouraged to create questions that naturally extended the preceding conversation while varying in answerability types, question types, and multi-turn patterns (Section 3.1), which were noted during annotation time. Once annotators wrote a question, the chat application queried the retriever for potentially relevant passages, ran the generator to produce a model response based on the retrieved passages, and presented everything to the user for repair. Annotators were guided to ensure passage diversity and answer quality (Section 3.2).

To ensure diversity, they had to make sure to ask questions that were not relevant to the same few passages throughout the conversation while still maintaining a naturally flowing conversation. To assist with this, they used the application’s custom passage search to increase the coverage of relevant passages. The annotators were allowed to use any query terms via the custom passage search to improve coverage and diversity, such as rewriting the query and searching for words from the answer.

To maintain high answer quality, the annotators were instructed to repair the answer to ensure it satisfied our properties (Section 3.2) using the edit functionality. In the final benchmark, the simi-

ilarity between the original and repaired response based on Rouge-L (denoted as the *edit score*) is 60.7 on average, indicating significant amount of repair. This novel process simulates real-time conversations, which is missing in prior work (Section 2).

(iii) **Review:** The resulting 126 conversations were then reviewed. Annotators could accept or reject conversations, and repair responses, passage relevance, and dimensions as needed. They were not allowed to edit the questions or passages as such changes could negatively affect the conversation flow. During this phase, most conversations were kept (see review process details in Appendix B.3). Our human evaluation (Section 7) shows that the reference responses are preferred by humans.

4.4 Data Statistics

This process yielded a benchmark of 110 conversations (29 ClapNQ, 27 FiQA, 28 Govt, 26 Cloud) with an average 7.7 turns per conversation (min = 3 turns, max = 12 turns), leading to 842 conversational turns in total. For evaluation purposes, for each conversational turn we create a corresponding *evaluation task* (or simply *task*). The task, which corresponds to a turn, contains the entire conversation so far together with the last user question (e.g., the task created for turn k includes all user and assistant questions/responses for the first $k - 1$ turns plus the user question for turn k). All evaluations described next are performed at the task level. Figure 2 shows the distribution of the 842 tasks in mTRAG on the benchmark dimensions.

5 Retrieval

We use mTRAG to evaluate both the retrieval and generation components of RAG systems. Retrieval evaluation is described in this section,

		Recall				nDCG			
		@1	@3	@5	@10	@1	@3	@5	@10
BM25	Last Turn	0.08	0.15	0.20	0.27	0.17	0.16	0.18	0.21
	Query Rewrite	0.09	0.18	0.25	0.33	0.20	0.19	0.22	0.25
BGE-base 1.5	Last Turn	0.13	0.24	0.30	0.38	0.26	0.25	0.27	0.30
	Query Rewrite	0.17	0.30	0.37	0.47	0.34	0.31	0.34	0.38
Elser	Last Turn	0.18	0.39	0.49	0.58	0.42	0.41	0.45	0.49
	Query Rewrite	0.20	0.43	0.52	0.64	0.46	0.45	0.48	0.54

Table 3: Retrieval performance of models on our benchmark using Recall and nDCG metrics.

while generator evaluation is presented in the next section. Note that this is not intended to be a comprehensive evaluation of all retrievers and generators. Instead our goal is to demonstrate the challenging nature of MTRAG and how it can help evaluate and surface issues in state-of-the-art RAG systems.

5.1 Experimental Setup

We use MTRAG with the indexed corpora from the creation stage (Section 4.2) to evaluate lexical (BM25), dense (BGE-base 1.5) (Xiao et al., 2024), and sparse (Elser) retrieval. The competitiveness of these models is shown on the MTEB leaderboard.³

We report the commonly used Recall and nDCG metrics (Thakur et al., 2021) @1, 3, 5, 10. Since we use Elser for retrieval during data creation, there may be some biases towards Elser.

The retrieval task is performed on the reference passages (i.e., those marked as relevant during the creation process⁴). It is therefore only computed on answerable and partially answerable tasks where reference passages exist.

5.2 Retrieval Strategies

We experimented with several strategies to query the retriever for relevant passages, including sending the full conversation up to the current user turn, all the user turns without the responses, subsets of the conversation, and only the last user turn. Using the full conversation, or even just a

few turns from it, consistently under-performed, often causing the retriever to bring back the same passages over and over. The most effective strategy was just using the last user turn. A particular challenge in multi-turn conversations is that a user turn may be non-standalone, employing shortcuts to express intent, or referencing entities or concepts from earlier in the conversation. To alleviate this we used a *query rewrite* strategy, also known as contextual query rewriting (Zhou et al., 2023; Sun et al., 2023) to rewrite the user turn using an LLM so that it incorporates all necessary parts from the context into an unambiguous, standalone question (see Appendix C.1 for implementation details). An example of query rewrite is shown below:

User: Who is the CEO of Apple Inc.?
Assistant: The CEO of Apple Inc. is TIM COOK.
User: its address?

[Rewriting] What is the address of Apple Inc?

5.3 Retrieval Results

We highlight the results in Tables 3 and 4. Table 3 shows that the query rewriting strategy consistently outperforms using only the last turn (without rewriting), across all metrics, for all models. Elser outperforms BM25 and BGE-base 1.5. Table 4 shows the Elser results with query rewriting broken down by domain, first turn vs later turn, and whether the question is standalone. Retrieval performance is significantly lower for later turns than for the first turn, and non-standalone questions continue to pose a challenge (though query rewriting helps). These results

³<https://www.elastic.co/search-labs/blog/elasticsearch-elser-relevance-mteb-comparison>.

⁴We do not distinguish the ranking for the relevant passages; all are considered equally ranked as 1 while all other irrelevant passages are ranked 0. It is also very likely that there is selection bias and some relevant passages are missed, a challenge in all retrieval datasets (Thakur et al., 2021).

	Subset	R@5
By Turn	Turn 1 (102)	0.89
	> Turn 1 (675)	0.47
By Standalone (> Turn 1)	Standalone (555)	0.48
	Non-Standalone (120)	0.42
By Domain	CLAPNQ (208)	0.56
	FiQA (180)	0.50
	Govt (201)	0.56
	Cloud (188)	0.47

Table 4: Elser retrieval results with query rewrite on subsets of the data to highlight multi-turn properties. Numbers in parentheses denote size of each subset.

highlight two key areas of improvement for retrieval components: 1) multi-turn retrieval, and 2) non-standalone questions.

6 Generation

We next present the generator experiments. We start with the experimental setup, followed by the results, using automated metrics including LLM judges. Section 7 will complement this with a human evaluation on a subset of MTRAG. Given a task, we send to the model the following information: the question, preceding turns, N passages, and instructions. We choose $N = 5$ passages because it achieves considerable improvement compared to top 3, while remaining a manageable amount of passages (Section 5). For more generation format details, see Appendix D.2.

6.1 Retrieval Settings

We evaluate how LLMs perform under three retrieval settings, simulating ideal/noisy retrieval.

Reference (●): Generation using reference passages or no passages if unanswerable/conversational. No retrieval is performed in this setting; it simulates a perfect retriever.

Reference + RAG (○): Partial retrieval followed by generation, where the reference passages are supplemented by the top retrieved passages (using Elser with rewrite) to yield a total of 5 passages. We restrict this to the 426 tasks that have at most two reference passages to ensure all passages needed for the reference are included. This can be

considered an upper bound, where the retrieval is successful but there are additional noisy passages.

Full RAG (○): Retrieval using Elser with rewrite followed by generation, where the top $N = 5$ passages are retrieved (the standard RAG setting).

6.2 Models

We evaluate the following auto-regressive models.

Llama 3.1 Models (Dubey et al., 2024): The Llama 3.1 family of models are instruction-tuned models that support up to 128K tokens. We evaluate the 8B, 70B, and 405B models.

Mixtral Mixture-of-Expert Models (Jiang et al., 2024): The instruction fine-tuned Mixtral $8 \times 22B$ model that supports up to 32K tokens.⁵

GPT-4o Models:⁶ We use the GPT-4o and GPT-4o-mini model in our experiments. These support context lengths of up to 128K tokens.

Command R+⁷: This is a 104B parameter multi-lingual model optimized for RAG and tool use.

Qwen 2.5 models (Team, 2024): We use the instruct versions of the 7B and 72B models. They support a context length of up to 128K tokens.

6.3 Metrics

We use three metrics, described below, to evaluate the quality of RAG systems and understand if a model response exhibits the desirable *FANC* properties outlined in Section 3.2. We use these metrics to provide insights into the trends of LLMs on multi-turn RAG rather than declare winners. See Section 8 for a discussion on how we picked the metrics and the challenges associated with finding good evaluation metrics for multi-turn RAG. Moreover, for implementation details, see Appendix F.

6.3.1 Aggregate Performance Metrics

The first two metrics, RB_{alg} and RB_{llm} , are both reference-based (*RB*) metrics. They both measure the aggregate performance of model responses by comparing the model response (*MR*)

⁵We use the 22B and not the 7B variant of Mixtral, due to the fact that the latter was used during conversation generation.

⁶<https://openai.com/index/gpt-4o-system-card/>.

⁷<https://huggingface.co/CohereForAI/c4ai-command-r-plus>.

to the reference answer (RA) utilizing different techniques:

\mathbf{RB}_{alg} is the harmonic mean of three algorithmic metrics (Adlakha et al., 2024): Bert-Recall (*Bert-Rec*), Bert-K-Precision (*Bert-K-Prec*), and *Rouge-L*. The intuition is as follows: Bert-Rec is an approximation for completeness as it measures the semantic overlap between MR and RA , tending to prefer longer answers. Bert-K-Prec compares MR to the passages, P , and is an approximation for faithfulness and completeness. Rouge-L measures whether phrases from RA are in MR and is an approximation for appropriateness.

\mathbf{RB}_{llm} is an LLM judge inspired by RAD-Bench (Kuo et al., 2024). We adapt RAD-Bench’s approach of comparing MR to RA but modify the prompt to add P and anchor the evaluation on the metrics of faithfulness, appropriateness, and completeness. To minimize model biases and improve evaluation reliability, we use several models as judges and use the median as the final score.

6.3.2 Faithfulness Metric

We also use a metric specifically for faithfulness, which is important for RAG applications and a challenge for LLMs (see Section 7). \mathbf{RL}_F , the Faithfulness LLM judge from RAGAS (Es et al., 2024), appears to be a good judge for faithfulness (see Section 8). In contrast to the other two metrics, this is a reference-less (RL) metric, as it does not rely on the reference answer.

6.3.3 Conditioning Metrics on Answerability

Prior to computing the metrics we employ an IDK (“I Don’t Know”) judge to detect whether the response has a full or partial answer. It is important to first determine whether a response is IDK because intuitively, words used to indicate not knowing the answer may not match the context; this is also reflected in the metrics, which were not designed to measure IDK correctly. Our IDK judge achieves an accuracy of over 97% (see Appendix F.3). We condition the metric score, $\phi \in \{\mathbf{RB}_{\text{alg}}, \mathbf{RB}_{\text{llm}}, \mathbf{RL}_F\}$, on answerability value A and IDK value IDK as shown on Table 5.

We define *answerability accuracy* (*Ans. Acc.*) of a model as the accuracy of IDK correctly predicting A . Finally, conversational questions also require special handling for evaluation. Since only 10 tasks fall in this category, we exclude

	$IDK = \text{no, partial}$	$IDK = \text{yes}$
$A = \text{yes, partial}$	ϕ	0
$A = \text{no}$	0	1

Table 5: Table showing how a metric score ϕ is conditioned on answerability value A and IDK value IDK .

them from the experiments, leaving their study as future work.

6.4 Generation: Evaluation Results

The overall results across the different retrieval settings are shown in Table 6.⁸ Within model families, the larger model does as well or better than its smaller counterparts across metrics. In general, GPT-4o and Llama 3.1 405B Instruct perform the best across metrics and settings. All models score significantly lower than the reference answer, indicating that there is still room for improvement in mRAG for all LLMs.

Comparing performance across retrieval settings, we see that the results degrade for the Reference-Based judges as the setting gets more challenging: $\bullet > \circ > \circ$, indicating the noise has an impact on generation. Exceptions to this are Llama 3.1 70B and 8B, where answerability accuracy improves in non-perfect settings. These models tend to respond with “I don’t know” too often in the perfect retrieval setting and they do a better job of providing an answer when there are more passages (even though that answer may be incorrect). Interestingly, Qwen 2.5 72B and Command-R+ are more competitive and closer to GPT-4o and Llama models in noisy settings.

Also note, that faithfulness, \mathbf{RL}_F , increases in the Reference+RAG (\circ) setting when there are more passages than the Reference (\bullet) setting; however this does not mean that the answer is more correct, but rather that it is copying from the additional passages. Finally, the faithfulness, \mathbf{RL}_F , decreases considerably for the reference answer in the RAG setting (\circ). This shows that the retrieval has missed relevant passages and that the reference answer is no longer a reference for comparing answers in this setting.

Focusing next on the Reference (\bullet) retrieval setting, we further explore how LLMs perform

⁸These results are limited to the 426 Reference + RAG (\circ) tasks to be consistent across retrieval settings. Appendix I shows that the trends persist throughout the benchmark.

	Ans. Acc.			RL _F			RB _{llm}			RB _{alg}		
	●	◐	○	●	◐	○	●	◐	○	●	◐	○
Reference	0.98	0.97	0.98	0.86	0.87	0.67	0.94	0.94	0.94	0.88	0.88	0.86
Command-R+ (104B)	0.86	<u>0.86</u>	0.87	0.69	0.71	0.66	0.66	0.62	0.59	0.43	0.40	0.38
GPT-4o	0.89	<u>0.86</u>	<u>0.86</u>	0.69	0.69	<u>0.65</u>	0.73	<u>0.68</u>	0.66	<u>0.46</u>	<u>0.40</u>	<u>0.38</u>
GPT-4o-mini	0.87	<u>0.86</u>	0.84	<u>0.66</u>	0.69	0.64	<u>0.72</u>	<u>0.68</u>	0.64	0.43	<u>0.40</u>	0.37
Llama 3.1 405B Instruct	0.87	<u>0.86</u>	0.85	0.69	<u>0.70</u>	<u>0.65</u>	0.70	<u>0.68</u>	0.63	0.47	0.42	0.39
Llama 3.1 70B Instruct	0.78	0.83	0.81	0.63	0.66	0.64	0.62	0.64	0.59	0.43	0.42	0.39
Llama 3.1 8B Instruct	0.71	0.75	0.74	0.50	0.51	0.53	0.54	0.56	0.54	0.36	0.33	0.34
Mixtral 8×22B Instruct	0.86	0.87	<u>0.86</u>	0.54	0.61	0.56	0.66	0.64	0.61	0.39	0.38	0.35
Qwen 2.5 (72B)	0.87	0.87	0.87	0.65	0.71	0.64	0.71	0.69	<u>0.65</u>	0.43	<u>0.40</u>	0.37
Qwen 2.5 (7B)	<u>0.88</u>	<u>0.86</u>	0.87	0.62	0.66	0.62	0.68	0.65	0.63	0.42	0.38	0.37

Table 6: Generation results by retrieval setting: Reference (●), Reference+RAG (◐), and RAG (○), w/ IDK conditioned metrics ($n = 426$). Per column, the best result is in **bold** and second best is underlined.

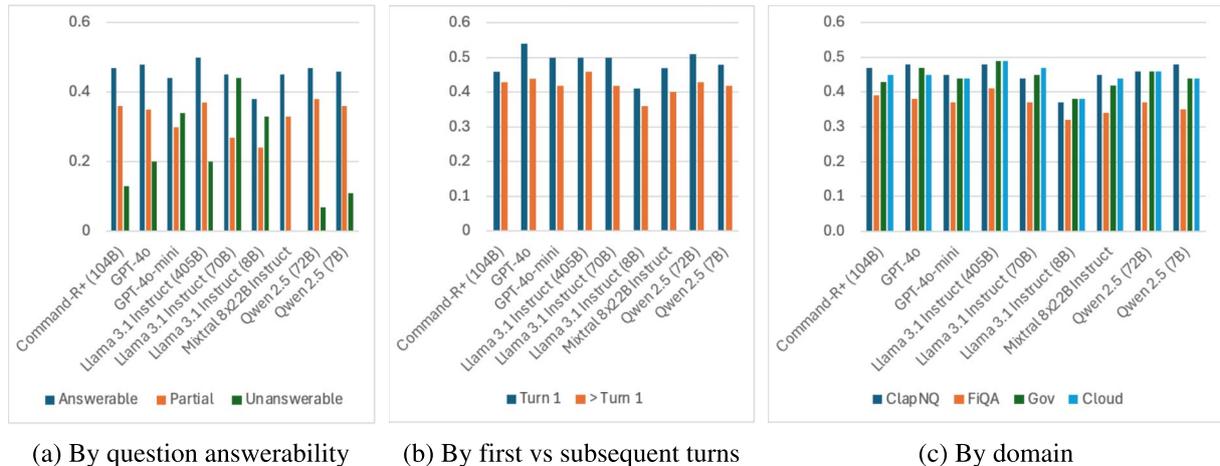


Figure 3: Generation results in the Reference (●) retrieval setting using a single metric, RB_{alg}, on three different dimensions: (a) answerability, (b) turns, and (c) domains.

on the different dimensions of mTRAG. Figure 3 shows the results. For space reasons we only report the RB_{alg} metric. We also leave out the breakdown by question types and multi-turn types, where we did not find interesting patterns. The additional metrics and breakdowns can be found in Appendix I.

Answerable vs Unanswerable Questions. As shown in Figure 3a, model performance generally drops when questions become partially answerable or unanswerable. In particular, models experience a dramatic drop in performance on unanswerables, struggling to declare that they do not know the answer when the answer is not included in the input passages. Interestingly, performance on unanswerables differs widely between model families. While GPT-4o and Llama 405B score low for unanswerables, they still perform much better than models in other families.

It is interesting to note that Llama 70B and 8B perform better on the unanswerables, because they say ‘‘I don’t know’’ too often, as discussed above.

First Turn vs Subsequent Turns. As shown in Figure 3b, in almost all cases the models perform better on first turn vs subsequent turn questions. This proves our conjecture that answering questions in a multi-turn setting is more challenging, as a model has to interpret a question in the context of the preceding conversation, which was one of the main motivations for this work.

Domains. Figure 3c shows that model performance is similar across domains except for FiQA where the results tend to be lower. We suspect this is due to the nature of the corpus; it contains posts from a financial discussion forum, which are typically short, very informal in style, and often subjective.

	Human Evaluation						Metrics		
	WR	F	A	N	C	All	RL _F	RB _{llm}	RB _{alg}
Ref.	59.3	3.8	3.8	4.0	3.9	3.8	0.86	0.97	0.87
GPT-4o	47.8	3.5	3.8	4.0	3.7	3.6	0.80	0.80	0.46
Llama3.1	47.4	3.5	3.9	3.9	3.7	3.6	0.78	0.78	0.49

Table 7: Results on the human evaluation for win-rate (WR), Faithfulness, Appropriateness, Naturalness, Completeness, and All (harmonic mean of FANC) compared to the IDK conditioned metrics. All results are reported on the Answerable subset.

7 Human Evaluation

We also performed a human evaluation on a subset of the benchmark. This serves several purposes: 1) Verifying that the reference responses are of high quality, 2) correlating our metrics with human judgment, and 3) analyzing frontier models.

We perform the human evaluation on two frontier models, GPT-4o and Llama 3.1 405B Instruct, and compare them to the reference answers. We select 5 conversations per domain for a total of 159 evaluation tasks. We ask annotators to measure the quality across the desirable response properties, FANC (Section 3.2), on a scale of 1 (Low) - 4 (High). We also ask them to perform a pair-wise comparison of the models, from which we calculate the *Win-Rate (WR)*. Ties are allowed but we ask them to do so sparingly.

The annotators completed the task using the Appen platform.⁹ Each evaluation task comprised the conversation thus far up to the last user question, the relevant passages, and the corresponding model responses (anonymized and in random order). Each task was performed by three skilled annotators. The annotator agreement was very high: {F: 89.6, A: 92.1, N: 95.6, C: 90.3, WR: 87.4} (see Appendix E).

Table 7 shows the evaluation results using median to consolidate the scores of the three annotators on the answerable/partial subset (152 tasks) of the human evaluation. The reference answer is exceedingly preferred by annotators over the model responses, as evidenced by the win-rate. It also exhibits the highest score for most individual properties. This highlights the quality of the human-generated reference answers and shows that even frontier LLMs still have room for improvement. The results show that LLMs do a

⁹<https://www.appen.com/>.

	RB _{alg}	RB _{llm}	RL _F	RL _R	RL _{MTB}
WR	0.24	0.33	0.01	-0.03	-0.17

Table 8: Weighted Win-Rate Spearman correlation with reference-based (RB) and reference-less (RL) metrics.

good job of providing natural and appropriate answers, but they still struggle with faithfulness and completeness, where they receive lower scores. Both frontier models are equally preferred for answerable questions.

We also explore the few unanswerable questions, and observe that Llama 3.1 405B Instruct is the least preferred by win-rate. Llama still provides answers (hallucinations) to unanswerable questions.

8 Automatic Evaluation

Human evaluation is not feasible as a long-term solution for evaluating models as it does not scale easily. We explore reference-less and reference-based automated evaluation via algorithmic metrics and LLM judges. Table 8 shows the correlation of these metrics with the human evaluation’s win-rate and Figure 4 shows the correlation of metrics that correlate positively with win-rate for the individual properties of FANC. We employed a weighted correlation on human judgements to correct for data imbalance. We exclude the reference responses from Figure 4 because by design the ranking of the reference is 1 for the Reference-Based metrics. These findings drive our decision to report RB_{alg}, RB_{llm}, and RL_F as our main metrics in Section 6.

Reference-based Metrics. We consider an adapted version of RAD-Bench (Kuo et al., 2024), denoted as RB_{llm}, as well as the harmonic mean of algorithmic metrics, denoted as RB_{alg}, as described in Section 6.3. As shown in Table 8, these metrics correlate reasonably well with win-rate; in particular RB_{llm}. Figure 4 shows that they also tend to correlate well with our desired properties; in particular RB_{alg}.

Reference-less Metrics. We investigate popular metrics in literature: the Faithfulness (RL_F) and Answer Relevance (RL_R) metrics from RAGAS (Es et al., 2024); an adapted version of MT-Bench (RL_{MTB} as described in Appendix F.2); Zheng

F	0.48	0.69	0.33	0.53	0.68	0.27
A	0.45	0.28	0.49	0.45	0.32	0.46
N	0.11	0.36	0.21	0.42	0.34	0.30
C	0.45	0.59	0.35	0.49	0.63	0.33
FANC	0.46	0.64	0.39	0.53	0.63	0.32
	Rougel	Bert-KPrec	Bert-Rec	RB _{alg}	RL _F	RB _{llm}

Figure 4: Weighted Spearman correlation of human evaluation with the automated metrics on the answerable subset for the GPT-4o and Llama 3.1 405B Inst. models.

et al., 2023; as well as Bert-K-Prec. RAGAS Answer Relevance (RL_R), evaluates the model response by asking the LLM to determine the question associated with the response. This question is then compared with the actual question. This is not suitable in a multi-turn setup where the question is often non-standalone. Similarly, RL_{MTB} , while adapted to fit multi-turn conversations, has a very low correlation with win-rate, favoring machine-generated text. In contrast, we find that the two metrics focusing on faithfulness, Bert-K-Prec and RAGAS Faithfulness (RL_F), correlate well with the human faithfulness scores and are therefore useful in settings where no reference is available, although it is clear that evaluating only the faithfulness aspect of a response is incomplete.

9 Synthetic Conversations

Manually creating data is an expensive and time-consuming process that does not scale well. Automating this process has become popular via synthetic data generation (Soudani et al., 2024) and can serve as useful evaluation. To explore this direction, we construct a companion benchmark, $MTRAG-S$, of synthetically generated conversations. We extend the recently proposed framework of Lee et al. (2024) to automatically generate multi-turn conversations (see Appendix G.1 for details). To ensure that the two benchmarks are comparable, we utilize the same corpora, question types, multi-turn patterns, and answerability types (Section 3.1).

In Table 9, we see that synthetic conversations are typically shorter (averaging 5.9 vs 7.7 turns)

	$MTRAG$	$MTRAG-S$
Avg # turns per conversation	7.7	5.9
Avg # unique passages per conv.	16.9	4.6
Avg # edited responses per conv.	7.3	–
Avg # of question words	8.6	13.7
Avg # of response words *	97.2	72.4
Avg # of relevant passages *	2.8	4.1
Total # of Conversations	110	200
Total # of Tasks	842	1,181

Table 9: Comparing human-generated $MTRAG$ and synthetically generated $MTRAG-S$. (Note: Properties with * computed on answerable+partial subset.)

and exhibit a lower passage diversity (4.6 vs 16.9 unique passages per conversation). Question and response lengths also differ with synthetic data having longer (potentially more detailed) questions but shorter answers. We also found that mimicking several important characteristics posed challenges: Our attempts to synthetically generate unanswerable questions were not very successful as the model would often create questions with at least a partial answer. Moreover, increasing the number of turns tended to lead to repetitive user questions and a higher likelihood of hallucinated agent responses.

Since the aggregate metrics (RB_{alg} and RB_{llm}) rely on a reference answer, which does not exist for synthetic data, we employ the reference-less RL_F and Bert-K-Prec to evaluate faithfulness on $MTRAG-S$. Using either metric, we see that models across the board receive a higher faithfulness score on the synthetic than on the human-generated data (see Appendix G.2). There are multiple potential explanations for this, ranging from potential idiosyncrasies of synthetic data generation approaches, to the reliability of the presence of desired characteristics, to the quality and biases of automatic evaluation metrics. More work is needed to compare human-generated and synthetic data and we hope that our companion synthetic benchmark serves as a valuable asset towards that goal.

10 Conclusions and Future Work

We present $MTRAG$, a comprehensive and diverse benchmark of 110 multi-turn human-generated

conversations averaging 7.7 turns for a total of 842 tasks.¹⁰ These tasks are used to test the full RAG pipeline. mTRAG is the first end-to-end human-generated multi-turn RAG benchmark that reflects real-world properties of multi-turn conversations. Our experimental results, employing both automated metrics and a human evaluation, highlight the quality of our benchmark and outline several trends and challenges related to multi-turn RAG systems that state-of-the-art retrievers and LLMs face during retrieval and generation. Our findings encourage future research on improving retrieval and generation performance, especially in longer multi-turn conversations, unanswerable questions, and non-standalone user questions. In addition, the ability to scale indicates a clear need for i) more accurate reference-less automatic evaluation metrics, which align more closely with human judgment and can better differentiate model performance; and ii) synthetic data to obtain more conversations. We are also motivated to extend mTRAG in the future to include adversarial turns, additional domains, and multilingual conversations.

Acknowledgments

We would like to thank our annotators for their high-quality work in generating and evaluating this dataset: Mohamed Nasr, Joeeki Gurski, Tamara Henderson, Hee Dong Lee, Roxana Passaro, Chie Ugumori, Marina Variano, and Eva-Maria Wolfe. We would like to also thank Lihong He and Krishnateja Killamsetty for their help in preparing the two new document corpora (Govt and Cloud), Krishnateja Killamsetty for developing and running the classifier for identifying question types in mTRAG-S, and Hans Florian for contributing the code for the corpora ingestion. Finally, we thank our TACL action editor, Hai Zhao, and the anonymous reviewers for their suggestions that helped us further improve this manuscript.

References

Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2024. Evaluating correctness and

¹⁰mTRAG is publicly available at <https://github.com/ibm/mt-rag-benchmark>.

faithfulness of instruction-following models for question answering. *Transactions of the Association for Computational Linguistics*, 12:681–699. https://doi.org/10.1162/tacl_a_00667

Mohammad Aliannejadi, Zahra Abbasiantaeb, Shubham Chatterjee, Jeffrey Dalton, and Leif Azzopardi. 2024. TREC iKAT 2023: A test collection for evaluating conversational and interactive knowledge assistants. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, pages 819–829, New York, NY, USA. Association for Computing Machinery. <https://doi.org/10.1145/3626772.3657860>

Yang Bai and Daisy Zhe Wang. 2022. More than reading comprehension: A survey on datasets and metrics of textual question answering. *Preprint*, arXiv:2109.12264v2.

B. Barla Cambazoglu, Mark Sanderson, Falk Scholer, and Bruce Croft. 2021. A review of public datasets in question answering research. *SIGIR Forum*, 54(2). <https://doi.org/10.1145/3483382.3483389>

Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking large language models in retrieval-augmented generation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17754–17762. <https://doi.org/10.1609/aaai.v38i16.29728>

Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wen-tau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. QuAC: Question answering in context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2174–2184, Brussels, Belgium. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D18-1241>

Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. Wizard of wikipedia: Knowledge-powered conversational agents. In *International Conference on Learning Representations*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra,

Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie

Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Barambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute,

- Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U., Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A., Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojuan Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu Wang, Yuchen Hao, Yundi Qian, Yuze He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The Llama 3 herd of models. *Preprint*, arXiv:2407.21783v3.
- Nouha Dziri, Ehsan Kamaloo, Sivan Milton, Osmar Zaiane, Mo Yu, Edoardo M. Ponti, and Siva Reddy. 2022. FaithDial: A faithful benchmark for information-seeking dialogue. *Transactions of the Association for Computational Linguistics*, 10:1473–1490. <https://doi.org/10.1162/tacl-a-00529>
- 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*, pages 150–158, St. Julians, Malta. Association for Computational Linguistics.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. ELI5: Long form question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*,

- pages 3558–3567, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1346>
- Song Feng, Siva Sankalp Patel, Hui Wan, and Sachindra Joshi. 2021. MultiDoc2Dial: Modeling dialogues grounded in multiple documents. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6162–6176, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.emnlp-main.498>
- Robert Friel, Masha Belyi, and Atindriyo Sanyal. 2024. RAGBench: Explainable benchmark for retrieval-augmented generation systems. *Preprint*, arXiv:2407.11005v2.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th  ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2024. Mixtral of experts. *Preprint*, arXiv:2401.04088v1.
- Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7969–7992, Singapore. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.emnlp-main.495>
- Tzu-Lin Kuo, Feng-Ting Liao, Mu-Wei Hsieh, Fu-Chieh Chang, Po-Chun Hsu, and Da-Shan Shiu. 2024. RAD-Bench: Evaluating large language models capabilities in retrieval augmented dialogues. *Preprint*, arXiv:2409.12558v2.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466. https://doi.org/10.1162/tacl_a_00276
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent retrieval for weakly supervised open domain question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6086–6096, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1612>
- Young-Suk Lee, Chulaka Gunasekara, Danish Contractor, Ram  n Fernandez Astudillo, and Radu Florian. 2024. Multi-document grounded multi-turn synthetic dialog generation. *Preprint*, arXiv:2409.11500v1.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich K  ttler, Mike Lewis, Wen-tau Yih, Tim Rockt  schel, et al. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Chien-Chang Lin, Anna Y. Q. Huang, and Stephen J. H. Yang. 2023. A review of AI-driven conversational chatbots implementation methodologies and challenges (1999–2022). *Sustainability*, 15(5). <https://doi.org/10.3390/su15054012>
- Macedo Maia, Siegfried Handschuh, Andr   Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. WWW’18 open challenge: Financial opinion mining and question answering. In *Companion Proceedings of the The Web Conference 2018*, WWW ’18, pages 1941–1942, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee. <https://doi.org/10.1145/3184558.3192301>
- Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, KaShun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2024. RAGTruth: A hallucination corpus for developing trustworthy retrieval-augmented language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10862–10878,

- Bangkok, Thailand. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2024.acl-long.585>
- Chen Qu, Liu Yang, Cen Chen, Minghui Qiu, W. Bruce Croft, and Mohit Iyyer. 2020. Open-retrieval conversational question answering. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '20*, pages 539–548, New York, NY, USA. Association for Computing Machinery. <https://doi.org/10.1145/3397271.3401110>
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P18-2124>
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D16-1264>
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. CoQA: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266. https://doi.org/10.1162/tacl_a_00266
- Anna Rogers, Matt Gardner, and Isabelle Augenstein. 2023. QA dataset explosion: A taxonomy of NLP resources for question answering and reading comprehension. *ACM Computing Surveys*, 55(10). <https://doi.org/10.1145/3560260>
- Sara Rosenthal, Avirup Sil, Radu Florian, and Salim Roukos. 2025. CLAPNQ: Cohesive long-form answers from passages in natural questions for RAG systems. *Transactions of the Association for Computational Linguistics*, 13:53–72. https://doi.org/10.1162/tacl_a_00729
- Marzieh Saeidi, Max Bartolo, Patrick Lewis, Sameer Singh, Tim Rocktäschel, Mike Sheldon, Guillaume Bouchard, and Sebastian Riedel. 2018. Interpretation of natural language rules in conversational machine reading. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2087–2097, Brussels, Belgium. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D18-1233>
- Maojia Song, Shang Hong Sim, Rishabh Bhardwaj, Hai Leong Chieu, Navonil Majumder, and Soujanya Poria. 2025. Measuring and enhancing trustworthiness of LLMs in RAG through grounded attributions and learning to refuse. In *The Thirteenth International Conference on Learning Representations*.
- Heydar Soudani, Roxana Petcu, Evangelos Kanoulas, and Faegheh Hasibi. 2024. A survey on recent advances in conversational data generation. *Preprint*, arXiv:2405.13003v1.
- Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022. ASQA: Factoid questions meet long-form answers. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8273–8288, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.emnlp-main.566>
- YuHong Sun, Zhangyue Yin, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Hui Zhao. 2024. Benchmarking hallucination in large language models based on unanswerable math word problem. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2178–2188, Torino, Italia. ELRA and ICCL.
- Zhongkai Sun, Yingxue Zhou, Jie Hao, Xing Fan, Yanbin Lu, Chengyuan Ma, Wei (Sawyer) Shen, and Chenlei (Edward) Guo. 2023. Improving contextual query rewrite for conversational AI agents through user-preference feedback learning. In *EMNLP 2023*. <https://doi.org/10.18653/v1/2023.emnlp-industry.41>
- Qwen Team. 2024. Qwen2.5: A party of foundation models.

Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.

Zhen Wang. 2022. Modern question answering datasets and benchmarks: A survey. *Preprint*, arXiv:2206.15030v1.

Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. 2024. C-Pack: Packed resources for general Chinese embeddings. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, pages 641–649, New York, NY, USA. Association for Computing Machinery. <https://doi.org/10.1145/3626772.3657878>

Xiao Yang, Kai Sun, Hao Xin, Yushi Sun, Nikita Bhalla, Xiangsen Chen, Sajal Choudhary, Rongze Gui, Ziran 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. In *The Thirty-eighth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric 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*, volume 36, pages 46595–46623. Curran Associates, Inc.

Yingxue Zhou, Jie Hao, Mukund Rungta, Yang Liu, Eunah Cho, Xing Fan, Yanbin Lu, Vishal Vasudevan, Kellen Gillespie, and Zeynab Raeesy. 2023. Unified contextual query rewriting. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 608–615, Toronto, Canada. Association for Computa-

tional Linguistics. <https://doi.org/10.18653/v1/2023.acl-industry.58>

A Dimension Definitions

This section provides a detailed list of the mTRAG’s dimensions, introduced in Section 3.1. In particular, Tables 10, 11, and 12 show the definitions of the question types, multi-turn types, and answerability types included in the benchmark.

B Benchmark Creation Details

B.1 Creating the Govt and Cloud Corpora

The Cloud corpus was created by crawling the public technical documentation of a major cloud provider, thus containing inter-connected pages covering various aspects of each cloud offering.

Similarly, the Govt corpus was created starting from 100 seed web pages from the .gov and .mil domains, selected to be suitable to write conversations on (e.g., excluding very short pages or navigation pages) and diverse in terms of topics (e.g., covering parks, NASA, Veteran Affairs, city web-sites, DMV). For each seed page, we crawled web pages appearing in its neighborhood and thus expected to cover related topics. For each neighborhood the top 150 pages were selected based on their page rank score and additionally filtered to remove duplicates and other low-quality pages, leading to the final set of 7,661 pages/documents.

B.2 Creating Conversations

The following instructions were given to the annotators for creating conversations:

- Given a seed question, continue a conversation between you and the assistant. Continue the conversation to generate a total of 6–9 turns (turn = question + response).
- Aim to have around 25% unanswerable questions. There should be no more than 2 unanswerable questions in one conversation. Most conversations should be completely answerable.
- We encourage diversity of questions across types.
- The questions in your conversation should be connected, but diverse enough to include different passages. All questions shouldn’t have the same relevant passages.
- The assistant is a bot and the answers may be wrong. In addition to creating questions, you will also edit the answers if they are wrong. Answer correctness is based on the passages provided with the answer. Good conversations will have questions that require edits to repair the conversation.

Question type	Definition	Example questions
Comparative	Asking for comparison. This can be comparison (a) of multiple entities/concepts, (b) of characteristics of a single entity, or (c) comparison with decision.	(a) ‘‘What’s the difference between effective and marginal tax’’, (b) ‘‘pros and cons of credit cards’’, (c) ‘‘is X better than Y’’
Composite	Comprises several questions. They may be related or dependent.	‘‘Am I eligible for a driver’s license and how do I apply?’’
Explanation	Explain the reason behind something.	‘‘Why do I have to do/have X. . . ?’’
Factoid	Asking for a specific piece of information, such as a date, quantity, name, yes/no answer, or other singular fact. It can be answered directly and concisely, and does not require an explanation, opinion, interpretation, or subjective judgment to answer.	‘‘Am I eligible for a driver’s license?’’, ‘‘What is the link to the application portal?’’
How-To	Instructions describing how to perform a task.	‘‘How should I do X . . . ’’, ‘‘How do I apply for social security disability benefits?’’, ‘‘I need a license. What do I need to do?’’
Keyword	Asking using keywords (not full sentence/phrase). This may be ambiguous.	‘‘vacation days’’, ‘‘ios 17 upgrade’’
Non-question	Not asking a question but instead answering a question or providing information asked by the model.	‘‘I am in Sacramento’’, in response to the model saying: ‘‘The procedure to file a restraining order depends on the type of restraining order and the court where you want to file it in.’’
Opinion	Asking model’s opinion on something. The question could also be phrased as leading.	‘‘Don’t you think iphone is better than samsung?’’, ‘‘Which color car is the best?’’
Summarization	Asking to summarize a process or a policy.	‘‘What’s the policy on vacation days?’’
Troubleshooting	Finding solutions to issues, problems, challenges.	‘‘I have the error X. . . what should I do?’’

Table 10: Definitions of question types.

Multi-turn type	Definition	Example
Follow-up	Ask a question that requests more information or related information to continue the conversation.	<ul style="list-style-type: none"> • User: Can you tell me about the responsibilities of SSA? • Assistant: Social Security Administration (SSA) assigns Social Security numbers and runs Social Security retirement and disability insurance programs. • User [Follow-up]: <i>I have forgotten my social security number. Should I contact SSA regarding that?</i>
Clarification	Clarify user’s intent or model’s answer: <p>(a) Write a statement clarifying the user’s intent (typically used when the assistant misinterpreted one of the prior questions).</p> <p>(b) Ask a question to clarify the model’s previous answer. Clarifications are typically asked when something is unclear or hard to understand.</p>	<p>(a) Clarification of user’s intent:</p> <ul style="list-style-type: none"> • User: graphql • Assistant: GraphQL is an open-source data query and manipulation language for APIs and a query runtime engine. • User [Clarification]: <i>No, I meant, how do I set it up.</i> <p>(b) Clarification of model’s answer:</p> <ul style="list-style-type: none"> • User: Can you tell me about the responsibilities of SSA? • Assistant: Social Security Administration (SSA) assigns Social Security numbers and runs Social Security retirement and disability insurance programs. • User [Clarification]: <i>Can you explain what you meant by ‘‘Social Security retirement and disability insurance programs’’?</i>

Table 11: Definitions of multi-turn types.

Answerability type	Definition
Answerable	The question can be fully answered from the passages.
Partially answerable	Only part of the question can be answered from the passages.
Unanswerable	The question cannot be answered neither fully nor partially from the passages.
Conversational	The user turn does not contain a question but is a conversational statement (e.g., ‘‘Hello’’, ‘‘Hi, I had a question’’, ‘‘Cool’’, ‘‘That’s interesting’’, ‘‘That was all’’, ‘‘Thank you’’)

Table 12: Definitions of answerability types.

B.3 Reviewing Conversations

As described in Section 4.3, the conversations initially generated by annotators went through a

round of review before being finalized. During the review process, 16 out of the 126 conversations were rejected due to lack of turn-wise and topic-level coherency as well as passage diversity,

Given the following conversation, please reword the final utterance from the user into a single utterance that does not need the history to understand the user's intent. Output in proper JSON format indicating the "class" (standalone or non-standalone) and the "reworded version" of the last utterance. Use this format: {"class": "type of last utterance", "reworded version": "the last utterance rewritten into a standalone question, IF NEEDED"}.

In your rewording of the last utterance, do not do any unnecessary rephrasing or introduction of new terms or concepts that were not mentioned in the prior part of the conversation. Be minimal, by staying as close as possible to the shape and meaning of the last user utterance. If the last user utterance is already clear and standalone, the reworded version should be THE SAME as the last user utterance, and the class should be 'standalone'.

```

{{conversation}}
ASSISTANT:

```

Figure 5: Query rewrite prompt.

leading to 110 accepted conversations. Additionally, during review, 264 responses (31.3%) were repaired, leading to a slight increase on the number of repaired responses in the final benchmark (from 778 repaired responses before review to 799 responses after review, respectively). The average *edit score* on these responses also increased from 59.5 to 69.2.

C Retrieval Experiment Details

C.1 Query Rewriting

For the retrieval experiments, we implemented query rewriting by sending the prompt of Figure 5 to Mixtral 8×7B Instruct. This is only a reference implementation, and better ones could be obtained with larger models (Mixtral 8×22B Instruct, Llama 3.1 405B, etc.), few-shot prompts with in-context-learning examples, or via fine-tuning. However, this implementation is sufficient to make the point in Section 5 that query rewriting is both necessary and effective in mitigating the challenges of non-standalone questions in multi-turn data.

D Generation Experiment Details

D.1 Model Choices

To pick the models used in the generation experiments, we selected a suite of state-of-the-art models of varying sizes. It is not an exhaustive list, but shows the common trends in frontier models. For OpenAI models, we also experimented with the latest models in the o1¹¹ family. However, initial experiments showed that o1 did not perform as well as GPT-4o. We suspect that this is due to o1’s focus on reasoning and multimodal rather than on RAG.

¹¹<https://platform.openai.com/docs/models#o1>.

D.2 Model Invocation

To evaluate LLMs on mTRAG, for each task we sent to the LLM the following information: question, preceding turns, passages, and instruction. This was composed into the prompt below, adapted to different models using HuggingFace’s ChatTemplate.¹² In all cases the prompt fit within the model’s context length without need for truncation:

```

Given one or more documents and a user query, generate a response to the query using less than 150 words that is grounded in the provided documents. If no answer can be found in the documents, say, "I do not have specific information"

PASSAGE 1
...
PASSAGE M

User turn 1

Agent Turn 1
...
User Turn N

```

D.3 Example Output

Table 13 shows an example of an mTRAG task using the Reference (●) retrieval setting together with the Reference response as well as responses from several LLMs.

E Annotator Agreement

We next describe the computation of the annotator agreement reported in the human evaluation (Section 7). For each property, the annotator agreement was computed as the percentage of tasks on which the 3 annotators had either an *absolute agreement* (i.e., they all agreed on the score) or *high agreement* (i.e., two of the annotators agreed on the score and the third annotator gave a score that was a single point away from the score given by the majority).

F Metrics Details

F.1 RB_{llm} Implementation

To adapt the RAD-Bench judge (Kuo et al., 2024) to our RAG setting, we make two modifications: (1) We include the retrieved passages as additional input to the evaluation, and (2) we anchor

¹²https://huggingface.co/docs/transformers/main/en/chat_templating.

Conversation
<p>User: where does doctor strange get his powers from</p> <p>Assistant: Doctor Strange’s powers come from mystical entities such as Agamotto, Cyttorak, Ikonn, Oshtur, Raggadorr, and Watoomb, who lend their energies for spells. He also wields mystical artifacts including the Cloak of Levitation, the Eye of Agamotto, the Book of the Vishanti, and the Orb of Agamotto which give him additional powers.</p> <p>User: did he always have powers? what was he before?</p> <p>Assistant: No, Doctor Strange did not always have powers. He was a brilliant but egotistical neurosurgeon. He was in a car accident that ended his surgical career. He searched the globe for a way to repair his hands and encountered the Ancient One, eventually becoming one of her students.</p> <p>User: how many films does he appear in</p>
Passages
<p>Passage 1</p> <p>The character appears in the Marvel Cinematic Universe. He is first referenced as merely “Stephen Strange” in dialogue by Agent Sitwell in the 2014 film Captain America: The Winter Soldier, and first appears in the 2016 film Doctor Strange, in which he is portrayed by Benedict Cumberbatch . He reprises the role in the films Thor: Ragnarok and Avengers: Infinity War, and will appear again in the untitled sequel to Avengers: Infinity War. Doctor Strange appears in the 2016 animated film Hulk: Where Monsters Dwell, voiced by Liam O’Brien.</p>
<p>Passage 2</p> <p>In 2008, Doctor Strange was ranked 83rd in Wizard ’s “200 Greatest Comic Book Characters of All Time” list, and in 2012 was ranked 33rd in IGN ’s list of “The Top 50 Avengers” . He was also ranked 38th on IGN ’s list of “Top 100 Comic Book Heroes” . The character was first portrayed in live - action by Peter Hooten in the 1978 television film Dr. Strange . Benedict Cumberbatch stars as the character in the Marvel Cinematic Universe, first appearing in the role in the 2016 film Doctor Strange . He reprised the role in the 2017 film Thor: Ragnarok, the 2018 film Avengers: Infinity War, and will return in the fourth Avengers film in 2019 .</p>
<p>Passage 3</p> <p>Strange appeared as a supporting character for the bulk of the 2000s . He appeared regularly in The Amazing Spider - Man under J. Michael Straczynski, before being cast into a time loop by Baron Mordo . He later appeared on and off in The New Avengers, where he was stated as being part of the secret group known as the Illuminati to deal with future threats to Earth . Ultimately Strange joined the team and allowed them to use his home as a base after the events of Civil War, which he sat out . Doctor Strange was critical of the federal Superhuman Registration Act and aided the anti-registration Avengers team led by Luke Cage .</p>
Responses
<p>Reference</p> <p>I do not have an exact count of films in which Doctor Strange appears, but here are some: Doctor Strange, Thor: Ragnarok, Avengers: Infinity War, and The Amazing Spider-Man. Doctor Strange, voiced by Liam O’Brien, appears in the 2016 animated film Hulk: Where Monsters Dwell and will appear again in the untitled sequel to Avengers: Infinity War, which will return in the fourth Avengers film in 2019.</p>
<p>Llama 3.1 405B</p> <p>According to the text, Doctor Strange appears in the following films in the Marvel Cinematic Universe: 1. Doctor Strange (2016) 2. Thor: Ragnarok (2017) 3. Avengers: Infinity War (2018) 4. The untitled sequel to Avengers: Infinity War (Avengers: Endgame, released in 2019, although not explicitly mentioned in the text) Additionally, he appears in the 2016 animated film Hulk: Where Monsters Dwell and the 1978 television film Dr. Strange, played by Peter Hooten.</p>
<p>GPT-4.0</p> <p>Doctor Strange, portrayed by Benedict Cumberbatch, appears in the following Marvel Cinematic Universe films: 1. Doctor Strange (2016) 2. Thor: Ragnarok (2017) 3. Avengers: Infinity War (2018) 4. Avengers: Endgame (2019, referenced as the fourth Avengers film) 5. Doctor Strange in the Multiverse of Madness (subsequent appearances expected, although not specified in the provided documents) He also appears in the 2016 animated film Hulk: Where Monsters Dwell, voiced by Liam O’Brien.</p>
<p>Command-R+ (104b)</p> <p>According to the information provided, Benedict Cumberbatch first portrayed the character of Doctor Strange in the Marvel Cinematic Universe in 2016. Cumberbatch reprised the role in two other films, “Thor: Ragnarok” and “Avengers: Infinity War”, and will appear again in the untitled sequel to “Avengers: Infinity War”. Doctor Strange also appears in the 2016 animated film “Hulk: Where Monsters Dwell”, voiced by Liam O’Brien. Doctor Strange was first portrayed in live action by Peter Hooten in a 1978 television film.</p>
<p>Qwen 2.5 (72b)</p> <p>Doctor Strange, portrayed by Benedict Cumberbatch, appears in at least four films within the Marvel Cinematic Universe: “Captain America: The Winter Soldier” (referenced), “Doctor Strange” (2016), “Thor: Ragnarok” (2017), and “Avengers: Infinity War” (2018). He is also set to appear in the 2019 sequel to “Avengers: Infinity War”. Additionally, he appears in the 2016 animated film “Hulk: Where Monsters Dwell”, voiced by Liam O’Brien.</p>

Table 13: Example of an mTRAG task using the Reference (●) retrieval setting (this task is part of the conversation shown in Figure 1) with the Reference response as well as responses from several LLMs. It is clear that several models hallucinate.

the evaluation on our desired properties of faithfulness, appropriateness, and completeness. Figure 6 shows the final prompt that we used to implement RB_{llm} . Finally, in order to minimize model biases and improve evaluation reliability, we use four models as judges: GPT-4o-mini (2024-07-18), Qwen 2.5 (72B), Mixtral 8×22B Instruct, and

Llama 3.1 405B Instruct, taking the median as the final score.

F.2 RL_{MTB} Implementation

The MT-Bench judge was originally designed for evaluating conversations containing exactly two turns (i.e., user question/assistant response pairs),

```

[Instruction]
Please act as an impartial judge and evaluate the quality of the
response provided by an AI assistant to the user question given the
provided document and a reference answer.

Your evaluation should assess the faithfulness, appropriateness,
and completeness. Your evaluation should focus on the assistant's
answer to the question of the current turn. You will be given the
assistant's answer and a reference answer. You will also be given
the user questions and assistant's answers of the previous turns
of the conversation. You should consider how well the assistant's
answer captures the key information, knowledge points mentioned
in the reference answer, and how it respects or builds upon the focus
and knowledge points from the previous turns.

[Faithfulness]: You are given the full conversation, the question of
the current turn, the assistant's answer, and documents. You should
evaluate how faithful is the assistant's answer to the information in
the document and previous conversation.
[Appropriateness]: You should evaluate if the assistant's answer is
relevant to the question of the current turn and if it addresses all the
issues raised by the question without adding extra information.
[Completeness]: You should evaluate whether the assistant's answer
is complete with information from the documents.

Begin your evaluation by comparing the assistant's answer against
the reference answer in this turn. Be as objective as possible, and
provide a detailed justification for your rating. After providing your
explanation, you must rate the response on a scale of 1 to 10, strictly
following this format: "Rating: [[rating]]", for example: "Rating:
[[5]]".

[The Start of Previous Conversation]
{previous_conversation}
[The End of Previous Conversation]

[The Start of Current Turn Question]
{current_question}
[The End of Current Turn Question]

[The Start of Reference Answer]
{reference_answer}
[The End of Reference Answer]

[The Start of Assistant's Answer]
{response}
[The End of Assistant's Answer]

[The Start of Document]
{passages}
[The End of Document]

```

Figure 6: Prompt used for the RB_{llm} judge.

with a focus on properties such as helpfulness, depth, and creativity (Zheng et al., 2023). To use it for assessing general multi-turn conversations with relevant passages, we make the following changes: (1) We focus the evaluation on the last user turn of a multi-turn conversation, (2) we include retrieved passages as additional input, and (3) we modify the instruction to include desired properties, such as faithfulness and completeness. Figure 7 shows the final prompt that we used to implement RL_{MTB}.

F.3 IDK Correction to the Metrics

To determine whether a model response corresponds to an IDK (I-Don't-Know) answer, we use

```

[System]
Please act as an impartial judge and evaluate the quality of the
response provided by an AI assistant to the user question given
the provided document. Your evaluation should consider factors
such as the faithfulness, appropriateness, conciseness, relevance,
accuracy, completeness, and correctness of their responses. You
evaluation should focus on the assistant's answer to the final user
question in the conversation. Begin your evaluation by providing a
short explanation. Be as objective as possible. After providing your
explanation, you must rate the response on a scale of 1 to 10 by
strictly following this format: [[rating]], for example: Rating: [[5]].

<[The Start of Assistant A's Conversation with User]>
### User:
{question_turn_1}

### Assistant A:
{answer_turn_1}
...

### User:
{question_turn_last}

### Assistant A:
{answer_turn_last}

### Document:
{passages}

<[The End of Assistant A's Conversation with User]>

```

Figure 7: Prompt used for the RL_{MTB} judge

```

Read an inquiry and a response to it. Decide if the response con-
tains, in part or in whole, an indication that there is not sufficient
information to fully answer the inquiry. Use only one of the follow-
ing labels: yes, no, partial. Answer "yes" if the response indicates
that the inquiry cannot be answered, for example, due to lack of
information (in a reference document or text). Answer "partial"
if the response indicates insufficient information (in a reference
document or text) to answer part of the inquiry, but still attempts
to answer other part of the inquiry that can be answered. Answer
"no" if the response does not indicate a lack of information prevent-
ing answering part of the inquiry, even if some reasoning steps or
general knowledge are needed to address the inquiry.

[Inquiry]
How can I learn SQL without taking a course?

[Response]
I'm unable to provide an answer regarding how you can learn SQL
without taking a course. However, I can assist with the No SQL
databases topic.

[Output]
yes

[Additional ICL examples]

[Inquiry]
{inquiry}

[Response]
{response}

[Output]

```

Figure 8: IDK judge prompt.

an LLM judge that we developed for this purpose. The judge sends the prompt of Figure 8 to GPT-4o-mini (2024-07-18). To measure the judge's performance, we evaluated its predictions

Model	RL_F	RL_F	RB_{llm}	RB_{llm}	RB_{alg}	RB_{alg}
Reference	0.18	0.87	0.89	0.87	0.75	0.87
GPT-4o	0.19	0.20	0.66	0.20	0.26	0.20
Llama 3.1	0.21	0.20	0.69	0.20	0.27	0.20

Table 14: Generation results in the Reference (●) retrieval setting before and after conditioning with IDK judge on 55 unanswerable questions.

on the human evaluation results across model responses and the reference answer, where the judge achieved over 97% accuracy. We assess the impact of the IDK judge on the unanswerable questions using our reference answers and responses from frontier models (GPT-4o, Llama 3.1 405B Instruct) in Table 14. These observations show that the IDK judge is needed to correct the scores.

G Synthetic Conversation Details

We next provide additional details on the mTRAG-S benchmark of synthetic multi-turn conversations.

G.1 Synthetic Conversation Generation

To create mTRAG-S, we extended the recently proposed framework of (Lee et al., 2024) to automatically generate user questions corresponding to the question types used in mTRAG. To classify question types within the synthetic conversations, we employed a question type classifier trained on metadata derived from the human-generated conversations in mTRAG. Finally, we used Mixtral 8×22B Instruct v0.1 as the LLM for the conversation generation and restricted the conversations to a maximum of 8 turns, as further increasing the number of turns tends to lead to repetitive user questions and a higher likelihood of hallucinated assistant responses.

G.2 Faithfulness of Models on Synthetic Data vs Human-generated Data

We next provide detailed evaluation results on the comparison of faithfulness of models on mTRAG and mTRAG-S, discussed in Section 9. Tables 15a and 15b show faithfulness on the two benchmarks based on the RL_F and BERT-K-Prec metrics, respectively. Using either metric, we see that models across the board receive a higher faithfulness score on the synthetic than on the human-generated data.

H Detailed Retrieval Results

Expanding on Table 3, which shows average retrieval results, in Table 16 we provide detailed retrieval results per domain for Elser, the best performing retriever.

I Detailed Generation Results

Expanding on the generation experiment results of Section 6 we provide two additional sets of results.

I.1 Generation Results for All Three Metrics

Expanding on Figure 3, which shows generation results in the Reference Generation (●) setting using a single metric, RB_{alg} , we now show generation results for the same setting using all three metrics: RB_{alg} , RB_{llm} , and RL_F . Tables 17a, 17b, and 17c present the results broken down by answerability, turns, and domain, respectively.

I.2 Generation Results by Question Type and Multi-turn Type

Using the same three metrics, we also include generation results broken down by question and multi-turn type, shown in Figures 9 and 10, respectively.

	MT-RAG	MT-RAG-S		MT-RAG	MT-RAG-S
Command-R+ (104B)	0.76	0.83	Command-R+ (104B)	0.33	0.38
GPT-4o-mini	0.71	0.81	GPT-4o-mini	0.27	0.34
Llama 3.1 405B Instruct	0.75	0.85	Llama 3.1 405B Instruct	0.33	0.4
Mixtral 8×22B Instruct	0.61	0.79	Mixtral 8×22B Instruct	0.29	0.41
Qwen 2.5 (72B)	0.72	0.82	Qwen 2.5 (72B)	0.3	0.38

(a) Based on RL_F (b) Based on BERT-K-Prec

Table 15: Comparing faithfulness of models on the human-generated MT-RAG and synthetic MT-RAG-S using two metrics: RL_F and BERT-K-Prec.

		Recall				nDCG			
		@1	@3	@5	@10	@1	@3	@5	@10
CLAP _{NQ} (208)	Last Turn	0.20	0.43	0.53	0.65	0.48	0.46	0.49	0.54
	Query Rewrite	0.22	0.45	0.56	0.70	0.54	0.50	0.53	0.59
FiQA (180)	Last Turn	0.15	0.34	0.44	0.55	0.38	0.36	0.40	0.45
	Query Rewrite	0.18	0.39	0.50	0.63	0.43	0.41	0.46	0.52
Govt (203)	Last Turn	0.18	0.42	0.50	0.58	0.42	0.44	0.46	0.49
	Query Rewrite	0.21	0.47	0.56	0.67	0.47	0.48	0.51	0.56
Cloud (189)	Last Turn	0.19	0.38	0.48	0.56	0.40	0.39	0.42	0.47
	Query Rewrite	0.20	0.40	0.47	0.57	0.42	0.41	0.43	0.48
All (780)	Last Turn	0.18	0.39	0.49	0.58	0.42	0.41	0.45	0.49
	Query Rewrite	0.20	0.43	0.52	0.64	0.46	0.45	0.48	0.54

Table 16: Elser retrieval performance of models on our benchmark using Recall and nDCG metrics per domain.

	Overall			Answerable			Partial			Unans.
	RL _F	RB _{llm}	RB _{alg}	RL _F	RB _{llm}	RB _{alg}	RL _F	RB _{llm}	RB _{alg}	
Reference	0.87	0.95	0.88	0.88	0.96	0.88	0.71	0.88	0.83	0.87
Command-R+ (104B)	0.76	0.69	0.44	0.82	0.74	<u>0.47</u>	0.59	0.63	0.36	0.13
GPT-4o	<u>0.75</u>	0.76	<u>0.45</u>	0.82	0.81	0.48	0.53	<u>0.71</u>	0.35	0.20
GPT-4o-mini	0.71	<u>0.75</u>	0.43	0.77	0.79	0.44	0.39	0.62	0.30	<u>0.34</u>
Llama 3.1 405B Instruct	<u>0.75</u>	0.74	0.47	<u>0.81</u>	0.79	0.50	<u>0.58</u>	0.66	<u>0.37</u>	0.20
Llama 3.1 70B Instruct	0.69	0.66	0.44	0.74	0.69	0.45	0.42	0.47	0.27	0.44
Llama 3.1 8B Instruct	0.55	0.59	0.36	0.59	0.62	0.38	0.34	0.47	0.24	0.33
Mixtral 8×22B Instruct	0.61	0.69	0.41	0.68	0.75	0.45	0.41	0.68	0.33	0.00
Qwen 2.5 (72B)	0.72	0.74	0.44	0.79	<u>0.80</u>	<u>0.47</u>	0.53	0.72	0.38	0.07
Qwen 2.5 (7B)	0.68	0.72	0.43	0.74	0.77	0.46	0.44	0.67	0.36	0.11

(a) By question answerability

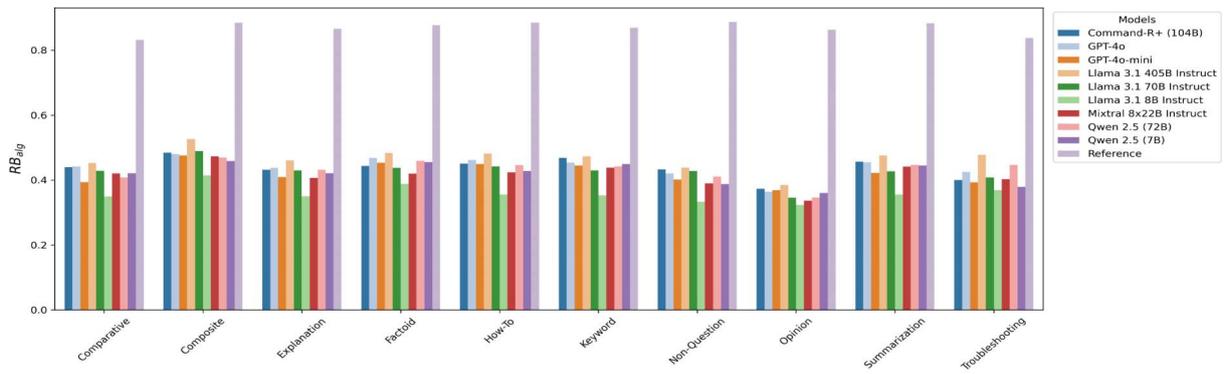
	RL _F		RB _{llm}		RB _{alg}	
	TURN 1	> TURN 1	TURN 1	> TURN 1	TURN 1	> TURN 1
Reference	0.89	0.86	0.97	0.95	0.89	0.87
Command-R+ (104B)	0.83	0.75	0.70	0.69	0.46	0.43
GPT-4o	0.86	<u>0.74</u>	<u>0.78</u>	0.76	0.54	<u>0.44</u>
GPT-4o-mini	<u>0.84</u>	0.69	0.79	<u>0.74</u>	<u>0.50</u>	0.42
Llama 3.1 405B Instruct	0.81	<u>0.74</u>	0.74	<u>0.74</u>	<u>0.50</u>	0.46
Llama 3.1 70B Instruct	0.80	0.68	0.69	0.66	<u>0.50</u>	0.42
Llama 3.1 8B Instruct	0.66	0.53	0.56	0.59	0.41	0.36
Mixtral 8×22B Instruct	0.82	0.58	0.73	0.69	0.47	0.40
Qwen 2.5 (72B)	<u>0.84</u>	0.70	0.77	<u>0.74</u>	0.51	0.43
Qwen 2.5 (7B)	0.82	0.65	0.71	0.72	0.48	0.42

(b) By first vs subsequent turns

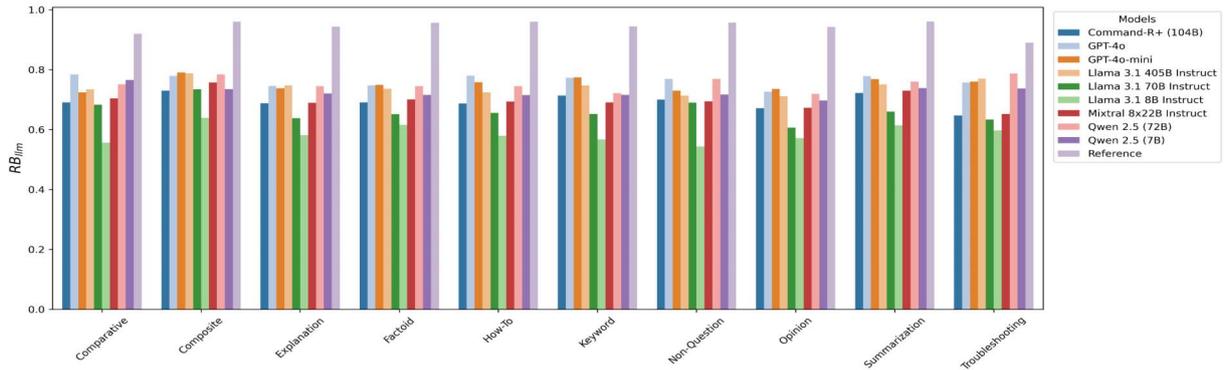
	RL _F				RB _{llm}				RB _{alg}			
	CLAPNQ	FiQA	Govt	Cloud	CLAPNQ	FiQA	Govt	Cloud	CLAPNQ	FiQA	Govt	Cloud
Reference	0.86	0.89	0.85	0.87	0.96	0.95	0.95	0.95	0.88	0.88	0.88	0.87
Command-R+	0.78	0.76	<u>0.76</u>	0.73	0.72	0.70	0.68	0.66	0.47	<u>0.39</u>	0.43	0.45
GPT-4o	<u>0.73</u>	<u>0.74</u>	0.77	0.78	0.77	0.74	0.78	0.74	0.48	0.38	<u>0.47</u>	0.45
GPT-4o-mini	0.69	0.71	0.71	0.75	<u>0.76</u>	0.74	0.75	<u>0.73</u>	0.45	0.37	0.44	0.44
Llama 3.1 405B Inst.	0.72	0.76	0.77	<u>0.76</u>	0.75	<u>0.73</u>	0.75	<u>0.72</u>	0.48	0.41	0.49	0.49
Llama 3.1 70B Inst.	0.66	0.67	0.73	0.71	0.66	0.62	0.69	0.65	0.44	0.37	0.45	<u>0.47</u>
Llama 3.1 8B Inst.	0.52	0.56	0.56	0.56	0.62	0.55	0.60	0.57	0.37	0.32	0.38	0.38
Mixtral 8×22B Inst.	0.60	0.60	0.61	0.64	0.72	0.68	0.70	0.68	0.45	0.34	0.42	0.44
Qwen 2.5 (72B)	0.68	0.70	0.74	<u>0.76</u>	0.75	<u>0.73</u>	<u>0.76</u>	<u>0.73</u>	<u>0.46</u>	0.37	0.46	0.46
Qwen 2.5 (7B)	0.70	0.63	0.68	0.69	0.72	0.69	0.74	0.71	0.48	0.35	0.44	0.44

(c) By domain

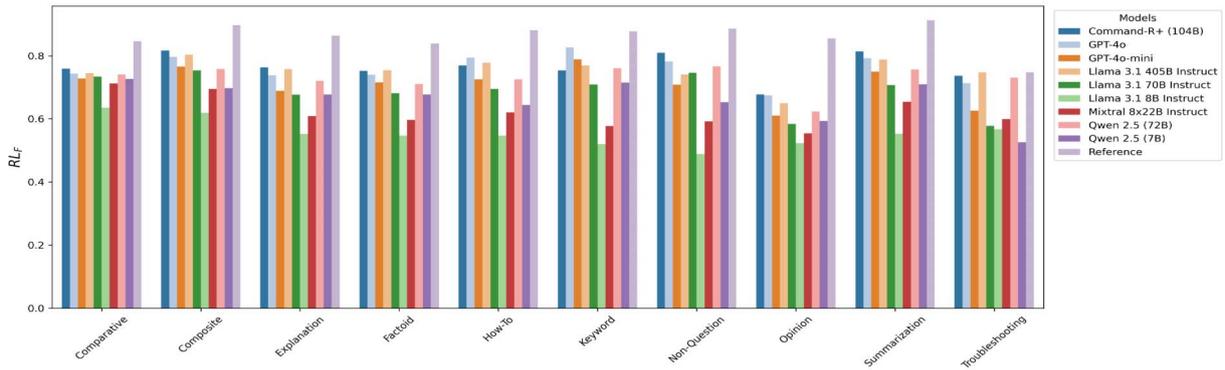
Table 17: Detailed generation results in the Reference (●) retrieval setting using three metrics, RL_F, RB_{llm}, and RB_{alg}, on three different dimensions: (a) answerability, (b) turns, and (c) domains.



(a) Based on RB_{alg}

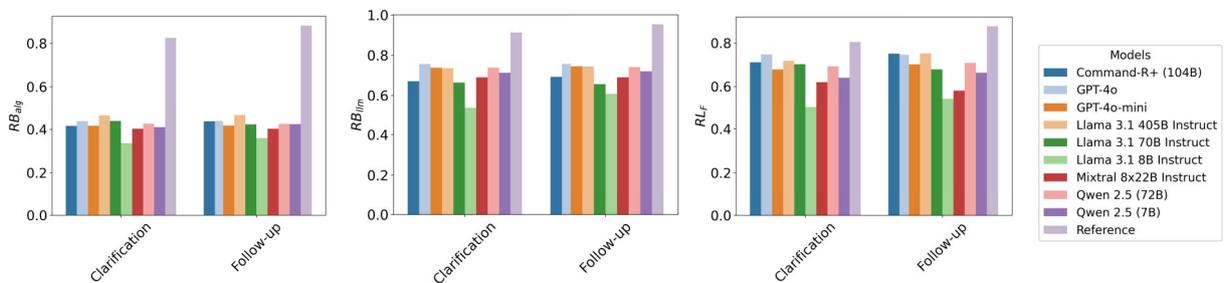


(b) Based on RB_{llm}



(c) Based on RL_F

Figure 9: Generation results in the Reference (●) retrieval setting by question type based on three metrics: RB_{alg} , RB_{llm} , and RL_F .



(a) Based on RB_{alg}

(b) Based on RB_{llm}

(c) Based on RL_F

Figure 10: Generation results in the Reference (●) retrieval setting by multi-turn type based on three metrics: RB_{alg} , RB_{llm} , and RL_F .