Vendi-RAG: Adaptively Trading-Off Diversity And Quality Significantly Improves Retrieval Augmented Generation With LLMs

Anonymous ACL submission

Abstract

Retrieval-augmented generation (RAG) enhances large language models (LLMs) for domain-specific question-answering (QA) tasks by leveraging external knowledge sources. However, traditional RAG systems primarily focus on relevance-based retrieval and often struggle with redundancy, especially when reasoning requires connecting information from multiple sources. This paper introduces Vendi-RAG, a framework based on an iterative 011 process that jointly optimizes retrieval diversity and answer quality. This joint optimization leads to significantly higher accuracy for multi-hop QA tasks. Vendi-RAG leverages the Vendi Score (VS), a flexible similarity-based diversity metric, to promote semantic diversity in document retrieval. It then uses an LLM 017 judge that evaluates candidate answers, gener-019 ated after a reasoning step, and outputs a score that the retriever uses to balance relevance and diversity among the retrieved documents 021 during each iteration. Experiments on three challenging datasets-HotpotQA, MuSiQue, and 2WikiMultiHopQA-demonstrate Vendi-RAG's effectiveness in multi-hop reasoning The framework achieves significant tasks. accuracy improvements over traditional single-step or multi-step RAG approaches, with accuracy increases reaching +4.2% on HotpotQA, +4.1% on 2WikiMultiHopQA, and +1.3% on MuSiQue compared to Adaptive-RAG, the current best baseline. The benefits of Vendi-RAG are even more pronounced as the number of retrieved documents increases. Finally, we evaluated Vendi-RAG across different LLM backbones, including GPT-3.5, GPT-4, and GPT-4o-mini, and observed consistent improvements, demonstrating that the framework's advantages are model-agnostic.

1 Introduction

043

Retrieval-augmented generation (RAG) has emerged as a transformative framework for enhancing the performance of large language models (LLMs) in domain-specific tasks such as question-answering (QA). By retrieving relevant information from external sources beyond the training set, RAG enables LLMs to answer specialized queries more effectively (Achiam et al., 2023; Team et al., 2023; Jiang et al., 2024). This approach has been particularly successful in single-hop QA, where a question can be answered using information from a single document (Raiaan et al., 2024; Kwiatkowski et al., 2019). For instance, answering a question such as "Who wrote the novel Frankenstein?" only requires retrieving relevant information from a single document containing this fact. 044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

081

084

However, multi-hop QA introduces significantly greater complexity. Finding the correct answer to queries in multi-hop QA requires reasoning across multiple sources (Press et al., 2022; Tang and Yang, 2024). For instance, answering "Which city is the capital of the African country where Mount Kilimanjaro is located?" necessitates first identifying that Mount Kilimanjaro is in Tanzania, and then determining that Dodoma is the capital of Tanzania. This process involves not only retrieving information from multiple documents but also synthesizing these different sources effectively to form an accurate answer, which greatly increases the complexity of both retrieval and reasoning and leads to redundancy.

To address these challenges, iterative RAG pipelines have been developed. These pipelines refine the retrieval process through repeated modifications and re-querying of retrieved documents, aiming to resolve ambiguities and improve relevance. Notable examples include Adaptive-RAG (Lewis et al., 2020), which controls the number of iterations of the pipeline including the retrieval process and modifying the queries based on a classification model's assessment of the input query, Self-RAG(Asai et al., 2023), which incorporates iterative self-reasoning, and IROC (Trivedi et al., 2022),



Figure 1: The process begins with an initial retrieval step, where a diverse set of documents is retrieved using the Vendi Score, ensuring broad semantic coverage. Next, leveraging a reasoning step to construct a coherent path to the final answer, the LLM generates an answer, which then undergoes quality assessment by an LLM judge. Based on the answer quality, the retriever is adjusted to balance diversity and relevance: high-quality answers limit the emphasis on diversity, while low-quality answers prompt the retriever to prioritize diversity more heavily. This adjustment is controlled by an adaptive parameter, *s*, which is updated over iterations. The process continues until the answer quality reaches an optimal threshold, denoted by Thr. Finally, the highest-quality responses and documents are selected, ensuring both diversity and accuracy.

which progressively refines retrieval to optimize the final answer (Wei et al., 2022; Wang and Zhou, 2024).

Despite their success, iterative RAG methods typically rely solely on relevance-based retrieval, which focuses on the similarity between the query and dataset entries. This approach presents a fundamental limitation: it does not actively manage the diversity and quality of the retrieved information to properly address the query. More complex queries require diverse retrieval. We therefore propose a novel retrieval method called *Vendi retrieval* to address the limitation of existing retrieval pipelines. Vendi retrieval leverages the Vendi Score (VS) to enhance the diversity of retrieved documents while accounting for retrieval quality through a simple weighting mechanism.

Building on Vendi retrieval, we propose an iterative RAG pipeline called Vendi-RAG that balances diversity and quality. More specifically, the pipeline is as follows: an initial set of candidate documents is retrieved. Based on these retrieved documents, the system generates chain-of-thought (CoT) reasoning steps. Using these reasoning steps and retrieved documents, the LLM then generates candidate answers. An LLM-based evaluator then assesses these candidates for relevance, coherence, and completeness. The highest-scoring answer is selected as the final response. If the answer does not meet the quality threshold, the Vendi retrieval process dynamically adjusts the balance between diversity and relevance in document selection, ensuring broader semantic exploration or increased specificity as needed. This iterative refinement continues until a high-quality response is achieved. Figure 1 provides a detailed overview of the Vendi-RAG framework.

117

118

119

120

121

We evaluated the Vendi retrieval process and Vendi-RAG on three challenging multi-hop QA 123 datasets, HotpotQA (Yang et al., 2018), MuSiQue 124 (Trivedi et al., 2022), and 2WikiMultiHopQA (Ho 125 et al., 2020). To assess the Vendi retrieval method 126 we measured the diversity of retrieved documents 127 on the three datasets using two different diver-128 sity metrics, the VS and the max pairwise dis-129 tance (MPD). We found that the Vendi retrieval 130 process yields more diverse documents compared 131 to the baselines according to both metrics. Sec-132 ond, we evaluated Vendi-RAG in terms of sev-133 eral performance metrics, looking at both accu-134 racy and diversity. The results showed that Vendi-135 RAG substantially improves response accuracy, 136 outperforming existing RAG approaches. Us-137 ing GPT-3.5 as the LLM backbone, Vendi-RAG 138 demonstrated significant accuracy gains across all datasets, with accuracy increases reaching +4.2% 140 on HotpotQA, +4.1% on 2WikiMultiHopQA, and 141 +1.3% on MuSiQue compared to Adaptive-RAG, 142 the best baseline. Notably, the accuracy improve-143 ment remained consistent across different LLM 144 backbones-GPT-40, GPT-40-mini, and GPT-3.5-145 indicating that Vendi-RAG's advantages are model-146 Additionally, our experiments with agnostic. 147 varying numbers of retrieved documents-beyond 148

113

114

115

116

the standard two-document setting—showed that Vendi-RAG maintained its superior performance, especially as the number of retrieved documents increased. This underscores the critical role of the Vendi retrieval process in handling complex retrieval scenarios. For instance, when retrieving ten documents from HotpotQA, Vendi-RAG outperformed Adaptive-RAG by 7.8% in accuracy using GPT-40-mini as the backbone LLM.

This work introduces a diversity-guided retrieval approach that optimizes both diversity and quality to address the challenges of multi-step reasoning in multi-hop QA. Our experimental results highlight the effectiveness of Vendi-RAG in enhancing retrieval diversity and response accuracy, underscoring its potential as a robust solution for complex multi-hop QA tasks.

2 Related Work

149

150

151

152

153

154

155

156

157

158

160

161

162

163

164

165

166

168

169

170

171

172

173

174

176

178

179

180

181

182

183

184

185

189

There are three main approaches to QA: nonretrieval-based methods (Petroni et al., 2019), single-step RAG (Lewis et al., 2020), and multistep RAG (Asai et al., 2023). Non-retrieval-based QA methods pass queries directly to an LLM and use its generated output as the answer, without consulting external sources. While efficient, these methods struggle with queries requiring external or up-to-date information and suffer from hallucinations on out-of-distribution queries (Shuster et al., 2021). Single-step RAG methods integrate external knowledge retrieved from a knowledge base (e.g., Wikipedia). These methods improve factual accuracy but are limited by retrieval noise and perform poorly in complex reasoning tasks (Trivedi et al., 2022). Multi-step RAG methods are designed for complex multi-hop queries (Jeong et al., 2024; Asai et al., 2023; Tang and Yang, 2024). They iteratively retrieve documents and refine answers until they converge on a final response. This iterative refinement approach enables reasoning across multiple sources but introduces computational overhead and is prone to error accumulation (Jeong et al., 2024).

Advances in multi-hop QA. Recent improve-190 ments in multi-hop QA focus on question decom-191 position (Radhakrishnan et al., 2023), chain-of-192 thought reasoning (Wei et al., 2022; Liu et al., 193 194 2024a), and iterative retrieval (Jeong et al., 2024; Shao et al., 2023; Yu et al., 2024). Methods like 195 ReCite (Sun et al., 2022) and IRCoT (Trivedi et al., 196 2022) refine retrieval with progressive reasoning, while Self-RAG (Asai et al., 2023) adapts retrieval 198

strategies based on query complexity. Decomposed prompting (Khot et al., 2022) further enhances retrieval for complex queries (Zhang et al., 2024). MultiHop-RAG (Tang and Yang, 2024) integrates decomposition and retrieval pipelines but remains constrained by relevance-based retrieval, leading to redundancy and limited synthesis of diverse information.

199

200

201

202

203

204

205

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

224

225

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

Vendi scoring. The Vendi Score (VS) (Friedman and Dieng, 2023) is a similarity-based diversity metric applied in machine learning (Berns et al., 2023; Pasarkar and Dieng, 2023; Mousavi and Khalili, 2025; Nguyen and Dieng, 2024; Kannen et al., 2024; Jalali et al., 2024; Askari Hemmat et al., 2024; Rezaei and Dieng, 2025; Bhardwaj et al., 2025), chemistry (Pasarkar et al., 2023), materials science (Liu et al., 2024b), and biology (Pasarkar and Dieng, 2025). Vendi-RAG integrates VS into retrieval, balancing diversity and quality beyond conventional ranking systems (Carbonell and Goldstein, 1998; Slivkins et al., 2010). Unlike standard relevance-based retrieval (Guu et al., 2020), this approach enhances robustness and accuracy in multi-hop QA by incorporating semantic diversity into document retrieval.

3 Method

We now describe Vendi-RAG, including the novel retrieval process it uses.

3.1 Vendi Retrieval

Diversity in retrieved documents is essential for multi-hop QA, as it ensures broad semantic coverage, reduces redundancy, and incorporates multiple perspectives (Sun et al., 2022; Carbonell and Goldstein, 1998; Thakur et al., 2021). The most used methods for diverse retrieval are similarity search (SS) (Thakur et al., 2021) and maximal marginal relevance (MMR) (Carbonell and Goldstein, 1998). SS maximizes relevance to the query but retrieves highly similar documents, leading to redundancy. MMR balances relevance and novelty using pairwise comparisons but also struggles to account for global semantic diversity.

To overcome these limitations, we propose a novel retrieval method that leverages the VS (Friedman and Dieng, 2023) to explicitly optimize retrieval diversity. Let $\mathcal{D} = \{d_1, \ldots, d_n\}$ be a set of retrieved documents and $k(\cdot, \cdot)$ a positive semidefinite similarity kernel such $k(d_i, d_i) = 1$ for all *i*. Denote by K the corresponding similarity ma-

295

296

311 312 313

325

326

327

329

330

331

332

333

334

335

336

337

339

trix that is such that $K_{ij} = k(d_i, d_j)$. The VS is defined as

248

249

251

253

258

262

263

264

265

270

272

273

274

275

276

281

284

287

$$\operatorname{VS}_k(\mathcal{D}) = \exp\left(-\sum_{i=1}^n \lambda_i \log \lambda_i\right),$$
 (1)

where $\lambda_1, \ldots, \lambda_n$ are the eigenvalues of the normalized kernel matrix K/n. As argued by Friedman and Dieng (2023), the VS is the effective number of unique documents in D, reaching its maximum value n when all the documents are distinct and its minimal value 1 when all the documents are the same.

While accounting for diversity is good for retrieval, especially for complex queries, it shouldn't be the only criterion. Quality also matters. To balance these two criteria, the Vendi retrieval process uses a convex combination of the two,

$$VRS = s \cdot VS_k(q, \mathcal{D}) + (1 - s) \cdot SS(q, \mathcal{D}), \quad (2)$$

where VRS stands for *Vendi retrieval score* and $s \in [0, 1]$ is a tunable parameter controlling the trade-off between diversity and similarity. When handling complex queries, such as those with multiple possible answers, a higher diversity weight *s* promotes the selection of a semantically diverse set of documents. In contrast, for simpler or more specific queries that require precise information, a smaller value of *s* prioritizes similarity-based retrieval.

The similarity score SS(q, D) is computed using dense vector representations of both the query and the documents. The document representations are used to provide a meaningful comparison between queries and candidate documents in a high-dimensional semantic space, ensuring that retrieval is based on conceptual similarity. Incorporating contextual understanding through transformerbased embeddings ensures semantic matching beyond simple lexical overlap.

3.2 Vendi-RAG

We integrate the Vendi retrieval process into a flexible RAG pipeline that balances diversity and relevance for improved performance on multi-hop QA.

2881. Initial retrieval. The process begins by re-
trieving a set of documents using Vendi retrieval.290This first step prioritizes broad semantic coverage
(we set s = 0.8 initially in all our experiments),
ensuring that the retrieved documents capture mul-
tiple perspectives and to prevent recovering seman-
tically redundant documents. This initial diversity

is particularly critical for multi-hop QA, where synthesizing information from varied sources is essential to accurately answering the query.

2. Reasoning generation. Based on the retrieved documents, the system generates CoT reasoning steps. These intermediate reasoning steps help contextualize the retrieved information, building a coherent pathway to the final answer.

3. Candidate answer generation. Using the reasoning steps and retrieved documents, the LLM generates candidate answers. These proposed answers are evaluated to determine their quality and completeness.

4. Quality evaluation. An LLM judge assesses the candidate answers. This evaluation considers factors such as coherence, relevance, and alignment with the query. A quality score Q_t is produced at the end of this quality-check. Here t is used to indicate the iteration step.

5. Dynamic adjustment of the VRS. Based on the quality score Q_t , the parameter s is adjusted dynamically. We denote by s_t the value of the parameter s at the t^{th} iteration. It controls the trade-off between diversity (via VS) and relevance (via similarity search). If Q_t is low, s_t should be increased, to prioritize greater diversity in the subsequent retrievals. This ensures broader semantic exploration, which is beneficial for refining answers in cases where the retrieved information is already relevant but lacks coverage. Conversely, if Q_t is high, s_t should be decreased to focus more on relevance, retrieving documents that are closely aligned with the query to address potential gaps in specificity. We therefore define s_t as

$$s_t = f(Q_{t-1}) = 1 - \frac{Q_{t-1}}{\max(Q_{t-1})},$$
 (3)

where f is a simple linear function that maps Q_{t-1} to the interval [0, 1], ensuring that higher quality scores correspond to lower diversity scores.

6. Iterative refinement. The retrieval and reasoning steps are repeated iteratively, with adjustments to *s* dynamically balancing diversity and relevance at each stage. This process continues until the desired answer quality is reached, ensuring that the system converges on an optimal set of documents and reasoning steps.

340

34

347

3

333

365 366 367

363

3

371

372

373 374

376 377

378

3

3

7. Final answer selection. Once the iterative refinement process is complete, the final set of documents and answers are selected based on their quality scores. This ensures that the output reflects both broad semantic coverage and high-quality, relevant information. Algorithm 1 summarizes the procedure.

Why Adjusting *s* Matters: The dynamic adjustment of *s* is critical for striking the right balance between diversity and relevance during the retrieval process. High diversity is essential for exploring various facets of a complex query, especially in multi-hop QA, where information from disparate sources must be synthesized. However, excessive diversity can dilute the relevance of retrieved documents, potentially introducing noise and reducing the quality of generated answers. On the other hand, overemphasizing relevance can lead to redundancy and failure to capture the breadth of information needed for comprehensive reasoning.

By reducing *s* when the quality score is high, the Vendi-RAG pipeline encourages exploration of less-redundant, semantically diverse documents. This ensures that even if the current answer is sufficient, the model explores additional perspectives that may enhance the depth and breadth of the final response. Conversely, increasing *s* when quality is low allows the system to focus on retrieving documents that are more closely aligned with the query, addressing gaps in specificity or relevance.

The strength of Vendi-RAG lies in this adaptive approach to document retrieval. Unlike traditional RAG systems that use fixed retrieval strategies, Vendi-RAG's dynamic adjustment of the diversityrelevance trade-off (the parameter *s*) allows it to respond to the specific requirements of each query and intermediate reasoning step. When the system detects that current retrievals are yielding highquality but potentially narrow responses, it automatically shifts toward greater diversity, exploring complementary perspectives. Conversely, when responses lack precision, the system can focus on more closely related documents to improve specificity.

Performance characteristics. In practice, Vendi RAG exhibits distinctive performance patterns that
 reflect its sophisticated design. The system natu rally adapts its computational effort to query com plexity, requiring more iterations for intricate multi hop queries while converging quickly for simpler
 ones. Though the computational overhead exceeds

Algorithm 1 Vendi-RAG Inference Pipeline

Require: Query q , Knowledge base \mathcal{D} , Max iterations N , Ouality threshold τ
Ensure: Final answer \hat{a}^*
Initialize context: $q_1 \leftarrow q$, set initial parameter: $s_1 \leftarrow 0.8$
for $i = 1$ to N do
$\operatorname{VRS}_i \leftarrow s_i \cdot \operatorname{VS}_k(q_i, \mathcal{D}) + (1 - s_i) \cdot \operatorname{SS}(q_i, \mathcal{D})$
$D_i \leftarrow \text{Vendi-Retrieval}(q_i, \text{scores}_i; \mathcal{D})$
Generate reasoning steps: $r_i \leftarrow CoT(q_i, D_i)$
Produce answer: $\hat{a}_i \leftarrow \text{LLM}(q, D_i, r_{1:i})$
if LLM-Judge $(\hat{a}_i) \geq \tau$ then
return \hat{a}_i
end if
Update query: $q_{i+1} \leftarrow \text{RewriteQuery}(q_i, \hat{a}_i, r_i)$
Update weight parameter: $s_{i+1} \leftarrow f(Q_i)$
end for
return \hat{a}_N

that of basic RAG systems, the improved retrieval quality often results in better final answers. The system maintains reasonable scalability with document corpus size, as the primary computational bottleneck—eigenvalue computation—depends on the number of retrieved documents rather than the total corpus size. These characteristics make Vendi-RAG particularly suitable for complex tasks such as multi-hop QA.

391

392

393

394

395

396

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

494

425

4 Experiments

This section presents a comprehensive evaluation of Vendi-RAG on multi-hop QA datasets. First, we investigate the effectiveness of the Vendi retrieval process in enhancing retrieval diversity. Next, we evaluate the Vendi-RAG pipeline, demonstrating its ability to handle complex queries requiring multistep reasoning compared to the baselines.

Datasets. Our experiments are conducted on three challenging benchmark multi-hop QA datasets: MuSiQue (Trivedi et al., 2022), HotpotQA (Yang et al., 2018), and 2WikiMultiHopQA (Ho et al., 2020).

MuSiQue evaluates a model's ability to synthesize facts spread across across multiple document sources. It includes questions spanning diverse domains such as history, science, and culture, requiring logical reasoning and synthesis of interdependent information. Given its emphasis on multi-step comprehension, this dataset challenges models to accurately identify and integrate relevant information to generate correct answers to queries. This is the most challenging dataset among the three.

HotpotQA assesses reasoning and fact verification across various domains, including geography, entertainment, and history. Its questions

Table 1: Retrieval diversity as measured by the Vendi Score (VS) and Max Pairwise Distance (MPD) for Vendi-RAG, Adaptive-RAG, and Adaptive Retrieval across different datasets. Vendi-RAG achieves higher diversity than the baselines across all datasets and both metrics.

Dataset	Adapti	ve Retrieval	Adapt	ive-RAG	Vendi-RAG		
Dutuset	VS MPD		VS	MPD	VS	MPD	
MuSiQue	6.13	1.25	6.55	1.42	7.12	1.95	
HotpotQA	4.95	1.10	5.21	1.31	6.82	1.89	
2WikiMultiHopQA	5.34	1.32	5.81	1.45	6.68	1.78	

necessitate reasoning over two or more interconnected documents linked via hyperlinks. Additionally, the dataset includes "comparison" questions that require juxtaposing information from multiple sources, making it a challenging benchmark for evaluating both retrieval quality and reasoning ability.

426

427

428

429

430

431

432

433

434 435

436

437

438

439

440

441

2WikiMultiHopQA leverages Wikipedia's complex structure to pose complex reasoning challenges. Questions are derived from real-world knowledge graphs and require navigating reasoning paths across multiple documents. Topics span science, politics, and sports, emphasizing logical relationships such as cause-effect dependencies, making it an essential tool for evaluating structured knowledge reasoning.

442 Vendi retrieval improves document retrieval diversity. To assess the impact of the Vendi re-443 trieval process on retrieval diversity, we compared 444 the diversity of outputs from Vendi-RAG against 445 Adaptive-RAG and Adaptive Retrieval. We mea-446 sured diversity using two different metrics, the 447 VS and the max pairwise distance (MPD). Ta-448 ble 1 summarizes the results. Vendi-RAG achieves 449 higher diversity compared to Adaptive Retrieval 450 and Adaptive-RAG on all dataset, demonstrating 451 its ability to retrieve documents that capture mul-452 tiple perspectives relevant to the query. This is a 453 crucial advancement, as increased diversity in re-454 trieval directly correlates with improved robustness 455 in multi-hop reasoning (see Table 2). Adaptive-456 RAG, which incorporates iterative refinement but 457 lacks explicit diversity control, shows moderate 458 retrieval diversity improvement over Adaptive Re-459 460 trieval.

Accuracy on multi-hop QA tasks. We further evaluated the performance of the Vendi-RAG
pipeline, to assess its ability to reason across multiple documents. The results in Table 2 indicate that

Vendi-RAG consistently outperforms other methods in response accuracy across all datasets, showcasing the efficacy of balancing retrieval diversity with quality. Additionally, Vendi-RAG achieves competitive performance in Exact Match (EM) and F1 score. These findings highlight Vendi-RAG's capability to enhance answer correctness for complex reasoning tasks through improved document retrieval. 465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

Impact of the number of retrieved documents **on performance.** To further examine the impact of document size on retrieval effectiveness, we analyze the performance of Vendi-RAG and Adaptive-RAG across varying document sizes. Figure 2 illustrates the relationship between document size and performance on the HotPotQA dataset. Vendi-RAG consistently outperforms Adaptive-RAG in accuracy for document sizes greater than two. As document size increases, accuracy improves for both methods, but the gain is notably higher for Vendi-RAG. Similar to accuracy, EM and F1 scores exhibit an increasing trend as document size grows. Vendi-RAG shows a more pronounced improvement, underscoring its capacity to retrieve more informative and relevant documents, thereby enhancing answer quality.

The VS also increases with document size. This is evidence that Vendi-RAG alleviates redundancy since the VS is known to be invariant under duplication (Friedman and Dieng, 2023). An increasing VS indicates less redundancy in the retrieved documents. By leveraging the VS in its retrieval process, Vendi-RAG avoids the redundancy issue that often plagues RAG pipelines.

These results indicate that increasing document size enhances both retrieval diversity and answer correctness. However, the degree of improvement varies across methods, with Vendi-RAG achieving superior gains in all metrics. However, we are computationally bottlenecked primarily by the LLM's

Table 2: Performance on multi-hop QA datasets using GPT-3.5 (Turbo). Here we use three different flavors of accuracy: exact match (EM), F1-score (F1), and traditional accuracy (Acc). Vendi-RAG outperforms the baselines in all 3 datasets, except in terms of F1-score, where it performs similarly to Adaptive-RAG.

			MuSiQue		HotpotQA		2WikiMultiHopQA				
Steps	Types	Methods	EM	F1	Acc	EM	F1	Acc	EM	F1	Acc
Single-step	Simple	No Retrieval Single-step Approach	20.40 16.40	31.30 26.70	24.40 23.60	37.40 39.60	51.04 50.44	43.20 45.60	37.00 46.80	48.50 57.69	43.40 52.60
Multi-step	Adaptive	Adaptive Retrieval Adaptive-RAG Vendi-RAG	18.80 21.80 24.4	30.30 32.60 32.52	24.80 29.60 30.4	38.60 40.40 42.2	50.70 52.56 57.04	43.20 47.00 58.4	44.20 46.60 47.2	55.11 60.09 58.9	50.60 56.80 61.4



Figure 2: Performance comparison of Vendi-RAG and Adaptive-RAG across different document sizes in terms of Exact Match, F1-score, Accuracy, and Vendi Score on HotPotQA. The backbone LLM here is GPT-4o-mini. Vendi-RAG consistently outperforms Adaptvie-RAG on all the metrics. In particular, performance improves as the number of documents increases.

context window limitation and processing time. As the number of retrieved documents increases, we must either truncate documents to fit within the model's maximum context length or process documents in multiple batches, both of which have significant computational overhead. The bottleneck occurs specifically in the final stage of the pipeline where the LLM processes the retrieved documents to generate answers.

505

506

507

508

510

511

512

513

Performance for different LLM Backbones and 514 retrieval strategies. To evaluate the impact of 515 different LLM backbones and retrieval strategies on 516 the performance of the Vendi-RAG framework, we 517 conducted experiments using three LLMs: GPT-40, GPT-4o-mini, and GPT-3.5, across all the multi-519 hop QA datasets described above. The results, shown in Figure 3, highlight the effectiveness of 521 Vendi-RAG compared to Adaptive-RAG, the best 523 baseline, across all datasets and backbones, except for F1-score on the 2WikiMultiHopQA dataset. In 524 general, larger LLM backbones, such as GPT-40, 525 achieve superior performance across all datasets, particularly for tasks requiring complex reasoning 527

and synthesis across multiple documents. However, even with smaller models like GPT-40-mini, the Vendi-RAG model maintains competitive performance, demonstrating its effectiveness.

528

529

530

531

532

5 Conclusion

While retrieval-augmented generation (RAG) has 533 proven effective in enhancing large language model 534 (LLM) performance for domain-specific question-535 answering (QA) tasks, traditional RAG frameworks 536 often struggle with redundancy, particularly in 537 multi-hop reasoning tasks. To address this short-538 coming, we introduce Vendi-RAG, a novel frame-539 work that jointly optimizes retrieval diversity and 540 answer quality through an iterative refinement pro-541 cess. Vendi-RAG leverages the Vendi Score and 542 an LLM judge to promote semantic diversity while 543 maintaining relevance during retrieval. Our experi-544 ments on HotpotQA, MuSiQue, and 2WikiMulti-545 HopQA demonstrate Vendi-RAG's effectiveness. 546 Specifically, Vendi-RAG outperforms the best base-547 line by +4.2% on HotpotQA, +4.1% on 2Wiki-548 MultiHopQA, and +1.3% on MuSiQue. These 549



Figure 3: Performance comparison of Vendi-RAG and Adaptive-RAG variants across the three datasets using three evaluation metrics: F1-score, Exact Match, and Accuracy. Results show that Vendi-RAG-40 consistently outperforms other variants across all datasets and metrics, with a particularly strong performance on HotpotQA.

gains become even more pronounced as the number of retrieved documents increases, highlighting the importance of retrieval diversity in complex reasoning tasks. Furthermore, we evaluated Vendi-RAG across multiple LLM backbones, including GPT-3.5, GPT-4, and GPT-4o-mini, and observed consistent performance improvements, demonstrating that the framework is model-agnostic. These findings establish Vendi-RAG as an effective and adaptable solution for multi-hop QA.

6 Limitations

Vendi-RAG introduces computational overhead due to LLM-based quality scoring, which may limit scalability in real-time applications. Additionally, like all RAG approaches, its performance depends

on the quality and completeness of external knowledge sources, making it susceptible to biases or gaps in the retrieved information.

565

566

567

568

569

570

571

572

573

574

575

576

Ethics Statement 7

The deployment of LLMs, including their use in Vendi-RAG, necessitates careful ethical consideration. Since the model relies on external knowledge sources, concerns arise regarding the credibility and accuracy of retrieved content. Ensuring the reliability and factual integrity of information is crucial to mitigating risks related to bias and misinformation.

551

553

555

557

560

561

564

References

580

585

586

588

596

598

610

611

613

614

615

617

618

619

621

622

624

625

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*.
- Reyhane Askari Hemmat, Melissa Hall, Alicia Sun, Candace Ross, Michal Drozdzal, and Adriana Romero-Soriano. 2024. Improving geo-diversity of generated images with contextualized vendi score guidance. In *European Conference on Computer Vision*, pages 213–229. Springer.
 - Sebastian Berns, Simon Colton, and Christian Guckelsberger. 2023. Towards Mode Balancing of Generative Models via Diversity Weights. *arXiv preprint*. ArXiv:2304.11961 [cs.LG].
 - Utkarsh Bhardwaj, Vinayak Mishra, Suman Mondal, and Manoj Warrier. 2025. A robust machine learned interatomic potential for nb: Collision cascade simulations with accurate defect configurations. *arXiv preprint arXiv*:2502.03126.
- Jaime Carbonell and Jade Goldstein. 1998. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In *Proceedings* of the 21st annual international ACM SIGIR conference on Research and development in information retrieval, pages 335–336.
- Dan Friedman and Adji Bousso Dieng. 2023. The Vendi Score: A Diversity Evaluation Metric for Machine Learning. *Transactions on Machine Learning Research*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*.
- Mohammad Jalali, Azim Ospanov, Amin Gohari, and Farzan Farnia. 2024. Conditional vendi score: An information-theoretic approach to diversity evaluation of prompt-based generative models. *arXiv preprint arXiv:2411.02817*.
- Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong C Park. 2024. Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity. *arXiv preprint arXiv:2403.14403*.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*. 630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

- Nithish Kannen, Arif Ahmad, Marco Andreetto, Vinodkumar Prabhakaran, Utsav Prabhu, Adji Bousso Dieng, Pushpak Bhattacharyya, and Shachi Dave. 2024. Beyond aesthetics: Cultural competence in text-toimage models. *arXiv preprint arXiv:2407.06863*.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453– 466.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Jingyu Liu, Jiaen Lin, and Yong Liu. 2024a. How much can rag help the reasoning of llm? *arXiv preprint arXiv:2410.02338*.
- Tsung-Wei Liu, Quan Nguyen, Adji Bousso Dieng, and Diego A Gómez-Gualdrón. 2024b. Diversitydriven, efficient exploration of a mof design space to optimize mof properties. *Chemical Science*, 15(45):18903–18919.
- Mohsen Mousavi and Nasser Khalili. 2025. Vsi: An interpretable bayesian feature ranking method based on vendi score. *Available at SSRN 4924208*.
- Quan Nguyen and Adji Bousso Dieng. 2024. Quality-Weighted Vendi Scores And Their Application To Diverse Experimental Design. In *International Conference on Machine Learning*.
- Amey P Pasarkar, Gianluca M Bencomo, Simon Olsson, and Adji Bousso Dieng. 2023. Vendi sampling for molecular simulations: Diversity as a force for faster convergence and better exploration. *The Journal of chemical physics*, 159(14).
- Amey P Pasarkar and Adji Bousso Dieng. 2023. Cousins of the vendi score: A family of similaritybased diversity metrics for science and machine learning. *arXiv preprint arXiv:2310.12952*.
- Amey P. Pasarkar and Adji Bousso Dieng. 2025. The vendiscope: An algorithmic microscope for data collections. *arXiv preprint arXiv:2502.04593*.

- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2019. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*.
- Ansh Radhakrishnan, Karina Nguyen, Anna Chen, Carol Chen, Carson Denison, Danny Hernandez, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamilė Lukošiūtė, et al. 2023. Question decomposition improves the faithfulness of model-generated reasoning. *arXiv preprint arXiv:2307.11768*.
- Mohaimenul Azam Khan Raiaan, Md Saddam Hossain Mukta, Kaniz Fatema, Nur Mohammad Fahad, Sadman Sakib, Most Marufatul Jannat Mim, Jubaer Ahmad, Mohammed Eunus Ali, and Sami Azam. 2024. A review on large language models: Architectures, applications, taxonomies, open issues and challenges. *IEEE Access*.

701

704

710

712

713

714

715

716

718

719

720 721

722

725

727

728

730

731

732

733

734 735

736

737

738

- Mohammad Reza Rezaei and Adji Bousso Dieng. 2025. The *alpha*-alternator: Dynamic adaptation to varying noise levels in sequences using the vendi score for improved robustness and performance. *arXiv preprint arXiv:2502.04593*.
- Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy. *arXiv preprint arXiv*:2305.15294.
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation. *arXiv preprint arXiv:2104.07567*.
- Aleksandrs Slivkins, Filip Radlinski, and Sreenivas Gollapudi. 2010. Learning optimally diverse rankings over large document collections. In *Proc. of the* 27th International Conference on Machine Learning (ICML 2010).
- Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. 2022. Recitation-augmented language models. *arXiv preprint arXiv:2210.01296*.
- Yixuan Tang and Yi Yang. 2024. Multihop-rag: Benchmarking retrieval-augmented generation for multihop queries. *arXiv preprint arXiv:2401.15391*.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. Beir:

A heterogenous benchmark for zero-shot evaluation of information retrieval models. *arXiv preprint arXiv:2104.08663*.

- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509*.
- Xuezhi Wang and Denny Zhou. 2024. Chain-ofthought reasoning without prompting. *arXiv preprint arXiv:2402.10200*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Chunliang Yang, Rosalind Potts, and David R Shanks. 2018. Enhancing learning and retrieval of new information: a review of the forward testing effect. *NPJ science of learning*, 3(1):8.
- Tian Yu, Shaolei Zhang, and Yang Feng. 2024. Auto-rag: Autonomous retrieval-augmented generation for large language models. *arXiv preprint arXiv:2411.19443*.
- Zhihao Zhang, Alan Zhu, Lijie Yang, Yihua Xu, Lanting Li, Phitchaya Mangpo Phothilimthana, and Zhihao Jia. 2024. Accelerating retrieval-augmented language model serving with speculation. *arXiv preprint arXiv:2401.14021*.

769 770

- 771
- 773
- 775

779

- 781

- 784

790

- 793 794

796

797

Evaluation Metrics Α

To compare model performance across different datasets, we employ the following key evaluation metrics:

> • Exact Match (EM): This metric calculates the percentage of predictions that exactly match the ground truth answers. It is defined as:

$$EM = \frac{\text{Number of exact matches}}{\text{Total number of queries}} \times 100$$
(4)

EM is a strict metric that grants credit only when the predicted answer matches the ground truth exactly in both content and format. It is particularly useful for assessing a model's precision in generating accurate responses.

• F1 Score (F1): The F1 score captures the harmonic mean of precision and recall at the token level, providing a balanced measure of correctness. It is defined as:

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$
(5)

where precision is the fraction of retrieved tokens that are relevant, and recall is the fraction of relevant tokens that are retrieved. The F1 score is particularly relevant for multi-hop QA tasks, where partial correctness (e.g., retrieving some but not all supporting evidence) is informative.

• Accuracy (Acc): Accuracy measures the proportion of correct predictions over all evaluated queries. It is defined as:

$$Acc = \frac{\text{Number of correct predictions}}{\text{Total number of queries}} \times 100$$
(6)

Unlike EM, which requires exact matches, accuracy provides a broader assessment by capturing overall correctness, including responses that convey the intended meaning.

• Max Pairwise Distance (MPD): This metric evaluates the maximum Euclidean distance between pairs of retrieved data points, measuring diversity. It is defined as:

$$MPD = \max_{i,j} \|x_i - x_j\|_2, \quad i \neq j \quad (7)$$

where x_i and x_j represent document embed-808 dings in the feature space. Higher values in-809 dicate greater diversity among retrieved docu-810 ments. 811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

Each of these metrics offers a unique perspective on model performance. EM is a stringent measure of precision, F1 balances precision and recall, and accuracy provides an overall correctness measure. Meanwhile, MPD and diversity-based metrics assess the variety and independence of retrieved documents-critical for multi-hop QA tasks requiring integration of diverse information.

B **Implementation Details for Dataset Ingestion and Vector Database**

Preparing datasets for question-answering requires transforming data into a searchable vector database to enable efficient retrieval. This workflow includes document chunking and semantic embedding to optimize performance.

B.1 Dataset Processing and Chunking

The dataset, provided in JSON format with context paragraphs and metadata, is processed by splitting each document into smaller chunks. Each chunk has a maximum size of 512 tokens, with a 50-token overlap to preserve context across chunk boundaries and facilitate multi-hop reasoning in long documents.

B.2 Embedding Model and Vector Database

We use the SentenceTransformer model, specifically all-mpnet-base-v2, to generate dense vector representations for documents and queries. These embeddings are stored locally to avoid redundant downloads and improve reusability. The Chroma vector database efficiently stores and retrieves these vectorized documents along with metadata, such as document titles and chunk IDs.

B.3 Batch Processing and Database Population

To efficiently populate the vector database, document chunks are processed in batches of up to 10,000. This approach optimizes memory usage while ensuring completeness in the ingestion process. The total number of processed chunks is logged for verification.

854 855

856

857

858

859

860

861

862

863 864

865

866

867

868

869

B.4 Query Answering Workflow

For queries such as "Who is the father-in-law of Queen Hyojeong?", relevant chunks are retrieved using **Chroma**'s similarity-based search mechanism. The system ranks the top 10 chunks based on their semantic similarity to the query, leveraging embeddings generated by all-mpnet-base-v2 to ensure precise and relevant results.

B.5 Key Configuration Details

The system is configured with the following parameters:

- **Embedding Model:** all-mpnet-base-v2, optimized for semantic similarity tasks.
- Vector Database: Chroma, persisted to disk for efficient reuse.
- **Chunk Size:** 512 tokens per chunk, with a 50-token overlap.
- Batch Size: Up to 10,000 chunks per batch.