# GRAMMAR: Grounded and Modular Methodology for Assessment of Closed-Domain Retrieval-Augmented Language Models

**Anonymous ACL submission** 

### Abstract

001 Retrieval-Augmented Generation (RAG) systems are widely used across various indus-003 tries for querying closed-domain and in-house knowledge bases. However, evaluating these systems presents significant challenges due to the private nature of closed-domain data and a scarcity of queries with verifiable ground truths. 007 800 Moreover, there is a lack of analytical methods to diagnose problematic modules and identify types of failure, such as those caused by knowledge deficits or issues with robustness. To address these challenges, we introduce GRAM-MAR (GRounded And Modular Methodology for Assessment of RAG), an evaluation frame-014 work comprising a grounded data generation process and an evaluation protocol that effectively pinpoints defective modules. Our val-017 idation experiments reveal that GRAMMAR provides a reliable approach for identifying vulnerable modules and supports hypothesis testing for textual form vulnerabilities. An 022 open-source tool accompanying this framework will be released to easily reproduce our results and enable reliable and modular evaluation in closed-domain settings.

### 1 Introduction

The emergent capabilities of Large Language Models (LLMs) have driven significant research and the widespread deployment of Retrieval-Augmented Generation (RAG) systems for closed-domain settings. This growing trend highlights the critical need for reliable evaluation methods tailored to RAG systems.

A key challenge in evaluating RAG systems is the collection of domain-specific data that includes accurate ground truths. Researchers often rely on data originally designed for human assessment, as seen in studies such as Santurkar et al. (2023); Wang et al. (2022); Zhong et al. (2022, 2023); Hendrycks et al. (2021); Choi et al. (2023). However, these data sources typically reflect open-domain commonsense and world knowledge, whereas industrial in-house RAG systems require data that captures closed-domain knowledge, such as details about company projects and employees. Although user queries can be gathered during system deployment, obtaining accurate ground truths remains a significant challenge. To address this, recent studies have explored referencefree LLM evaluators for cases where ground truths are unavailable (Chern et al., 2023; Min et al., 2023; Es et al., 2023). However, the reliability of these reference-free evaluation methods is still questionable. 041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

Another critical challenge is assessing robustness in RAG systems, where a system may possess adequate knowledge for the semantics of a query but respond inconsistently to different query forms, as highlighted in Shen et al. (2023). Developers and researchers often hypothesize about factors that could impact robustness, such as vulnerable modules (e.g., retrieval mechanisms or language models) or specific input attributes. However, existing research, e.g., studies on adversarial robustness (Alzantot et al., 2018; Li et al., 2020, 2019), lacks effective tools for identifying non-robust modules and analytically testing these hypotheses.

Recognizing these challenges, the paper makes the following contributions:

**Proposing a grounded, controllable data generation process (GRAMMAR-Gen).** To address the issue of ground truth scarcity, we propose GRAMMAR-Gen, a grounded data-generation process designed to ensure reliable evaluation. This process leverages relational databases and LLMs to extract ground truths via SQL queries. Additionally, the template-based query generation enables scalable data creation. By leveraging the capabilities of LLMs, it allows the controlled generation of data with diverse linguistic attributes. Notably, a similar data-generation process for evaluating do-



Figure 1: An Example of Applying the GRAMMAR Framework for Modular Evaluation and Hypothesis Testing. The upper section demonstrates the data generation process for creating sets of hypothetically robust (Drobust) and non-robust (Dnon-robust) data. The lower section depicts the evaluation protocol that utilizes the generated data to identify defective modules and facilitate hypothesis testing.

main knowledge was proposed concurrently by Tu et al.  $(2024)^1$ .

**Proposing a protocol for modular evaluation and hypothesis testing (GRAMMAR-Eval).** GRAMMAR-Eval combines a grouping and tagging mechanism based on GRAMMAR-Gen with effective strategies to both identify knowledge gaps in the retrieval database and determine whether robustness issues stem from the retrieval mechanism or the language models (LMs). Our empirical results demonstrate the effectiveness of GRAMMAR-Eval in identifying deliberately bugged retrieval systems and accurately testing hypotheses about which input attributes contribute to robustness issues. An example of this process is illustrated in Figure 1.

### 2 Background

Evaluating the performance of Retrieval-Augmented Generation (RAG) systems, especially in closed-domain settings, poses unique challenges, leading to an absence of reliable reference-based evaluation protocols. This section outlines existing approaches and highlights the gaps addressed by our proposed framework, GRAMMAR.

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

**Retrieval-augmented Generation (RAG)** RAG systems augment LLMs with retrieval for domainspecific use where extensive databases can potentially provide necessary information. A retrieval can be used for retrieving relevant passages or documents. A dense retrieval (Mialon et al., 2023; Lewis et al., 2020; Borgeaud et al., 2022) consists of an embedding model to produce a vector  $v_q$  for any given query q. This vector  $v_q$  is then used to

<sup>&</sup>lt;sup>1</sup>It was published after the initial preprint/submission of this work.

205

206

207

208

209

210

211

166

167

identify pertinent text segments normally via dot 115 product calculations with the vectors of text chunks. 116 The most relevant chunks are aggregated to form 117 an additional context c, which serves as the factual 118 basis for the response generation, until adding an-119 other chunk would surpass the maximum allowable 120 context length. Subsequently, an LLM will process 121 this aggregated context alongside the original query 122 q to generate a predicted answer  $\hat{a}$  for it. 123

124

125

126

127

128

129

130

131

132

In contrast, sparse retrieval computes similarity without parametric embedding models. For instance, keyword-matching approaches (Manning, 2009) select documents that share the highest number of common words for context retrieval. Additionally, methods like TF-IDF (Chen et al., 2017) factor in the inverse document frequency, emphasizing the significance of less common words which are likely to be more indicative of pertinent content.

Reference-Free Evaluation Reference-free eval-133 uation methods (Es et al., 2023) assess model 134 performance without requiring predefined ground 135 truths. These methods can be categorized accord-136 ing to their applicable systems: 1) Protocols for LLMs: FactScore (Min et al., 2023) and FacTool 138 (Chern et al., 2023) use retrieval combined with 139 LLMs as validators, but these are not directly ap-140 141 plicable to RAG systems as retrieval is reserved for evaluation. 2) Protocols for RAG: ARES and 142 RAGAS focus on evaluating the faithfulness of 143 generated answers against retrieved contexts, with 144 LLMs scoring the relevance and correctness of the 145 context (Saad-Falcon et al., 2023; Es et al., 2023). 146 3) Protocols for Both: SelfCheck (Manakul et al., 147 2023) utilizes the stochastic nature of LLMs for 148 self-validation. While these methods offer some 149 insights, their reliability is often questioned, especially for closed-domain applications where spe-151 cific, accurate answers are crucial. 152

Reference-Based Evaluation While reference-153 free evaluation can be applied for evaluating 154 both open-domain and closed-domain queries, 155 reference-based evaluation typically assesses opendomain knowledge. The reason is that open-157 domain queries are widely available with ground-158 truth answers from public resources, exams, and 159 surveys (Santurkar et al., 2023; Wang et al., 2022; 160 161 Zhong et al., 2022, 2023; Hendrycks et al., 2021; Choi et al., 2023). These evaluations compare gen-162 erated answers to known correct answers, calculat-163 ing metrics that reflect the model's accuracy. However, closed-domain scenarios lack reference-based 165

evaluation due to the unavailability of ground-truth answers. GRAMMAR-Gen is proposed to solve this issue.

### **3** Is Reference-free Evaluation Reliable?

This section analyzes the reliability of referencefree evaluation methods using 198 examples from an engineering company. The RAG system for evaluation is implemented with dense retrieval and GPT- $3.5^{2}$ .

This preliminary study briefly highlights the potential benefits of reference-based evaluation, as our primary contribution focuses on proposing a reference-based evaluation framework. A detailed analysis of reference-free evaluation on state-ofthe-art models is left for future work.

**Two Evaluation Perspectives: Optimism and Cynicism** To structure this analysis, referencefree evaluation is treated as a binary classification task. The ground truth reflects the actual correctness of RAG responses, while the evaluation models' judgments serve as the task's predictions. Thus, low precision indicates an optimistic bias (where incorrect RAG responses are often judged as correct), and low recall indicates a cynical bias (where correct RAG responses are often judged as incorrect).

**Two Evaluation Protocols** We assess two reference-free protocols.

- RAGAS-Fact (Es et al., 2023): This protocol utilizes the context-query-response triplets to assess the veracity of responses. It evaluates the faithfulness of a response by calculating the ratio of claims grounded on the context to the total claims made. This process involves identifying statements that hold atomic facts, following the methodology outlined by Chern et al. (2023).
- SelfCheck (Manakul et al., 2023): This protocol relies on the stochastic generation of responses, based on the premise that incorrect answers are unlikely to be consistently produced. This principle, initially applied to LLMs, is adapted in our analysis of RAG systems. In a departure from its original application, we enhance the evaluation prompt to include not only stochastic responses but

<sup>&</sup>lt;sup>2</sup>An OpenAI embedding model for document embedding and *gpt-3.5-turbo-16k* 

212also the query itself. Refer to Appendix A213for details. To generate four stochastic sam-214ples, we adjust the temperature setting to 1.0,215contrasting with a temperature of 0.0 used to216generate the primary response. A response is217deemed correct if it aligns consistently across218all stochastic samples.

219

224

227

231

237

238

241

242

243

245

246

247

248

**Results:** Both are Extremely Optimistic on Wrong Predictions While SelfCheck Becomes Too Cynical on Correct Predictions As shown in Table 1, both RAGAS and SelfCheck achieve low precision (19% and 15%, respectively), highlighting a deficiency in correctly identifying erroneous predictions. For RAGAS, this may be because the RAG system produces responses that, while contextually relevant, fail to directly address the intended query. However, RAGAS demonstrates high recall, indicating accurate assessment of correct RAG responses. Overall, the preliminary insight suggests that reference-free evaluation methods are not reliable for evaluation, underscoring the critical need for a robust, reference-based evaluation framework to ensure reliable model assessment.

	Precision	Recall
RAGAS-Fact	19% (11-27%)	92% (86-97%)
SelfCheck	15% (0.08-0.22%)	50% (0.4-0.59%)

Table 1: Reliability of reference-free evaluation protocols. Taking the size of the sample into account, the Z-test with 95% confidence intervals is utilized.

### 4 GRAMMAR

This section introduces GRAMMAR-Gen and GRAMMAR-Eval for grounded and modular evaluation: 1) GRAMMAR-Gen begins by generating templates based on a database schema, enabling scalable data creation <sup>3</sup>. By leveraging LLMs, GRAMMAR-Gen offers controllability in generating text variations, which not only supports scalable data generation but also facilitates modular robustness evaluation and hypothesis testing. The process of SQL and text template generation is detailed in § 4.1. SQL queries, serving as intermediate representations, are used to create ground-truth answers by querying relational databases, as explained in § 4.2. The overall process is illustrated in Algorithm 1 and Figure 2. 2) GRAMMAR-Eval, an

### Algorithm 1 GRAMMAR-Gen

**Input:** SQL Template Generator  $g_{sql}$ , Text Template Generator  $g_t$ , Semantic criteria for SQL template generation  $C_{sql}$ , Linguistic criteria for text template generation  $C_t$ , Database  $\mathcal{D}$ , Database Schema  $\mathcal{S} = \{S_1, S_2, \ldots, S_n\}$ , where each  $S_i$  is a schema for a table.

**Output**: Final evaluation data Q

$Q \leftarrow \emptyset$
$\mathcal{S}_{ ext{target}} \in \mathcal{S}$
$\{tpl_{sql}\} \leftarrow g_{sql}(\mathcal{S}_{target}, \mathcal{C}_{sql})$
for $tpl_{sql} \in \{tpl_{sql}\}$ do
$\{tpl_{t}\} \leftarrow g_{t}(tpl_{sql}, \mathcal{C}_{t})$
$P \leftarrow \text{extract placeholders from } tpl_{sql}$
for each $p \in P$ do
$C \leftarrow \text{extract column from } p$
$V_p \leftarrow$ query $\mathcal{D}$ for distinct values of C
end for
$Comb \leftarrow Cartesian \text{ product of } \{V_p : p \in P\}$
for each $comb \in Comb$ do
$q_{sql} \leftarrow substitute \ comb \ into \ tpl_{sql}$
$a \leftarrow $ query $\mathcal{D}(q_{sql})$
$\{q_t\} \leftarrow \text{substitute } comb \text{ into } tpl_t \in \{tpl_t\}$
$Q \leftarrow Q \cup \{(\{q_t\}, a)\}$
end for
end for
return Q

evaluation protocol that effectively pinpoints defective modules by assessing the knowledge deficit of the retrieval database and detecting robustness issues originating from either the retrieval mechanisms or the language models. Further details can be found in § 4.3. 252

253

254

255

256

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

280

### 4.1 Generating Query Templates

The GRAMMAR-Gen process first utilizes database schemas and LLMs to generate SQL query templates and their corresponding textual forms.

**Database Schema** Database schema serves as the blueprint that defines how classes of entities are organized and the relations among them. Appendix B demonstrates the reasons why we use database schema for knowledge representation. Let a database schema with a collection of N tables  $S = \{S_1, S_2, \ldots, S_N\}$ , where each  $S_i$  is a schema for a database table. A table schema  $S_i$  consists of the table name  $T_i$ , a set of attributes  $\{A_1, A_2, \ldots, A_M\}$  and a set of constraints, i.e.,  $S_i = (T_i, \{A_1, A_2, \dots, A_m\}, C),$ where the constraint C demonstrates the primary key PK and foreign keys FK. Specifically,  $FK(T_i.A_k \to T_j.A_h)$  indicates a foreign key  $A_k$ in table  $T_i$  referencing the primary key  $A_h$  in table  $T_i$ . The schema defines the knowledge structure of an entity, which can be enough to generate query semantics that is not tied to any specific linguistic

<sup>&</sup>lt;sup>3</sup>Further details can be found in Appendices D.2 and E

284

286

287

290

291

296

297

300

expressions and sensitive data.

**Generating SQL Templates** A SQL template  $tpl_{sql}$  essentially represents the semantics of a query. Various query semantics are formulated through the application of relational algebra operators. In our study, we focus on three fundamental operators: Select  $(\sigma)$ , Project  $(\pi)$ , and Join  $(\bowtie)$ . An automated SQL template generator  $g_{sql}$  is powered by a generative LLM. When given a target schema Starget and specified SQL criteria Csql,  $g_{sql}$  generates a set of SQL query templates.

**Generating Text Templates**  $tpl_{sql}$  can be instantiated into various text templates, each denoted as  $tpl_t$ . This process, driven by the text template generator  $g_t$ , leverages linguistic criteria  $C_t$  to transform  $tpl_{sql}$  into natural language forms that align with specified linguistic characteristics, e.g., complexity, length and stylistic nuances <sup>4</sup>. The text template generator  $g_t$  is also operationalized using a sophisticated language model, such as GPT-4, which is adept at producing diverse linguistic variations of the same semantic content.

$$\{tpl_{\mathsf{t}}\} = g_{\mathsf{t}}(tpl_{\mathsf{sql}}, \mathcal{C}_{\mathsf{t}}),$$

where  $\{tpl_t\}$  signifies the resulting list of text templates.

### 4.2 From Templates To Evaluation Data

Utilizing a database  $\mathcal{D}$  that aligns with the predefined schema, the framework samples a diverse range of queries along with their corresponding ground-truth answers. This dataset, derived from SQL and text templates filled with database content, provides a rich source for evaluation.

Generating Queries Via Placeholder Fill-in The placeholder fill-in step populates the templates with actual data from the rows of tables in  $\mathcal{D}$ . Specifically, the approach employs a SE-LECT query format: "SELECT DISTINCT {column name} FROM {table name};" to ensure variety in the data points. It handles SQL query templates with multiple placeholders by employing 309 combinations of placeholders and using Cartesian products to generate multiple query permutations. 310 This approach leads to a comprehensive set of SQL 311 queries  $(q_{sql})$  and their natural language equiva-312 lents, the text queries  $(q_t)$ . 313

Generating AnswersAnswer generation is in-<br/>tegral to completing our dataset. Each SQL query314 $(q_{sql})$  is executed against  $\mathcal{D}$  to match with a factual<br/>answer a. Each answer is then paired with the rel-<br/>evant text queries  $\{q_t\}$ , forming the basis of our<br/>evaluation data.314

321

322

323

324

325

326

327

328

329

330

332

333

334

335

336

### 4.3 Modular Evaluation

Evaluating the overall performance of an RAG model can obscure the specific weaknesses or robustness of its constituent modules. This section introduces a modular evaluation protocol to address this.

**Query Grouping and Tagging** The process involves grouping and tagging queries based on their semantics. Specifically, with GRAMMAR-Gen, textual queries Q generated from a particular SQL query  $q_{sql}$  are organized into a semantic group. Depending on the model M's performance,  $q_{sql}$  (and consequently, the group Q) are tagged into three principal categories, as shown in Figure 1.

• Gap Groups emerge when the model M consistently provides incorrect answers for every text query within a subset  $S \subseteq Q$  produced through the data generation process. This reflects a deficiency in M 's knowledge or capability. Formally, this scenario is described as:

$$\forall q_t \in S, \neg M_e(M(q_t)),$$

where  $M_e$  is the evaluation model, which returns "TRUE" if the model's response  $M(q_t)$  is correct.

• **Robust Groups** are established when model *M* accurately responds to all text queries within the subset *S*, evidencing a robust comprehension of the SQL logic. Formally, this is articulated as:

$$\forall q_t \in S, M_e(M(q_t)).$$

• Non-robust Groups are characterized by model M's ability to potentially answer the query semantics correctl, but with at least one query  $q_t$  within S that M fails to predict correctly. This is formally denoted as:

$$\exists q_t \in S, M_e(M(q_t)) \land \exists q_t \in S, \neg M_e(M(q_t)).$$

<sup>&</sup>lt;sup>4</sup>Examples of  $C_{sql}$  and  $C_t$  are specified in Appendix C.





Assessing Retrieval Database via Gap Groups
Gap examples, which are incorrectly predicted instances within gap groups, highlight errors in the
retrieval database that arise due to its limited coverage. The adequacy of the knowledge within the
database can be quantified as follows:

$$Acc_{retrieval\_db} = 1 - \frac{\text{Number of Gap Groups}}{\text{Total Number of Semantic Groups}}$$
(1)

343

347

352

359

A low Acc<sub>retrieval\_db</sub> suggests the need for expanding the database.

**Isolating Retrieval Database Errors** To ensure accurate evaluation of the retrieval and language models without interference from retrieval database errors (i.e., distinguishing knowledge gaps from robustness performance), two strategies have been developed:

1. Removing Gap Examples: To assess the other two modules, we calculate refined accuracy, denoted as *R*, by excluding gap examples:

$$R = \frac{\text{\# with Correct Predictions}}{\text{Total \# - Total \# in Gap Groups}}, \quad (2)$$

where "#" represents "the number of instances".

2. Balancing Gap Groups: To understand the advantage of this refined accuracy over baseline accuracy Acc, consider the following relationship:

360

361

362

363

364

365

367

369

370

371

372

373

375

376

378

379

381

382

$$Acc = R \times \frac{\text{Total } \# - \text{Total } \# \text{ in Gap Groups}}{\text{Total } \#}$$
$$= R \times (1 - \frac{\text{Total } \# \text{ in Gap Groups}}{\text{Total } \#})$$
(3)

This equation shows that a smaller normal-Total # in Gap Groups, enhances ized term,  $\lambda =$ Total # evaluation accuracy. Intuitively, a smaller  $\lambda$  implies fewer gap examples (errors propagated from the initial module), leading to more accurate results. Maintaining a consistent  $\lambda$  across evaluation models or datasets ensures accurate assessments, such as for the hypothetically robust dataset  $\mathcal{D}_{robust}$  and the non-robust dataset  $\mathcal{D}_{non-robust}$  discussed in §5. With GRAMMAR-Gen, an equal number of text queries can be generated within each semantic group by adjusting the text template generator.

Identifying Non-Robust Retrieval or Non-Robust LM: A "Context Comparison" Approach To determine whether robustness issues stem from the retrieval module or the language model, we focus on incorrectly predicted queries within the non-robust group, each referred to as

	SQL	SQL	Text	Text	Queries	Text Queries (	Balanced)
	Templates	Queries	Templates	Short	Long	Short	Long
Aurp	11	157	30	314	257	471	471
Spider-Open	5	57	10	397	436	570	570
Spider-Closed	5	57	10	426	430	570	570

Isolating Errors from		Aurp		Spider-Closed		Spider-Open	
Retrieval Database	LLM	Robust	Non-Robust	Robust	Non-Robust	Robust	Non-Robust
Baseline: No Action	No Action	0.27	0.51	0.26	0.28	0.67	<b>0.87</b>
	Context Comparison	0.29	0.51	0.31	0.34	<b>0.37</b>	0.35
Remove Gap Groups	No Action	0.95	<b>0.96</b>	0.58	0.44	0.98	0.96
	Context Comparison	1	0.98	0.62	0.49	0.51	0.39
Balance Gap Examples	No Action	0.27	0.27	0.28	0.21	0.69	0.68
	Context Comparison	0.29	0.28	0.32	0.28	0.36	0.29

Table 3: Validating GRAMMAR's tricks in in detecting non-robust retrieval module under a deliberately constructed RAG pipeline. Accuracy is calculated on selected examples across various types of semantic groups.

 $q_{\text{target}}$ . Specifically, we pose the following question:

### Does the retrieval module provide sufficient context for $q_{target}$ ?

By leveraging other semantically identical but correctly answered queries within the group, we can assess whether the context provided is sufficient. Specifically, we compare the document indices idxtarget associated with qtarget to the indices of correctly answered queries, denoted as I = idx1, idx2, ... If idxtarget matches any index in I, the context provided by idxtarget is considered sufficient.

### 5 A Toy Case for Validation: Detecting A Vulnerable Retrieval

This section demonstrates the effectiveness of GRAMMAR on question answering tasks.

### 5.1 Experiment Setup

383

389

393

397

398

400

401

402

403

404

405

406

407

408

409

**Evaluation Target** To validate our framework, we deliberately introduce a vulnerability in the retrieval component by using a basic keywordmatching retrieval approach. We use GPT3.5 for language generation. While this method is less common in practical applications, it is inherently weak in handling lengthy queries, making it an ideal candidate for testing our evaluation framework.

410 **Evaluation Benchmarks and Datasets** Closed-411 domain RAG systems lack established benchmarks, especially because questions must be answered using documents containing knowledge that is not publicly available or known to public LLMs. To address this gap, we modify the Spider benchmark (Yu et al., 2018) and create two synthetic benchmarks with fictitious data. 1) Context-Augmented Spider (Open-Domain): We use the "company\_employee" relational database from Spider to simulate a general real-world scenario. Documents for retrieval are synthesized using GPT-4 based on the database. 2) Context-Augmented Spider (Closed-Domain): This is similar to the first dataset, but with company names modified to fictitious entities that are either unused by specific companies or generic according to the latest GPT-4 models. 3) Aurp: To create a strictly closeddomain evaluation, we synthesize fictitious facts about a fictional company called Aurp. A relational database and corresponding retrieval documents are generated so that questions derived from the database can only be answered by the provided documents. Details on the data synthesis process are provided in Appendix G, and the database schemas are outlined in Appendix F. Using these three settings, datasets are generated by GRAMMAR-Gen to validate the design of GRAMMAR-Eval. The statistics of the generated datasets are summarized in Table 2.

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

### 5.2 Results

This section validates the effectiveness of modular evaluation within the proposed framework in

identifying defective modules. 443

**GRAMMAR-Eval's Effectiveness on Modular** 444 **Evaluation** The results in Table 3 highlight the 445 value of modular evaluation in identifying the de-446 fective retrieval. Regardless of the methods used to 447 isolate errors from the retrieval database and LLMs, 448 the accuracy consistently identifies the vulnerable 449 retrieval module. The next three paragraphs will 450 provide a more detailed analysis. 451

Preventing Forward Error Propagation Pre-452 venting the propagation of knowledge deficits from 453 retrieval databases is crucial. The two proposed 454 strategies effectively address this issue. In all 455 benchmarks, after either removing gap groups or 456 balancing gap examples, accuracy on the robust 457 text form consistently outperforms the non-robust 458 text form (with the exception of Row 3 in Table 3, 459 where the results are very close). 460

Impact of Backward Errors from LLMs on Re-461 trieval Assessment In some cases, such as Row 462 3 in Table 3, language model inaccuracies reduce 463 the accuracy, masking the true performance of re-464 trieval. While these instances are not prevalent in 465 our dataset, they still distort the results. This dis-466 tortion may be due to the stochastic nature of the 467 468 LM's performance, or, in a less favorable scenario, because the non-robust text forms used in retrieval 469 inadvertently benefit the LM. Based on the results 470 from the other two datasets, it is more likely the 471 former. 472

Applicability of GRAMMAR to Open-Domain 473 **RAG** Table 3 also demonstrates the effectiveness 474 475 of GRAMMAR for performing modular evaluation on the open-domain Spider dataset. However, in 476 open-domain settings, the GPT model may generate accurate answers from its internal knowledge. 478 This can lead to inaccuracies in assessing both the 479 retrieval system's knowledge gaps and the robust-480 ness of the retrieval mechanism. Specifically, the potential limitations in applying GRAMMAR-Eval 482 to open-domain settings include: 1) Misjudgment 483 of the Retrieval Database: Some gap examples in the retrieval database may go undetected, leading 485 to an inflated Accretrieval\_db. This issue is em-486 pirically demonstrated in Table 4, which compares Accretrieval\_db before and after filtering out opendomain queries for Spider-Open. Closed-domain 489 queries are identified by applying GRAMMAR to 490 detect gap groups using answers from the LLMonly system. 2) Limitations of Context Com-492

477

481

484

487 488

491

parison: Context comparison may fail when cor-493 rectly answered examples do not actually receive 494 sufficient context. This issue is evident in Table 495 3, where the accuracy decreased under "Context 496 Comparison" compared to "No Action." The de-497 crease in accuracy is likely due to the introduction 498 of incorrect context, which adds noise to the lan-499 guage model's generation, thereby hindering its 500 performance. In contrast, accuracy consistently 501 improves for the two closed-domain benchmarks 502 when context comparison is applied.

	$Acc_{retrieval\_db}$
With Open-Domain Queries	0.53
Without Open-Domain Queries	0

Table 4: Wrong judgement on retrieval database with open-domain queries

#### 6 Conclusions

In this work, we have introduced the GRAMMAR framework specifically designed for Retrieval-Augmented Generation (RAG) systems in closeddomain settings. The data generation module (GRAMMAR-Gen) improves the reliability of evaluation, while the evaluation protocol (GRAMMAR-Eval) enables detailed modular analysis.

**Reproducibility** Our open-source Python package enables easy reproduction of all results, including the regeneration of experimental data and outcomes from our validation study. Detailed instructions on using the package are provided in Appendix H. Due to data privacy policies, the 198 examples and the RAG system from the sponsoring company are not included.

#### 7 Limitations

The data generation process has certain limitations, such as its inability to handle scenarios with multiple correct answers (refer to Appendix D.1). Additionally, the expressiveness of the generated queries is constrained by the database schema and SQL (refer to Appendix D.3), limiting the ability to generate queries that require multi-step reasoning or free-form responses. However, as demonstrated in the case study, simple queries are sufficient for identifying and addressing intrinsic robustness issues within the models.

Besides, while the simple context comparison method performs well for closed-domain evalua504

505

507

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

tion, it may need refinement for more complex
retrieval tasks. This method was primarily used
to validate the importance of modular evaluation
and to establish a foundation, along with a reusable
code base, for future development.

### References

539

540

541

542

543

544

545

546

547

548

550

552

553

554

555

557

558

559

560

561

562 563

564

567

568

571

573 574

577

578

579

581

585

- Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. 2018. Generating natural language adversarial examples. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 2890–2896, Brussels, Belgium. Association for Computational Linguistics.
  - Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022.
     Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.
  - Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to answer opendomain questions. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1870–1879, Vancouver, Canada. Association for Computational Linguistics.
  - I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, Pengfei Liu, et al. 2023. Factool: Factuality detection in generative ai–a tool augmented framework for multi-task and multi-domain scenarios. *arXiv preprint arXiv:2307.13528*.
  - Jonathan H Choi, Kristin E Hickman, Amy Monahan, and Daniel Schwarcz. 2023. Chatgpt goes to law school. *Available at SSRN*.
  - Shahul Es, Jithin James, Luis Espinosa-Anke, and Steven Schockaert. 2023. Ragas: Automated evaluation of retrieval augmented generation.
  - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
     2021. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledgeintensive nlp tasks. In Advances in Neural Information Processing Systems, volume 33, pages 9459– 9474. Curran Associates, Inc.

Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2019. Textbugger: Generating adversarial text against real-world applications. In 26th Annual Network and Distributed System Security Symposium, NDSS 2019, San Diego, California, USA, February 24-27, 2019. The Internet Society. 586

587

589

590

592

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. BERT-ATTACK: Adversarial attack against BERT using BERT. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6193–6202, Online. Association for Computational Linguistics.
- Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models.
- Christopher D Manning. 2009. An introduction to information retrieval. Cambridge university press.
- Grégoire Mialon, Roberto Dessi, Maria Lomeli, Christoforos Nalmpantis, Ramakanth Pasunuru, Roberta Raileanu, Baptiste Roziere, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. 2023. Augmented language models: a survey. *Transactions on Machine Learning Research*. Survey Certification.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. *arXiv preprint arXiv:2305.14251*.
- Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. 2023. Ares: An automated evaluation framework for retrieval-augmented generation systems.
- Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect? *arXiv* preprint arXiv:2303.17548.
- Xinyue Shen, Zeyuan Chen, Michael Backes, and Yang Zhang. 2023. In chatgpt we trust? measuring and characterizing the reliability of chatgpt. *arXiv preprint arXiv:2304.08979*.
- Shangqing Tu, Yuanchun Wang, Jifan Yu, Yuyang Xie, Yaran Shi, Xiaozhi Wang, Jing Zhang, Lei Hou, and Juanzi Li. 2024. R-eval: A unified toolkit for evaluating domain knowledge of retrieval augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD24-ADS).*
- Siyuan Wang, Zhongkun Liu, Wanjun Zhong, Ming Zhou, Zhongyu Wei, Zhumin Chen, and Nan Duan. 2022. From lsat: The progress and challenges of complex reasoning. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 30:2201–2216.

741

694

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3911–3921, Brussels, Belgium. Association for Computational Linguistics.

642

643

645

646

647

651

653

658

671

672

673

677

678

685

- Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. Agieval: A human-centric benchmark for evaluating foundation models. arXiv preprint arXiv:2304.06364.
- Wanjun Zhong, Siyuan Wang, Duyu Tang, Zenan Xu, Daya Guo, Yining Chen, Jiahai Wang, Jian Yin, Ming Zhou, and Nan Duan. 2022. Analytical reasoning of text. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2306–2319, Seattle, United States. Association for Computational Linguistics.

### A Prompt Templates for LLM-Based Evaluation

For SelfCheck, given that our responses are typically concise, we opt for assessing the overall response rather than rather than performing a sentence-by-sentence validation as in the original implementation, as shown in Table 5.

# B Knowledge Structure and Representation

Commonly, the knowledge is formalized by entities along with their relations, e.g., a knowledge graph, ontology and a database schema. Instead of using knowledge graphs containing concrete entities and their relations, we use the schema-based definition, where the schema is characterized by entity types (classes of entities). An entity type is defined by its name, its attributes and relations with other entity types that will be concretized as components in a database, i.e., the table name, columns and foreign keys.

The definition of entity types can provide combinatorial expansion of data generation, protecting the leak of private information (normally stored as rows of tables) during the data generation process and the ability of retrieving ground-truths via Structured Query Language (SQL). The three advantages are the reasons why we represent knowledge in database schema rather than knowledge graphs. It disentangles the private data stored in the database schema. Since schemas and templates contain only meta-data for structuring the real data, there is no private issue to use commercial large language models for this process.

## C Generator Criteria

Each criterion for the SQL generator and text generator ( $C_{sql}$  and  $C_t$ , respectively) consists of brief instructions provided to the LLM, specifying the type of SQL or text that needs to be generated.

Table 6 provides examples:

### **D** Details of SQL Template Generation

### D.1 Constraints/Guidelines for Creating SQL Templates

The constraints limit the types of templates that can be generated. This section demonstrates two important constraints, which are verbalized as the prompt for LLMs.

Selecting Attributes That Are Understandable to Humans The SQL template generator is required to follow the constraint: "The selected and condition columns in the query MUST BE MEAN-INGFUL and DESCRIPTIVE to ensure the queries are easily understood by non-technical users.".

Avoiding Answer Multiplicity In the evaluation of question-answering (QA) models, a unique challenge arises from the existence of multiple valid answers to a single query, which necessitates a nuanced approach to assessing model performance. Consider the question: "Get the name of the client associated with the project named Innovation Precinct" For such a question, a set of correct responses could include any combination of names from a predefined list, such as {Apple, Amazon, Meta, Facebook }. This multiplicity of correct answers underscores the complexity of evaluating QA models, as it requires the assessment mechanism to recognize and validate the full spectrum of possible correct answers rather than comparing the model's output against a single 'gold standard' answer. This scenario demands a more flexible method for answer validation that can accommodate the variability in correct responses. Also, evaluation metrics may be required to effectively measure the performance of QA models in handling diverse and equally valid answers. Such metrics must account for the exhaustive set of correct answers and evaluate the model's ability to retrieve any or all valid responses within the context of the query, thereby ensuring a comprehensive assessment of the model's

SelfCheck from (Manakul et al., 2023) Context: {context} Sentence: {sentence} Is the sentence supported by the context above? Answer Yes or No: SelfCheck-QA Query: {query} Answer A: {answer} Answer B: {stochastic\_answer} Do both answers address the query with equivalent meaning? Use only "Yes" or "No" for your evaluation:

**Ragas** from (Es et al., 2023) Natural language inference. Use only 'Yes' (1), 'No' (0) and 'Null' (-1) as verdict. context: {context} statement: {statement} verdict:

### Ours

Evaluate the accuracy of the given response in relation to the true answer for the specified query. After evaluating, provide a judgement as either "Correct" or "Incorrect" based on whether the ##Given Response## accurately matches the ##True Answer##. ##Query##: {query} ##True Answer##: {true\_answer} ##Given Response##: {given\_response} ##Judgement##:

Table 5: Prompt templates for LLM-based Evaluation.

 $C_{sql}$ : Generating SQL queries with one placeholder

Each query must contain at least one parameter placeholder in the WHERE clause.

 $C_t$ : Generating short queries

Short and Clear: Keep your queries short and straightforward. Cut down on words and skip parts of speech, such as conjunctions and articles. It's okay to use fragmented phrases as long as they still convey the full meaning. Valid examples: "client of '[Project.Name]'" or "client for '[Project.Name]'"; Invalid Examples: "Find the client of a project named '[Project.Name]'.

 $\mathcal{C}_t$ : Generating long queries

Complex Sentence Structure: Ensure your queries are always in complete sentences. Opt for longer, more complex sentence structures, incorporating elements of speech like conjunctions and articles for fuller expression. Each query should be at least 30 words long. You can add context and background information to the query.

Table 6: Examples of criteria for the SQL generator and text generator.

understanding and response generation capabili-742 ties. Generally, encompassing all valid responses 743 within the database is not feasible. Hence, rather than let the non-robust results affect the accuracy 745 of evaluation, we look for an approach to avoid it, 746 i.e., only generating queries ensuring one ground-747 truth answer. Specifically, the evaluation of certain 748 SQL queries, such as "SELECT Name, Depart-749 ment, BusinessAddress FROM Employee WHERE 750 JobTitle = 'designer''', requires a complete and

thorough listing of multiple answers that can be dy-752 namically changed. To solve the issue, the solution 753 involves adding a specific criterion to the prompt 754 for generating SQL queries. This criterion, "Ensure 755 the query yields a specific and singular answer", 756 aims to produce queries that result in a single, clear 757 answer, thereby avoiding the complexities of mul-758 tiple possible correct answers. For example, with 759 this criterion, the query "SELECT Industry FROM 760 Company WHERE Name = '[Company.Name]';" 761

You are a SQL query Template Generator: Generate ACCEPTABLE SQL query templates with placeholders according to the give data schema and requirements. A simple example of an acceptable SQL query template is: SELECT Industry FROM Company WHERE Name = '[Company.Name]'; You must follow the basic criteria below except for other requirements: ##CRTTFRTA## - The placeholder format should be a combination of a table name and a column name, enclosed within square brackets, e.g., '[User.Name]'. - Use only 'SELECT' queries. - Select specific column(s) instead of using '\*'. Avoid projecting attributes that appear in the predicate. - The selected and condition columns in the query MUST BE MEANINGFUL and DESCRIPTIVE to ensure the queries are easily understood by non-technical users. - Avoid using technical column names that don't clearly signify the nature of the entities or objects involved, e.g., column for semantically void record identifiers. - Do not create redundant or semantically duplicated queries when translated into natural language. - Each query must contain at least one parameter placeholder in the WHERE clause Ensure the query yields a specific and singular answer to avoid multiplicity issues, thus facilitating accurate chatbot evaluation. {SPECIFIC\_REQUIREMENTS} - If no acceptable SQL template can be generated with the given table and column information, do not generate any text. ##RESPONSE FORMAT## - Output each SQL template as a single line, without any prefix or suffix. - Do not include any other text in your response, even something like ##RESPONSE\_END##. ##DATA SCHEMA## {GIVEN\_SCHEMA} ##RESPONSE\_START##

Table 7: A prompt template designed to guide the LLM in functioning as a controllable SQL template generator

is preferred as it's likely to yield a singular answer
about a company's industry based on a specific
company name. In contrast, without this criterion, queries like "SELECT Name FROM Company WHERE Industry = '[Company.Industry]';"
are acceptable but may result in multiple names,
leading to evaluation difficulties due to data completeness and query multiplicity issues.

### D.2 Possibilities Of SQL Templates

779

**One Table/Entity** Given a schema with only one table, the possibilities for SQL templates can be analyzed by considering the basic SQL operators like SELECT, WHERE. Here's a breakdown:

• Variation in Selected Attributes (SELECT): The number of SQL templates varies based on the combination of columns selected. If the "Company" table has n columns, then theoretically, there are  $2^n - 1$  possible combinations of columns for selection (excluding the case where no column is selected).

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

• Conditions in Queries (WHERE): Each SQL query can include zero or more conditions in the WHERE clause. The number of possible conditions is determined by the number of columns, the type of each column (text, numeric, date, etc.), and the range of operators applicable to these types (like =, <, >, LIKE, IN for text columns; =, !=, <, >, BETWEEN for numeric columns). The complexity increases combinatorially with multiple conditions combined using AND/OR.

**Two Entities** The relations between tables further amplifies the number of possible templates. For example, with two entities, "Company" and "Project," and their associative table "Company\_Project", the possibilities for SQL templates

- 799
- 8
- 8
- 8(
- 8(
- 804
- 8
- 8
- 810 811
- 812
- 813 814
- 815 816
- 817
- 818 819
- 820
- 821 822

- 824 825
- 8
- 827 828

829

8

833

- 835
- 836 837
- 83 83 84

expand significantly due to the introduction of joins
and more complex WHERE clauses. Let's break
down the possibilities.

- Selection Variations (SELECT): The number of SQL templates grows with the combination of columns selected across the three tables: "Company", "Project" and "Company\_Project". If "Company" has n columns, "Project" has m columns, and "Company\_Project" has p columns, the possible combinations for selection are  $(2^n 1) \times (2^m 1) \times (2^p 1)$ .
- Join Conditions (JOIN): The introduction of the associative table "Company\_Project" allows for meaningful JOIN operations between "Company" and "Project." Templates can include joins like Company JOIN Company\_Project ON condition and Project JOIN Company\_Project ON condition, or a multitable join linking all three. The variety of JOIN conditions adds another layer of complexity to the possible templates.

• WHERE Clause Complexity: With more tables, the WHERE clause can include a wider range of conditions, potentially involving attributes from any of the three tables. The complexity increases with the number of columns and their types across all tables, and combinations of these conditions.

# D.3 Limitations

**Scope of Using Automated SQL Generators** The effectiveness of queries generated by LLMs hinges on the meaningfulness of table and column names to accurately reflect the essence of the entities or concepts they represent. If these names lack contextual clarity, the resulting queries may be impractical. To mitigate this, practitioners can conduct an informal assessment or consult established benchmarks that evaluate LLMs' proficiency in domain knowledge and commonsense reasoning, for instance, (Zhong et al., 2023; Hendrycks et al., 2021).

840Queries Beyond SQL ExpressivenessWhile841SQL and relational algebra offer a wide range of842operations enabling the formulation of numerous843user queries, certain semantic nuances exceed the844expressive capabilities of SQL. The following ex-845amples illustrate such limitations:

• Real-world queries often contain ambiguity and subjective interpretations that SQL struggles to accommodate. For instance, the term "major" in the query "What are some major rail projects we've been involved with?" does not directly translate into SQL criteria without additional interpretative steps.

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

• SQL queries are typically structured to elicit specific, predefined responses, contrasting with the open-ended nature of many real-world inquiries that seek exploratory or comprehensive answers.

# E Breakdown of Scalable Data Generation

The scalability of the proposed data generation process is achieved through two key factors: 1) the creation of templates using various SQL formulations from a modest, manually crafted database with a limited number of attributes and rows, and 2) the generation of text queries with LLMs. The potential number of query permutations can be expressed succinctly by the formula (see Figure 2 for a visual representation):

Total Query Variations =  $M \times N \times Q$ , (4)

where "M" denotes the array of SQL templates that can be extracted from a given schema, "N" refers to the vast number of text templates that can be produced from a single SQL template, showcasing the flexibility of natural language, and "Q" accounts for the broad spectrum of text queries that can be generated from a single text template, with this diversity arising from different combinations of unique values filling the placeholders.

Let's break down each step of the proposed data generation process:

1 schema ⇒ M SQL Templates: The number of possible SQL templates (M) that can be generated from a defined schema is influenced by the SQL operators used, the number of tables, and the columns in the database schema. Each combination of tables and columns, along with different SQL operators (like SELECT, WHERE, JOIN), can lead to a unique SQL query template, e.g., the exponential complexity of the SELECT operator, and the combinatorial increase of possible predicates (See Appendix D.2 for details). Note that while the theoretical maximum is

982

983

984

985

986

987

940

941

942

943

944

945

946

947

high, practical and meaningful queries will be a smaller subset, as specified in the case study and Appendix D.

895

926

927

928

929

930

931

932

- 897• 1 SQL Template  $\Rightarrow$  N Text Templates: For898each SQL template, there can be an arbitrary899number (N) of textual expressions. This varia-900tion arises from the different ways to linguis-901tically express the same SQL query due to the902flexibility and richness of natural language,903e.g., the variation of synonyms and sentence904structures.
- 1 Text Template  $\Rightarrow$  Q Text Queries: Each text 905 template can lead to a number of text queries 906 (Q), depending on the unique values avail-907 able for each placeholder. With more than 908 one placeholder, the potential for growth in 909 the number of text queries is combinatorial, 910 barring semantic conflicts (context conflict), 911 where certain combinations of column values 912 may not be semantically valid. Overall, this 913 framework demonstrates a combinatorial ex-914 pansion at each transition stage, especially 915 notable in steps involving natural language 916 due to its inherent variability. However, this 917 growth is tempered by practical constraints 918 such as the meaningfulness of queries (seman-919 tic validity) and the actual data distribution 920 in the database. The exponential increase is 921 most pronounced in the transition from text templates to text queries, where the permuta-923 tions of placeholders can lead to a vast array 924 of unique query possibilities. 925

# F Database Schema Overview

Figures 3a and 3b illustrate the database schemas used for an actual industrial context and a fabricated scenario, respectively, within our research.

# G Synthetic Processes for Data Generation

# G.1 Aurp Setup

The generation starts from a company profile, i.e., a fictitious company named Aurp, and then goes to structural knowledge to create a relational database, including organizational structures, employees, clients and projects. Finally, all the information above is used to generate synthetic project documents. Below is a detailed process.

- 1. Generating Company Profile: Initially, a company profile for Aurp is created, detailing its foundation year, headquarters location, CEO, number of employees, and the services it offers, such as bespoke architectural solutions, sustainable urban planning, and structural health monitoring.
- 2. Generating Organizational Structure: Next, a project-oriented organizational structure is established, naming key positions and employees within the company, similar to real-world firms. This includes a wide range of roles from executive positions to specialized engineers and support staff.
- 3. Generating Employee Information: For each employee listed in the organizational structure, detailed job titles, departments, and direct supervisors or managers are fabricated, creating a network of relationships and reporting lines within the company.
- 4. Generating Client Information: The process creates a list of ten clients across various industries—ranging from technology and real estate to hospitality and healthcare—each with specified locations, thereby illustrating the company's diverse portfolio.
- 5. **Generating Projects**: Specific projects are then devised, including names, locations, start and end dates, clients, project directors, and project managers. This step integrates previously generated data (client and employee information) to create realistic project scenarios.
- 6. Generating Project Reports: This step involves synthesizing data from the previous steps to produce detailed analyses, updates, and outcomes of the various projects, which may or may not contain detailed information for clients and employees.

# G.2 Context-Augmented Spider

Since relational databases have been constructed in Spider, simple company and employee profiles with factual knowledge are manually crafted using the stringified templates in Table 8. The simple profiles are then used to prompt GPT-4 to generate Wikipedia-style pages containing 200-400 words. Knowledge required for queries is sparse in all documents.



(a) Spider.

(b) Aurp.

Figure 3: Entity-Relationship Diagrams for Data Generation

Company Profile \*\*Name\*\*: {company\_name} \*\*Headquarters\*\*: {headquarter} \*\*Industry\*\*: {industry} \*\*Sales in Billion\*\*: {sales} \*\*Profits in Billion\*\*: {profits} \*\*Assets in Billion\*\*: {assets} \*\*Market Value in Billion\*\*: {market\_value} Employee Profile

\*\*Name\*\*: {name}
\*\*Age\*\*: {age}
\*\*Nationality\*: {nationality}
\*\*Graduation College\*\*: {graduation\_college}

Table 8: Templates for company and employee profiles. Placeholders in curly brackets will be replaced by factual information in Spider databases.

**Evaluation of Factuality in Generated Documents** To evaluate the factuality of generated documents, a Machine Reading Comprehension (MRC) task is setup, where the generated document is given as context to answer the questions encoding the facts it generate on. Note that the correct answer can only source from the generated document.

### H Reproducing Data and Results

### H.1 Reproducing Data

989

990

991

993

995

997

1002

1003

1004

1005

1006

1007

1008

The following steps illustrate how to use our opensource tool to reproduce the data. The datasets are also available along with the tool.

```
from grammar.db_tool import DBTool
from grammar.llm import AnyOpenAILLM
from grammar.sql_template_generator
    import SQLTemplateGenerator
from grammar.text_template_generator
    import TextTemplateGenerator
from grammar.qa_generator import
    QADataGenerator
```

```
1010
llm = AnyOpenAILLM(model_name = "gpt4-
                                                   1011
    short")
                                                   1012
                                                   1013
setup_env = "aurp"
                                                   1014
if setup_env == "spider" or setup_env ==
                                                   1015
     "spider_closed":
                                                   1016
    database_name = 'spider'
                                                   1017
    connection_string = f'sqlite:///{
                                                   1018
    database_name }/rel_database/
                                                   1019
    company_employee.sqlite
                                                   1020
    schemas = [('company',)
                                                   1021
                               ('people',)
      ('company', 'people')]
                                                   1022
elif setup_env == "aurp":
                                                   1023
    database_name = 'Aurp
                                                   1024
    connection_string = "mysql+pymysql
    ://root:password@localhost:3306/Aurp
                                                   1026
                                                   1027
    schemas = [('client',), ('employee'
                                                   1028
    ,), ('project', )]
db_tool = DBTool(connection_string)
                                                   1030
# Step 1: Generate SQL Query Templates
                                                   1031
file_path = f"{setup_env}/
                                                   1032
    SQLTemplateGenerator/sql_templates.
                                                   1034
    ison'
sql_template_generator =
                                                   1035
    SQLTemplateGenerator.from_file(
                                                   1036
    file_path, sql_connection=
                                                   1037
                                                   1038
    connection_string, llm=llm)
entities_to_sql_templates =
                                                   1039
    sql_template_generator.
                                                   1040
    generate_batch(schemas, override=
    False, verbose=True)
sql_templates = [tpl for entity, tpls in
                                                   1043
     entities_to_sql_templates.items()
                                                   1044
    for tpl in tpls]
                                                   1045
```

1009

1046

1047

1048

Below are examples of generated SQL templates in JSON format, where the templates are keyed by table names.

	1043
	105
"('client',)": [	105
"SELECT Location FROM	105
Client WHERE Name = '[Client.	1053
Name]';",	1054

```
"SELECT Industry FROM
   Client WHERE Name = '[Client.
   Name]';"
    ],
    "('employee',)": [
         "SELECT JobTitle FROM
   Employee WHERE Name = '[
   Employee.Name]';",
         "SELECT Department FROM
   Employee WHERE Name = '[
   Employee.Name]';",
         " SELECT
   SupervisorOrManager FROM
   Employee WHERE Name = '[
   Employee.Name]';"
    ],
    . . .
3
# Step 2: Generate Text Query Templates
linguistic_attr = "long"
file_path = f'{setup_env}/
   TextTemplateGenerator/{
   linguistic_attr}.json
```

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1873

1074

1075

1076 1077

1078

1079

1081

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1188

1109

1110

1111

{

```
linguistic_attr}.json'
text_template_generator =
    TextTemplateGenerator.from_file(
    file_path=file_path, verbalize_attrs
    =linguistic_attr, llm=llm) # Load
    existing generations to avoid re-
    generation
sql_to_text_templates =
    text_template_generator.
    generate_batch(sql_templates,
    verbose=True, num_generations=3,
    override=False)
text_template_generator.save(file_path=
    file_path, override=True)
```

Below are examples of generated text templates saved in JSON format.

```
"SELECT StartDate FROM
Project WHERE Name = '[Project
.Name]';": [
    "Start date for project '
[Project.Name]'",
    "Look up start date of '[
Project.Name]'",
    "Get '[Project.Name]'
start date"
],
...
```

# Step 3: Generate Evaluation Data (Text Queries and Answers) save\_file = f"{linguistic\_attr}.json"

<pre>qa_generator = QADataGenerator(db_tool)</pre>	11
all_answers_to_text_queries =	11
qa_generator.generate(	11
<pre>sql_to_text_templates)</pre>	11
qa_generator.save(	11
all_answers_to_text_queries,	11
<pre>database_name, save_file, overwrite=</pre>	11
True)	11

12 13

14

15 16

17

18

19

1120

1121

1122

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

Below are examples of generated query-answer pairs saved in JSON format.

Γ

٦

```
1123
 Ε
                                            1124
      "[('Maldives',)]",
                                            1125
      Γ
                                            1126
           "Get 'Blue Horizon
                                            1127
Hotels' location details",
                                            1128
           "Find which location
                                            1129
'Blue Horizon Hotels' is
                                            1130
located in",
                                            1131
           "Determine the
                                            1132
location of 'Blue Horizon
                                            1133
Hotels'"
                                            1134
      ]
                                            1135
 ],
                                            1136
                                            1137
                                            1138
```

### H.2 Reproducing Experiment Results

The tool aslo provides two high-level Python objects to reproduce results from our experiments: TaggedGroup to perform tagging on each semantic groups and generate metrics for modular evaluation. The code below demonstrates an example to reproduce the results in Table 3. The input to get\_eval\_results and TaggedGroup is a list of RAGResult objects, which is a data class including a query, the ground-truth answer and the RAG response.

```
import json
                                                    1151
from grammar.eval.result import
                                                    1152
                                                    1153
   RAGResult
from grammar.eval.tag_group import
                                                    1154
   TaggedGroup
                                                    1155
from grammar.eval.match import
                                                    1156
                                                    1157
   SemanticsMatch
                                                    1158
def get_eval_results(eval_results,
                                                    1159
   linguistic_attr, root_dir, file_path
                                                    1160
                                                    1161
   ):
    tagged_group = TaggedGroup(
                                                    1162
    eval_results)
                                                    1163
    semnatics_match = SemanticsMatch.
                                                    1164
    from_file(root_dir=root_dir,
                                                    1165
                                                    1166
   verbalize_attrs=linguistic_attr)
                                                    1167
    for eval_result in eval_results:
                                                    1168
```

```
# sleep for 20 seconds after 9
   examples
      # if results.index(result) % 9
    == 0 and results.index(result) != 0:
       #
              print("Sleeping for 20
    seconds")
       #
              time.sleep(20)
        #
              print("Waking up")
        eval_result.
    judge_retrieval_response(
    tagged_group=tagged_group, method='
   use_exist')
        eval_result.judge_rag_response(
    semnatics_match)
    num_retrieval_failure = sum([result.
    retrieval_judgement==0 for result in
    eval_results])
    print(f"Retrieval failed in {
    num_retrieval_failure} out of {len(
    eval_results)} examples")
    num_rag_failure = sum([result.
    judgement=="Incorrect" for result in
    eval_results])
    print(f"RAG failed in {
    num_rag_failure} out of {len(
    eval_results)} examples")
    semnatics_match.save(root_dir=f'{
    root_dir}', override=True)
    # semnatics_match.llm.
    gpt_usage_record.write_usage(
   model_name='chatgptk' )
    # save results
    results = [result.asdict() for
    result in eval_results]
    # ensure json serializable
    for result in results:
        result['true_document_ids'] =
    list(result['true_document_ids'])
        result['retrieved_document_ids']
    = list(result['
    retrieved_document_ids'])
    with open(file_path, 'w') as f:
        json.dump(results, f, indent=4)
    return eval_results, tagged_group
root_dir = 'aurp'
closed_domain = True
results, metric = get_eval_results( '
   short', root_dir, file_path=f'{
    root_dir}/eval_results/
   results_short_balanced.json')
# re-produce metrics in Table 3
print('Baseline Accuracy: ', metric.
   get_accuracy())
print('Accuracy (Remove LLM Errors): ',
    metric.get_accuracy(for_retrieval=
   True))
print('Removing Gap Examples: ', metric.
   get_robustness())
print('Removing LLM Errors & Gap
   Examples: ', metric.get_robustness(
   for_retrieval=True))
```

1170

1171

1172

1173

1174

1175

1176 1177

1178

1179

1180

1181

1182

1183 1184

1185

1186

1187

1188

1189

1190

1191 1192

1193

1194

1195 1196

1197

1198

1199

1200

1201 1202

1203

1204

1205

1206

1207 1208

1209

1210

1211

1212

1213 1214 1215

1216

1217

1218

1219

1220 1221

1222

1223 1224

1225

1226

1227

1228 1229

1230

1231

1232

1233