

Knowledge Graph Prompting for Multi-Document Question Answering

Yu Wang,¹ Nedim Lipka,² Ryan A. Rossi,² Alexa Siu,² Ruiyi Zhang,² Tyler Derr¹

¹ Vanderbilt University, Nashville, USA

² Adobe Research, San Jose, USA

yu.wang.1@vanderbilt.edu, {lipka, ryrossi, asiu, ruizhang}@adobe.com, tyler.derr@vanderbilt.edu

Abstract

The ‘pre-train, prompt, predict’ paradigm of large language models (LLMs) has achieved remarkable success in open-domain question answering (OD-QA). However, few works explore this paradigm in multi-document question answering (MD-QA), a task demanding a thorough understanding of the logical associations among the contents and structures of documents. To fill this crucial gap, we propose a Knowledge Graph Prompting (KGP) method to formulate the right context in prompting LLMs for MD-QA, which consists of a graph construction module and a graph traversal module. For graph construction, we create a knowledge graph (KG) over multiple documents with nodes symbolizing passages or document structures (e.g., pages/tables), and edges denoting the semantic/lexical similarity between passages or document structural relations. For graph traversal, we design an LLM-based graph traversal agent that navigates across nodes and gathers supporting passages assisting LLMs in MD-QA. The constructed graph serves as the global ruler that regulates the transitional space among passages and reduces retrieval latency. Concurrently, the graph traversal agent acts as a local navigator that gathers pertinent context to progressively approach the question and guarantee retrieval quality. Extensive experiments underscore the efficacy of KGP for MD-QA, signifying the potential of leveraging graphs in enhancing the prompt design and retrieval augmented generation for LLMs. Our code: <https://github.com/YuWVandy/KG-LLM-MDQA>.

1 Introduction

Due to the emergence of large language models (LLMs), the ‘pre-train, prompt, and predict’ paradigm has revolutionized natural language processing (NLP) in real-world applications, such as open-domain question answering, fact-checking, and arithmetic reasoning (Chen et al. 2017; Thorne et al. 2018; Asai et al. 2019; Karpukhin et al. 2020; Aly et al. 2021; Qin et al. 2023; Zou and Caragea 2023; Liu, Dong, and Zhang 2023). However, no significant efforts have investigated this framework in the scenario of multi-documental question answering (MD-QA), which enjoys practical usage in academic research, customer support, and financial/legal inquiries that require deriving insightful analysis from multiple documents (Tessuto 2011; Bolino, Long, and Turnley 2016).

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

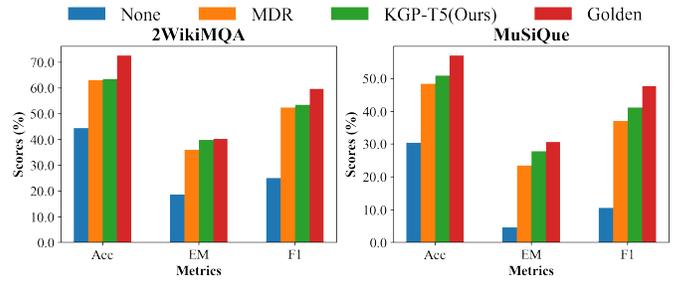


Figure 1: MD-QA performance when prompting ChatGPT with the context retrieved using different strategies.

To investigate the capability of LLMs for MD-QA, we randomly sample multi-document questions from the development set of 2WikiMQA (Ho et al. 2020) and MuSiQue (Trivedi et al. 2022b), and then prompt LLMs in four different strategies for the answer¹. Successfully answering these questions requires knowledge from multiple Wikipedia documents. As shown in Figure 1, on 2WikiMQA and MuSiQue, directly prompting LLMs without providing any context, i.e., None, achieves only 25.07%/10.58% F1 and 18.60%/4.60% EM on 2WikiMQA/MuSiQue, which is far less than 59.69%/47.75% F1 and 40.20%/30.60% EM when prompting with supporting facts² provided as contexts, i.e., the Golden one. This demonstrates the limitation of fulfilling MD-QA using solely the knowledge encoded in LLMs. One common solution to overcome this limitation in conventional OD-QA and single document question-answering (D-QA) (Xu et al. 2020; Mathew, Karatzas, and Jawahar 2021) is to retrieve grounding contexts and derive faithful answers from the contexts, i.e., retrieve-and-read (Zhu et al. 2021; Ju et al. 2022). However, unlike OD-QA and D-QA, the primary challenge of MD-QA roots in its demands for alternatively retrieving and reasoning knowledge across different documents (Pereira et al. 2023; Caciularu et al. 2023). For example, successfully answering questions in Figure 2(a)-(b) requires reasoning over distinct passages from two different documents (in these two cases, Wikipedia pages). Moreover, each document is essentially a

¹Detailed experimental setting is presented in Section 5.

²Supporting facts: passages that are assumed to contain the answer to the question.

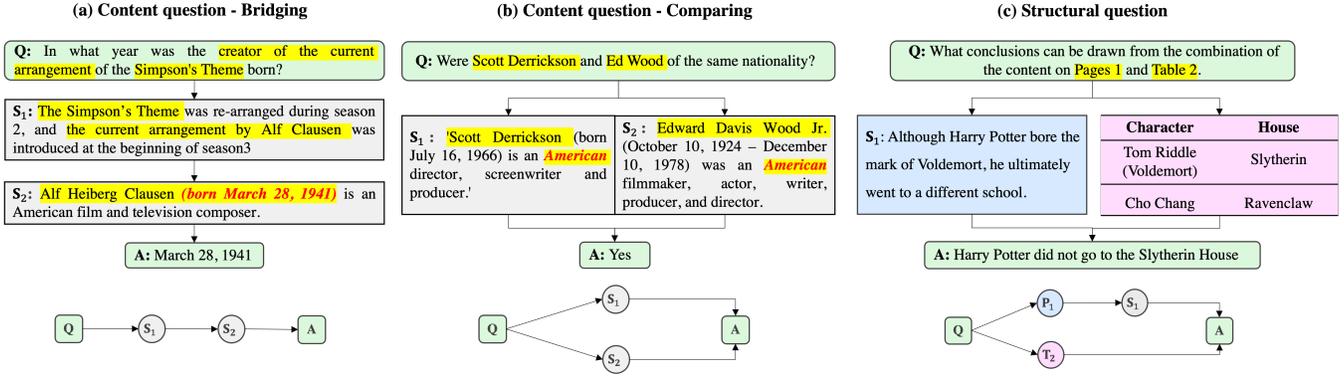


Figure 2: Three popular questions that require reasoning and retrieving over passages/pages/tables from multiple documents. (a) **Bridging questions** rely on sequential reasoning while (b) **Comparing questions** rely on parallel reasoning over different passages. (c) **Structural questions** rely on fetching contents in the corresponding document structures.

compilation of multi-modality structured data (e.g., pages, sections, paragraphs, tables, and figures) and some questions may specifically ask for the content in certain structures, which necessitates a comprehensive grasp of these complex document structures. For example, the question in Figure 2(c) asks about the difference between Page 1 and Table 2, which is unanswerable if leveraging heuristic methods like BM25 or deep-learning ones like DPR (Karpukhin et al. 2020). Building on top of previous challenges, the advent of LLMs introduces new complexities.

For the challenge of alternatively retrieving and reasoning knowledge across different documents, although previous works train a multi-hop retriever (Xiong et al. 2020; Yavuz et al. 2022) to imitate such process by sequentially fetching the next passage based on the already-retrieved ones, none of them explore the potential of engaging LLMs into this process. More recent works design different prompting strategies such as Chain/Tree/Graph-of-thought (Trivedi et al. 2022a; Wei et al. 2022; Yao et al. 2023; Yao, Li, and Zhao 2023) to guide LLMs approaching answers progressively. However, prompting non-open-sourced LLMs back and forth incurs forbiddable latency as well as unaffordable consumption. In addition, how to integrate different document structures into the prompt design so that LLMs can understand them is still an open-ended question.

Given the above challenges, we propose a knowledge graph prompting (KGP) method for enhancing LLMs in MD-QA. Specifically, we construct a KG over the given documents with nodes symbolizing passages or document structures and edges denoting their lexical/semantic similarity between passages or intra-document structural relations. Then for the first challenge of alternative reasoning and retrieving knowledge across different documents, we design an LLM-based KG traversal agent, which can alternatively generate the next evidence to approach the question, i.e., reasoning, and select the most promising neighbor to visit from the constructed KG based on the generated evidence, i.e., retrieval. Moreover, we apply the instruction fine-tuning strategy to augment the reasoning capability of the LLM-based KG traversal agent and hence refrain from repeatedly prompting non-open-sourced LLMs for evidence generation. For the

multi-modality challenge, we add different types of nodes to the KG characterizing different document structures and hence enabling content retrieval within those specific structures. We highlight our contributions as follows:

- **Generally-applicable KG Construction.** We propose three KG construction methods over documents, with passages or document structures as nodes and their lexical/semantic similarity or structural relations as edges. Then we empirically evaluate the quality of the constructed KGs in MD-QA by checking the level of overlap between the neighborhood and the supporting facts for each question (Figure 5). Additionally, we provide a comprehensive summary of both our proposed and existing KG construction methods in Table 5 in Supplementary.
- **Engaging KG for Prompt Formulation.** We design a Knowledge Graph Prompting (KGP) method, which leverages the LLM-based KG traversal agent to retrieve the question-relevant contexts by traversing the constructed KG. Moreover, we fine-tune this agent to adaptively traverse the most promising neighbors for approaching the question based on the visited nodes (retrieved passages).
- **Case Studies Verifying MD-QA Framework.** We compare the performance of MD-QA when using different types of LLM agents in graph traversal (Table 2) on the KGs constructed over different numbers of documents (Figure 7(c)). We conduct case studies on visualizing KGP for MD-QA in Section 8.7 in Supplementary.

2 Notations

Following (Tian et al. 2023a), let $G = (\mathcal{V}, \mathcal{E})$ be a knowledge graph constructed from a set of documents \mathcal{D} , where the node set $\mathcal{V} = \{v_i\}_{i=1}^n$ representing document structures (e.g., passages/pages/tables, etc.) and the edge set $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$ representing the connections among different nodes (e.g., semantic/lexical similarity and belonging relations among document structures, etc.). Let $\mathcal{X} = \{\mathcal{X}_i\}_{i=1}^n$ be node features and \mathcal{X}_i corresponds to the feature of node v_i , the form of which could be the text for the passage, the markdown for the table and the page number for the page.

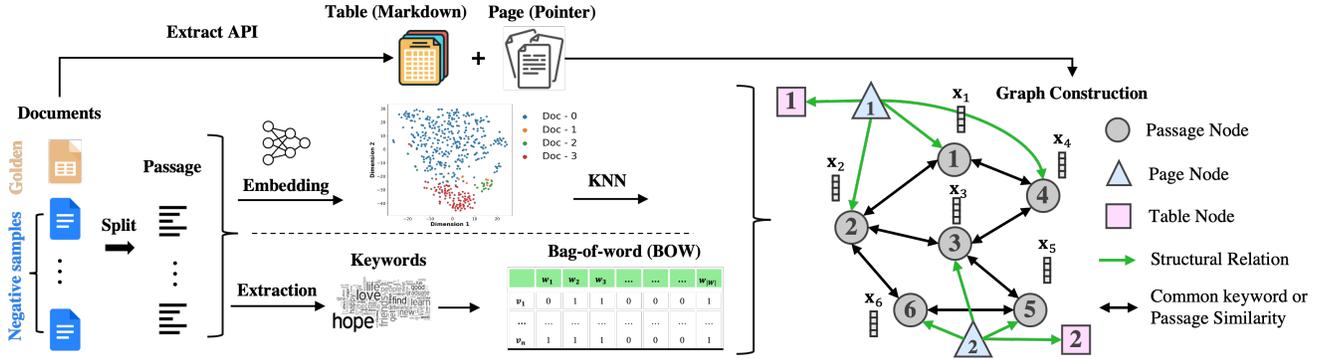


Figure 3: Knowledge Graph Construction. We split each document in the document collection into passages. For each passage, we either directly obtain their embeddings via pre-trained encoders or extract their keywords to build bag-of-word (BOW) features. Then we connect two passages based on their embedding similarity or whether they share common keywords. Additionally, we extract tables/pages via Extract-PDF API and add them as structural nodes to the KG. If pages include passages and tables, we add a directed edge to denote the belonging relations. The table nodes include the markdown formatted content of that table as Figure 8 in Supplementary has empirically shown that LLMs are able to understand tables in this format.

3 Knowledge Graph Construction

Despite numerous well-established KGs (Hoffart et al. 2013; Tian et al. 2023b), they treat nodes/edges as entities/relations, which necessitates sophisticated relational extraction techniques and thereby limits their applicability in general domains (Huang et al. 2021). Additionally, their primary focus on the Wikipedia domain also restricts their usage for answering non-Wikipedia questions such as ones over legal or financial documents. To remedy this issue, we propose generally applicable KG construction methods.

We first analyze two representative questions in Figure 2(a)-(b) to motivate our KG construction. Answering these two questions necessitates the deduction of logical associations among different passages. These associations are encoded either through 1) lexical similarity: common keywords shared among different passages, e.g., ‘Alf Clausen’ bridges passage S_1 and passage S_2 in Figure 2(a), or 2) semantic similarity: syntactic elements that convey semantic relations, e.g., ‘nationality’ and ‘American director’ in Figure 2(b). This motivates us to construct the graph by modeling passages as nodes and their lexical/semantic similarity as edges. More specifically in Figure 3, we split each document into individual passages, and for each passage S_i , we add a node v_i to the KG with its feature being the text of that passage \mathcal{X}_i . Then we add edges by checking the lexical/semantic similarity between pairs of passage nodes.

TF-IDF KG Construction For adding edges according to lexical similarity, we first apply TF-IDF keyword extraction (Ramos et al. 2003) over each document to filter out meaningless words such as supporting verbs and articles, which also reduces the dimension of bag-of-word (BOW) features, sparsifies the constructed graph and increases the graph traversal efficiency. In addition, we add the document title into the extracted keyword set since some questions focus on title entities. We collect the extracted keywords from all documents to form the keyword space \mathcal{W} and then connect two passages if they share any common keyword in \mathcal{W} .

KNN-ST/MDR KG Construction For adding edges according to semantic similarity, we can readily employ pre-existing models such as sentence transformers to generate passage embedding \mathbf{x}_i for each node v_i and subsequently compute pairwise similarity matrix to construct the K-nearest neighbor (KNN) graph. However, these off-the-shelf models, typically trained on tasks not so-related to MD-QA, may not adequately encapsulate necessary logical associations in their embedding similarity demanded by the question. To overcome this problem, we follow the training strategy of MDR (Xiong et al. 2020) and train a sentence encoder by predicting the subsequent supporting facts based on previously supporting facts, thereby endowing the encoder with reasoning capability. Consequently, the embedding similarity and the corresponding constructed KNN graph fundamentally encapsulate the necessary logical associations between different passages.

TAGME Moreover, we employ TAGME (Min et al. 2019) to extract Wikipedia entities from each passage and construct the graph based on whether two passage nodes share common Wikipedia entities.

In addition to passage nodes, we further add structural nodes into the graph by extracting document structures via Extract-PDF³. In this paper, we only consider adding pages and tables but the constructed KG can include more different types of document structures. The feature of table nodes is the markdown since LLMs can understand this as demonstrated in Figure 8 in Supplementary. The feature of page nodes is the page number and we add directed edges from it to sentence/table nodes in that page. *Note that we do not aim to propose a one-size-fits-all KG construction method. Instead, we seek to compare the merits and limitations of various methods in Table 5, offering guidance on which KGs are best suited for specific scenarios.*

³<https://developer.adobe.com/document-services/docs/overview/pdf-extract-api/>

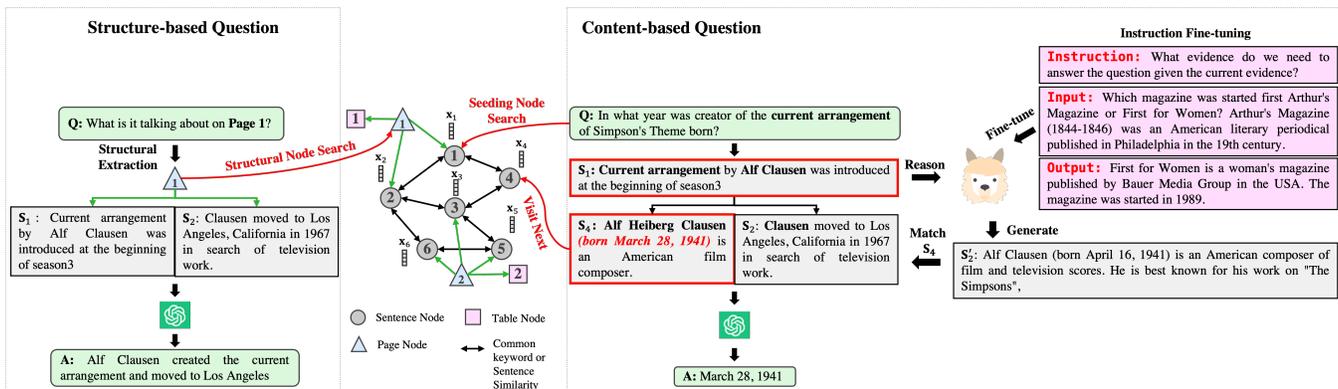


Figure 4: LLM-based KG traversal agent for context retrieval. For questions on document structures (left), we employ LLM to extract structures and retrieve their corresponding contents (the content of pages are passages belonging to that page and the content of tables is the markdown-formatted text). For questions on document content, we concatenate it with the currently retrieved context and prompt the LLM to generate the next evidence to answer the question. By comparing the similarity between the candidate neighboring sentences and the generated passage, we determine the next passage node to traverse. Correspondingly, the candidate neighbors are updated for the next round of traversal.

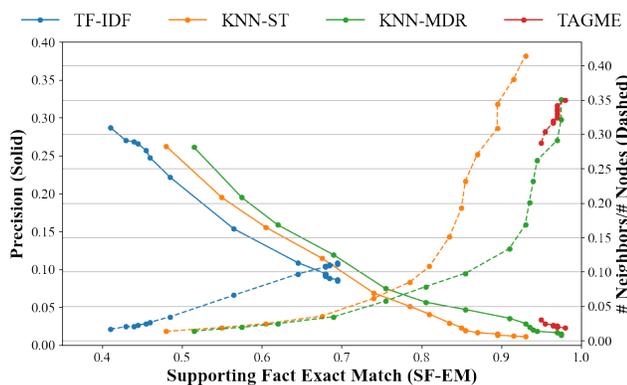


Figure 5: Quality of KGs on HotpotQA. For each KG Construction method, as the average number of neighbors increases (KG becomes denser) in the right y-axis, the SF-EM increases while the precision decreases. KNN-MDR achieves a better trade-off than TF-IDF and KNN-ST. KGs constructed by TAGME are denser than others.

To verify the constructed KGs indeed encode the necessary information for MD-QA, we randomly sample questions from HotpotQA and construct KGs over the set of documents for each of these questions using our proposed methods. We vary the hyperparameters to control the sparsity of the constructed KG and measure how much percentage of the supporting facts are covered by neighbors of the seeding passages initially searched by TF-IDF. More details about the four construction methods and their hyperparameters are included in Section 8.5 in Supplementary. As shown in Figure 5, as the constructed graph becomes denser, the chance that the neighboring node passages hit the supporting facts increases (i.e., SF-EM increases) although the redundant information also increases (i.e., the precision decreases). Given the common keywords shared between one passage to all other passages are typically far less than the total number of passages across all documents, the density of

the constructed graph by TF-IDF would be upper-bounded, causing lower SF-EM (evidenced by SF-EM below 0.7 in Figure 5 for TF-IDF curve). For TAGME, we empirically find it identifies a larger quantity of entities mentioned in a single passage, which leads to a denser graph and causes the starting SF-EM of TAGME to be already around 0.95. In addition, since KNN-MDR is pre-trained by predicting the next supporting facts (Xiong et al. 2020) on HotpotQA, it achieves better trade-off than KNN-ST where the embeddings are directly obtained from the sentence transformer without dataset-specific pre-training.

To summarize, although high SF-EM indicates that the supporting facts for most questions are fully covered by the neighbors of seeding passages, low precision signifies that most of these neighboring passages are irrelevant to the question. Therefore, if we blindly perform graph traversal without any question-tailored adaptation, our retrieved contexts would include redundant passages and compromise the capability of LLMs in MD-QA (which is also verified by the lower performance of KGP-Random in Table 2). To remedy this issue, in the next section, we introduce an LLM-based KG traversal agent to adaptively visit neighboring passages that are most conducive to answering the given question.

4 LLM-based KG Traversal Agent

A natural solution to enable adaptive knowledge graph traversal is to rank the candidate nodes, i.e., the neighbors of the already-visited nodes in our case, thereby determining which ones to visit next. The most straightforward way is to apply heuristic-based fuzzy matching or embedding-based similarity ranking, which cannot capture the intrinsic logical relations between the already traversed paths and the nodes to visit next. Instead, we fine-tune a large language model (LLM) to guide the knowledge graph traversal towards the next most promising passages in approaching the question based on the visited passages, which we term as the LLM-based KG traversal agent.

Given a question q asking about the document content, the LLM-based graph traversal agent reasons over previously visited nodes/retrieved passages $\{s_k\}_{k=0}^j$ and then generates the next passage s_{j+1} as follows:

$$s_{j+1} = \arg \max_{v \in \mathcal{N}_j} \phi(g(\mathcal{X}_v), f(\|\|_{k=0}^j \mathcal{X}_k)), \quad (1)$$

where $\|\|_{k=0}^j \mathcal{X}_k$ concatenates the textual information of previously retrieved passages/visited nodes. For the choice of f , one way is to employ encoder-only models like Roberta-base (Asai et al. 2019; Xiong et al. 2020; Yavuz et al. 2022) and correspondingly g would be another encoder model with $\phi(\cdot)$ being the inner product measuring the embedding similarity. Another way is to employ encoder-decoder models such as T5 (Brown et al. 2020; Touvron et al. 2023) and correspondingly g would be an identity function with $\phi(\cdot)$ measuring the textual similarity. To mitigate the hallucination issue and enhance the reasoning capability (Wei et al. 2022; Ji et al. 2023) of the LLM traversal agent, we further instruction fine-tune f (Chung et al. 2022) by predicting the next supporting facts based on previous supporting facts, thereby integrating commonsense knowledge encoded originally in their pre-trained parameters with the enhanced reasoning capability inherited from the instruction fine-tuning. After visiting the top-scoring nodes selected from the candidate neighbor queue by Eq (1), the candidate neighbor queue is updated by adding neighbors of these newly visited nodes. We iteratively apply this process until hit the preset budget. Next, we illustrate the above process with an example in Figure 4 and present the algorithm thereafter.

Figure 4 presents the content-based question asking ‘In what year was the creator of the current arrangement of Simpson’s Theme born?’. We use TF-IDF search to initialize the seeding passage Node 1, which reads: ‘Alf Heiberg Clausen (born March 28, 1941) is an American film composer’. Subsequently, we prefix the currently retrieved-context (Node 1) with the question and prompt the LLM to generate the next evidence required to approach the question one step closer. Because we augment the reasoning capability of the LLM by instruction fine-tuning, it is expected to recognize the logical association between the question and the currently retrieved context. Consequently, it can predict the subsequent passage that *maintains logical coherence, albeit may contain factual mistakes*, i.e., ‘Alf Clausen (born April 16, 1941) is an American composer of film and television scores.’ To rectify this potential factual mistake, we select nodes from the candidate neighbors that match the most with the LLM-generated passage, in this case, Node 4 ‘Alf Heiberg Clausen (born March 28, 1941) is an American film composer’. Since this passage sources directly from documents, it inherently ensures the validity of the information. Then we prompt LLMs along with the retrieved context Node 1 and 4 for the answer.

Additionally, for questions asking about document structures, we extract the document structure names and locate their corresponding structural nodes in the KG. For the table node, we retrieve its markdown formatted content while for the page node, we traverse its one-hop neighbor and obtain passages belonging to that page.

Algorithm 1: LLM-based KG Traversal Algorithm to Retrieve Relevant Context for Content-based Question.

Input: A question q over a set of documents \mathcal{D} , the constructed KG $G = \{\mathcal{V}, \mathcal{E}, \mathcal{X}\}$ over \mathcal{D} , the fine-tuned LLM-guided graph traversal f_{GT} , the preset context budget K , the TF-IDF search function g .

```

1 Initialize seed passages  $\mathcal{V}^s = g(\mathcal{V}, \mathcal{X}, q)$ 
2 Initialize the retrieved passage queue  $\mathcal{P} = [\{v_i\} | v_i \in \mathcal{V}^s]$ 
3 Initialize the candidate neighbor queue  $\mathcal{C} = [\mathcal{N}_i | v_i \in \mathcal{V}^s]$ 
4 Initialize the retrieved passage counter  $k = \sum_{\mathcal{P}_i \in \mathcal{P}} |\mathcal{P}_i|$ 
5 while queue  $\mathcal{P}$  and queue  $\mathcal{C}$  are not empty do
6    $\mathcal{P}_i \leftarrow \mathcal{P}.dequeue(), \mathcal{C}_i \leftarrow \mathcal{C}.dequeue()$ 
7    $\mathcal{V}'_i = \text{Graph Traversal}(\{q\} \cup \mathcal{P}_i, \mathcal{C}_i, k)$  by Eq (1)
8   for  $v \in \mathcal{V}'_i$  do
9      $\mathcal{P}.enqueue(\mathcal{P}_i \cup \{v\}), \mathcal{C}.enqueue(\mathcal{N}_v)$ 
10     $k \leftarrow k + 1$ 
11    if  $k > K$  then
12      Terminate
13 return Retrieved Passage Queue  $\mathcal{P}$ 

```

Here we present the algorithm for our proposed KGP method for MD-QA. Given a question, we first apply LLM to classify whether the question is asking about the document structure or content. If the question focuses on the document structure, we extract the structural keywords such as Page or Table, and retrieve the content in the corresponding structural nodes in KG. If the question focuses on the document content, we follow the step according to Algorithm 1. Specifically, we first initialize seeding passages \mathcal{V}^s and the reasoning path queue \mathcal{P} by TF-IDF search. Then for each seeding passage $v_i \in \mathcal{V}^s$, we add its neighboring passage nodes \mathcal{N}_i into the candidate neighbor queue \mathcal{C} (lines 1-4). After that, we iteratively dequeue the earliest enqueued reasoning path/candidate neighborhood $\mathcal{P}_i/\mathcal{C}_i$ from \mathcal{P}/\mathcal{C} and employ the fine-tuned LLM-based graph traversal agent to rank the dequeued neighbors in \mathcal{C}_i by Eq. (1) (lines 5-7). Last, we select top-k passage nodes \mathcal{V}'_i from \mathcal{C}_i to visit next based on their rank and correspondingly update the candidate neighbor queue and reasoning path queue (lines 8-13). The above process terminates when either the candidate neighbor queue becomes empty or the prefixed budget K for the retrieved passages is met. The time and space complexity are thoroughly analyzed in Section 8.3 in Supplementary.

5 Experiment

In this section, we conduct experiments to verify the proposed knowledge graph prompting method (KGP) for MD-QA. In particular, we answer the following questions:

- **Q1 - Section 5.1:** How well does KGP perform MD-QA compared with existing baselines?
- **Q2 - Section 5.2-5.3:** How do the quality of the constructed KG and the LLM-based graph traversal agent impact the MD-QA performance?

Due to the space limitation, we comprehensively introduce our experimental setting, including dataset collection, baselines, and evaluation criteria, in Supplementary 8.1-8.2.

Table 1: MD-QA Performance (%) of different baselines. The best and runner-up are in **bold** and underlined. None: no passages but only the question is provided. Golden: supporting facts are provided along with the question.

Method	HotpotQA			IIRC			2WikiMQA			MuSiQue			PDFTriage	Rank	
	Acc	EM	F1	Struct-EM	w PDFTriage	w/o PDFTriage									
None	41.80	19.00	30.50	19.50	8.60	13.17	44.40	18.60	25.07	30.40	4.60	10.58	0.00	8.53	9.00
KNN	71.57	40.73	57.97	43.82	25.15	37.24	52.40	31.20	42.13	44.70	18.86	30.04	-	7.00	7.33
TF-IDF	76.64	<u>45.97</u>	64.64	47.47	27.22	40.80	58.40	34.60	44.50	44.40	21.59	32.50	-	4.85	5.00
BM25	71.95	41.46	59.73	41.93	23.48	35.55	55.80	30.80	40.55	44.47	21.11	31.15	-	6.92	7.25
DPR	73.43	43.61	62.11	48.11	26.89	<u>41.85</u>	62.40	35.60	51.10	44.27	20.32	31.64	-	5.31	5.50
MDR	75.30	45.55	<u>65.16</u>	50.84	<u>27.52</u>	43.47	<u>63.00</u>	36.00	<u>52.44</u>	<u>48.39</u>	<u>23.49</u>	<u>37.03</u>	-	<u>3.07</u>	<u>3.08</u>
IRCoT	74.36	45.29	64.12	<u>49.78</u>	27.73	41.65	61.81	<u>37.75</u>	50.17	45.14	22.46	34.21	-	4.00	4.08
KGP-T5	<u>76.53</u>	46.51	66.77	48.28	26.94	41.54	63.50	39.80	53.50	50.92	27.90	41.19	67.00	2.69	2.75
Golden	82.19	50.20	71.06	62.68	35.64	54.76	72.60	40.20	59.69	57.00	30.60	47.75	100.00	1.00	1.00

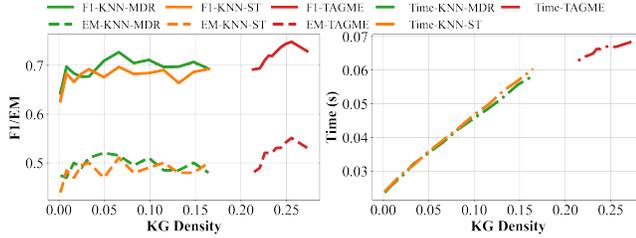


Figure 6: The performance/latency increases as the KG density increases. The results are averaged across 100 randomly sampled questions on HotpotQA.

5.1 Performance Comparison on MD-QA

We compare the MD-QA performance of the proposed KGP-T5 and other baselines in Table 1. Firstly, the baselines ‘None/Golden’ achieve the worst/best performance because one provides no context and the other provides the golden context. All other baselines achieve the performance in-between because the retrieved context only covers the partial of the golden supporting facts. Our proposed methods KGP-T5 rank at the Top-1 except for the Golden baseline. The 2nd-performing baseline MDR fine-tunes a RoBERTa-base encoder by predicting the next supporting fact based on the question and the already retrieved contexts (Xiong et al. 2020). This next-passage prediction pretext task equips the model with the reasoning capability of the knowledge across different passages and hence increases the quality of the retrieved contexts. The other deep-learning-based retriever DPR achieves much worse performance than MDR because it only fine-tunes the encoder by maximizing the similarity between the query and its supporting facts regardless of their sequential order, demonstrating the importance of understanding the logical order of different knowledge when solving MD-QA (Xiong et al. 2020). By comparing the MD-QA performance across different datasets, we find that all baselines perform better on HotpotQA than on IIRC. This is because questions in HotpotQA are generally simpler than in IIRC. Existing works (Jiang and Bansal 2019) have shown that some questions in HotpotQA can be easily answered following shortcuts while questions in IIRC sometimes necessitate arithmetic skills to derive answers numerically, e.g., ‘How many years did the event last when Wingfield lost his fortune?’, which poses unique difficulty due to LLMs’ inferior arithmetic capability (Yuan et al. 2023).

Table 2: Comparing different LLM-based KG Traversal Agents, including off-the-shelf ChatGPT equipped with few-shot demonstration with fine-tuned LLaMA/T5/MDR on TAGME-constructed KG.

Traversal Agent	HotpotQA			IIRC			2WikiMQA			MuSiQue		
	Acc	EM	F1									
TF-IDF	73.52	43.79	63.14	46.30	27.70	41.43	58.12	35.07	45.95	44.67	21.93	32.90
MDR	75.72	46.09	65.77	49.58	29.32	43.21	60.94	37.22	51.29	51.22	27.76	41.11
ChatGPT	77.80	46.03	66.57	46.27	26.01	39.35	61.62	36.16	49.39	50.61	26.92	38.66
LLaMA	75.66	<u>46.22</u>	66.31	<u>49.57</u>	<u>28.09</u>	<u>42.56</u>	62.45	<u>37.55</u>	<u>52.45</u>	50.81	26.72	40.01
T5	<u>76.53</u>	46.51	66.77	48.28	26.94	41.54	63.50	39.80	53.50	<u>50.92</u>	27.90	41.19

Moreover, without any particular design for document structures, no existing baselines can handle structural questions in PDFTriage, e.g. ‘What is the difference between Page 1 and Page 2’ or ‘In Table 3, which station has the highest average flow rate?’. Fortunately, with the constructed KG incorporating structural nodes and our designed traversal algorithm retrieving structural contexts, our proposed method achieves 67% Struct-EM.

5.2 Impact of the Constructed Graph

We construct KGs with varying densities by varying the hyperparameters of TF-IDF/KNN-ST/KNN-MDR/TAGME, and studying its impact on the performance and the neighbor matching time of MD-QA using KGP-T5. Since the LLM-based graph traversal agent selects the next node to visit from neighbors of already visited nodes, the chance that it hits the supporting facts increases as neighbors increase. In contrast, the neighbor matching efficiency decreases as the candidate pool, i.e., \mathcal{N}_j in Eq (1), increases. As evidenced in Figure 6, we observe a similar trend, i.e., as KG density increases, the F1/EM increases and stays stable while the latency for selecting the most promising neighbors to visit next also increases. KNN-MDR achieves better performance than KNN-ST when the density of the two constructed KGs is the same. This is because the encoder in KNN-ST is pre-trained on wide-spectrum datasets while the encoder in MDR is specifically pre-trained on the HotpotQA by the pretext task of predicting the next supporting facts. Therefore, the embedding similarity and the corresponding neighbor relations better reflect the logical associations among different passages, which aligns with the better constructed KG by KNN-MDR than KG by KNN-ST in Figure 5. Compared with KNN-MDR/ST, TAGME delivers superior performance at the cost of increasing latency since the generated KG by TAGME is denser than KGs by KNN-ST/MDR.

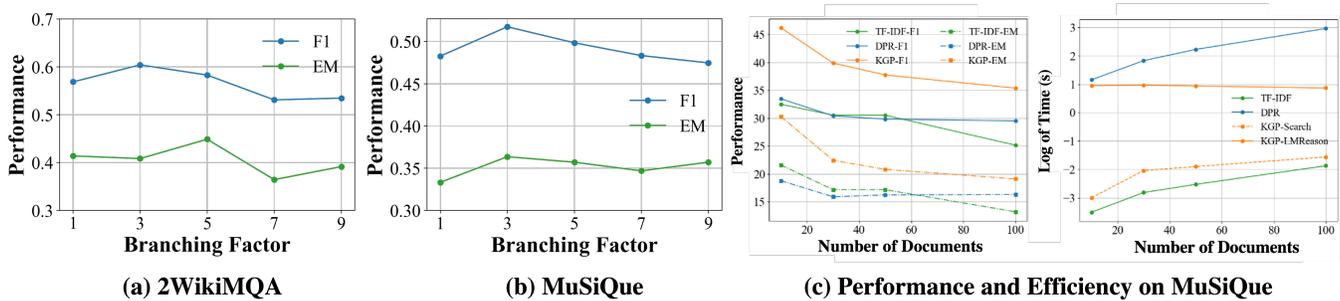


Figure 7: (a)-(b): Performance first increases and then decreases as the branching factor increases. The results are averaged across 100 sampled questions on 2WikiMQA and MuSiQue. (c): Performance/Efficiency increases/decreases as the number of documents increases on MuSiQue. KGP-T5 achieves higher performance/efficiency than DPR.

5.3 Impact of Graph Traversal Agent

Here we study the influence of using different LLM agents to traverse over TAGME-constructed KG on MD-QA. Specifically, we compare agents that select the next neighbor to visit randomly or intelligently via guidance from ChatGPT, LLaMA, T5, and MDR in Table 2. Because the random agent only blindly traverses the KG without any guidance from LLM, it unavoidably collects irrelevant passages and hence achieves the worst performance than others under LLMs’ guidance. This aligns with our previous observation on the low precision in Figure 5 and further demonstrates the necessity of using LLMs to guide the graph traversal. Interestingly, we find that KGP-T5 performs better than LLaMA even though the parameters of LLaMA-7B are more than the ones with T5-0.7B. We hypothesize this is because LLaMA-7B requires more data to fine-tune than T5-0.7B.

5.4 Sensitivity Analysis

Here we perform the sensitivity analysis of the branching factor (the number of nodes selected from candidate neighbors to visit next). In Figure 7(a)-(b), the performance first increases as the branching factor increases because more passage nodes selected from the candidate neighbors lead to more reasoning paths to reach the final answer. However, as we fix the context budget to ensure fair comparison (i.e., the total number of passages we are allowed to retrieve for each question is the same across all baselines), the performance declines as the branching factor increases because the number of initial seeding nodes diminishes, leading to reduced coverage of the KG. Furthermore, we compare the efficiency of KGP when the constructed KG includes different numbers of documents in Figure 7(c). KGP consistently achieves higher performance than other baselines and higher efficiency than embedding-based DPR. TF-IDF is slightly faster than KGP because it is a purely heuristic-based method.

6 Related Work

Question answering Question Answering (QA) aims to provide answers to users’ questions in natural language (Zhu et al. 2021; Pandya and Bhatt 2021), and most QA systems are composed of information retrieval (IR) and answer extraction (AE) (Mao et al. 2021; Ju et al. 2022; Liu and Qin 2022). In IR, the system searches for query-relevant factual passages using heuristic methods

(BM25) (Robertson, Zaragoza et al. 2009) or neural-ranking ones (DPR) (Karpukhin et al. 2020). In AE, the final answer is extracted usually as a textual span from related passages. Although this framework has been broadly applied in O-QA (Mao et al. 2021) and D-QA (Xu et al. 2020; Mathew, Karatzas, and Jawahar 2021), no previous work focus on MD-QA, which demands alternatively reasoning and retrieving knowledge from multiple documents. To tackle this issue, we construct KGs to encode logical associations among different passages across documents and design an LLM-based graph traversal agent to alternatively generate the reason and visit the most matching passage node.

Pre-train, Prompt, and Predict with LLMs With the emergence of LLMs, the paradigm of ‘pre-train, prompt, predict’ has gained magnificent popularity in handling a wide spectrum of tasks (Gururangan et al. 2020; Liu et al. 2023; Yu et al. 2023). This approach begins with pre-training LLMs by pretext tasks to encode world knowledge into tremendous parameters (Wu et al. 2023) followed by a prompting function to extract pertinent knowledge for downstream tasks (Yang et al. 2023). Recent advancements explore different prompting strategies to enhance LLMs’ reasoning capabilities (Wei et al. 2022; Jin et al. 2023). In contrast to that, our work offers a novel perspective by transforming the prompt formulation into the KG traversal.

7 Conclusion

Answering multi-document questions demands knowledge reasoning and retrieving from different documents across various modalities, presenting challenges for applying the paradigm of ‘pre-train, prompt and predict’ with LLMs. Recognizing the logical associations among passages and structural relations within documents, we propose a Knowledge Graph Prompting method (KGP) for aiding LLMs in MD-QA. The KGP constructs KGs from documents with nodes as sentences or document structures, and edges as their lexical/semantic similarity/structural relations. Since constructed KGs may contain irrelevant neighbor information, we further design an LLM-based graph traversal agent that selectively visits the most promising node in approaching the question. In the future, we plan to investigate the capability of LLMs in understanding graph topology and explore the potential of fine-tuning/prompting LLMs to encode complex topological signals hidden in the graph.

Ethics Statement

Due to page limitation, the supplementary material and reproducing details are publically available at https://github.com/YuWVandy/KG-LLM-MDQA/blob/main/AAAI24_KGPrompting_Supplementary.pdf.

Acknowledgements

This research is supported by Adobe Research and the National Science Foundation (NSF) under grant number IIS2239881.

References

- Aly, R.; Guo, Z.; Schlichtkrull, M.; Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; Cocarascu, O.; and Mittal, A. 2021. FEVEROUS: Fact extraction and VERification over unstructured and structured information. *arXiv preprint arXiv:2106.05707*.
- Asai, A.; Hashimoto, K.; Hajishirzi, H.; Socher, R.; and Xiong, C. 2019. Learning to retrieve reasoning paths over wikipedia graph for question answering. *arXiv preprint arXiv:1911.10470*.
- Bolino, M.; Long, D.; and Turnley, W. 2016. Impression management in organizations: Critical questions, answers, and areas for future research. *Annual Review of Organizational Psychology and Organizational Behavior*, 3: 377–406.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.
- Caciularu, A.; Peters, M. E.; Goldberger, J.; Dagan, I.; and Cohan, A. 2023. Peek Across: Improving Multi-Document Modeling via Cross-Document Question-Answering. *arXiv preprint arXiv:2305.15387*.
- Chen, D.; Fisch, A.; Weston, J.; and Bordes, A. 2017. Reading wikipedia to answer open-domain questions. *arXiv preprint arXiv:1704.00051*.
- Chung, H. W.; Hou, L.; Longpre, S.; Zoph, B.; Tay, Y.; Fedus, W.; Li, E.; Wang, X.; Dehghani, M.; Brahma, S.; et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Dong, J.; Zhang, Q.; Huang, X.; Duan, K.; Tan, Q.; and Jiang, Z. 2023. Hierarchy-Aware Multi-Hop Question Answering over Knowledge Graphs. In *Proceedings of the ACM Web Conference 2023*, 2519–2527.
- Ferragina, P.; and Scaiella, U. 2010. Tagme: on-the-fly annotation of short text fragments (by wikipedia entities). In *Proceedings of the 19th ACM international conference on Information and knowledge management*, 1625–1628.
- Gionis, A.; Indyk, P.; Motwani, R.; et al. 1999. Similarity search in high dimensions via hashing. In *Vldb*, volume 99, 518–529.
- Gururangan, S.; Marasović, A.; Swayamdipta, S.; Lo, K.; Beltagy, I.; Downey, D.; and Smith, N. A. 2020. Don’t stop pretraining: Adapt language models to domains and tasks. *arXiv preprint arXiv:2004.10964*.
- Ho, X.; Nguyen, A.-K. D.; Sugawara, S.; and Aizawa, A. 2020. Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*.
- Hoffart, J.; Suchanek, F. M.; Berberich, K.; and Weikum, G. 2013. YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia. *Artificial intelligence*, 194: 28–61.
- Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Huang, X.; Zhang, J.; Xu, Z.; Ou, L.; and Tong, J. 2021. A knowledge graph based question answering method for medical domain. *PeerJ Computer Science*, 7: e667.
- Ji, Z.; Lee, N.; Frieske, R.; Yu, T.; Su, D.; Xu, Y.; Ishii, E.; Bang, Y. J.; Madotto, A.; and Fung, P. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12): 1–38.
- Jiang, Y.; and Bansal, M. 2019. Avoiding reasoning shortcuts: Adversarial evaluation, training, and model development for multi-hop QA. *arXiv preprint arXiv:1906.07132*.
- Jin, B.; Liu, G.; Han, C.; Jiang, M.; Ji, H.; and Han, J. 2023. Large Language Models on Graphs: A Comprehensive Survey. *arXiv preprint arXiv:2312.02783*.
- Ju, M.; Yu, W.; Zhao, T.; Zhang, C.; and Ye, Y. 2022. Grape: Knowledge graph enhanced passage reader for open-domain question answering. *arXiv preprint arXiv:2210.02933*.
- Karpukhin, V.; Oğuz, B.; Min, S.; Lewis, P.; Wu, L.; Edunov, S.; Chen, D.; and Yih, W.-t. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.
- Liu, H.; and Qin, Y. 2022. Heterogeneous graph prompt for community question answering. *Concurrency and Computation: Practice and Experience*, e7156.
- Liu, P.; Yuan, W.; Fu, J.; Jiang, Z.; Hayashi, H.; and Neubig, G. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9): 1–35.
- Liu, X.; Dong, Z.; and Zhang, P. 2023. Tackling Data Bias in MUSIC-AVQA: Crafting a Balanced Dataset for Unbiased Question-Answering. *arXiv preprint arXiv:2310.06238*.
- Mao, Y.; He, P.; Liu, X.; Shen, Y.; Gao, J.; Han, J.; and Chen, W. 2021. RIDER: Reader-Guided Passage Reranking for Open-Domain Question Answering. *arXiv preprint arXiv:2101.00294*.
- Mathew, M.; Karatzas, D.; and Jawahar, C. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2200–2209.
- Min, S.; Chen, D.; Zettlemoyer, L.; and Hajishirzi, H. 2019. Knowledge guided text retrieval and reading for open domain question answering. *arXiv preprint arXiv:1911.03868*.
- Pandya, H. A.; and Bhatt, B. S. 2021. Question answering survey: Directions, challenges, datasets, evaluation matrices. *arXiv preprint arXiv:2112.03572*.

- Pereira, J.; Fidalgo, R.; Lotufo, R.; and Nogueira, R. 2023. Visconde: Multi-document QA with GPT-3 and Neural Reranking. In *European Conference on Information Retrieval*, 534–543. Springer.
- Qin, C.; Zhang, A.; Zhang, Z.; Chen, J.; Yasunaga, M.; and Yang, D. 2023. Is ChatGPT a general-purpose natural language processing task solver? *arXiv preprint arXiv:2302.06476*.
- Qu, A.; Wang, Y.; Hu, Y.; Wang, Y.; and Baroud, H. 2020. A data-integration analysis on road emissions and traffic patterns. In *Driving Scientific and Engineering Discoveries Through the Convergence of HPC, Big Data and AI: 17th Smoky Mountains Computational Sciences and Engineering Conference, SMC 2020, Oak Ridge, TN, USA, August 26-28, 2020, Revised Selected Papers 17*, 503–517. Springer.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1): 5485–5551.
- Ramos, J.; et al. 2003. Using tf-idf to determine word relevance in document queries. In *Proceedings of the first instructional conference on machine learning*, volume 242, 29–48. Citeseer.
- Robertson, S.; Zaragoza, H.; et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4): 333–389.
- Saad-Falcon, J.; Barrow, J.; Siu, A.; Nenkova, A.; Rossi, R. A.; and Dernoncourt, F. 2023. PDFTriage: Question Answering over Long, Structured Documents. *arXiv preprint arXiv:2309.08872*.
- Tessuto, G. 2011. Legal problem question answer genre across jurisdictions and cultures. *English for Specific Purposes*, 30(4): 298–309.
- Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; and Mittal, A. 2018. FEVER: a large-scale dataset for fact extraction and VERification. *arXiv preprint arXiv:1803.05355*.
- Tian, Y.; Dong, K.; Zhang, C.; Zhang, C.; and Chawla, N. V. 2023a. Heterogeneous graph masked autoencoders. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 9997–10005.
- Tian, Y.; Song, H.; Wang, Z.; Wang, H.; Hu, Z.; Wang, F.; Chawla, N. V.; and Xu, P. 2023b. Graph neural prompting with large language models. *arXiv preprint arXiv:2309.15427*.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Trivedi, H.; Balasubramanian, N.; Khot, T.; and Sabharwal, A. 2022a. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509*.
- Trivedi, H.; Balasubramanian, N.; Khot, T.; and Sabharwal, A. 2022b. MuSiQue: Multihop Questions via Single-hop Question Composition. *Transactions of the Association for Computational Linguistics*, 10: 539–554.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35: 24824–24837.
- Wu, X.; Zhou, K.; Sun, M.; Wang, X.; and Liu, N. 2023. A survey of graph prompting methods: techniques, applications, and challenges. *arXiv preprint arXiv:2303.07275*.
- Xiong, W.; Li, X. L.; Iyer, S.; Du, J.; Lewis, P.; Wang, W. Y.; Mehdad, Y.; Yih, W.-t.; Riedel, S.; Kiela, D.; et al. 2020. Answering complex open-domain questions with multi-hop dense retrieval. *arXiv preprint arXiv:2009.12756*.
- Xu, Y.; Xu, Y.; Lv, T.; Cui, L.; Wei, F.; Wang, G.; Lu, Y.; Florencio, D.; Zhang, C.; Che, W.; et al. 2020. Layoutlmv2: Multi-modal pre-training for visually-rich document understanding. *arXiv preprint arXiv:2012.14740*.
- Yang, J.; Jin, H.; Tang, R.; Han, X.; Feng, Q.; Jiang, H.; Yin, B.; and Hu, X. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *arXiv preprint arXiv:2304.13712*.
- Yao, S.; Yu, D.; Zhao, J.; Shafraan, I.; Griffiths, T. L.; Cao, Y.; and Narasimhan, K. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*.
- Yao, Y.; Li, Z.; and Zhao, H. 2023. Beyond Chain-of-Thought, Effective Graph-of-Thought Reasoning in Large Language Models. *arXiv preprint arXiv:2305.16582*.
- Yasunaga, M.; Bosselut, A.; Ren, H.; Zhang, X.; Manning, C. D.; Liang, P. S.; and Leskovec, J. 2022. Deep bidirectional language-knowledge graph pretraining. *Advances in Neural Information Processing Systems*, 35: 37309–37323.
- Yasunaga, M.; Ren, H.; Bosselut, A.; Liang, P.; and Leskovec, J. 2021. QA-GNN: Reasoning with language models and knowledge graphs for question answering. *arXiv preprint arXiv:2104.06378*.
- Yavuz, S.; Hashimoto, K.; Zhou, Y.; Keskar, N. S.; and Xiong, C. 2022. Modeling multi-hop question answering as single sequence prediction. *arXiv preprint arXiv:2205.09226*.
- Yu, Y.; Zhuang, Y.; Zhang, J.; Meng, Y.; Ratner, A.; Krishna, R.; Shen, J.; and Zhang, C. 2023. Large Language Model as Attributed Training Data Generator: A Tale of Diversity and Bias. *arXiv preprint arXiv:2306.15895*.
- Yuan, Z.; Yuan, H.; Tan, C.; Wang, W.; and Huang, S. 2023. How well do Large Language Models perform in Arithmetic tasks? *arXiv preprint arXiv:2304.02015*.
- Zhu, F.; Lei, W.; Wang, C.; Zheng, J.; Poria, S.; and Chua, T.-S. 2021. Retrieving and reading: A comprehensive survey on open-domain question answering. *arXiv preprint arXiv:2101.00774*.
- Zou, H.; and Caragea, C. 2023. JointMatch: A Unified Approach for Diverse and Collaborative Pseudo-Labeling to Semi-Supervised Text Classification. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 7290–7301.

8 Supplementary

8.1 Dataset Collection

This section introduces the collection of datasets used for the experiments conducted in this paper.

Document Set Collection and Precession As no previous works focus on MD-QA, we create our own datasets to simulate real-world scenarios where users maintain folders containing various documents and pose questions to which the answers are only from certain parts of these documents. To imitate this scenario, we randomly sample questions from the development set of existing datasets: HotpotQA/IIRC/2WikiMQA/MuSiQue, and then for each specific question, we fetch documents from Wikipedia that encompass supporting facts pertaining to the question⁴ and term these documents as golden documents. Then we randomly sample negative documents from Wikipedia and pair them with golden documents to constitute the document collection. For each document in the collected document set, we split it into multiple passages with the default passage length being 250 as it empirically yields superior performance. As questions from these existing datasets only focus on document contents, we additionally incorporate the ‘PDFTriage’ dataset, an internal company collection of real-world questions focusing on document structures. We refer readers to the paper (Saad-Falcon et al. 2023) for more details.

Knowledge Graph Construction We construct a knowledge graph for each question and its corresponding collection of documents. For datasets where the questions are from Wikipedia: HotpotQA, IIRC, WikiMHop, and Musique, we only have passage nodes since answering questions in these datasets does not require information about document structures. For the PDFTriage dataset, in addition to passage nodes, we apply ExtractAPI to obtain the page and table information so that the constructed KG also has pages/tables as nodes. For all of these datasets, we add edges following Section 3. Table 3 summarizes the average statistics of the document collections across all questions with their corresponding KGs. The code for the dataset collection and preprocessing is publically available at <https://github.com/YuWVandy/KG-LLM-MDQA>.

Table 3: Statistics of document collections and their corresponding knowledge graph used in Table 1 and 2 average across all questions.

Dataset	#Docs	#Questions	#Passages	#Edges	Passage Avg. Length	KG Density
HotpotQA	12	500	715.22	70420.68	37.55	0.23
IIRC	12	477	1120.55	143136.17	37.24	0.20
WikiMHop	12	500	294.19	19235.15	37.24	0.27
MuSiQue	12	500	748.04	97931.28	38.56	0.29

More details about the PDFTriage dataset can be found at PDFTriage (Saad-Falcon et al. 2023).

⁴The HotpotQA/IIRC/2WikiMQA/Musique datasets already have the supporting facts for each question.

Sequential Data Collection Training MDR (Xiong et al. 2020) requires rearranging supporting facts into the sequential order that progressively approaches the answer. To fulfill this requirement, we directly follow MDR and use the pre-processed HotpotQA data from the GitHub Repository⁵ to train the encoder and apply it to other datasets that do not provide the sequential order of supporting facts. For instruction fine-tuning LLaMA, we still use the above HotpotQA data and rearrange it into the instruction-input-output format and use the instruction ‘What evidence do we need to answer the question given the current evidence’. We present one example in Listing 1. For T5-large, we use the same input-output but prefix the reasoning instruction to the input following the original T5 input format (Raffel et al. 2020).

8.2 Experiment Details

Training DPR and MDR For training DPR (Karpukhin et al. 2020), we pair each question with its supporting facts as its positive passages, and some randomly sampled passages as its negative passages. For training MDR (Xiong et al. 2020), as each question in HotpotQA only requires 2 supporting facts to derive the answer, we set the first supporting fact as the positive pair for each question. Further, we concatenate this question and the first supporting fact to form a new question and for this newly-formed question, we set the second supporting fact as its positive pair. For both the original question and the concatenated one, we randomly sample other passages as the negative pair. Following (Xiong et al. 2020; Karpukhin et al. 2020), we use RoBERTa-base as the default encoder. The search space of hyperparameters is summarized in Table 4.

Table 4: Hyperparameters used for tuning DPR and MDR. The value of most of them are directly taken from their original GitHub Repository.

Hyperparameter	Search Space
Encoder	RoBERTa-base
Hidden Dimension	768
Max Context Length	{128, 256, 350}
Batch Size	{128, 256, 512}
Epoch	50
Warmup Steps	300
Learning Rate	2e-5
Gradient Clipping Range	2

Instruction Fine-tuning LLaMA⁶ and T5-Large⁷ We fine-tune LLaMA using instruction data in Listing 1. Due to the computational limitation, we choose LLaMA-7B and use LoRA (Hu et al. 2021). For fine-tuning T5-Large, we use the same instruction data except that we remove the instruction but only prefix the reasoning instruction to the input (Raffel et al. 2020). We use the default hyperparameters from their original GitHub repository to fine-tune these two LLMs.

⁵https://github.com/facebookresearch/multihop_dense_retrieval/tree/main

⁶<https://github.com/Lightning-AI/lit-llama>

⁷<https://shivanandroy.com/fine-tune-t5-transformer-with-pytorch/>

Prompting LLMs for MD-QA - Table 1 and 2 Following (Trivedi et al. 2022a), we randomly select questions from the development set for reporting the performance. To ensure a fair comparison, we set the number of retrieved passages to 30 across all baselines and use ChatGPT as the downstream LLM for reading the retrieved passages and generating the answer. We summarize the key implementation details for each baseline as follows:

- **KNN**: We employ the sentence-transformer variant ‘multi-qa-MiniLM-L6-cos-v1’ to obtain passage embeddings as it has been trained on 215M (question, answer) pairs from diverse sources. Then we select the top-15 passages according to the embedding similarity and the top-15 passages according to the fuzzy matching⁸.
- **MDR**: We use beam search with the inner product as the scoring function to rank passages. We limit the search depth to 2 as answering questions in HotpotQA requires at most 2-hop reasoning steps (Xiong et al. 2020). We set the number of passages to be 15 in the first-hop retrieval and for each of these passages, we further retrieve 3 more passages in the second round, which in total generates 45 passage pairs. Then we rank these 45 passage pairs by the product of the scores between the first-hop and the second-hop retrieval and select the top 30 ones as the final context.
- **IRCoT**: Instead of directly employing the original IR-CoT code (Trivedi et al. 2022a), we modify it based on our problem setting. The first reason is that passages to be retrieved in IRCoT (Trivedi et al. 2022a) are the pre-processed Wikipedia Corpus and do not cover the whole contents of Wikipedia documents, which thereby is not aligned with our MD-QA setting. The second reason is that the question-answering reader employed in IRCoT requires running on A100-80G GPU, which is not affordable on our side. Therefore, we modify the IRCoT by replacing the question reader with the ChatGPT and using our pre-processed Wikipedia document collections as introduced in Section 8.1. For the prompt used in the reasoning step, we select 2 examples from ‘gold_with_2_distractors_context’ for the demonstration purpose. We iteratively select top-5 passages based on the generated reason from LLM along with their document titles and add them to the retrieved context until hitting the prefix budget. For the prompt used in the reading step, we use exactly the same prompt as other baselines as we find it empirically leads to better performance than the original one used in IRCoT (Trivedi et al. 2022a).
- **KGP-T5/LLaMA/MDR/ChatGPT**: We use T5-large/LLaMA-7B/MDR/ChatGPT as the LLM to guide the graph traversal respectively. For content-based questions, similar to MDR, we perform a 2-hop retrieval but for each hop, we only search the node to visit next from neighbor candidates. In the 1st-hop retrieval, we select 10 passages and in 2nd-hop retrieval, we select 3 passages, which totally forms 30 reasoning paths. Note that passages in the 1st-hop retrieval are allowed to overlap

⁸We use Levenshtein-distance to measure the lexical distance between two passages.

with the ones in the 2nd-hop retrieval. For structural-based questions, we first use ChatGPT to extract page/table structures and then fetch relevant contents in those structures. Future work could explore how to pre-train a structural extraction model to obtain document structures.

- **KGP-TF-IDF**: We remove the LLM-guided graph traversal but select passage nodes based on their TF-IDF similarity to the given question.

Note that we put the prompt template for running all the above baselines in Section 8.9.

8.3 Complexity Analysis for KGP

Algorithm 2: LLM-based KG Traversal Algorithm to Retrieve Relevant Context for Content-based Question.

Input: A question q over a set of documents \mathcal{D} , the constructed KG $G = \{\mathcal{V}, \mathcal{E}, \mathcal{X}\}$ over \mathcal{D} , the fine-tuned LLM-guided graph traversal f_{GT} , the preset context budget K , the TF-IDF search function g .

- 1 Initialize seed passages $\mathcal{V}^s = g(\mathcal{V}, \mathcal{X}, q)$
- 2 Initialize the retrieved passage queue $\mathcal{P} = [\{v_i\} | v_i \in \mathcal{V}^s]$
- 3 Initialize the candidate neighbor queue $\mathcal{C} = [\mathcal{N}_i | v_i \in \mathcal{V}^s]$
- 4 Initialize the retrieved passage counter $k = \sum_{\mathcal{P}_i \in \mathcal{P}} |\mathcal{P}_i|$
- 5 **while** queue \mathcal{P} and queue \mathcal{C} are not empty **do**
- 6 $\mathcal{P}_i \leftarrow \mathcal{P}.dequeue(), \mathcal{C}_i \leftarrow \mathcal{C}.dequeue()$
- 7 $\mathcal{V}'_i = \text{Graph Traversal}(\{q\} \cup \mathcal{P}_i, \mathcal{C}_i, k)$ by Eq (1)
- 8 **for** $v \in \mathcal{V}'_i$ **do**
- 9 $\mathcal{P}.enqueue(\mathcal{P}_i \cup \{v\}), \mathcal{C}.enqueue(\mathcal{N}_v)$
- 10 $k \leftarrow k + 1$
- 11 **if** $k > K$ **then**
- 12 **Terminate**
- 13 **return** Retrieved Passage Queue \mathcal{P}

Since our algorithm can be essentially deemed as the combination of the neighborhood ranking by Eq. (1) and the breadth-first-search. The time complexity would be the multiplication between the time of bread-first-search $\mathcal{O}(|\mathcal{V}| + |\mathcal{E}|)$ and the time of neighborhood ranking $\mathcal{O}(|\mathcal{N}|\gamma) = \mathcal{O}(\hat{d}\gamma)$ where γ is the time for computing the embedding similarity between a specific neighbor passage and the retrieved reasoning path and \hat{d} is the average degree of the KG. Therefore the final time complexity would be $\mathcal{O}((|\mathcal{V}| + |\mathcal{E}|)\hat{d}\gamma)$, which is in-between the linear and quadratic to the size of the graph. As users typically maintain 10-100 documents, correspondingly the number of nodes in the constructed KG would be around 1,000-10,000 (according to Table 3, a collection of 12 documents have roughly 200-1000 passage nodes), which is affordable even with the quadratic time complexity. Moreover, we can apply advanced techniques to further reduce the time complexity for neighborhood ranking, such as LSH (Gionis et al. 1999) and KD-tree (Qu et al. 2020).

In addition, whenever there are some changes over the document set (e.g., the user adds a new document into the folder or removes an existing document), we can remove/add all sentence nodes from/to the graph. To guarantee the linear time complexity for removing sentences from one

document, we need to maintain a pointer from the document to its sentence nodes. For adding sentence nodes of one document, we need to first apply the KG construction method to compute the lexical/semantic similarity between each of the newly added sentence nodes and the existing nodes in KG, and then add corresponding edges connecting them, which is also linear to the size of the current graph.

For space complexity, it takes $\mathcal{O}(|\mathcal{V}|(\alpha + \beta))$ to maintain the constructed KG on the fly where α is the average space for saving the passage embedding vector while β is the average space for saving the textual information of that passage. Although our constructed KG treats passages as nodes, which cannot scale very well when the graph is extremely large, the total number of documents a user maintains in a folder is typically around 10-100, which is still affordable.

8.4 Markdown-Formatted Table

Figure 8 demonstrates that by sending Tables in the markdown format, ChatGPT can successfully understand their content and perform information retrieval based on the given questions. However, we do observe that such a markdown-formatted solution is not feasible for the long table due to the input token limitation of ChatGPT, we plan to explore the solution using SQL as the prompt content or modeling the Table as the grid graph to solve the issue in the future.

Here is a table, I will give you some questions based on it. What is the first row talking about?

Item	In Stock	Price
Python Hat.	True	\$23.99
SQL Hat	True	\$23.99
Codecademy Tee	False	\$19.99
Codecademy Hoodie	False	\$42.99

The first row in the table is talking about the "Python Hat." It states that the Python Hat is currently in stock (In Stock: True) and its price is \$23.99.

What is the price of Codecademy Tee?

The price of the Codecademy Tee is \$19.99, as indicated in the table.

Which one has higher price, Python Hat or Codecademy Hoodies?

The Python Hat has a price of \$23.99, and the Codecademy Hoodie has a price of \$42.99. Therefore, the Codecademy Hoodie has a higher price than the Python Hat.

Figure 8: An example demonstrating that ChatGPT can understand table in the markdown format.

8.5 Knowledge Graph Construction Comparison

Table 5 compares different knowledge graph construction methods and their pros and cons.

- **TAGME**: TAGME (Ferragina and Scaiella 2010) is very effective in extracting Wikipedia Entities from a passage despite the low efficiency. In our graph construction, it usually takes more than 8 hours to extract entities of all passages for even just 12 Wikipedia documents. Even after we apply parallel processing, it still takes more than 2

hours. In addition, it can only handle entities mentioned in the existing Wikipedia system and hence cannot generalize to documents from other domains.

- **TF-IDF and KNN-ST**: Although there is no domain limitation, it is hard to guarantee the extracted keywords or the embedding semantic similarity can precisely encode the relationships that are desired for answering the given question between any two passages. We empirically find TF-IDF is more likely to extract meaningless keywords even after removing supporting verbs and articles.
- **KNN-MDR**: Since KNN-MDR pre-trains the sentence encoder by predicting the next supporting passage given already-retrieved passages, the embedding similarity between two passages is more likely to encode necessary logical associations required for MD-QA. However, the main bottleneck here is how to obtain the logically ordered supporting facts that can progressively reach the answer. Obtaining these sequential data is non-trivial and usually requires a large number of human resources for well-curated annotation.
- **Existing Knowledge Base**: One common approach in the literature is to use existing knowledge bases or extract subgraphs from them for specific tasks (Yasunaga et al. 2022; Dong et al. 2023; Yasunaga et al. 2021). Because the factual information is characterized as a triplet consisting of two entity nodes and their relationship, it is very powerful in encoding factual information/commonsense knowledge and also avoids the scalability issue (since two different passages might share the same entity). Despite its potency and ease of use, constructing this type of KGs demands meticulously designed relation extractors, which is still deemed a challenging task in the literature. Recent research has explored using LLMs for relation extraction. However, with increasing document numbers, using non-open-sourced LLMs can become prohibitively expensive. A potential solution is fine-tuning an open-sourced LLM specifically for relation extraction. Detailed discussion on this is beyond the scope of this study and is thus omitted.

To put it in a nutshell, there's no one-size-fits-all method for KG construction. Our paper offers an in-depth analysis of the proposed KG construction methods alongside other existing ones. The best approach often depends on the specific use case. For broad domains containing general factual information, tools like 'TAGME' or 'Knowledge Base' might be apt. However, for more niche or sensitive areas, methods like TF-IDF/KNN-ST are more appropriate. In certain situations, gathering domain-specific data and pre-training encoders is the most effective way to build the KG.

Table 5: Systematically Comparison among existing and our proposed Knowledge Graphs.

KG	Node	Edge	Domain	Constructor	Scalability	Hyperparameters	Advantage	Disadvantage
TAGME	Passage	Common Wikipedia Entity	Wikipedia	/	No	Prior Threshold	Effectively Identify Wikipedia Entities	Low efficiency for Entity Identification Narrow Domain Application
TF-IDF	Passage	Common Keyword	General	/	No	# Keywords	No Domain Limitation	Common keywords irrelevant to question
KNN-ST	Passage	Semantic Similarity	General	Sentence Transformer	No	# Neighbors	No Domain Limitation	Semantic Similarity irrelevant to question
KNN-MDR	Passage	Semantic Similarity	General	MDR	No	# Neighbors	Encoding the logical association for QA	Require logically ordered supporting facts to pre-train the model
Knowledge Base	Entity	Relationship	Specific	Human	Yes	/	Powerful in encoding factual information	Relation Extraction is non-trivial Domain Specific

8.6 Additional Results and Discussions

Quality of KG on MuSiQue Similar to the setting used for Figure 5, we change the hyperparameters to construct KGs for each question in MuSiQue with varying levels of sparsity and measure how much percentage of the supporting facts are covered by neighbors of the seeding passages that are initially retrieved by TF-IDF. The general trend is similar to the one in Figure 5, i.e., as the graph becomes denser, the precision decreases while the SF-EM increases. However, on MuSiQue, KNN-MDR achieves the worst trade-off between Precision and SF-EM compared with KNN-ST and TF-IDF. This is because our KNN-MDR is pre-trained on HotpotQA and due to the distribution shift from HotpotQA to MuSiQue, it is expected for the graph constructed with KNN-MDR to have less quality. Note that although here KNN-ST leads to a better KG than KNN-MDR, it does not mean the KNN baseline in Table 1 should perform better than MDR because the baseline name only refers to the retrieval method while the name in this figure refers to the KG construction method.

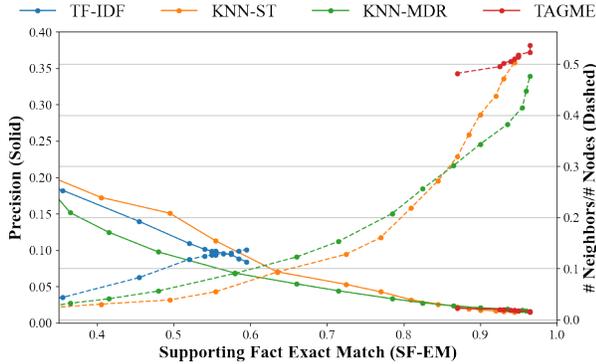


Figure 9: Quality of constructed KGs with different methods on MuSiQue. **TF-IDF**: lexical similarity based on common keywords extracted by TF-IDF. **KNN-ST**: KNN graph constructed based semantic similarity of embeddings from sentence-transformer; **KNN-MDR**: KNN graph constructed based on semantic similarity of embeddings from the pre-trained MDR (Xiong et al. 2020); **TAGME**: graph constructed based on whether two passages share common Wikipedia entity mentions

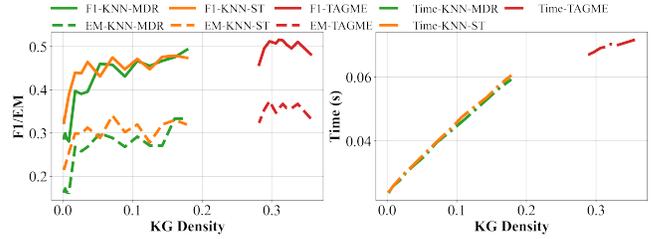


Figure 10: The performance/latency increases as the KG density increases. The results are averaged across 100 randomly sampled questions on MuSiQue.

The impact of KG on MuSiQue Similar to the setting used for Figure 6, we compare the MD-QA performance for KGP-T5 using TAGME-based KG with different levels of density. Similar to Figure 6, here we also observe that as the KG becomes denser, the MD-QA performance increases while the time for the next node search increases. However, on MuSiQue, in most cases, KNN-ST achieves better F1/EM than KNN-MDR, which exactly aligns with the constructed KG quality observed in Figure 9, i.e., KNN-ST achieves better Precision/SF-EM trade-off than KNN-MDR on MuSiQue.

8.7 Case study on Structural/Content Questions

In this section, we conduct six MD-QA case studies using our self-designed user interface coupled with the proposed method on the backend. Examples include two table-based QA (Figure 11-12), one page-based QA (Figure 13), one single-document content-based QA (Figure 14) and two multi-document content-based QA (Figure 15-18). In our designed interface, we can upload documents we are interested in reading and the model on the backend will split each of them into multiple passages. In addition, on the left side, we can ask questions related to the currently uploaded documents. By clicking the button ‘SUBMIT’, the question would be sent to the model on the backend and it retrieves relevant context and arranges them as the prompt to get the answer from ChatGPT. In the figures below, we can see our system can understand the Table/Page questions and also questions requiring knowledge across multiple documents.

Document question

Based on Table on Page 4, how many people belong to the membership grade 'Fellow'?

Answer

There are 24 people who belong to the membership grade 'Fellow'.

Evidence

Page 4/ Table 1:

Membership Grade	Number of People	Section Annual Dues
Fellow	24	\$25
Member	153	\$25
Institute Affiliate	8	\$25
Student Member	47	\$0*
Esteemed Colleague	6	\$0
TOTAL	238	

SUBMIT

Oregon Section 2007/2008 Annual Report

DR--182866691

2. MEMBERSHIP

The following table provides a summary of the Oregon Section's membership, including grade and section dues.

Membership Grade	Number of People	Section Annual Dues
Fellow	24	\$25
Member	153	\$25
Institute Affiliate	8	\$25
Student Member	47	\$0*
Esteemed Colleague	6	\$0
TOTAL	238	

*Note: Students also receive a substantial discount on meeting registration fees.

Membership Committee:

The membership committee continued its organizational captain membership program over the past year. The program was designed to identify specific individuals within public and private organizations who might be interested in joining ITE or attending ITE functions. The captains also assist with advertising meetings, recruiting new members, and promoting ITE.

Institute of Transportation Engineers Page 4 of 13

UPLOAD PDF

Figure 11: Table QA asking for the number of people belonging to the membership grade 'Fellow'. It is shown that ChatGPT can understand table structure in the format of markdown and successfully fetch the number of people belonging to membership 'Fellow'.

Document question

According to Table on Page 7, where will the event occur on Date 5-18-07?

Answer

The event on Date 5-18-07 will occur at the Crowne Plaza Hotel in Lake Oswego, OR.

Evidence

Page 7/ Table 1:

Date	Meeting	Attendees	Subject
4-10-07	April - Joint meeting with WTS	60	Kris Strickler and Danielle Cogan gave a presentation on the Columbia River Crossing project on the I-5 bridge that connects Oregon and Washington.
5-18-07	May - Joint meeting with Oregon Traffic Control Devices Committee (OTCDC)	45	Sections newly elected officers were announced and a presentation on the Highway 25 emergency reconstruction efforts that followed a large storm on Mt. Hood early 2007.
7-13 to 18-07	2007 ITE Western District (District 6) Annual Meeting & Quad Conference	200	Oregon ITE section hosted a very successful joint 2007 ITE District 6 Annual Meeting & Quad Conference in Portland. This year's annual meeting broke many previous records with over 500 meeting participants at the meeting.
9-21-07	2007 Golf Tournament	57	Due to timing of ITE District 6 conference in Portland, this year's golf tournament was held in late summer. We had good turn out and good corporate sponsorship donations.
9-25-07	September Meeting	56	Jim Peters (DKS Associates) and Jay McCoy (City of Gresham) gave a presentation on the use of recent SCA's Adaptive Traffic Signal System in City of Gresham, Oregon.
10-23-07	October Meeting	73	Senator Rick Meninger, Chairman of the Senate Transportation Committee, discussed about congestion pricing and his "Moving Oregon" statewide transportation tour, and effort to build support for major investment in Oregon's transportation system.
11-15-07	2007 Student Traffic Bowl	123 including 52 students	Oregon ITE 10 th Annual Student Traffic Bowl competition featured six universities from around the northwest. This year's 1 st place prize went to University of Portland, while University of Washington and Oregon Institute of Technologies (OIT) both tied for 2 nd place.

SUBMIT

Oregon Section 2007/2008 Annual Report

DR--182866691

5. MEETINGS

The Oregon Section conducted six (6) general meetings and the summer golf tournament over the past year, as well as hosting Joint 2007 ITE District 6 annual meeting. The general meetings included luncheons with speakers, joint meetings with other professional societies, the annual traffic bowl, and a technical workshop. The table below summarizes the general meetings conducted over the past year and those scheduled for the remainder of 2008.

Date	Meeting	Subject	Location	Attendees
4-10-07	April - Joint meeting with WTS	Kris Strickler and Danielle Cogan gave a presentation on the Columbia River Crossing project on the I-5 bridge that connects Oregon and Washington.	Embassy Suites Hotel Portland, OR	60
5-18-07	May - Joint meeting with Oregon Traffic Control Devices Committee (OTCDC)	Sections newly elected officers were announced and a presentation on the Highway 25 emergency reconstruction efforts that followed a large storm on Mt. Hood early 2007.	Crowne Plaza Hotel, Lake Oswego, OR	45
7-13 to 18-07	2007 ITE Western District (District 6) Annual Meeting & Quad Conference	Oregon ITE section hosted a very successful joint 2007 ITE District 6 Annual Meeting & Quad Conference in Portland. This year's annual meeting broke many previous records with over 500 meeting participants at the meeting.	Hyatt Hotel, Portland, OR	200
9-21-07	2007 Golf Tournament	Due to timing of ITE District 6 conference in Portland, this year's golf tournament was held in late summer. We had good turn out and good corporate sponsorship donations.	Oregon Golf Association (OGA) Golf Course, Washburn, OR	57
9-25-07	September Meeting	Jim Peters (DKS Associates) and Jay McCoy (City of Gresham) gave a presentation on the use of recent SCA's Adaptive Traffic Signal System in City of Gresham, Oregon.	Kelly's Irish Pub Portland, OR	56
10-23-07	October Meeting	Senator Rick Meninger, Chairman of the Senate Transportation Committee, discussed about congestion pricing and his "Moving Oregon" statewide transportation tour, and effort to build support for major investment in Oregon's transportation system.	Hotel Monaco, Portland, OR	73
11-15-07	2007 Student Traffic Bowl	Oregon ITE 10 th Annual Student Traffic Bowl competition featured six universities from around the northwest. This year's 1 st place prize went to University of Portland, while University of Washington and Oregon Institute of Technologies (OIT) both tied for 2 nd place.	McMenamins Edgefield Troutdale, OR	123 including 52 students

Institute of Transportation Engineers Page 7 of 13

UPLOAD PDF

Figure 12: Table QA asking for the place where the event on Date 5-18-07 will occur.

Document question

What is main content on Page 2?

Answer

The main content on Page 2 is a description of the participants and methods used in the study, including details about the chairs and equipment used, the protocol followed, and the measurement of energy expenditure and heart rate.

Evidence

Page 2:
[PARTICIPANTS AND METHODS Participants']
'Participants provided informed written consent, and the Mayo Clinic Institutional Review Board approved the protocol. Sixteen participants (nine women and seven men) were included with a mean (±SD) age, 23±5 years and body mass index (BMI), 26±5.5 kg/m² (↔)table 1.'
'Description of chairs Standard office chair (control chair): The criterion model chair is a standard office chair (Steelcase; Grand Rapids, Michigan, USA).'
'FootFidget, (<http://footfidget.com>)http://footfidget.com (FootFidget, Lake Zurich, Illinois, USA) (↔)figure 1), is an under-desk elasticated footrest that encourages leg activity while seated. It comprises of a steel 17"×10"×10" frame support base. The elasticised central footpad consists of a 7" foam-covered cylindrical rigid tube centred on two 17" flex-ible resistance cords that run through the tube and attach to the four upright legs on the stand. The user repeatedly 'bounces' their foot on the cylindrical tube that encounters resistance.'
'CoreChair, (<https://www.corechair.com>)https://www.corechair.com (CoreChair, Aurora, Ontario, Canada) (↔)figure 1), is a chair designed to promote activity while a person stays seated. It is a modified five-wheel office chair. It has a low, 9" backrest with adjustable depth and is without armrests. The seat is sculpted and covered in 2" thick foam padding. The main feature of the CoreChair is the mechanical core that allows for lateral movement while seated. Severity of seat tilt is adjustable and has a range of motion up to 14° in all directions. For the study trial, tilt severity was set to allow for the greatest range of motion. Interchangeable

2

Kepp GA, et al. *BMJ Open Sport Exerc Med* 2016;2:e000152. doi:10.1136/bmjsem-2016-000152

SUBMIT

Open Access

PARTICIPANTS AND METHODS

Participants

Participants provided informed written consent, and the Mayo Clinic Institutional Review Board approved the protocol. Sixteen participants (nine women and seven men) were included with a mean (±SD) age, 23±5 years and body mass index (BMI), 26±5.5 kg/m² (table 1).

Description of chairs

Standard office chair (control chair): The criterion model chair is a standard office chair (Steelcase; Grand Rapids, Michigan, USA).
FootFidget, http://footfidget.com (FootFidget, Lake Zurich, Illinois, USA) (figure 1), is an under-desk elasticated footrest that encourages leg activity while seated. It comprises of a steel 17"×10"×10" frame support base. The elasticised central footpad consists of a 7" foam-covered cylindrical rigid tube centred on two 17" flex-ible resistance cords that run through the tube and attach to the four upright legs on the stand. The user repeatedly 'bounces' their foot on the cylindrical tube that encounters resistance.
CoreChair, https://www.corechair.com (CoreChair, Aurora, Ontario, Canada) (figure 1), is a chair designed to promote activity while a person stays seated. It is a modified five-wheel office chair. It has a low, 9" backrest with adjustable depth and is without armrests. The seat is sculpted and covered in 2" thick foam padding. The main feature of the CoreChair is the mechanical core that allows for lateral movement while seated. Severity of seat tilt is adjustable and has a range of motion up to 14° in all directions. For the study trial, tilt severity was set to allow for the greatest range of motion. Interchangeable center columns of different lengths (tall and short) make the chair height adjustable.

Protocol

Participants were tested in thermal comfort, 2 hours after eating and after 30 min of rest. Prior to testing, patients were shown the equipment and the experimen-tal protocol was explained. Participants were weighed on a calibrated Seca 569 scale in light clothing (ath-letics shorts and vobars), and height was measured using a Seca 219 stadiometer without shoes (Seca, Chino, California, USA).
Participants were asked to sit on a standard office chair (Steelcase Criterion) for 20 min while working, checking email or using the internet. During this time, energy expenditure and heart rate were monitored.

Participants then continued their work-like activities and were provided with the FootFidget. Energy expenditure and heart rate were measured for 20 min. Participants then sat on the CoreChair during which time energy expenditure and heart rate were measured for 20 min. Participants remained sitting in the CoreChair, stopped their work-like activities and followed a 7 min chair-based exercise video. The participant was given a 5 min break. Subsequently, the participant followed the 7 min video for a second sample. Finally, participants walked at 2 mph for 20 min on a calibrated treadmill (PaceMaster Bruner; Aerobic, West Caldwell, New Jersey, USA).

Methods

Energy expenditure

Energy expenditure was measured using indirect calorimetry¹¹ (Metamax 3B; Cortex, Leipzig, Germany). The calorimeter was calibrated using 5.0% CO₂, 15.0% O₂, balance nitrogen (Praxair, Danbury, Connecticut, USA) and ambient air according to the manufacturer's specifi-cations. In addition, it was volume calibrated before each participant using a 3 L syringe. The calorimeter collects health-by-breath CO₂ and O₂ production and consumption, respectively, and energy expenditure is cal-culated using standard formulae.¹²

Heart rate monitoring

Participants were also fitted with a Polar Heart Rate Monitor H7 (Polar, Lake Success, New York, USA). Heart rate samples were recorded and synchronized for each beach.

RESULTS

Energy expenditure of the four seated conditions and slow walking (2 mph) are shown in figure 2. While sitting in the standard office chair, as expected, resting energy expenditure (sitting in a standard chair) showed a positive correlation with body weight (r=0.55, p=0.03). The relationship was described by the equation, resting energy expenditure (kcal/hour)=0.976×weight (kg) +0.917.
Energy expenditure increased significantly while using the FootFidget (~30%) when compared to the standard office chair. Energy expenditure increased in all partici-pants from a mean of 76±5 to 98±12 kcal/hour (p<0.001). Heart rate did not increase significantly, however (75±6, cf. 78±11 bpm). Similarly, energy expenditure increased significantly while using the

Table 1 Demographic information for 16 volunteers engaged in the study

	Age (years)	Height (cm)	Weight (kg)	Systolic blood pressure (mm hg)	Diastolic blood pressure (mm hg)
9 women;	23±4.7	170±10	77±17	116±18	78±19
7 men					

Data are shown as mean±SD.

2

Kepp GA, et al. *BMJ Open Sport Exerc Med* 2016;2:e000152. doi:10.1136/bmjsem-2016-000152

UPLOAD PDF

DR--18286691
DR--1058108

Figure 13: Page QA asking the main content on Page 2. The answer provides a high-level summarization of Page 2, covering the title of each section.

Document question

What is associated with chronic health conditions and impair cognitive function and obesity?

Answer

Sedentariness is associated with chronic health conditions, impaired cognitive function, and obesity. Breaking up sitting time with standing or walking while working can effectively decrease sedentariness and improve insulin sensitivity and lipids. Solutions to promote physical activity are necessary to reverse sedentariness and prevent chronic diseases.

Evidence

1: Levine JA. Health-chair reform: your chair: comfortable but deadly. (<http://dx.doi.org/10.2337/db10-1042>)Diabetes 2010;59:2715–16. Ng SW, Popkin BM. Time use and physical activity: a shift away from movement across the globe.

2: 2–5 6 7 8 9 10 Sedentariness is associated with a myriad of chronic diseases, impaired cognition and obesity. (↔)The mechanism by which sitting excessively causes disease is not well understood, but it is known that breaking up sitting improves insulin sensitivity and lipids. Several studies have examined the effectiveness of programmes to displace sitting with standing or walking while working. These measures can effectively decrease sitting time and improve productivity although their

3: 2 7 20–23 24 25 The importance of sedentariness in chronic disease and obesity is established. Sedentary behaviours occur during work and while at home. Many people spend the majority of their weekly waking hours at work and so solutions to reverse sedentariness and promote physical activity are necessary.

4: (<http://dx.doi.org/10.1016/j.genhosppsych.2012.08.001>)Gen Hosp Psychiatry 2013;35:16–22. Tudor-Locke C, Schuna JM, Jr. Steps to preventing type 2 diabetes: exercise, walk more, or sit less?

2

Kepp GA, et al. *BMJ Open Sport Exerc Med* 2016;2:e000152. doi:10.1136/bmjsem-2016-000152

SUBMIT

Open Access

Original research

BMJ Open Sport Exercise Medicine

Chair-based fidgeting and energy expenditure

Gabriel A Koepf,¹ Graham K Moore,¹ James A Levine^{1,2}

To cite: Koepf GA, Moore GK, Levine JA. Chair-based fidgeting and energy expenditure. *BMJ Open Sport Exerc Med* 2016;2:e000152. doi:10.1136/bmjsem-2016-000152

► Realisation history for this paper is available online. To view these files please visit the journal online (http://dx.doi.org/10.1136/bmjsem-2016-000152)

Accepted 14 August 2016

ABSTRACT

Introduction: Sedentariness is associated with chronic health conditions, impaired cognitive function and obesity. Work contributes significantly to sedentariness because many work tasks necessitate sitting. Few sustained solutions exist to reverse workplace sedentariness. Here, we evaluated a chair and an under-table device that were designed to promote fidgeting while seated, that hypothesis was that an under-table leg fidget bar and/or a fidget-promoting chair significantly increased energy expenditure. We compared these devices with chair-based exercise and walking.

Materials and methods: We measured energy expenditure and heart rate in 16 people while they sat and walked using a standard chair, an under-desk device that encourages leg fidgeting and a fidget-promoting chair. We compared outcomes with chair-based exercise and walking.

Results: Energy expenditure increased significantly while using either an under-table leg fidget bar or a fidget-promoting chair, when compared to the standard office chair (standard chair, 76±5 kcal/hour; leg fidget bar, 88±12 kcal/hour (p<0.001); fidget chair, 89±10 kcal/hour (p<0.03)). However, heart rate did not increase significantly in either case. Routs of exercise performed while seated provided energetic and heart rate equivalency to walking at 2 mph.

Conclusions: Chairs and devices that promote fidgeting can increase energy expenditure by ~20–30% but not increase heart rate. Dynamic sitting may be among a host of options to help people move more while at work.

INTRODUCTION

Sedentariness is associated with a myriad of chronic diseases, impaired cognition¹ and obesity.² The mechanism by which sitting excessively causes disease is not well understood, but it is known that breaking up sitting improves insulin sensitivity and lipids.³ Several studies have examined the effectiveness of programmes to displace sitting with standing or walking while working.⁴ These measures can effectively decrease sitting time and improve productivity^{5–7} although their

New findings

It is clear that chair-based fidgeting can increase energy expenditure.

How might it impact on clinical practice in the near future?

Chair-based fidgeting may be part of corporate weight maintenance or weight loss programmes. Chair-based fidgeting may be recommended by the medical community to reduce sedentary time at work.

long-term health benefits have not been proven.

When a person walks at even 1 mph, energy expenditure doubles when compared to basal metabolic rate.⁸ Sitting, however, even while fidgeting, is not substantially exothermic (5–10% increase above basal metabolic rate).^{9–11} We wanted to assess whether chairs that promote 'dynamic sitting'¹² can increase energy expenditure significantly above resting values. To this end, we examined the thermogenic impact of a chair designed specifically to encourage body fidgeting and an underdesk device that encourages leg fidgeting. In both cases, we measured the changes in energy expenditure and heart rate that accompanied their use. In order to better understand how we can encourage movement in people who need to work seated, we also examined the thermogenic and heart rate responses to chair-based exercises delivered through video. We compared these conditions to a 2 mph walk since walking is known to improve health.¹³ Our hypothesis was that chairs and gadgets designed to promote fidgeting while seated are associated with an increase in energy expenditure. The null hypothesis was that chair-promoted fidgeting does not increase energy expenditure above resting values. Secondary hypotheses examined how 'dynamic sitting' influenced heart rate.

2

Kepp GA, et al. *BMJ Open Sport Exerc Med* 2016;2:e000152. doi:10.1136/bmjsem-2016-000152

UPLOAD PDF

DR--18286691
DR--1058108

Figure 14: Single Document Content QA asking Sedentariness. The 2nd retrieved sentence includes the answer and corresponds to the first sentence in the abstract of the paper.

Document question: Who borns later, Michael Jordan or LeBron James?

Answer: LeBron James was born later than Michael Jordan, on December 30, 1984, while Jordan was born on February 17, 1963.

Evidence:

- Former NBA players, including Michael Jordan and Magic Johnson, were also critical of James, condemning him for joining with Bosh and Wade in Miami and not trying to win a championship as "the man".
- "James is also discussed within the context of being the greatest basketball player of all-time, which has resulted in frequent comparisons to Michael Jordan.[a]
- Michael Jeffrey Jordan (born February 17, 1963), also known by his initials MJ, is an American former professional basketball player and businessman.
- In a 2016 interview with Sports Illustrated, James acknowledged that his motivation was surpassing Jordan as the greatest.
- He also joined Michael Jordan as the only players to win an NBA MVP award, NBA championship, NBA Finals MVP, and Olympic gold medal in the same year.
- LeBron Raymone James Sr. (/læˈbrɒn/ lə-brɒn; born December 30, 1984), also known as LBJ, is an American professional basketball player for the Los Angeles Lakers of the National Basketball Association (NBA).

WIKIPEDIA The Free Encyclopedia
Michael Jordan
 Michael Jeffrey Jordan (born February 17, 1963), also known by his initials MJ, is an American former professional basketball player and businessman. The official National Basketball Association (NBA) website states: "By acclamation, Michael Jordan is the greatest basketball player of all time." He played fifteen seasons in the NBA, winning six NBA championships with the Chicago Bulls. He was integral in popularizing the sport of basketball and the NBA around the world in the 1980s and 1990s, becoming a global cultural icon.

Jordan played college basketball for three seasons under coach Dean Smith with the North Carolina Tar Heels. As a freshman, he was a member of the Tar Heels' national championship team in 1982. Jordan joined the Bulls in 1984 as the third overall draft pick and quickly emerged as a league star, entertaining crowds with his prolific scoring while gaining a reputation as one of the game's best defensive players. His leaping ability, demonstrated by performing slam dunks from the free-throw line in Slam Dunk Contests, earned him the nicknames "Air Jordan" and "His Airness". Jordan won his first NBA title with the Bulls in 1991 and followed that achievement with titles in 1992 and 1993, securing a three-peat. Jordan abruptly retired from basketball before the 1993–94 NBA season to play Minor League Baseball but returned to the Bulls in March 1995 and led them to three more championships in 1996, 1997, and 1998, as well as a then-record 72 regular season wins in the 1995–96 NBA season. He retired for the second time in January 1999 but returned for two more NBA seasons from 2001 to 2003 as a member of the Washington Wizards. During the course of his professional career, he was also selected to play for the United States national team, winning four gold medals—at the 1985 Pan American Games, 1984 Summer Olympics, 1992 Tournament of the Americas and 1992 Summer Olympics—while also being undefeated.

Jordan's individual accolades and accomplishments include six NBA Finals Most Valuable Player (MVP) awards, ten NBA scoring titles (both all-time records), five NBA MVP awards, ten All-NBA First Team designations, nine Defensive First Team honors, fourteen NBA All-Star Game selections, three NBA All-Star Game MVP awards, three NBA steals titles, and the 1988 NBA Defensive Player of the Year Award. He holds the NBA records for career regular season scoring average (30.1 points per game) and career playoff scoring average (32.4 points per game). In 1999, he was named the 20th century's greatest North American athlete by ESPN and was second to Babe Ruth on the Associated Press' list of athletes of the century. Jordan was twice inducted into the Naismith Memorial Basketball Hall of Fame, once in 2000 for his individual career, and again in 2010 as part of the 1992

Michael Jordan - Wikipedia

Jordan in April 2014

Personal information	
Born	February 17, 1963 New York City, New York, U.S.
Listed height	6 ft 6 in (1.98 m)
Listed weight	216 lb (98 kg)
Career information	
High school	Emley A. Lansy (Wilmington, North Carolina)
College	North Carolina (1981–1984)
NBA draft	1984: 1st round, 3rd overall pick
Playing career	Selected by the Chicago Bulls 1984–1989, 1995–1998, 2001–2003
Position	Shooting guard / small forward

https://en.wikipedia.org/wiki/Michael_Jordan

DR--182866691
 DR--1058108
 LeBron James
 Ohio
 Michael Jordan

SUBMIT

UPLOAD PDF

Figure 17: Multi-document Comparing Question comparing Lebron James and Michael Jordan. It requires the birthday information of Lebron and Jordan.

Document question: Who borns later, Michael Jordan or LeBron James?

Answer: LeBron James was born later than Michael Jordan, on December 30, 1984, while Jordan was born on February 17, 1963.

Evidence:

- Former NBA players, including Michael Jordan and Magic Johnson, were also critical of James, condemning him for joining with Bosh and Wade in Miami and not trying to win a championship as "the man".
- "James is also discussed within the context of being the greatest basketball player of all-time, which has resulted in frequent comparisons to Michael Jordan.[a]
- Michael Jeffrey Jordan (born February 17, 1963), also known by his initials MJ, is an American former professional basketball player and businessman.
- In a 2016 interview with Sports Illustrated, James acknowledged that his motivation was surpassing Jordan as the greatest.
- He also joined Michael Jordan as the only players to win an NBA MVP award, NBA championship, NBA Finals MVP, and Olympic gold medal in the same year.
- LeBron Raymone James Sr. (/læˈbrɒn/ lə-brɒn; born December 30, 1984), also known as LBJ, is an American professional basketball player for the Los Angeles Lakers of the National Basketball Association (NBA).

WIKIPEDIA The Free Encyclopedia
LeBron James
 LeBron Raymone James Sr. (/læˈbrɒn/ lə-brɒn; born December 30, 1984), also known as LBJ, is an American professional basketball player for the Los Angeles Lakers of the National Basketball Association (NBA). Nicknamed "King James", he is widely regarded as one of the greatest players in the history of the sport and is often compared to Michael Jordan in debates over the greatest basketball player of all time. James is the all-time leading scorer in NBA history and ranks fourth in career assists. He has won four NBA championships (two with the Miami Heat, one each with the Lakers and Cleveland Cavaliers), and has competed in 10 NBA Finals. He has also won four Most Valuable Player (MVP) Awards, four Finals MVP Awards, and two Olympic gold medals, and has been named an All-Star 19 times, selected to the All-NBA Team 19 times (including 13 First Team selections) and the All-Defensive Team six times, and was a runner-up for the NBA Defensive Player of the Year Award twice in his career.

James grew up playing basketball for St. Vincent–St. Mary High School in his hometown of Akron, Ohio. He was heavily touted by the national media as a future NBA superstar for his all-around scoring, passing, athleticism and playmaking abilities. A prep-to-pro, he was selected by the Cleveland Cavaliers with the first overall pick of the 2003 NBA draft. Named the 2004 NBA Rookie of the Year, he soon established himself as one of the league's premier players, leading the Cavaliers to their first NBA Finals appearance in 2007 and winning the NBA MVP award in 2009 and 2010. After failing to win a championship with Cleveland, James left in 2010 as a free agent to join the Miami Heat; this was announced in a nationally televised special titled *The Decision* and is among the most controversial free agency moves in sports history.

James won his first two NBA championships while playing for the Heat in 2012 and 2013; in both of these years, he also earned the league's MVP and Finals MVP awards. After his fourth season with the Heat in 2014, James opted out of his contract and re-signed with the Cavaliers. In 2016, he led the Cavaliers to victory over the Golden State Warriors in the Finals by coming back from a 3–1 deficit, delivering the team's first championship and ending the Cleveland sports curse. In 2020, James exercised his contract option to leave the Cavaliers and signed with the Lakers, where he won the 2020 NBA championship and his fourth Finals MVP. James is the first player in NBA history to accumulate \$1 billion in

LeBron James - Wikipedia

LeBron James

James with the Los Angeles Lakers in 2022

No. 23 – Los Angeles Lakers	
Position	Small forward / power forward
League	NBA
Personal information	
Born	December 30, 1984 Akron, Ohio, U.S.
Listed height	6 ft 9 in (2.06 m)
Listed weight	250 lb (113 kg)
Career information	
High school	St. Vincent–St. Mary (Akron, Ohio)
NBA draft	2003: 1st round, 1st overall pick
Playing career	Selected by the Cleveland Cavaliers 2003–present
Career history	

https://en.wikipedia.org/wiki/LeBron_James

DR--182866691
 DR--1058108
 LeBron James
 Ohio
 Michael Jordan

SUBMIT

UPLOAD PDF

Figure 18: Multi-document Comparing Question comparing Lebron James and Michael Jordan. It requires the birthday information of Lebron and Jordan.

8.8 Visualizing the Reasoning-and-Retrieving Process of LM-guided Graph Traverser

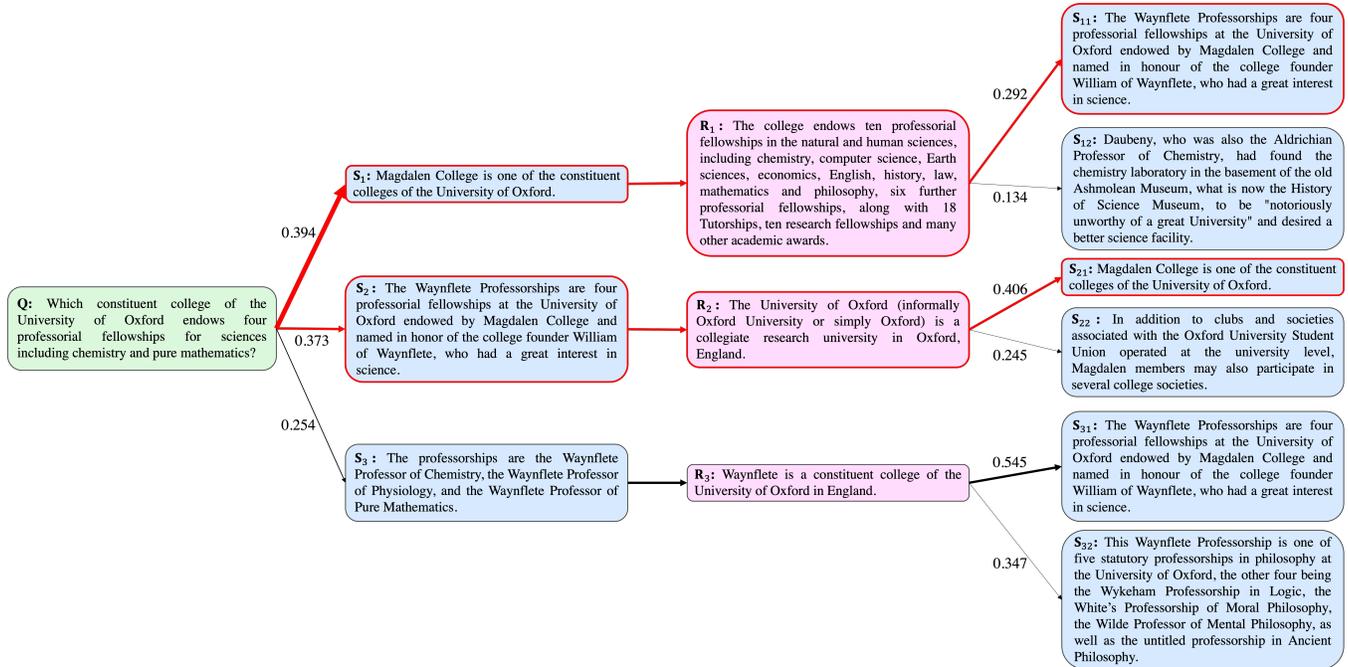


Figure 19: Visualizing the graph traversal over MD-QA-Example 1.

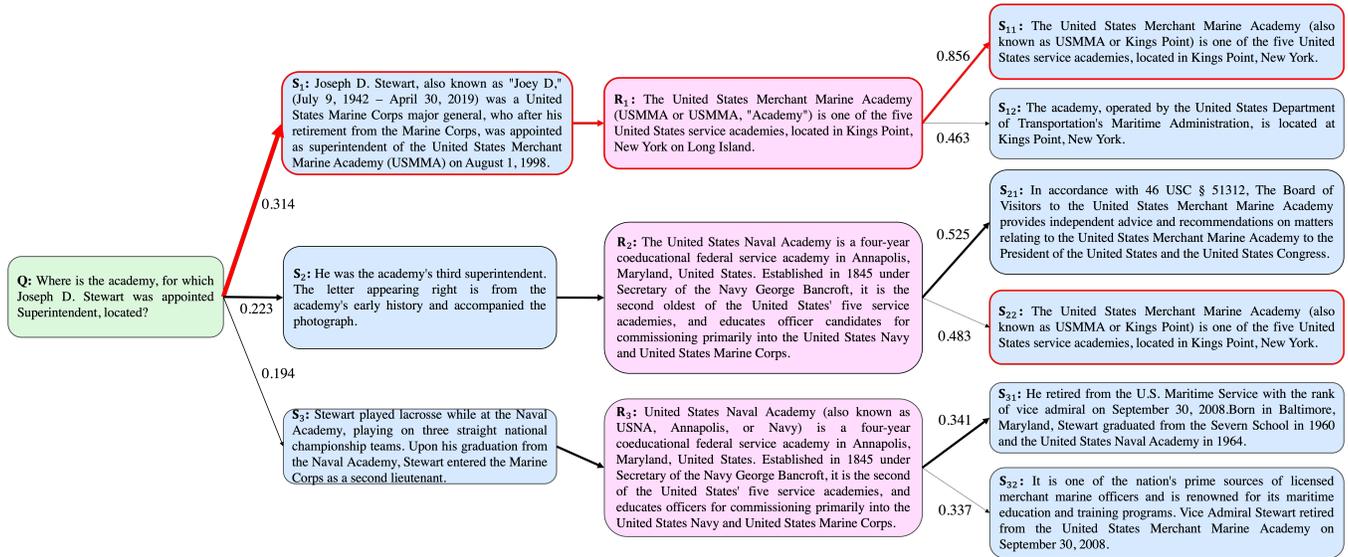


Figure 20: Visualizing the graph traversal over MD-QA-Example 2.

8.9 Prompt template used throughout this work

Listing 1: Examples of the Instruction Data for Fine-tuning LLaMA.

Question: Which magazine was started first Arthur's Magazine or First for Women?

Answer: Arthur's Magazine

Supporting Facts:

- (1) Arthur's Magazine (1844–1846) was an American literary periodical published in Philadelphia in the 19th century.
- (2) First for Women is a woman's magazine published by Bauer Media Group in the USA. The magazine was started in 1989.

Instruction: What evidence do we need to answer the question given the current evidence?

Input: Which magazine was started first Arthur's Magazine or First for Women? Arthur's Magazine (1844–1846) was an American literary periodical published in Philadelphia in the 19th century.

Output: First for Women is a woman's magazine published by Bauer Media Group in the USA. The magazine was started in 1989.

=====

Question: In what year was the creator of the current arrangement of Simpson's Theme born?

Answer: March 28, 1941

Supporting Facts:

- (1) The theme was re-arranged during season 2, and the current arrangement by Alf Clausen was introduced at the beginning of season 3.
- (2) Alf Heiberg Clausen (born March 28, 1941) is an American film and television composer.

Instruction: What evidence do we need to answer the question given the current evidence?

Input: In what year was the creator of the current arrangement of Simpson's Theme born? The theme was re-arranged during season 2, and the current arrangement by Alf Clausen was introduced at beginning of season 3.

Output: Alf Heiberg Clausen (born March 28, 1941) is an American film and television composer.

Listing 2: Example of the Prompt for QA without Retrieved Contexts.

Given the following question, create a final answer to the question.

=====

QUESTION: What is the birthday of this Anglo-Irish actress, courtesan, and mistress, who was the mother to the illegitimate daughter of King William IV?

=====

ANSWER: Please answer in less than 6 words.

Listing 3: Example of the Prompt for QA with Retrieved Contexts.

Given the following question and contexts, create a final answer to the question.

=====

QUESTION: During which years was the model of car, featured on the cover of Earth's "Pentastar: In the Style of Demons" manufactured?

=====

CONTEXT:

- 1: Pentastar: In the Style of Demons is the third full-length studio album by the drone doom band Earth.
- 2: In 1957, he published The Interpersonal Diagnosis of Personality, which the Annual Review of Psychology called the "most important book on psychotherapy of the year".
- 3: During the evanescent heyday of the cyberdelic counterculture, he served as a consultant to Billy Idol in the production of the 1993 album Cyberpunk.
- 4: During the development of the Barracuda, one of the worst-kept secrets was Ford's plan to introduce a new sporty compact car based on the inexpensive Falcon chassis and running gear (which was eventually released as the Mustang in mid-model year 1964); the extent of the other changes was not known.
- 5: "Peace in Mississippi" is a cover of the Jimi Hendrix song. The original vinyl release of the album has an alternative take of "Peace in Mississippi".
- 6: A 1975 Barracuda had been planned before the end of the 1970–74 model cycle.
- 7: In the spring of 2021, when the third wave of the coronavirus epidemic arrived, Varadi called their airline one of the "rare rays of hope" for investors.
- 8: During this time the first U.S. Federal auto safety standards were phased in, and Chrysler's response a requirement for side-marker lights distinguishes each model year of the second-generation Barracuda:As the pony-car class became established and competition increased, Plymouth began to revise the Barracuda's engine options.
- 9: The Barracuda sold for a base price of US\$2,512 (\$24,000 today).The 1964 model year was the first for the Barracuda and also the last year for push-button control of the optional Torqueflite automatic transmission.
- 10: In the words of symbolist poet Stephane Mallarme:Languages are imperfect because multiple; the supreme language is missing...no one can utter words which would bear the miraculous stamp of Truth Herself Incarnate...how impossible it is for language to express things...in the Poet's hands...by the consistent virtue and necessity of an art which lives on fiction, it achieves its full efficacy.
- 11: In France, the heart of the Decadent movement was during the 1880s and 1890s, the time of fin de siecle, or end-of-the-century gloom.
- 12: Pentastar: In the Style of Demons is the third full-length studio album by the drone doom band Earth, released in 1996. It has a more rock-oriented sound than their earlier drone doom work, although in a very minimalist style.

- 13: The game was a rematch of the previous year's Russell Athletic Bowl, which Clemson won 40-6. The two participants for the game were two of the semifinalists which were the Clemson Tigers and Oklahoma Sooners.
- 14: The effect of the war on Ernst was devastating; in his autobiography, he wrote of his time in the army thus: "On the first of August 1914 M[ax].E[rnst]. died. He was resurrected on the eleventh of November 1918".
- 15: Plymouth's executives had wanted to name the new model Panda, an idea unpopular with its designers. In the end, John Samsen's suggestion of Barracuda prevailed. Based on Chrysler's A-body, the Barracuda debuted in fastback form on April 1, 1964.
- 16: The Scapigliati (literally meaning "unkempt" or "disheveled") were a group of writers and poets who shared a sentiment of intolerance for the suffocating intellectual atmosphere between the late Risorgimento (1860s) and the early years of unified Italy (1870s).
- 17: Recurrent themes in his literary works include the supremacy of the individual, the cult of beauty, exaggerated sophistication, the glorification of machines, the fusion of man with nature, and the exalted vitality coexisting with the triumph of death.
- 18: Disc brakes and factory-installed air conditioning became available after the start of the 1965 model year. For the 1966 model year, the Barracuda received new taillamps, new front sheet metal, and a new instrument panel.
- 19: "Perhaps the worst failing of the book is the omission of any kind of proof for the validity and reliability of the diagnostic system," Eysenck wrote.
- 20: Based on stretched underpinnings of the rear-drive Alfa Romeo Giulia, it was rumored to be powered by a turbocharged V6 and arrive within the 2019 model year.
- 21: Their investments are in fleet development and the construction of airports, the first of which will be opened in Brasov.
- 22: He broke the hill record and this innovation was widely copied in the years to come. [citation needed] Mays made his mark on the track in such events as the 1935 German Grand Prix (scene of a famous victory of Tazio Nuvolari), sharing his ERA with Ernst von Delius.
- 23: There is still a question about the truth of the disclosure. In the 1968 Dagnet episode "The Big Prophet", Liam Sullivan played Brother William Bentley, leader of the Temple of the Expanded Mind, a thinly fictionalized Leary.
- 24: The Belgian Felicien Rops was instrumental in the development of this early stage of the Decadent movement. A friend of Baudelaire, he was a frequent illustrator of Baudelaire's writing, at the request of the author himself.
- 25: After taking responsibility for the controlled substance, Leary was convicted of possession under the Marihuana Tax Act of 1937 on March 11, 1966, sentenced to 30 years in prison, fined \$30,000, and ordered to undergo psychiatric treatment.
- 26: The general court delegation from Sullivan County is made up of all of the members of the New Hampshire House of Representatives from the county. In total, there are 13 members from 11 different districts.
- 27: Both teams then exchanged field goals, which brought the score to 16-10 in favor of Clemson. With 2:17 remaining, Oklahoma drove down the length of the field to score a touchdown, which gave the Sooners a one-point lead.
- 28: The average household size was 2.41 and the average family size was 2.88. 23.90% of the population were under the age of 18, 6.40% from 18 to 24, 28.00% from 25 to 44, 25.90% from 45 to 64, and 15.80% who were 65 years of age or older.
- 29: The band announced the release of a deluxe version of the album "How It Feels To Be Lost", which came out on August 21, 2020. On June 2, 2021, the band released the single "Bloody Knuckles" from their upcoming album.
- 30: The 82nd Orange Bowl was a College Football Playoff semifinal with the winner of the game competing against the winner of the 2015 Cotton Bowl: Alabama Crimson Tide football in the 2016 College Football Playoff National Championship, which took place at the University of Phoenix Stadium in Glendale, Arizona.

=====
QUESTION: During which years was the model of car, featured on the cover of Earth's "Pentastar: In the Style of Demons" manufactured?
=====

ANSWER: Please answer in less than 6 words.

Listing 4: Example of the Prompt for QA with Retrieved Contexts for MDR, KGP-T5, KGP-LLaMA and KGP-MDR.

Given the following question and contexts, create a final answer to the question.

=====

QUESTION: Anthony Avent played basketball for a High School that is located in a city approximately 8 mi west of where?

=====

CONTEXT:

- 1: Newark is the second largest city in the New York metropolitan area, located approximately 8 mi west of lower Manhattan.\n Prior to Seton Hall, Avent played at Malcolm X Shabazz High School in Newark, New Jersey.
- 2: Newark is the second largest city in the New York metropolitan area, located approximately 8 mi west of lower Manhattan.\n The United States District Court for the District of New Jersey is also located in the city.
- 3: Newark is the second largest city in the New York metropolitan area, located approximately 8 mi west of lower Manhattan.\n Near Market Street and includes a dormitory for boarding students; and Saint Vincent Academy which is an all-girls Roman Catholic high school founded and sponsored by the Sisters of Charity of Saint Elizabeth and operated continuously since 1869.Link Community School is a non-denominational coeducational day school that serves approximately 128 students in seventh and eighth grades.
- 4: Prior to Seton Hall, Avent played at Malcolm X Shabazz High School in Newark, New Jersey.\n Newark is the second largest city in the New York metropolitan area, located approximately 8 mi west of lower Manhattan.
- 5: Prior to Seton Hall, Avent played at Malcolm X Shabazz High School in Newark, New Jersey.\n The United States District Court for the District of New Jersey is also located in the city.
- 6: Prior to Seton Hall, Avent played at Malcolm X Shabazz High School in Newark, New Jersey.\n On Newark Bay, it is run by the Port Authority of New York and New Jersey and serves as the principal container ship facility for goods entering and leaving the New York metropolitan area and the northeastern quadrant of North America.
- 7: He played collegiately at Seton Hall University where he played in the 1989 NCAA championship game. Prior to Seton Hall, Avent played at Malcolm X Shabazz High School in Newark, New Jersey.\n Prior to Seton Hall, Avent played at Malcolm X Shabazz High School in Newark, New Jersey.
- 8: He played collegiately at Seton Hall University where he played in the 1989 NCAA championship game. Prior to Seton Hall, Avent played at Malcolm X Shabazz High School in Newark, New Jersey.\n The United States District Court for the District of New Jersey is also located in the city.
- 9: He played collegiately at Seton Hall University where he played in the 1989 NCAA championship game. Prior to Seton Hall, Avent played at Malcolm X Shabazz High School in Newark, New Jersey.\n As of the 202021 school year, the district, comprises 65 schools , had an enrollment of 40,423 students and 2,886.5 classroom teachers (on an FTE basis), for a studentteacher ratio of 14.0:1.Science Park High School, which was the 69th-ranked public high school in New Jersey out of 322 schools statewide, in New Jersey Monthly magazine's September 2010 cover story on the state's "Top Public High Schools", after being ranked 50th in 2008 out of 316 schools.
- 10: Anthony Avent (born October 18, 1969) is an American former professional basketball player who was selected by the Atlanta Hawks in the first round (15th pick overall) of the 1991 NBA draft.\n Newark is the second largest city in the New York metropolitan area, located approximately 8 mi west of lower Manhattan.
- 11: Anthony Avent (born October 18, 1969) is an American former professional basketball player who was selected by the Atlanta Hawks in the first round (15th pick overall) of the 1991 NBA draft.\n The United States District Court for the District of New Jersey is also located in the city.
- 12: Anthony Avent (born October 18, 1969) is an American former professional basketball player who was selected by the Atlanta Hawks in the first round (15th pick overall) of the 1991 NBA draft.\n Atlanta United 1, New York Red Bulls 2 The first game in Atlanta United history was played before a sellout crowd of 55,297.
- 13: Anthony Avent (born October 18, 1969) is a retired American professional basketball player who was selected by the Atlanta Hawks in the first round (15th pick overall) of the 1991 NBA Draft.\n The total school enrollment in Newark was 77,097 in the 20132017 ACS, with nursery and preschool enrollment of 7,432, elementary/high school (K12) enrollment of 49,532, and total college/graduate school enrollment of 20,133. The Newark Public Schools, a state-operated school district, is the largest school system in New Jersey.
- 14: Anthony Avent (born October 18, 1969) is a retired American professional basketball player who was selected by the Atlanta Hawks in the first round (15th pick overall) of the 1991 NBA Draft.\n As of the 202021 school year, the district, comprises 65 schools, had an enrollment of 40,423 students and 2,886.5 classroom teachers (on an FTE basis), for a studentteacher ratio of 14.0:1.Science Park High School, which was the 69th-ranked public high school in New Jersey out of 322 schools statewide, in New Jersey Monthly magazine's September 2010 cover story on the state's "Top Public High Schools", after being ranked 50th in 2008 out of 316 schools.

15: Anthony Avent (born October 18, 1969) is a retired American professional basketball player who was selected by the Atlanta Hawks in the first round (15th pick overall) of the 1991 NBA Draft. In the 2013--2017 American Community Survey, 13.6% of Newark residents ages 25 and over had never attended high school and 12.5% didn't graduate from high school, while 74.1% had graduated from high school, including the 14.4% who had earned a bachelor's degree or higher.

QUESTION: Anthony Avent played basketball for a High School that is located in a city approximately 8 mi west of where?

ANSWER: Please answer in less than 6 words.

Listing 5: Example of the Prompt for Grading QA.

You are an expert professor specialized in grading whether the prediction to the question is correct or not according to the real answer.

For example:

Question: What company owns the property of Marvel Comics?

Answer: The Walt Disney Company

Prediction: The Walt Disney Company

Return: 1

Question: Which constituent college of the University of Oxford endows four professorial fellowships for sciences including chemistry and pure mathematics?

Answer: Magdalen College

Prediction: Magdalen College.

Return: 1

Question: Which year was Marvel started?

Answer: 1939

Prediction: 1200

Return: 0

You are grading the following question:

Question: Anthony Avent played basketball for a High School that is located in a city approximately 8 mi west of where?

Answer: lower Manhattan

Prediction: Newark

If the prediction is correct according to the answer, return 1. Otherwise, return 0.

Return: your reply can only be one number '0' or '1'