

Entity Retrieval for Answering Entity-Centric Questions

Anonymous ACL submission

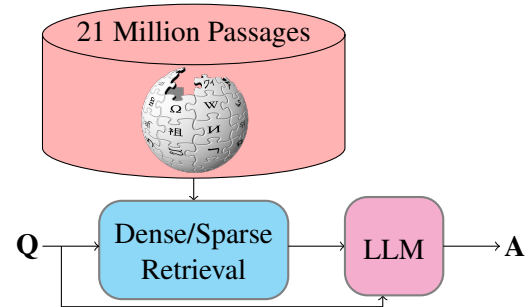
Abstract

The similarity between the question and indexed documents is a crucial factor in document retrieval for retrieval-augmented question answering. Although this is typically the only method for obtaining the relevant documents, it is not the sole approach when dealing with entity-centric questions. In this study, we propose *Entity Retrieval*, a novel retrieval method which rather than relying on question-document similarity, depends on the salient entities within the question to identify the retrieval documents. We conduct an in-depth analysis of the performance of both dense and sparse retrieval methods in comparison to *Entity Retrieval*. Our findings reveal that our method not only leads to more accurate answers to entity-centric questions but also operates more efficiently.

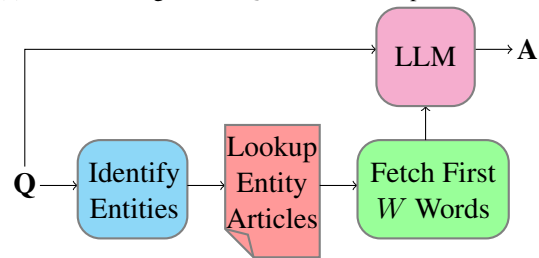
* We have included our source code implementation of the project, along with the generated model answers, in the Software section of our submission.

1 Introduction

Information retrieval has significantly enhanced the factual reliability of large language model (LLM) generated responses (Shuster et al., 2021) in question answering (Zhu et al., 2021; Zhang et al., 2023). This improvement is particularly notable in a research area known as retrieval-augmented generation (RAG; Lewis et al., 2020b; Izacard and Grave, 2021; Singh et al., 2021). RAG systems typically employ the *Retriever-Reader* architecture (Chen et al., 2017), with retrievers being either sparse (Peng et al., 2023), dense (Karpukhin et al., 2020), or a hybrid (Glass et al., 2022). The reader, which is a generative language model (e.g., BART; Lewis et al., 2020a, T5; Raffel et al., 2020, GPT-3; Brown et al., 2020), conditions its generated answers on the documents deemed relevant by the retriever. Recent RAG methodologies exploit the in-context learning capabilities of LLMs to incor-



(a) Retrieval-Augmented QA with Dense/Sparse Retrieval



(b) Retrieval-Augmented QA with *Entity Retrieval*

Figure 1: *Entity Retrieval* simplifies the process of obtaining augmentation documents by replacing the need to search through large indexed passages with a straightforward lookup. For **Q**: What is the capital of Seine-Saint-Denis? *Entity Retrieval* considers the first few sentences of Seine-Saint-Denis Wikipedia article which states “Its prefecture is **Bobigny**.” and returns **A = Bobigny** where the other retrieval methods return **A = Saint-Denis** or **A = Paris**.

porate the retrieved documents into the prompt (Shi et al., 2023; Peng et al., 2023; Yu et al., 2023).

Kandpal et al. (2023) demonstrate that retrieval-augmentation improves LLMs’ performance in answering entity-centric questions that seek factual information about real-world entities¹. They show that this technique is particularly helpful for questions about rare entities, which appear infrequently in LLM training and fine-tuning data.

¹Entity-centric questions typically have answers that are concise single words or short phrases. These answers often reference or directly stem from a knowledge base entity (Ranjan and Balabantaray, 2016).

050 But is there a correlation between the quality of
051 the retrieved documents and the generated response
052 quality? Sciavolino et al. (2021) demonstrate that
053 dense retrievers retrieve less relevant documents
054 for answering entity-centric questions than sim-
055 pler sparse retrievers. Additionally, Cuconasu et al.
056 (2024) show that the presence of irrelevant docu-
057 ments leads to worse answers.

058 In this paper, we propose *Entity Retrieval* (Figure
059 1b), which uses salient entities in the question to
060 lookup knowledge base (e.g., Wikipedia) articles
061 that correspond to each entity, and uses the first W
062 words of the articles as augmentation documents
063 for the question passed to the LLM.

064 Our contributions are as follows: (1) we pro-
065 pose *Entity Retrieval*, a novel method of acquiring
066 augmentation documents using salient entities in
067 the questions, (2) we compare the retrieval per-
068 formance quality of several retrieval techniques
069 (both dense and sparse) to *Entity Retrieval* for ques-
070 tions within two entity-centric question answering
071 datasets, (3) we study the retrieval-augmentation
072 quality of the compared techniques and *Entity Re-*
073 *trieval*, using salient entity annotations of the ques-
074 tions, and (4) we examine the application a recent
075 state-of-the-art entity linking method for *Entity Re-*
076 *trieval* in the absence of entity annotations in entity-
077 centric questions.

078 2 Retrieval for Retrieval-Augmentation

079 Retrieval-augmentation (Lewis et al., 2020b) is a
080 method of converting Closed-book question an-
081 swering² (Roberts et al., 2020) into extractive ques-
082 tion answering (Abney et al., 2000; Rajpurkar et al.,
083 2016), where the answers can be directly extracted
084 from the retrieved documents. Despite the abun-
085 dance of effective retrieval techniques for retrieval-
086 augmented question answering in existing literature
087 (Zhan et al., 2020a,b; Yamada et al., 2021; Izcard
088 et al., 2022; Santhanam et al., 2022; Ni et al., 2022,
089 *inter alia.*), this section will concentrate on a select
090 few methods³ utilized to study answering entity-
091 centric questions in this paper.

092 **BM25** (Robertson et al., 1994, 2009) is a prob-
093 abilistic retrieval method that ranks documents
094 based on the frequency of query terms appearing in
095 each document, adjusted by the length of the docu-

²Closed-book QA focuses on answering questions without additional context during inference.

³We selected the methods supported by `pyserini.io` for the similarity between the underlying modules, minimizing discrepancies across different implementations.

096 ment and overall term frequency in the collection.
097 It operates in the sparse vector space, relying on
098 precomputed term frequencies and inverse docu-
099 ment frequencies to retrieve documents based on
100 keyword matching.

101 **DPR** (Dense Passage Retrieval; Karpukhin et al.,
102 2020) leverages a bi-encoder architecture, wherein
103 the initial encoder processes the question and the
104 subsequent encoder handles the passages to be re-
105 trieved. The similarity scores between the two
106 encoded representations are computed using a dot
107 product. Typically, the encoded representations
108 of the second encoder are fixed and indexed in
109 FAISS (Johnson et al., 2019; Douze et al., 2024),
110 while the first encoder is optimized to maximize the
111 dot-product scores based on positive and negative
112 examples.

113 **ANCE** (Xiong et al., 2021) is another dense re-
114 trieval technique similar to DPR. It employs an en-
115 coder to transform both the questions and passages
116 into dense representations. These representations
117 are compared using dot product similarity. The key
118 distinction from DPR is that ANCE uses hard neg-
119 atives generated by periodically updating the pas-
120 sage embeddings during training, which helps the
121 model learn more discriminative features, thereby
122 enhancing retrieval performance over time.

123 3 Entity Retrieval for Question 124 Answering

125 While quite powerful, most retrieval-augmented
126 systems are notably time and resource-intensive,
127 necessitating the storage of extensive lookup in-
128 dices and the need to attend to all retrieved docu-
129 ments to generate a response (see Section 4.7).
130 This attribute renders such methods less desirable,
131 particularly given the drive to run LLMs locally
132 and on mobile phones (Alizadeh et al., 2023).

133 Entity recognition has been an integral com-
134 ponent of statistical question answering systems
135 (Aghaebrahimian and Jurčiček, 2016, *inter alia.*).
136 Additionally, the extensively studied field of
137 Knowledge Base Question Answering (Cui et al.,
138 2017, *inter alia.*) has underscored the significance
139 of entity information from knowledge bases in
140 question answering (Salnikov et al., 2023). A tra-
141 ditional neural question answering pipeline may
142 contain entity detection, entity linking, relation
143 prediction, and evidence integration (Mohammed
144 et al., 2018; Lukovnikov et al., 2019), where entity
145 detection can employ LSTM-based (Hochreiter and

Schmidhuber, 1997) or BERT-based (Devlin et al., 2019) encoders. Inspired by this body of work, we investigate the relevance of retrieval based on entity information as an alternative strategy to the proposed retrieval methods of Section 2, especially for answering entity-centric questions with LLMs.

Our proposed method *Entity Retrieval*, leverages the salient entities within the questions to identify and retrieve their corresponding knowledge base articles. We will then use the first W words of these articles as the documents augmenting entity-centric questions when prompting LLMs. Figure 1 presents a schematic comparison between *Entity Retrieval* and dense retrieval methods in identifying retrieval documents to enhance question answering with LLMs.

4 Experiments and Analysis

4.1 Setup

We focus on Wikipedia as the knowledge base and utilize the pre-existing BM25, DPR, and ANCE retrieval indexes in Pyserini (Lin et al., 2021). These indexes, follow established practices (Chen et al., 2017; Karpukhin et al., 2020) and segments the articles into non-overlapping text blocks of 100 words, resulting in 21,015,300 passages. For dense retrievers, the passages are processed with a pre-trained context encoder, generating fixed embedding vectors stored in a FAISS index (Douze et al., 2024). Our experimental entity-centric questions are encoded using the question encoder, and the top k relevant passages to the encoded question are retrieved from the FAISS index. For BM25 sparse retriever, the passages are stored in a Lucene index and the questions are keyword matched to this index.

As outlined in Section 3, the document retrieval process will require loading the entire index (as well as the question encoder for dense retrieval) into memory which entails significant time and memory consumption. To address this challenge, following Ram et al. (2023), we treat document retrieval as a pre-processing step, caching the most relevant passages for each question before conducting the question answering experiments.

For *Entity Retrieval*, similar to BM25, DPR, and ANCE, we maintain document lengths at 100 words. However, our approach diverges in sourcing documents: rather than drawing from a large index of 21 million passages, we employ the salient entities within the question and retrieve their corre-

sponding Wikipedia articles, which we then truncate to the initial 100 words. Nonetheless, to explore the impact of document size, beyond the standard 100-word segment aligned with comparable methods, we investigate *Entity Retrieval* across varied lengths, including the first 50, 300, and 1000 words from the retrieved Wikipedia articles.

We conduct our retrieval-augmented question answering experiments using LLaMA 3 model⁴, and in all such experiments⁵, we prevent it from generating sequences longer than 10 subwords.

We do not use any training question-answer pairs in the prompts of our models. In other words, aside from a simple instruction for answering the question, in the Closed-book setting, the prompt solely comprises the question, while in the retrieval-augmented settings using BM25, DPR, and ANCE, it includes the pre-fetched retrieved documents from the corresponding retrieval index along with the question. Similarly, for the *Entity Retrieval* settings, the prompt consists of the first W words of the Wikipedia pages corresponding to the salient entities in the question. We follow Ram et al. (2023) for question normalization and prompt formulation.

4.2 Data

We use the following datasets in our experiments:

EntityQuestions (Sciavolino et al., 2021) is created by collecting 24 common relations (e.g., ‘author of’ and ‘located in’) and transforming fact triples (subject, relation, object) that contain these relations, into natural language questions using pre-defined templates. The dataset comprises 176,560 train, 22,068 dev, and 22,075 test question-answer pairs. To expedite our analytical experiments in this paper, given the extensive size of the dev and test sets, we constrain the question-answer pairs in these subsets to those featuring salient entities within the top 500K most linked Wikipedia pages, as suggested by (Shavarani and Sarkar, 2023). Thus, the dev and test subsets of EntityQuestions considered in our experiments consist of 4,710 and 4,741 questions, respectively.

FactoidQA⁶ (Smith et al., 2008) contains 2,203 hand crafted question-answer pairs derived from Wikipedia articles, with each pair accompanied by

⁴<https://llama.meta.com/llama3/>.

⁵We run our experiments on one server containing 2 RTX A6000s with 49GB GPU memory each.

⁶https://www.cs.cmu.edu/~ark/QA-data/data/Question_Answer_Dataset_v1.2.tar.gz.

242 its corresponding Wikipedia source article included
243 in the dataset.

244 **StrategyQA**⁷ (Geva et al., 2021) is a complex
245 boolean question answering dataset, constructed
246 by presenting individual terms from Wikipedia
247 to annotators. Its questions contain references to
248 more than one Wikipedia entity, and necessitate
249 implicit reasoning for binary (Yes/No) responses.
250 The dataset comprises 5,111 answered questions
251 initially intended for training question answering
252 systems, with the system later tested on test set
253 questions with unreleased answers. This training
254 set is split into two subsets, based on the per-
255 ceived challenge of questions by adversarial anno-
256 tation models (Dua et al., 2019), resulting in train
257 and train_filtered subsets containing 2,290 and
258 2,821 questions, respectively.

259 4.3 Evaluation

260 We evaluate the performance of the retrieval meth-
261 ods using the following metrics:

- nDCG@ k (normalized Discounted Cumulative Gain at rank k ; Järvelin and Kekäläinen, 2002) evaluates the quality of a ranking system by considering both the relevance and the position of documents in the top k results. Mathematically, it is represented as

$$\text{nDCG}@k = \frac{\sum_{i=1}^k \frac{2^{r_i} - 1}{\log_2(i+1)}}{\sum_{i=1}^{|REL_k|} \frac{2^{r_i} - 1}{\log_2(i+1)}}$$

262 Where, r_i denotes the relevance score of a
263 document for a question, with relevance score
264 $r_i = 1$ if the document contains a normalized
265 form of the expected answer to the question,
266 and $r_i = 0$, otherwise. And, REL_k refers
267 to a subset of the retrieved documents that
268 contain a normalized form of the expected
269 answer. nDCG@ k scores range between 0 and 1,
270 where a score of 1 signifies an optimal ranking
271 with the most relevant documents positioned
272 at the top.

- MRR (Mean Reciprocal Rank; Voorhees and Harman, 1999) is the average of the reciprocal ranks of the first relevant document for each question. Mathematically, it is represented as

$$\text{MRR} = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{r_i}$$

⁷<https://allenai.org/data/strategyqa>.

273 where $|Q|$ represents the total number of ques-
274 tions and r_i denotes the rank of the first rele-
275 vant document for the i -th question.

- Top- k Retrieval Accuracy, as reported by Scialvolino et al. (2021), is calculated as the number of questions with at least one relevant document in the top k retrieved documents divided by the total number of questions in the dataset.

282 We evaluate the performance of the retrieval-
283 augmented question-answering models with each
284 retrieval method as follows:

- For FactoidQA and EntityQuestions datasets, we use OpenQA-eval (Kamalloo et al., 2023) scripts to evaluate model performance, and report exact match (EM) and F1 scores by comparing expected answers to normalized model responses.
- For StrategyQA, we present accuracy scores by comparing model responses to the expected boolean answers in the dataset. As well, to assess model comprehension of the task, we count the number of answers that deviate from Yes or No and report this count in a distinct column labeled “Inv #” for each experiment.

298 4.4 Entity Retrieval Performance using 299 Question Entity Annotations

300 We begin our analysis by comparing *Entity Re-*
301 *trieval* performance using BM25, DPR, and ANCE.
302 For this experiment, we calculate nDCG with var-
303 ious retrieved document sets of size $k = 1, 2, 3,$
304 $4, 5, 20,$ and 100 documents. We use the entity
305 annotations provided with the questions from Fac-
306 toidQA and the dev set of EntityQuestions to fetch
307 their corresponding Wikipedia articles, excluding
308 StrategyQA from our analysis as it does not include
309 entity annotations. On average, the FactoidQA and
310 EntityQuestions datasets contain one salient entity
311 per question.

312 To evaluate the effect of document length, we
313 compare *Entity Retrieval* with the first 100 words
314 (equivalent to the size of documents returned by
315 BM25, DPR, and ANCE; noted as $ER100w$) and
316 also consider the first 50, 300, and 1000 words of
317 the retrieved Wikipedia articles (noted as $ER50w,$
318 $ER300w,$ and $ER1000w$). An *Entity Retrieval* doc-
319 ument with 300 words has the same word count as
320 three documents returned by BM25 or DPR.

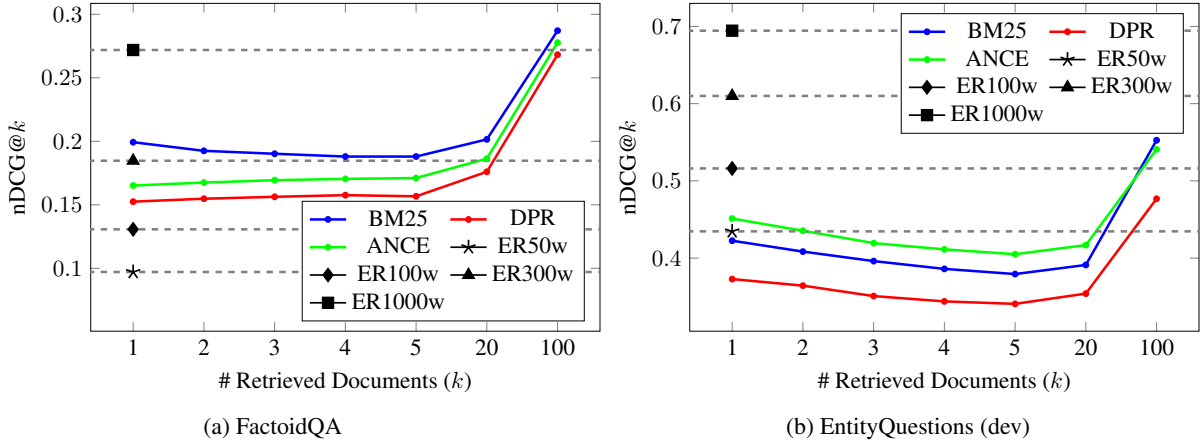


Figure 2: nDCG@ k scores comparing the quality of BM25, DPR, ANCE, and *Entity Retrieval* by considering both the relevance and the position of documents in the top k retrieved passages for each question.

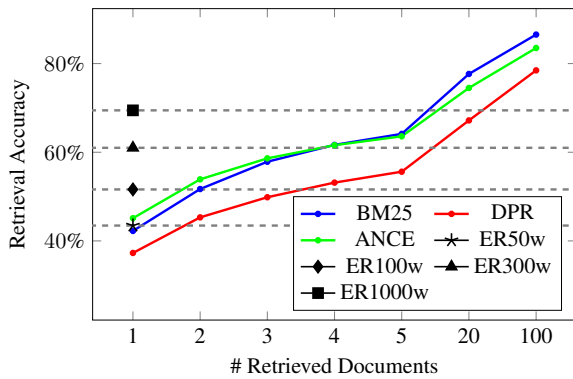


Figure 3: Retrieval Accuracy scores showcasing the correlation between the number of retrieved documents and the expected answers' coverage in EntityQuestions (dev) subset.

	FactoidQA	EntityQuestions (dev)
BM25	0.245	0.522
DPR	0.209	0.456
ANCE	0.222	0.536
ER50w	0.097	0.435
ER100w	0.131	0.516
ER300w	0.185	0.610
ER1000w	0.272	0.695

Table 1: MRR scores comparing the retrieval quality of BM25, DPR, ANCE, and *Entity Retrieval* through the average of the reciprocal ranks of the first relevant document for each question.

Figure 2 presents the computed nDCG@ k scores across varying document sizes, highlighting the superior performance of *Entity Retrieval* over other retrieval methods in the context of the entity-centric datasets under study. Notably, *ER1000w*, which corresponds to ten BM25 retrieved passages in terms of word count, exhibits a retrieval performance on par with 100 retrieved documents in FactoidQA and surpasses BM25, the top-performing retriever on EntityQuestions, by 25%. This impressive performance by *Entity Retrieval* can be attributed to its ability to retrieve fewer, yet more relevant, documents. This observation aligns with the conclusion drawn by Cuconasu et al. (2024), which emphasizes that the retrieval of irrelevant documents can negatively impact performance. *Entity Retrieval* effectively minimizes the retrieval of such documents. Further insights can be gleaned from the comparison of nDCG scores along the x-axis of the plots in Figure 2. As the number of retrieved documents increases, the likelihood of

retrieving irrelevant documents also rises, leading to a decline in retrieval performance when moving from 1 to 5 retrieved documents.

Table 1 showcases the calculated MRR scores, emphasizing the quicker attainment of relevant retrieval documents in *Entity Retrieval* compared to other retrieval methods. Concurrently, Figure 3 illustrates the impact of incrementing the number of retrieved documents on the expansion of the expected answers' coverage for the EntityQuestions dev subset.

While it may be appealing to consider 100 or more documents to simultaneously enhance both nDCG and Retrieval Accuracy, it is important to note that 100 retrieved documents would comprise 10,000 words. This could potentially overwhelm the model with excessive noise (irrelevant documents), and as well, could make it extremely costly to execute retrieval-augmented question answering, especially when the cost of API calls is calculated per token. We would need at least 10,000 tokens (optimistically, assuming each word equates to only

one token) in addition to the tokens in the question. These factors suggest that retrieving a few documents for each question is more beneficial.

Taking these considerations into account, along with the $nDCG@k$, MRR, and Retrieval Accuracy results from this section, we gain a comprehensive understanding of the trade-off between the quality of the retrieved documents, which diminishes as we consider more documents, and the answer coverage, which increases as the model has a higher chance of encountering the right document with the correct hint for the answer. Consequently, we opt for $k = 4$ as a default, and we will always retrieve the top-4 documents in our retrieval-augmented question answering experiments.

4.5 Retrieval-Augmented Question Answering

Next, we shift our focus to study the effectiveness of our proposed *Entity Retrieval* method compared to other retrieval methods in enhancing the quality of responses to entity-centric questions. In this section, we examine three distinct scenarios: (1) the Closed-book setting, where we use “Answer these questions:” as the task instruction, followed by the question, (2) the Retrieval-Augmented setting, where we use retrieved documents as a basis, followed by “Based on these texts, answer these questions:”, and then the question, and (3) the *Entity Retrieval* with question entity annotations, which uses the same prompt as the retrieval-augmented setting. The only difference lies in the documents retrieved, as we have previously discussed.

The initial eight rows of Table 2 present the results of our experiments using LLaMA 3 (8B) model. Upon examining these results, it is evident that *ER100w*, the most analogous *Entity Retrieval* setting to other retrieval methods, outperforms in terms of both EM and F1 scores. This setting returns identical 100-word documents as the other retrieval methods. Furthermore, our dense retrieval results align with the observations of Sciavolino et al. (2021), asserting that entity-centric questions indeed challenge dense retrievers. Although the BM25 method proves successful in enhancing the results compared to the Closed-book setting, it is noteworthy that even *Entity Retrieval* with the initial 50 words of the articles corresponding to the salient entities within questions yields superior results. This is particularly significant when compared to other retrieval methods which necessitate indexing the entire knowledge base on disk and

LLaMA3 (8B)	FactoidQA		EntityQuestions			
			dev		test	
	EM	F1	EM	F1	EM	F1
Closed-book	30.7	39.3	22.7	37.8	22.8	38.1
Retrieval-Augmented QA						
BM25	32.2	42.4	23.8	38.6	23.3	38.5
DPR	29.4	38.5	22.0	36.2	20.5	35.3
ANCE	30.5	40.0	23.1	37.9	22.7	37.9
<i>Entity Retrieval</i> w/ Question Entity Annotations						
ER50w	34.2	43.5	24.9	41.2	23.9	41.0
ER100w	33.6	42.8	26.2	42.8	25.7	42.4
ER300w	33.7	43.0	26.2	42.8	25.3	42.4
ER1000w	35.1	44.9	25.2	41.9	24.5	41.3
<i>Entity Retrieval</i> w/ SPEL Entity Annotations						
ERSp50w	29.7	38.6	24.3	39.2	24.0	39.7
ERSp100w	28.3	37.4	25.0	40.1	24.2	39.8
ERSp300w	26.8	35.6	24.4	39.7	24.6	40.2
ERSp1000w	21.3	30.4	24.4	39.7	23.0	39.2

Table 2: Question answering efficacy comparison between Closed-book and Retrieval-augmentation using BM25, DPR, ANCE, and *Entity Retrieval*. EM refers to the exact match between predicted and expected answers, disregarding punctuation and articles (a, an, the). Results represent the average of two runs with the margin of error values provided in Table 6 in the Appendix.

loading the index into memory; a process required in inference time where caching is not an option.

4.6 Entity Retrieval in absence of Question Entity Annotations

In this section, we concentrate on the most crucial component of the *Entity Retrieval* method: the salient entities within entity-centric questions. We explore a scenario where the entities are not explicitly provided in the question, suggesting the use of an entity linking method to extract these entities. Ideally, we would like to evaluate all recent entity linking methods to identify the most effective one. However, due to time and budget limitations, we depend on the recent benchmarking studies by Ong et al. (2024) to choose an entity linking method. They examine the latest entity linking methods in terms of performance against unseen data and endorse SPEL (Shavarani and Sarkar, 2023) as the top performer. Consequently, we investigate *Entity Retrieval* using entities identified with SPEL, while reserving the examination of other entity linking

Question	Who performed Alexis Colby?	What is the capital of Seine-Saint-Denis?
Answer	Joan Collins	Bobigny
Closed-Book	Diana Ross	Paris
BM25	Linda Evans	Saint-Denis
DPR	Alexis Cohen	Saint-Denis
ANCE	Nicollette Sheridan performed Alexis Colby.	Saint-Denis
ERSp100w	Joan Collins	Bobigny
Question	Where did John Snetzler die?	Where was Brigita Bukovec born?
Answer	Schaffhausen	Ljubljana
Closed-Book	He died in London, England, in 178	Brigita Bukovec was born in Slovenia
BM25	John Snetzler died in London.	Slovenia
DPR	John Snetzler died in London	in Slovakia
ANCE	in England	Ribniča
ERSp100w	Schaffhausen	Ljubljana

Table 3: Example questions from EntityQuestions (dev) to demonstrate the performance of *Entity Retrieval*.

techniques for *Entity Retrieval* for future research.

We maintain the *Entity Retrieval* settings as before, defining *ERSp50w*, *ERSp100w*, *ERSp300w*, and *ERSp1000w* for performing entity linking with SPEL, then retrieving the Wikipedia articles corresponding to the SPEL identified entities, and using the first 50, 100, 300, and 1000 words of these articles as documents to augment the question when prompting the LLM.

Passing the questions from our datasets to SPEL for analysis, we find that it generates a maximum of 8, 3, and 4 annotations for FactoidQA, EntityQuestions, and StrategyQA, respectively. On average, it produces 0.8, 0.7, and 1.1 annotations per question for these same datasets. SPEL successfully identifies and links entities in 56.5% of FactoidQA questions (1244/2203), 66.0% of EntityQuestions (dev) questions (3108/4710), 65.3% of EntityQuestions (test) questions (3095/4741), 75.8% of StrategyQA (train) questions (1735/2290), and 74.2% of StrategyQA (train_filtered) questions (2094/2821).

The final four rows of Table 2 showcase the comparative results of utilizing entities identified by SPEL for *Entity Retrieval*. Given that one-third of EntityQuestions and approximately half of FactoidQA lack identified annotations, the exact match scores reveal that *Entity Retrieval* performs robustly and surpasses BM25, the top-performing competitor retrieval method, for the entity-centric question-answering datasets under examination. This underscores the potential of *Entity Retrieval* within this paradigm. In addition, the disparity between the results with and without question entity annotations strongly indicates the necessity for further research in the Entity Linking domain, which could enhance entity-centric question answering as

LLaMA3 (8B)	train		train_filtered	
	Acc.	Inv #	Acc.	Inv #
BM25	43.8	601	49.1	679
ANCE	47.0	550	51.8	637
ERSp50w	49.7	378	56.2	417
ERSp100w	50.5	367	56.6	389
ERSp300w	46.2	508	53.9	538
ERSp1000w	40.2	778	43.2	924

Table 4: Comparison of *Entity Retrieval* using SPEL identified entities to the best-performing dense and sparse retrieval methods of Table 2 on the StrategyQA dataset. Given the expected boolean results for StrategyQA questions, we restricted LLaMA 3 to generate only one token. *Acc.* indicates the fraction of answers that correctly match the expected Yes or No responses in the dataset, while *Inv #* represents the count of labels that are neither Yes nor No, but another invalid answer. Results represent the average of two runs with the margin of error values provided in Table 7 in the Appendix.

a downstream task. Table 3 provides some example questions where *Entity Retrieval* has led to better answers.

Table 4 presents a comparison of the performance of *Entity Retrieval* using SPEL identified entities against other retrieval methods on the StrategyQA dataset. The results clearly demonstrate the superior performance of *Entity Retrieval* over the top-performing retrieval methods as shown in Table 2. It is important to note that the 100-word setting (*ERSp100w*) is the most analogous to other retrieval methods, given that the size of their retrieved documents is also 100 words. Interestingly, the results from the 1000-word setting suggest that longer documents do not necessarily enhance the

	Total Time	Disk Storage	Main Memory
BM25	45min	11GB	2.3GB
ANCE	960min	61.5GB	64.2GB
ERSp100w	34min	9.4GB	6.3GB

Table 5: Comparison of the required resources for each retrieval method in real-time execution. The reported total time values exclude the time taken to load the indexes and models, focusing solely on the time used to answer the questions.

model’s recall. In fact, beyond a certain length, the model may become overwhelmed by the sheer volume of noise, leading to confusion. Lastly, the invalid count values suggest that *Entity Retrieval* is more effective in assisting the model to comprehend the boolean nature of expected responses, eliminating the need to rely on retrieval from millions of passages.

4.7 Real-time Efficiency Analysis

Our analysis thus far has primarily focused on the retrieval performance, without consideration for the time and memory efficiency; crucial factors in retrieval method selection. In this section, we shift our focus to these aspects.

We begin by replacing the pre-built cache with the original retrieval modules that were used in creating the retrieval cache document sets. We load the indexes and the necessary models for fetching the retrieval documents. We then record the peak main memory requirement of each method during the experiment. It is important to note that all retrieval methods primarily rely on main memory, with minimal differences in GPU memory requirements. Therefore, we report an average GPU memory requirement of 35GB for the LLaMA 3 (8B) setting and exclude it from our results table. We then feed all 2,203 FactoidQA questions into the BM25, ANCE, and *Entity Retrieval* (using SPEL identified entities) to fetch the top-4 documents. We report the total time taken to generate answers to all the questions. Additionally, we keep track of all the pre-built models and indexes that each method requires for download and storage. We report the total size of all downloaded files to disk.

Table 5 presents our findings on time and memory requirements. It is evident that ANCE requires significantly more time to fetch and provide documents, six times more disk space to store its in-

dexes, and over ten times higher main memory demands to load its dense representations. In contrast, BM25 and *Entity Retrieval* are more resource-friendly. Notably, *Entity Retrieval* is 25% faster than BM25 in response generation while demanding the total memory and disk space of a standard personal computer. Future research can be directed towards reducing the memory requirements of *Entity Retrieval*; a direction which we find quite promising.

5 Related Work

Similar to our studies, [Kandpal et al. \(2023\)](#) investigate the impact of salient entities on question answering, and propose constructing oracle retrieval documents as the 300-word segment surrounding the ground-truth answer from the Wikipedia page that contains the answer (entity name). Our approach leverages salient entities from questions without directly involving answers. Additionally, they primarily use entities to classify questions into those concerning frequent knowledge base entries versus those about rare entries on the long-tail, whereas our approach assigns a more substantial role to entities, treating them as pointers guiding the retrieval of relevant documents to augment questions.

[Sciavolino et al. \(2021\)](#) compare DPR and BM25 retrievers for entity-centric questions, and demonstrate that DPR greatly underperforms BM25. They attribute this to dense retrievers’ difficulty with infrequent entities, which are less represented in training data. In contrast, BM25’s frequency-based retrieval is not sensitive to entity frequency. We take a parallel approach and propose a simple yet effective method that leverages salient entities in the question for identifying augmentation documents.

6 Conclusion

In this study, we focused on retrieval-augmented question answering, and explored various retrieval methods that rely on the similarity between the question and the content of the passages to be retrieved. We introduced a novel approach, *Entity Retrieval*, which deviates from the conventional similarity mechanism. Instead, it capitalizes on the salient entities within the question to identify retrieval documents. Our findings indicate that our proposed method is not only more accurate but also faster in the context of entity-centric question answering.

574 Limitations and Ethical Considerations

575 Our proposed *Entity Retrieval* method is specif- 624
576 ically tailored for answering entity-centric ques- 625
577 tions, with its performance heavily reliant on the 626
578 presence of question entities. In scenarios where 627
579 entity annotations are absent, the method’s effec- 628
580 tiveness is directly tied to the performance of exter- 629
581 nal entity linking methods. We acknowledge that 630
582 our exploration of potential entity linking methods 631
583 has not been exhaustive, and further investigation 632
584 may yield insights that could enhance the *Entity* 633
585 *Retrieval* method, even in the absence of question 634
586 entity annotations. 635

587 Furthermore, we recognize that entity linking 636
588 can occasionally result in ambiguous entities. Our 637
589 research has not delved into the impact of such am- 638
590 biguities on the *Entity Retrieval* method, and we 639
591 propose that future studies should focus on ensur- 640
592 ing the selection of the most contextually appropri- 641
593 ate entities for retrieval. 642

594 Our research is primarily centered on Wikipedia 643
595 as the knowledge base, a choice heavily influenced 644
596 by previous studies for the sake of comparability. 645
597 However, we acknowledge the importance of ex- 646
598 ploring other knowledge bases and ontologies, par- 647
599 ticularly in different domains, such as UMLS (Bo- 648
600 denreider, 2004) in the medical field. 649

601 In terms of benchmarking, we have compared 650
602 the *Entity Retrieval* method against a limited se- 651
603 lection of existing retrieval methods, guided by 652
604 our judgement, experience, and considerations of 653
605 implementation availability. We concede that our 654
606 comparison has not been exhaustive, and this rea- 655
607 soning extends to our comparison using different 656
608 LLMs and their available sizes. 657

609 Our research is on English only, and we acknowl- 658
610 edge that entity-centric question answering in other 659
611 languages is also relevant and important. We hope 660
612 to extend our work to cover multiple languages in 661
613 the future. We inherit the biases that exist in the 662
614 data used in this project, and we do not explicitly 663
615 de-bias the data. We are providing our code to the 664
616 research community and we trust that those who 665
617 use the model will do so ethically and responsibly. 666

618 References

619 Steven Abney, Michael Collins, and Amit Singhal. 2000. 667
620 *Answer extraction*. In *Sixth Applied Natural Lan- 668*
621 *guage Processing Conference*, pages 296–301, Seat- 669
622 tle, Washington, USA. Association for Computa- 670
623 tional Linguistics. 671

Ahmad Aghaebrahimian and Filip Jurčiček. 2016. 624
625 *Open-domain factoid question answering via knowl- 626*
627 *edge graph search*. In *Proceedings of the Workshop 628*
629 *on Human-Computer Question Answering*, pages 22– 630
28, San Diego, California. Association for Computa- 631
tional Linguistics. 632

Keivan Alizadeh, Iman Mirzadeh, Dmitry Belenko, 633
634 Karen Khatamifard, Minsik Cho, Carlo C Del Mundo, 635
636 Mohammad Rastegari, and Mehrdad Farajtabar. 2023. 637
638 *Llm in a flash: Efficient large language model 639*
640 *inference with limited memory*. *arXiv preprint 641*
642 *arXiv:2312.11514*. 643

Olivier Bodenreider. 2004. *The unified medical lan- 644*
645 *guage system (umls): integrating biomedical termi- 646*
647 *nology*. *Nucleic acids research*, 32(suppl_1):D267– 648
D270. 649

Tom Brown, Benjamin Mann, Nick Ryder, Melanie 650
651 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind 651
652 Neelakantan, Pranav Shyam, Girish Sastry, Amanda 652
653 Aske, et al. 2020. *Language models are few-shot 653*
654 *learners*. *Advances in neural information processing 654*
655 *systems*, 33:1877–1901. 655

Danqi Chen, Adam Fisch, Jason Weston, and Antoine 656
657 Bordes. 2017. *Reading Wikipedia to answer open- 657*
658 *domain questions*. In *Proceedings of the 55th Annual 658*
659 *Meeting of the Association for Computational Lin- 659*
660 *guistics (Volume 1: Long Papers)*, pages 1870–1879, 660
661 Vancouver, Canada. Association for Computational 661
662 Linguistics. 662

Florin Cuconasu, Giovanni Trappolini, Federico Sicil- 663
664 iano, Simone Filice, Cesare Campagnano, Yoelle 664
665 Maarek, Nicola Tonello, and Fabrizio Silvestri. 665
666 2024. *The power of noise: Redefining retrieval for 666*
667 *rag systems*. *arXiv preprint arXiv:2401.14887*. 667

Wanyun Cui, Yanghua Xiao, Haixun Wang, Yangqiu 668
669 Song, Seung-won Hwang, and Wei Wang. 2017. 669
670 *Kbqa: Learning question answering over qa corpora 670*
671 *and knowledge bases*. *Proceedings of the VLDB En- 671*
672 *dowment*, 10(5). 672

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 663
664 Kristina Toutanova. 2019. *BERT: Pre-training of 664*
665 *deep bidirectional transformers for language under- 665*
666 *standing*. In *Proceedings of the 2019 Conference of 666*
667 *the North American Chapter of the Association for 667*
668 *Computational Linguistics: Human Language Tech- 668*
669 *nologies, Volume 1 (Long and Short Papers)*, pages 669
670 4171–4186, Minneapolis, Minnesota. Association for 670
671 Computational Linguistics. 671

Matthijs Douze, Alexandr Guzhva, Chengqi Deng, 672
673 Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel 673
674 Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé 674
675 Jégou. 2024. *The faiss library*. *arXiv preprint 675*
676 *arXiv:2401.08281*. 676

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel 677
678 Stanovsky, Sameer Singh, and Matt Gardner. 2019. 677
678

679	DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.	
680		
681		
682		
683		
684		
685		
686		
687	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies . <i>Transactions of the Association for Computational Linguistics</i> , 9:346–361.	
688		
689		
690		
691		
692		
693	Michael Glass, Gaetano Rossiello, Md Faisal Mahbub Chowdhury, Ankita Naik, Pengshan Cai, and Alfio Gliozzo. 2022. Re2G: Retrieve, rerank, generate . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2701–2715, Seattle, United States. Association for Computational Linguistics.	
694		
695		
696		
697		
698		
699		
700		
701	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. <i>Neural computation</i> , 9(8):1735–1780.	
702		
703		
704	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning . <i>Transactions on Machine Learning Research</i> .	
705		
706		
707		
708		
709	Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 874–880, Online. Association for Computational Linguistics.	
710		
711		
712		
713		
714		
715		
716	Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of ir techniques . <i>ACM Transactions on Information Systems (TOIS)</i> , 20(4):422–446.	
717		
718		
719		
720	Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus . <i>IEEE Transactions on Big Data</i> , 7(3):535–547.	
721		
722		
723	Ehsan Kamaloo, Nouha Dziri, Charles Clarke, and Davood Rafiei. 2023. Evaluating open-domain question answering in the era of large language models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5591–5606, Toronto, Canada. Association for Computational Linguistics.	
724		
725		
726		
727		
728		
729		
730	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge . In <i>International Conference on Machine Learning</i> , pages 15696–15707. PMLR.	
731		
732		
733		
734		
	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6769–6781, Online. Association for Computational Linguistics.	735
		736
		737
		738
		739
		740
		741
	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	742
		743
		744
		745
		746
		747
		748
		749
		750
	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020b. Retrieval-augmented generation for knowledge-intensive nlp tasks . <i>Advances in Neural Information Processing Systems</i> , 33:9459–9474.	751
		752
		753
		754
		755
		756
	Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. Pyserini: A Python toolkit for reproducible information retrieval research with sparse and dense representations . In <i>Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)</i> , pages 2356–2362.	757
		758
		759
		760
		761
		762
		763
		764
	Denis Lukovnikov, Asja Fischer, and Jens Lehmann. 2019. Pretrained transformers for simple question answering over knowledge graphs . In <i>The Semantic Web–ISWC 2019: 18th International Semantic Web Conference, Auckland, New Zealand, October 26–30, 2019, Proceedings, Part I 18</i> , pages 470–486. Springer.	765
		766
		767
		768
		769
		770
		771
	Salman Mohammed, Peng Shi, and Jimmy Lin. 2018. Strong baselines for simple question answering over knowledge graphs with and without neural networks . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 291–296, New Orleans, Louisiana. Association for Computational Linguistics.	772
		773
		774
		775
		776
		777
		778
		779
		780
	Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. 2022. Large dual encoders are generalizable retrievers . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9844–9855, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	781
		782
		783
		784
		785
		786
		787
		788
	Nicolas Ong, Hassan S. Shavarani, and Anoop Sarkar. 2024. Unified examination of entity linking in absence of candidate sets . In <i>Proceedings of the 2024 Conference of the North American Chapter of the</i>	789
		790
		791
		792

907 Wenhao Yu, Zhihan Zhang, Zhenwen Liang, Meng
908 Jiang, and Ashish Sabharwal. 2023. [Improving lan-
909 guage models via plug-and-play retrieval feedback.](#)
910 *arXiv preprint arXiv:2305.14002*.

911 Jingtao Zhan, Jiabin Mao, Yiqun Liu, Min Zhang, and
912 Shaoping Ma. 2020a. [Learning to retrieve: How
913 to train a dense retrieval model effectively and effi-
914 ciently.](#) *arXiv preprint arXiv:2010.10469*.

915 Jingtao Zhan, Jiabin Mao, Yiqun Liu, Min Zhang, and
916 Shaoping Ma. 2020b. [Repbert: Contextualized text
917 embeddings for first-stage retrieval.](#) *arXiv preprint
918 arXiv:2006.15498*.

919 Qin Zhang, Shangsi Chen, Dongkuan Xu, Qingqing
920 Cao, Xiaojun Chen, Trevor Cohn, and Meng Fang.
921 2023. [A survey for efficient open domain question an-
922 swering.](#) In *Proceedings of the 61st Annual Meeting
923 of the Association for Computational Linguistics (Vol-
924 ume 1: Long Papers)*, pages 14447–14465, Toronto,
925 Canada. Association for Computational Linguistics.

926 Fengbin Zhu, Wenqiang Lei, Chao Wang, Jianming
927 Zheng, Soujanya Poria, and Tat-Seng Chua. 2021.
928 [Retrieving and reading: A comprehensive survey on
929 open-domain question answering.](#) *arXiv preprint
930 arXiv:2101.00774*.

931 **A Margin of Error Results**

LLaMA3 (8B)*	FactoidQA		EntityQuestions			
			dev		test	
	EM	F1	EM	F1	EM	F1
Closed-book	30.7±0.1	39.3±0.0	22.7±0.5	37.8±1.0	22.8±0.1	38.1±0.6
Retrieval-Augmented QA						
BM25	32.2±1.1	42.4±0.2	23.8±0.3	38.6±0.8	23.3±0.0	38.5±0.1
DPR	29.4±1.0	38.5±1.2	22.0±0.1	36.2±0.2	20.5±0.4	35.3±0.6
ANCE	30.5±0.4	40.0±0.4	23.1±0.7	37.9±0.6	22.7±0.7	37.9±0.9
<i>Entity Retrieval w/ Question Entity Annotations</i>						
ER50w	34.2±0.7	43.5±0.6	24.9±0.2	41.2±0.0	23.9±0.5	41.0±0.1
ER100w	33.6±0.5	42.8±0.5	26.2±0.0	42.8±0.1	25.7±0.1	42.4±0.0
ER300w	33.7±1.4	43.0±1.7	26.2±0.4	42.8±0.0	25.3±1.0	42.4±1.1
ER1000w	35.1±0.4	44.9±0.7	25.2±0.4	41.9±0.6	24.5±0.9	41.3±0.6
<i>Entity Retrieval w/ SPEL Entity Annotations</i>						
ERSp50w	29.7±0.3	38.6±0.7	24.3±0.2	39.2±0.1	24.0±0.1	39.7±0.0
ERSp100w	28.3±0.9	37.4±1.2	25.0±0.4	40.1±0.3	24.2±0.2	39.8±0.1
ERSp300w	26.8±0.6	35.6±0.7	24.4±0.0	39.7±0.1	24.6±0.3	40.2±0.5
ERSp1000w	21.3±0.5	30.4±0.8	24.4±0.1	39.7±0.1	23.0±0.7	39.2±0.7

Table 6: Question answering efficacy comparison between Closed-book and Retrieval-augmentation using BM25, DPR, ANCE, and *Entity Retrieval*. EM refers to the exact match between predicted and expected answers, disregarding punctuation and articles (a, an, the).

* Results represent the average of two runs, accompanied by a margin of error based on a 99% confidence interval.

LLaMA3 (8B)*	train		train_filtered	
	Acc.	Inv #	Acc.	Inv #
BM25	43.8±0.1	601±4	49.1±1.0	679±7
ANCE	47.0±1.2	550±15	51.8±1.0	637±42
ERSp50w	49.7±1.2	378±34	56.2±1.3	417±31
ERSp100w	50.5±2.0	367±21	56.6±0.5	389±1
ERSp300w	46.2±1.9	508±22	53.9±1.9	538±14
ERSp1000w	40.2±0.4	778±3	43.2±0.3	924±13

Table 7: Comparison of *Entity Retrieval* using SPEL identified entities to the best-performing dense and sparse retrieval methods on the StrategyQA dataset.

* Results represent the average of two runs, accompanied by a margin of error based on a 99% confidence interval.