SRAG: Structured Retrieval-Augmented Generation for Multi-Entity Question Answering over Wikipedia Graph

Anonymous ACL submission

Abstract

Multi-entity question answering (MEQA) 001 002 poses significant challenges for large language models (LLMs), which often struggle 003 to consolidate scattered information across 004 multiple documents. An example question 005 006 might be "What is the distribution of IEEE Fellows among various fields of study?", 007 which requires retrieving information from 008 diverse sources e.g., Wikipedia pages. The 009 effectiveness of current retrieval-augmented 010 011 generation (RAG) methods is limited by the LLMs' capacity to aggregate insights from 012 numerous pages. To address this gap, this 013 paper introduces a structured RAG (SRAG) 014 framework that systematically organizes ex-015 tracted entities into relational tables (e.g., 016 tabulating entities with schema columns like 017 "name" and "field of study") and then apply table-based reasoning techniques. Our 019 approach decouples retrieval and reasoning, enabling LLMs to focus on structured 021 data analysis rather than raw text aggregation. Extensive experiments on Wikipediabased multi-entity QA tasks demonstrate 024 that SRAG significantly outperforms stateof-the-art RAG solutions, achieving a 29.6% 027 improvement in accuracy. The results underscore the efficacy of structuring unstruc-028 tured data to enhance LLMs' reasoning capabilities.

1 Introduction

Recent progress in Retrieval-Augmented Gener-032 ation (RAG) has enhanced how language models 033 access external knowledge, improving applications like question answering and document integration (Fan et al., 2024). By merging advanced retrieval methods with powerful language mod-037 els, these systems have shown strong perfor-038 mance. However, challenges persist in accu-039 rately retrieving entity information from multidocument and heterogeneous knowledge bases. 041 This challenge becomes especially apparent in 042 Multi-Entity Question Answering, the challenge 043 lies not only in recognizing and extracting relevant entities precisely from data but also in 045 understanding the properties of these entities within the context of the query. Consider an-047 swering questions such as "What are the capitals of countries bordering France?" or "How many 049 Turing Award Winners are Canadian" (the query Q in Figure 1-a). Answering these questions requires the extraction of multiple documents, unless specific statistical analysis has been carried 053 out by hand in advance. Existing RAG methods 054 struggle with MEQA tasks, because useful infor-055 mation required to these tasks is badly scattered. 056 This characteristic makes it difficult for existing 057 RAG methods to accurately identify key information and perform global reasoning with noisy

031

retrieved content.

060

068

070

071

072

073

074

076

077

078

081

084

089

091

092

093

097

061To tackle the challenge, we propose an inno-062vative Structured RAG System (SRAG) which063includes two main parts: (1) Multi-entity Se-064mantic Retrieval as illustrated in Figure 1-065a, and (2) Structured Question Answering066(SQA) as illustrated in Figure 1-c, to effectively067address multi-entity QA.

Contributions Our notable contributions are summarized as follows.

• MEBench: Specialized Benchmark for MEQA. We curate MEBench, a standardized benchmark based on Wikipedia to evaluate the effectiveness of various approaches to address the complexities of information extraction and reasoning involving multiple entities.

 Multi-entity Semantic Retrieval. Multientity Semantic Retrieval innovatively enhances semantic retrieval by integrating language models with structured database. This integration enables the creation of precise SPARQL queries to effectively retrieve relevant entities and web pages from Wikipedia. Our key contribution lies in improving retrieval accuracy by leveraging contextual brilliance of language models with the rigor of structured data validation.

• SQA: Module for Structured Entity Information and Reasoning. We propose the SQA, a module for managing vast and unstructured data by extracting properties of entities and organizing information into structured tables. This module transforms textual information of entities into a format with a rigorous and accurate schema, which facilitates analysis. Our experiments demonstrate its remarkable performance, achieving SOTA results and outperforming098the strongest baselines by 29.6% in over-
all accuracy, while leading across all eight099subtasks.101

102

103

126

2 Related Work

2.1 Retrieval-Augmented Generation

Advanced RAG systems have evolved to incor-104 porate pre-retrieval, retrieval, and post-retrieval 105 strategies aimed at mitigating the limitations 106 of basic RAG methods. In parallel, Modular 107 RAG systems have introduced frameworks for 108 iterative and dynamic cycles of intertwined re-109 trieval and generation processes (Gao et al., 110 2023). Community summaries serve as a form 111 of self-memory (Cheng et al., 2023), enhanc-112 ing generation-augmented retrieval to support 113 future cycles of information generation. Addi-114 tionally, the parallel generation of community 115 answers from these summaries reflects an iter-116 ative (Shao et al., 2023) or federated (Wang 117 et al., 2024) retrieval-generation approach. Fur-118 thermore, using a hierarchical index and sum-119 marization aligns with alternative methods, such 120 as generating a hierarchical index of text chunks 121 through vector clustering (Sarthi et al., 2024). 122 These methods concentrate on semantic aggre-123 gation and overlook the structured views of mul-124 tiple entities. 125

2.2 Graph RAG

Recently, several studies have explored the in-127 corporation of graph structures into RAG sys-128 tems to enhance LLMs in addressing complex 129 question-answering tasks (Edge et al., 2024) 130 (Panda et al., 2024). One approach involves 131 utilizing pre-existing knowledge graphs, from 132 which subgraphs are extracted based on spe-133 cific queries. These subgraphs are then en-134 coded as soft prompts or converted into plain 135



Figure 1: An Overview of Multi-entity QA Solutions over Wikipedia Graph. (a) **Multi-Entity Retrieval**: In a1, a rough SPARQL query is generated using language model (GPT-4). In a2, we integrate the LLM's semantic parsing with Wikipedia's API and utilize verifiable query accuracy on structured Wikidata to accurately identify entities and properties. In step a3, we synthesize an exact SPARQL query. Finally, in a4, the refined SPARQL query is used to retrieve the relevant entities and web pages. (b) **Existing Reasoning Solutions**: b1 represents direct responses from LLMs, while b2 combines LLMs with RAG. (c) Our proposal: **Structured QA**. Initially, in step c1, a language model (GPT-4) is employed to analyze the question and determine the **table schema**. In c3, we utilize an **information extraction module** to populate the form. Finally, in step c4, the **TableQA** module is used to derive the final answer.

text for the generation module GraphGPT (Tang 136 et al., 2024) and other methods (He et al., 137 2024) (Guan et al., 2023). Another method fo-138 cuses on extracting entity-relation triples from 139 relevant text documents according to query 140 needs, allowing for the construction of graph 141 structures that facilitate knowledge augmenta-142 tion, like GraphRAG (Edge et al., 2024) and 143 RQ-RAG (Chan et al., 2024). These models 144 are inadequate for addressing multi-entity ques-145 tions. In response, we introduce a comprehen-146

sive Multi-Entities QA benchmark MEBench based on Wikipedia Graph and provide a firstof-its-kind original system SRAG.

3 Problem Statement

3.1 Wikipedia Graph

Wikipedia graph is represented as G=(V, E, 152) P), where $V = \{v_1, v_2, ..., v_n\}$ is the set of 153 nodes in the graph, with each node $(v_i \in 154)$ V) representing an entity; E is the set of 155 direct edges in the graph, with each edge 156

 $(e_i(v_i, v_k) \in E)$ representing a connection (or 157 relationships) between two nodes. An edge is 158 a tuple (v_a, v_b, r) , where $(v_a, v_b \in V)$ and r is 159 the type of relationship. For example, E =160 $\{(v_1, v_2, r_1), (v_2, v_3, r_2), \dots, (v_m, v_n, r_k)\}.$ P161 represents the set of properties associated with 162 both nodes and edges. Figure 2 is a sample of a 163 Wikipedia graph snippet. 164

3.2 Multi-Entities Question

165

166

167

168

169

170

171

172

173

174

175

188

189

190

191

192

193

194

A Multi-Entities Question can be formally defined as $Q = (t_Q, V_Q, P_Q)$, where $t_Q \in T$ denotes the query type, with $T = \{t_1, t_2, \ldots, t_8\}$ representing the set of eight predefined types. Details can be seen in Table 1. V_Q denotes the collection of entities directly associated with the question. P_Q represents the comprehensive set of properties pertinent to the question, encompassing both node and edge properties.

4 System Design

In our system, the Multi-entity Semantic Re-176 trieval involves conducting a SPARQL retrieval 177 across the Wikipedia graph to obtain relevant 178 Wikipedia pages. Secondly, in SQA module, 179 table generation begins with "guessing" a table 180 schema based on the given query, followed by 181 the extraction of information from the identified 182 entities to populate the table. Finally, we imple-183 ment an executor that processes the generated 184 table to respond to the query. Next, we will 185 elaborate on the details of each step. 186

4.1 Composite SPARQL Retrieval

Initially, we utilize GPT-4 to parse the question to construct rough SPARQL. The entities ID and properties ID contained in the rough SPARQL frequently turn out to be inaccurate. To make SPARQL valid, we deploy GPT-4 as Semantic Analysis Model to identify entities and properties. Integrating with the Wikipedia API, we get right (entityID, propertyID) pair to replacing 195 the wrong IDs in the rough SPARQL, as illus-196 trated in Figure 1-a2. Consequently, our sys-197 tem is capable of performing entity and property 198 identification without ambiguities. For multi-199 hop queries, the Semantic Analysis process initially deconstructs the queries into sub-queries, 201 allowing composite SPARQL retrieval step to be 202 applied sequentially to each sub-query to iden-203 tify named entities and extract properties until 204 all sub-queries have been processed. 205

4.2 Table Generation

Although SPARQL provides aggregation functions such as "SUM, AVG, COUNT", etc. they are insufficient for complex statistical problems. Therefore, we build tables instead of extended graphs to support more complex algorithms and analysis. The table generation consists of two steps: (1) Generation of schema; (2) Extracting entity information to fill the table.

207

209

210

211

212

213

214

233

Schema Generation. We employ GPT-4 to 215 systematically parsing the question to identify 216 critical entities, attributes, and their interrela-217 tionships, which are then formalized into a struc-218 tured schema. The generated schema may neces-219 sitate adjustments or refinements to align more 220 closely with the intent of the question. For example, for the question "How many Nobel Prizes in 222 Physics laureates have been awarded for discov-223 eries in Particle Physics?", the LLM produces a 224 schema (name, YearAwarded, field), there are 225 two issues, first is that (field) is oversimplified, 226 it should be (field in Physics). The second is-227 sue is that columns YearAwarded are redun-228 dant. Therefore, We prompt GPT-4 to critically 229 review content to minimize oversimplification, 230 omission of essential elements, and redundancy, 231 and the prompt is shown in Appendix A.2.1. 232

Entity Information Extraction. In our data



Figure 2: Illustration of a Wikipedia graph snippet.

processing workflow, we use the LLM (Mistral-7B) to extract information from the retrieved data, populating a table where each row represents a unique entity, as shown in Figure 1c3. This step transforms entities into structured tables for downstream table-based reasoning tasks.

4.3 Executor

241

242

243

244

245

246

247

248

249

250

We utilize GPT-4 to generate SQL according to the question. To increase accuracy, we include relevant information in the prompt, such as the table schema and data samples. The generated SQL are executed on the generated table to obtain results, which could be a single value or a subset of the table.

5 Experiment

5.1 Experiment Setup

MEBench Benchmark. It is a specialized 251 benchmark designed to evaluate systems ad-252 dressing multi-entity OA. The benchmark com-253 prises 4,780 methodically structured questions 254 partitioned into two subsets: a training set 255 256 (3,406 questions) for model fine-tuning and a test set (1,374 questions) for rigorous evaluation. 257 These questions are systematically categorized 258 into three primary categories, further divided 259 into eight distinct types (see Table 1), ensuring

broad coverage of real-world multi-entity reasoning scenarios. Table 2 details comprehensive statistics of the benchmark.

261

262

263

264

265

266

267

269

270

271

272

273

274

275

276

Baselines. For open-source LLMs, we conduct experiments using the representative Meta-Llama-3-8B-Instruct (Meta Llama3, 2024) and apply QLoRA (Dettmers et al., 2023) to finetune it with the training set of MEBench. For proprietary LLMs, we select the widely recognized GPT models, including GPT-3.5turbo (Ouyang et al., 2022) and GPT-4 (Achiam et al., 2023). Additionally, we incorporate RAG across all vanilla baseline models for comparative analysis and evaluation of the model's capacity to integrate and leverage external data sources.

Evaluation Metrics. We adopt Accuracy 277 (Acc) as the primary metric to assess the per-278 formance of LLMs on MEBench tasks. For the 279 subcategories of Variance Analysis, Correlation Analysis, and Distribution Compliance within 281 the Statistics tasks shown in Table 1, we focus solely on prompting LLMs to identify relevant columns and applicable methods, evalu-284 ating the accuracy of their selections instead of the computational results, as LLMs' abilities in precise calculations are not the central focus of 287 this study.

Categories	Types	Examples	
Comparison	Intercomparison	Which has more ACM fellow, UK or USA?	
companison	Superlative	Which city has the highest population?	
	Aggregation	How many ACM fellow are from MIT?	
Statistics	Distribution Compliance	Does the nationality of ACM fellows follow a normal distribution?	
	Correlation Analysis	Is there a linear relationship between number of events and records broken in Olympic Games?	
	Variance Analysis	Do the variances in the number of participat- ing countries and total events in the Summer Olympics differ significantly?	
Relationship	Descriptive Relationship	Is there a relationship between the year of ACM fellowship induction and the fellows' areas of expertise?	
	Hypothetical Scenarios	If China wins one more gold medal, will it over take the US in the gold medal tally at the 202 Olympics?	

Table	1.	Evami	hlee	of	multi_	entities	querie	20
rabic	1.	LAAIII	JIUS	01	munu-	cintuites	quein	<i>.</i> o.

Categories	MEBench-train	MEBench-test	MEBench-total
#-Queries	3406	1374	4780
#-one-hop Q	1406	606	2012
#-multi-hop Q	1322	768	2090
Ave. #-entities /Q	460	391	409
#-Topics	165	76	241
#-Comparison	1107	438	1545
#-Statistics	1440	585	2025
#-Relationship	859	351	1210

Table 2: Statistics of MEBench benchmark.

5.2 Results and Analysis

289

Various models exhibit notable variations in performance on MEBench. Table 3 presents experimental results alongside overall accuracy

on MEBench, and Figure 3 shows accuracy on
eight further-divided types.293294

Performance of SRAG and Baselines. Compared to baselines, our SRAG significantly im-295296296

Models	Accuracy				
	Comparison	Statistics	Relationship	Overall	
GPT-3.5-turbo	0.105	0.198	0.476	0.239	
GPT-3.5-turbo + RAG	0.605	0.260	0.476	0.425	
GPT-4	0.199	0.289	0.507	0.316	
GPT-4 + RAG	0.763	0.410	0.687	0.593	
Llama-3-Instruct	0.046	0.118	0.256	0.130	
Llama-3-Instruct + RAG	0.447	0.181	0.410	0.325	
FT Llama-3-Instruct	0.046	0.253	0.259	0.189	
FT Llama-3-Instruct + RAG	0.687	0.448	0.573	0.556	
SRAG (Ours)	0.934	0.908	0.803	0.889	

Table 3: Experimental results for MEBench.

proves overall accuracy, reaching 88.9% and 297 increasing the best baseline (GPT-4 + RAG) by 298 299 **29.6%**. Our approach outperforms the accuracy in the relational and comparative query types by 300 11.6% and 17.1%, respectively, while achieving 301 a remarkable improvement of 46% for statistical 302 query types. 303

Fine-grained Performance on Sub-tasks. Fig-304 ure 3 shows that vanilla LLMs perform well in 305 correlation analysis and descriptive relationship 306 sub-tasks, while RAG significantly improves 307 intercomparison and superlative tasks. How-308 ever, neither fine-tuning nor RAG overcomes 309 challenges in variance analysis and aggregation 310 311 tasks, while our proposed SRAG achieves superior accuracy of 87.3% and 97.9%. 312

Errors Analysis for SRAG. We sample and 313 analyze the output of the SRAG system. It faces 314 two challenges which are listed below. 315

316

317

318

• Relation semantic parsing. The semantic parsing model in SPARQL retrieval effectively recognizes entities but struggles with relationship identification, leading to 319 challenges in graph retrieval and negatively 320

affecting the performance of RAG-based 321 approaches, including SQA. For example, 322 in the query "How many US presidents 323 have served more than one term in office?" 324 The model incorrectly identifies the rela-325 tionship as "instance of" rather than the correct "position held", leading to erroneous 327 results.

• Insufficient information extraction. We 329 also identified errors in SRAG's information extraction during the table-filling 331 phase. An analysis of more than 2,000 ta-332 ble filling instances reveals that these errors 333 occur primarily as omissions (albeit with a low probability of approximately 0.1%). A 335 new challenge is the appearance of multiword synonyms within the same column, 337 such as "US" and "America", which negatively affects the accuracy of SQL execu-339 tion such as "SELECT".

Conclusion 6

Our research presents a novel framework, Struc-342 tured RAG system (SRAG), to address the com-343 plexities involved in multi-entity question an-

341



Figure 3: Experimental results for eight types queries of each model.

swering (QA) from Wikipedia. Existing meth-345 ods, particularly those employing RAG alongside LLMs, often fall short in effectively aggregating and reasoning over information scattered 349 across multiple Wikipedia pages. By leveraging the inherent structure of wiki-graph for multi-350 entity retrieval and introducing a system to or-351 ganize extracted entities into a relational table 352 format, SRAG significantly enhances the perfor-353 mance of multi-entity questions answering. The 354 exhaustive experiments conducted underscore the superior performance of SRAG in overcom-356 ing the limitations of traditional RAG-based solutions. This research not only presents a more 358 effective methodology for multi-entity OA but 359 also sets the stage for future explorations into 360 improving the accuracy and efficiency of infor-361 mation mining from large, unstructured knowl-362 edge bases. 363

7 Limitations

364

367

365 The limitations identified in the error analysis in Section 5.2 emphasize the need for better seman-366 tic understanding of queries and the use of large language models (LLMs) to improve data ex-368 traction accuracy. Another key limitation of the 369

system is its ability to handle ambiguous queries, 370 which in real-world scenarios can vary signif-371 icantly in complexity. This requires advanced 372 semantic inference to accurately grasp user in-373 tent, demanding the system to comprehend both 374 explicit content and implicit context. The chal-375 lenge is to balance precision and recall, allowing 376 diverse query handling while minimizing errors. 377 Also, we did not compare the SQL query-based 378 table method with simpler techniques such as 379 keyword extraction and summarization. This was based on the assumption that SQL's advan-381 tages in structured query handling are apparent. However, this overlooks the potential research value such a comparison could provide. 384

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-387 cia Leoni Aleman, Diogo Almeida, Janko Al-388 tenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. arXiv preprint 390 arXiv:2303.08774. 391
- Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yi-Ting Guo, and Jie Fu. 2024.

394 395	Rq-rag: Learning to refine queries for retrieval augmented generation. <i>ArXiv</i> , abs/2404.00610.	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong	436 437
	V' Chan D' La V' 'n Chan Land L'	Zhang, Sandhini Agarwal, Katarina Slama, Alex	438
396	Ain Cheng, Di Luo, Aluying Chen, Lemao Liu, Dongwon Zhoo, and Pui Yan. 2022. Lift your	Ray, et al. 2022. Training language models to	439
397	self up: Retrieval augmented text generation with	nonova in neural information processing systems	440
390	self memory	vances in neural information processing systems,	441
299	sen memory.	33.27730-27744.	442
400	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman,	Pranoy Panda, Ankush Agarwal, Chaitanya Devagup-	443
401	and Luke Zettlemoyer. 2023. Qlora: Efficient	tapu, Manohar Kaul, and Prathosh Ap. 2024.	444
402	finetuning of quantized llms. arXiv preprint	HOLMES: Hyper-relational knowledge graphs	445
403	arXiv:2305.14314.	for multi-hop question answering using LLMs. In	446
		Proceedings of the 62nd Annual Meeting of the As-	447
404	Darren Edge, Ha Trinh, Newman Cheng, Joshua	sociation for Computational Linguistics (Volume	448
405	Bradley, Alex Chao, Apurva Mody, Steven Truitt,	1: Long Papers), pages 13263–13282, Bangkok,	449
406	and Jonathan Larson. 2024. From local to global:	Thailand. Association for Computational Linguis-	450
407	A graph rag approach to query-focused summa-	tics.	451
408	rization. arXiv preprint arXiv:2404.16130.		
		Pratyush Patel, Esha Choukse, Chaojie Zhang, Inigo	452
409	Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie	Goiri, Aashaka Shah, Saeed Maleki, and Ricardo	453
410	Wang, Hengyun Li, Dawei Yin, Tat Seng Chua,	Bianchini. 2023. Splitwise: Efficient generative	454
411	and Qing Li. 2024. A survey on rag meeting llms:	Ilm inference using phase splitting. arXiv preprint	455
412	Towards retrieval-augmented large language mod-	arXiv:2311.186//.	456
413	els.	Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh	457
//1 /	Vunfan Gao, Vun Viong, Vinyu Gao, Kangyiang	Khanna, Anna Goldie, and Christopher D. Man-	458
415	Iia Iinliu Pan Yuxi Bi Yi Dai Jiawei Sun and	ning. 2024. Raptor: Recursive abstractive pro-	459
416	Haofen Wang 2023 Retrieval-augmented gener-	cessing for tree-organized retrieval.	460
417	ation for large language models: A survey arXiv		
418	nrenrint arXiv:2312 10997	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie	461
410		Huang, Nan Duan, and Weizhu Chen. 2023. En-	462
419	Xinyan Guan, Yanijang Liu, Hongyu Lin, Yaoije Lu,	hancing retrieval-augmented large language mod-	463
420	Ben He. Xianpei Han. and Le Sun. 2023. Mit-	els with iterative retrieval-generation synergy.	464
421	igating large language model hallucinations via	Lishin Tener, Maker Mara Wai Wai Lai Chi Linin	405
422	autonomous knowledge graph-based retrofitting.	Jiadin Tang, Yunao Yang, wei wei, Lei Shi, Lixin	465
		Su, Suqi Cheng, Dawei Till, and Chao Huang.	400
423	Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V. Chawla,	2024. Oraphypic Oraph Instruction tuning for	407
424	Thomas Laurent, Yann LeCun, Xavier Bresson,	A7th International ACM SIGIP Conference on Be	400
425	and Bryan Hooi. 2024. G-retriever: Retrieval-	4/In International ACM SIGIR Conference on Re-	469
426	augmented generation for textual graph under-	SICID '24 page 401 500 New York NY USA	470
427	standing and question answering.	Association for Computing Machinery.	471
428	Chunwei Liu, Matthew Russo, Michael Cafarella,	Shuai Wang Ekaterina Khramtooya Shengyao	/172
429	Lei Cao, Peter Baille Chen, Zui Chen, Michael	Zhuang and Guido Zuccon 2024 EebArag: Eval-	473
430	Franklin, Tim Kraska, Samuel Madden, and Ger-	uating federated search in the context of retrieval	475
431	ardo Vitagliano. 2024. A declarative system	augmented generation	476
432	for optimizing ai workloads. arXiv preprint	augmented generation.	770
433	arXiv:2405.14696.		
434	Meta Llama3, 2024, Meta llama3, https://llama		
435	meta.com/llama3/. Accessed: 2024-04-10.		
	C)	

Appendix 477 Α

478

481

482

483

485

486

487

488

489

491

492 493

A.1 Hops

In terms of Wikipedia graph systems, the term 479 'hop' refers to a step taken along the edges of a 480 graph from one node to another, so we consider 'hop' as a tuple (v_a, v_b, t) . For no hop, there is no relationship (edge) to track. Single hop track all entities have relationship t to v_a . Multi-hop 484 involves traversing multiple edges (hops) to find connections between nodes that are not directly linked, like track (v_a, v_b, t_1) (v_b, v_c, t_2) .

A.2 Prompt

A.2.1 Prompt for schema

Create a table schema that comprehensively captures information about {.....}. Ensure the schema is detailed and structured, avoiding over-simplification, missing elements, and redundancy. This schema should be structured so each row represents a unique instance, with each column capturing a distinct aspect of property details. Ensure there is no overlap in content between columns to avoid repetition.

A.3 Optimization 490

Two aspects of optimization are included in SRAG system to enhance the overall performance:

Model Selection. Model selection is straightfor-494 ward yet highly effective for optimization (Liu 495 et al., 2024). The SRAG system comprises mul-496 tiple tasks, necessitating the selection of the 497 most suitable model for different tasks. For ba-498 499 sic tasks, more affordable and faster LLMs can suffice, while utilization of the most advanced 500 LLMs is essential in more complex tasks to en-501 sure optimal performance. Specifically, SRAG 502 system employs powerful yet resource-intensive 503

GPT-4 for tasks such as semantic analysis or 504 generation of table schemas and SQL queries. 505 In contrast, for more basic information extrac-506 tion, we utilize open-source Mistral-7B, thereby 507 achieving a balance between cost efficiency and 508 functional performance. 509

LLM Input/Output Control. SplitWise (Patel 510 et al., 2023) shows that LLM inference time is 511 generally proportional to the size of input and 512 output tokens. Since GPT models decide the 513 cost based on the input token, we try to minimize 514 the input of large models. Meanwhile, we use 515 the instructive prompt to reduce the size of the 516 outputs generated by LLM without changing the 517 quality of these outputs. The example of prompt 518 is in Appendix A.3.1. 519

520

521

522

523

524

525

526

527

A.3.1 **Prompt for Output Control**

...review your output to ensure it meets all the above criteria. Your goal is to produce a clear, accurate, and well-structured output. Just output the {}, no other word or symbol.

A.4 Tables

Table 4 shows examples of topics and their entities' properties.

Table 5 shows examples of question templates to synthesize queries.

A.5 Automated QA Generation and Validation

We extract the introductory paragraph of textual 528 content for each entity from Wikipedia, akin to 529 an abstract of the entity's page, to derive relevant 530 property values. The preprocessing of graph 531 node properties is conducted using GPT-4. GPT-532 4 is deployed to generate the essential properties 533 of key entities for each topic, and subsequently, 534 property values are extracted from the respec-535 tive web pages of these entities. This process 536

Topics	Entities Properties	#-Entities	
ACM fellow	nationality, field of study, affiliation	1115	
Cities of the World	population, geographic coordinates, altitude,	7040	
	GDP		
Presidents of the US	term lengths, political parties, vice-presidents,	55	
Tresidents of the OS	birth states, previous occupations	55	
Chemical Elements	atomic number, atomic mass, boiling point,	166	
Chemical Elements	melting point, electron configuration	100	
Summer Olympic Comes	host cities, number of participating countries,		
Summer Orympic Games	total number of events, medal tally, records	35	
	broken		
Nobel Prize in Chemistry	categories, year of award, country of origin,	, 104	
Nobel I lize ill Chemisu y	field of contribution.	194	

Table 4: Example Topics and Their Entities Properties.

culminates in the formation of property tables.
An illustrative example of the topics and entities'
properties is provided in Appendix Table 4.

When questions or queries are posed, the 540 SRAG system efficiently navigates the graph by 541 utilizing both the connections (edges) and the 542 nodes along with the associated property tables 543 to retrieve relevant information. The property 544 545 tables, which contain attributes and values related to the entities within the graph, serve as a 546 comprehensive and structured data source that 547 can be queried alongside the graph structure. 548 This dual approach facilitates thorough analy-549 sis, as it takes into account both the relational 550 context (the connections among entities) and 551 the specific properties of the entities involved. 552 553 Moreover, such automated process benefits from low labor costs due to automation and optimiza-554 tion within the graph database system, reducing 555 the need for time-consuming and error-prone 556 manual data processing and analysis. 557

A.6 Quality Control of Questions

We devise several strategies to ensure the integrity and effectiveness of questions.

• Question Templates. The use of templates 561 ensures that every question is crafted with a clear structure, making it easier for respondents to understand and answer them 564 accurately. For relationship and complex 565 statistic questions, we turn the questions in a closed-ended style, as they require a 567 specific response of either "yes" or "no", which makes the answer in a standardized 569 format. We meticulously prepare all ques-570 tion templates, with examples in the Ap-571 pendix Table 5. 572

558

559

560

Question Refinement. After the initial development phase, each question undergoes a refinement process utilizing GPT-3.5-turbo. This stage is essential for improving the clarity, relevance, and neutrality of the questions. It also includes a thor-578

Categories	Types	Template Examples
Comparison	Intercomparison	Which has high [property], [entity A] or [entity B]?
	Superlative	Which [entity] has the highest/lowest [property]?
	Aggregation	How many [entities] have [specific property value]?
Statistics	Distribution Compliance Does [property] follow a normal distribution	
	Correlation Analysis	Is there a linear relationship between [property A] and [property B]?
	Variance Analysis	Are the variances in [property A] and [property B] significantly different?
	Descriptive Relationship	How is [entity A] related to [entity B]?
Relationship	Hypothetical Scenarios	What would be the impact if [entity A] collabo- rates with [entity B]?

Table 5: Template example for queries generated by the LLM (GPT-4).

ough review to identify and mitigate any potential bias, contributing to minimizing misunderstandings and elevating the overall quality of the questions.

• Manual review. We assess the questions for accuracy, ensuring they are factually correct and relevant to our purpose. Manual reviews can also provide insights into whether the questions are likely to effectively elicit the intended information from answers, thereby contributing to the reliability and validity of the benchmark.

A.7 Baseline Performance.

579

580

581

582

583

584

585

586

587

588

589

590

591

Introducing RAG significantly improves overall
performance, particularly in comparison tasks,
while fine-tuning LLaMA-3-Instruct alone does
not yield substantial gains without RAG. On
MEBench, open-source models like LLaMA-

3-Instruct, even with RAG, can't match proprietary models like GPT-4, which achieves a 59.3% accuracy compared to LLaMA-3-Instruct's 31.6%.