


# GRAPE : Knowledge Graph Enhanced Passage Reader for Open-domain Question Answering

Mingxuan Ju<sup>1\*</sup>, Wenhao Yu<sup>1\*</sup>, Tong Zhao<sup>2</sup>, Chuxu Zhang<sup>3</sup>, Yanfang Ye<sup>1</sup>

<sup>1</sup>University of Notre Dame, <sup>2</sup>Snap Inc., <sup>3</sup>Brandeis University  
<sup>1</sup>{mju2, wyu1, yye7}@nd.edu; <sup>2</sup>tzhao@snap.com; <sup>3</sup>chuxuzhang@brandeis.edu

## Abstract

A common thread of open-domain question answering (QA) models employs a retriever-reader pipeline that first retrieves a handful of relevant passages from Wikipedia and then peruses the passages to produce an answer. However, even state-of-the-art readers fail to capture the complex relationships between entities appearing in questions and retrieved passages, leading to answers that contradict the facts. In light of this, we propose a novel knowledge **Graph** enhanced **passage** reader, namely GRAPE , to improve the reader performance for open-domain QA. Specifically, for each pair of question and retrieved passage, we first construct a localized bipartite graph, attributed to entity embeddings extracted from the intermediate layer of the reader model. Then, a graph neural network learns relational knowledge while fusing graph and contextual representations into the hidden states of the reader model. Experiments on three open-domain QA benchmarks show GRAPE can improve the state-of-the-art performance by up to 2.2 exact match score with a negligible overhead increase, with the same retriever and retrieved passages. Our code is publicly available at <https://github.com/jumxglhf/GRAPE>.

## 1 Introduction

Open-domain question answering (QA) tasks aim to answer questions in natural language based on large-scale unstructured passages such as Wikipedia (Chen and Yih, 2020; Zhu et al., 2021). A common thread of modern open-domain QA models employs a *retriever-reader* pipeline, in which a retriever aims to retrieve a handful of relevant passages w.r.t. a given question, and a reader aims to infer a final answer from the received passages (Guu et al., 2020; Karpukhin et al., 2020; Lewis et al., 2020; Izacard and Grave, 2021). Although these methods have achieved remarkable

\* Equal contribution.

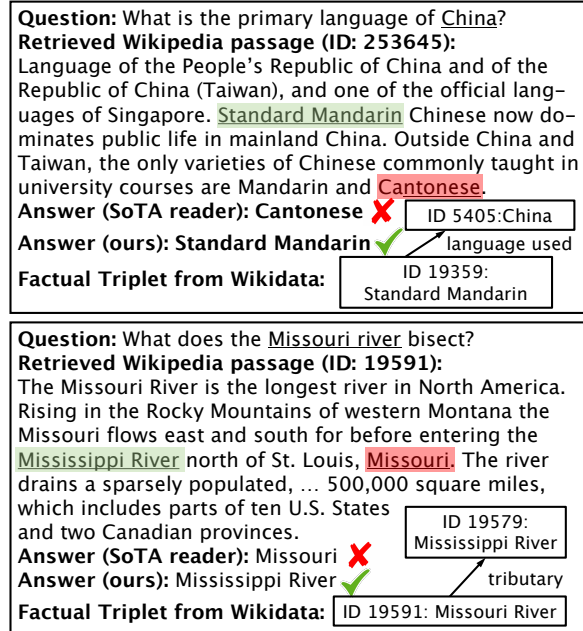


Figure 1: The answers produced by the SoTA reader FiD contradict the facts in the knowledge graph.

advances on various open-domain QA benchmarks, the state-of-the-art readers, such as FiD (Izacard and Grave, 2021), still often produce answers that contradict the facts. As shown in Figure 1, the FiD reader fails to produce correct answers due to inaccurate understanding of the factual evidence. Therefore, instead of improving the retrievers to saturate the readers with higher answer coverage in the retrieved passages (Yu et al., 2021; Oguz et al., 2022; Yu et al., 2022a), in this work, we aim at improving the readers by leveraging structured factual triples from the knowledge graph (KG).

A knowledge graph, such as Wikidata (Vrandečić and Krötzsch, 2014), contains rich relational information between entities, many of which can be further mapped to corresponding mentions in questions and retrieved passages. To verify the possible improvements brought by the KG, we conduct a simple analysis to examine the percentage of related fact triples present on the KG in the data,

Dataset	Fact-related examples	Error rate
NQ	736 (20.4%)	31.9%
TriviaQA	3,738 (33.0%)	17.4%
WebQ	1,181 (58.1%)	42.5%

Table 1: The error rate of state-of-the-art reader (i.e., FiD base) on the subset of data examples in the test set that have related fact triplets on the knowledge graph.

i.e., entities in questions are neighbors of answer entities in retrieved passages through any relation. We also wonder how many of the above examples are correctly answered by state-of-the-art readers. Table 1 shows that a great portion of examples (e.g., 58.1% in WebQ) can be matched to related fact triplets on the KG. However, without using the KG, FiD frequently produces incorrect answers to questions on these subsets, leaving us significant room for improvement. Therefore, a framework that leverages not only textual information in retrieved passages but also fact triplets from the KG is urgently desired to improve reader performance.

In this paper, we propose a novel knowledge **Graph enhanced passage** reader, namely GRAPE, to improve the reader performance for open-domain QA. Considering the enormous size of KGs and complex interweaving between entities (e.g., over 5 million entities and over 30 neighbors per entity on Wikidata), direct reasoning on the entire graph is intractable. Thus, we first construct a localized bipartite graph for each pair of question and passage, where nodes represent entities contained within them, and edges represent relationships between entities. Then, node representations are initialized with the hidden states of the corresponding entities, extracted from the intermediate layer of the reader model. Next, a graph neural network learns node representations with relational knowledge, and passes them back into the hidden states of the reader model. Through this carefully curated design, GRAPE takes into account both aspects of knowledge as a holistic framework.

To the best of our knowledge, we are the first work to leverage knowledge graphs to enhance the passage reader for open-domain QA. Our experiments demonstrate that, given the same retriever and the same set of retrieved passages, GRAPE can achieve superior performance on three open-domain QA benchmarks (i.e., NQ, TriviaQA, and WebQ) with up to 2.2 improvement on the exact match score over the state-of-the-art readers. In particular, our proposed GRAPE nearly doubles

the improvement gain on the subset that can be enhanced by fact triplets on the KG.

## 2 Related Work

**Text-based open-domain QA** Mainstream open-domain QA models employ a *retriever-reader* architecture, and recent follow-up work has mainly focused on improving the retriever or the reader (Chen and Yih, 2020; Zhu et al., 2021). For the retriever, most of them split text paragraphs on Wikipedia pages into over 20 million disjoint chunks of 100 words, each of which is called a passage. Traditional methods such as TF-IDF and BM25 explore sparse retrieval strategies by matching the overlapping contents between questions and passages (Chen et al., 2017; Yang et al., 2019). DPR (Karpukhin et al., 2020) revolutionized the field by utilizing dense contextualized vectors for passage indexing. Furthermore, other research improved the performance by better training strategies (Qu et al., 2021), passage re-ranking (Mao et al., 2021) or directly generating passages (Yu et al., 2022a). Whereas for the reader, extractive readers aimed to locate a span of words in the retrieved passages as answer (Karpukhin et al., 2020; Iyer et al., 2021; Guu et al., 2020). On the other hand, FiD and RAG, current state-of-the-art readers, leveraged encoder-decoder models such as T5 to generate answers (Lewis et al., 2020; Izacard and Grave, 2021). Nevertheless, these readers only used text corpus, failing to capture the complex relationships between entities, and hence resulting in produced answers contradicting the facts.

**KG-enhanced methods for open-domain QA** Recent work has explored incorporating knowledge graphs (KGs) into the *retriever-reader* pipeline for open-domain QA (Min et al., 2019; Zhou et al., 2020; Oguz et al., 2022; Yu et al., 2021; Hu et al., 2022; Yu et al., 2022b). For example, Unik-QA converted structured KG triples and merged unstructured text together into a unified index, so the retrieved evidence has more knowledge covered. Graph-Retriever (Min et al., 2019) and GNN-encoder (Liu et al., 2022) explored passage-level KG relations for better passage retrieval. KAQA (Zhou et al., 2020) improved passage retrieval by re-ranking according to KG relations between candidate passages. KG-FiD (Yu et al., 2021) utilized KG relations to re-rank retrieved passages by a KG fine-grained filter. However, all of these retriever-

enhanced methods focused on improving the quality of retrieved passages before passing them to the reader model. So, they still suffered from factual errors. Instead, our GRAPE is the first work to leverage knowledge graphs to enhance the reader, which is orthogonal to these existing KG-enhanced frameworks and our experiments demonstrate that with the same retriever and the same set of retrieved passages, GRAPE can outperform the state-of-the-art reader FiD by a large margin.

### 3 Proposed Method: GRAPE

In this section, we elaborate on the details of the proposed GRAPE. Figure 3 shows its overall architecture. GRAPE adopts a retriever-reader pipeline. Specifically, given a question, it first utilizes DPR to retrieve top- $k$  relevant passages from Wikipedia (§3.1). Then, to peruse the retrieved passages, it constructs a localized bipartite graph for each pair of question and passage (§3.2.1). The constructed graphs possess tractable yet rich knowledge about the facts among connected entities. Finally, with the curated graphs, structured facts are learned through a relation-aware graph neural network (GNN) and fused into token-level representations of entities in the passages (§3.2.2).

#### 3.1 Passage Retrieval

Given a collection of  $K$  passages, the goal of the retriever is to map all the passages in a low-dimensional vector, such that it can efficiently retrieve the top- $k$  passages relevant to the input question. Note that  $K$  can be very large (e.g., over 20 million in our experiments) and  $k$  is usually small (e.g., 100 in our experiments).

Following DPR (Karpukhin et al., 2020), we employ two independent BERT (Devlin et al., 2019) models to encode the question and the passage separately, and estimate their relevance by computing a single similarity score between their [CLS] token representations. Specifically, given a question  $q$  and a passage  $p_i \in \{p_1, p_2, \dots, p_K\}$ , we encode  $q$  by a question encoder  $E_Q(\cdot) : q \rightarrow \mathbb{R}^d$  and encodes  $p_i$  by a passage encoder  $E_P(\cdot) : p \rightarrow \mathbb{R}^d$ , where  $d$  is the hidden dimension of the used BERT. The ranking score  $r_q^i$  of  $p_i$  w.r.t  $q$  is calculated as:

$$r_q^i = E_Q(q)^\top \cdot E_P(p_i). \quad (1)$$

We select  $k$  passages whose ranking scores  $r_q$  are top- $k$  highest among all  $K$  passages. Before passing the retrieved passages into the reader model,

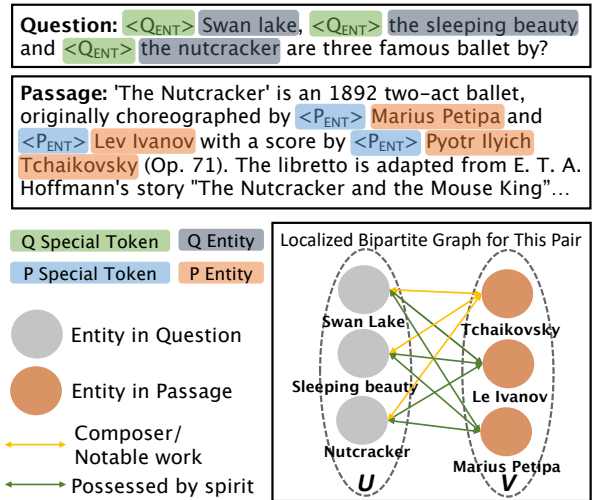


Figure 2: Given a pair of question and passage, the proposed GRAPE constructs a localized bipartite graph.

we process each question and passage by inserting special tokens before each entity. For entities in each passage, we use the special token  $\langle P_{ENT} \rangle$ ; for those in the question, we use another special token  $\langle Q_{ENT} \rangle$ , as shown in Figure 2. The special tokens play an important role in our proposed reader model, which is illustrated in more detail in §3.2.2.

#### 3.2 KG-enhanced Passage Reader

##### 3.2.1 Graph Construction

Given the retrieved and processed passages, our proposed GRAPE utilizes the factual triplets from KGs to construct localized bipartite graphs for each question-passage pair. A KG is defined as a set of triplets  $KG = \{(e_h, r, e_t)\}$ , where  $e_h$ ,  $e_t$ , and  $r$  refer to a head entity, a tail entity, and a corresponding relation between them, respectively. Knowledge graphs represent facts in the simple format of triplets, which can easily be leveraged to enrich our knowledge. Taking the question-passage pair in Figure 2 as an example, without any prior knowledge about the authorship of the ballets, the selection of answers between “Marius Petipa”, “Lev Ivanov” and “Tchaikovsky” is difficult. Nonetheless, factual triplets from the KG show that these three ballets are only “possessed by spirit” by “Marius Petipa” and “Lev Ivanov”. And their “composer” relations with “Tchaikovsky” make the answer obvious. By fusing such relational facts from KG triplets, the reader can better comprehend the concrete facts between involved entities and hence improve the performance for open-domain QA.

One naive solution could be fetching a sub-graph

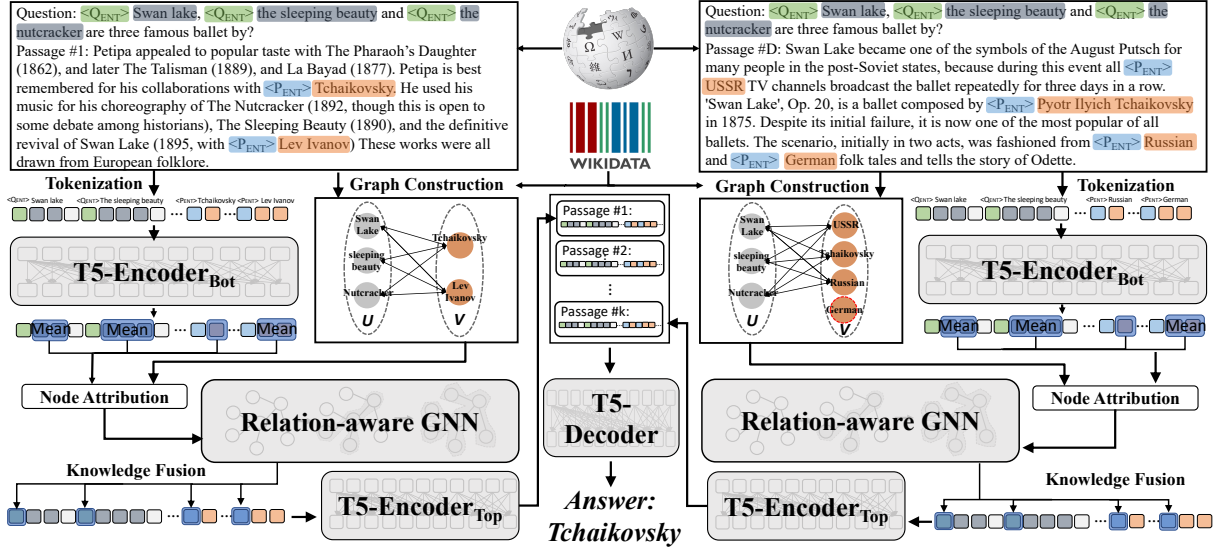


Figure 3: Two documents are independently encoded by our GRAPE with their corresponding localized bipartite graphs, leveraging both textual and structured information. The relation-aware GNN learns the structured knowledge from the localized bipartite graphs, attributed with entity representations extracted from the T5-Encoder<sub>Bot</sub>. The node representations are then fused into the T5-Encoder<sub>Top</sub>, which provides the hidden representations of the document. Finally, the T5-Decoder takes hidden states from all documents and generates the answer.

from the KG where all entities involved in the questions and the passages are included. While such design preserves all potentially relevant information, it suffers from dimensionality and noise issues. Therefore, we proposed to construct a localized bipartite graph for each question-passage pair, where only relational facts on relevant entities are kept. That is, in order to prune noisy peripheral relations, only the factual relations between question entities and passage entities are included in the localized bipartite graph. Let a bipartite graph be denoted as  $G = (\mathcal{U}, \mathcal{V}, \mathcal{E})$ , where  $\mathcal{U}$  and  $\mathcal{V}$  are two disjoint sets of nodes, and  $\mathcal{E}$  is the edge set containing edges that connect nodes from  $\mathcal{U}$  to  $\mathcal{V}$ , or vice versa. Specifically, in GRAPE,  $\mathcal{U}$  and  $\mathcal{V}$  are defined as the entity nodes in the question and the retrieved passage, respectively. There exists a bi-directional edge  $(e_h, e_t)$  between  $e_h \in \mathcal{U}$  and  $e_t \in \mathcal{V}$  if and only if  $\{(e_h, r_{h,t}, e_t) : r_{h,t} \in \mathcal{R}, (e_h, r_{h,t}, e_t) \text{ or } (e_t, r_{t,h}, e_h) \in KG\} \neq \emptyset$ , where  $\mathcal{R}$  denotes the set of all relation types on KG. Isolated nodes without any neighbors are removed from the graph. An example graph is shown in Figure 2, and Table 6 (in the appendix) shows the statistics of the constructed graphs.

### 3.2.2 Factual Relation Fusion

In this section, we illustrate how the proposed GRAPE fuses structured knowledge from our constructed localized bipartite graphs into the reader.

GRAPE uses FiD (Izcard and Grave, 2021) as the backbone architecture, which utilizes a T5 (Rafael et al., 2019) for encoding and decoding. To answer a question  $q$ , the input consists of  $k$  retrieved documents  $\{\text{doc}_1, \text{doc}_2, \dots, \text{doc}_k\}$ , where  $\text{doc}_i$  denotes to the concatenation of the token sequence of  $q$  and the token sequence of  $i$ -th retrieved passage  $p_i$ . Specifically,  $\text{doc}_i = \{q^1, q^2, \dots, q^t, p_i^1, p_i^2, \dots, p_i^o\}$ , where  $t$  and  $o$  are the length of the question and the passage sequence, respectively<sup>1</sup>. Given  $\text{doc}_i$ ,  $G_i$  denotes the localized bipartite graph constructed from it. And  $I_s(\text{doc}_i)$ ,  $I_e(\text{doc}_i)$ , and  $I_t(\text{doc}_i)$  denote the indices of the start, end, and special tokens of all entities in  $\text{doc}_i$ , respectively.

To fuse the relational knowledge from our constructed graphs, we split the encoder  $\text{Enc}(\cdot) : \text{doc} \rightarrow \mathbb{R}^{(t+o) \times d}$  of the reader (i.e., the encoder of T5) into two partitions  $\text{Enc}_{\text{top}}(\cdot)$  and  $\text{Enc}_{\text{bot}}(\cdot)$ . The bottom part  $\text{Enc}_{\text{bot}}(\cdot)$  contains the first  $L$  layers of  $\text{Enc}(\cdot)$  and the top part  $\text{Enc}_{\text{top}}(\cdot)$  contains the rest, where  $L$  is a hyper-parameter. Given  $\text{doc}_i$ ,  $\text{Enc}_{\text{bot}}(\cdot)$  delivers its encoded intermediate hidden states  $\mathbf{H}_i^b \in \mathbb{R}^{(t+o) \times d}$ , formulated as:

$$\mathbf{H}_i^b = \text{Enc}_{\text{bot}}(\text{doc}_i). \quad (2)$$

We then extract the node attributes  $\mathbf{X}_i^G \in \mathbb{R}^{|\mathcal{U} \cup \mathcal{V}| \times d}$  of its corresponding graph  $G_i$  from  $\mathbf{H}_i^b$  according

<sup>1</sup>For the simplicity of the notation, we assume all questions and passages have the equal length  $t$  and  $o$ , respectively.



to the span of its corresponding entity. For each entity node, its attribute vector is the average of the corresponding tokens’ representations. Formally,

$$\mathbf{X}_i^G = \bigoplus_{\substack{\text{StartIdx} \in I_s(\text{doc}_i) \\ \text{EndIdx} \in I_e(\text{doc}_i)}} \text{avg}(\mathbf{H}_i^b[\text{StartIdx}:\text{EndIdx}])^\top, \quad (3)$$

where  $\bigoplus$  denotes the vertical concatenation.

We use a relation-aware graph neural network (GNN) to conduct relation-aware message passing on the constructed graph  $G_i$  with attributes  $\mathbf{X}_i^G$ , denoted as  $\text{GNN}(\cdot, \cdot) : G \times \mathbb{R}^{|\mathcal{U} \cup \mathcal{V}| \times d} \rightarrow \mathbb{R}^{|\mathcal{U} \cup \mathcal{V}| \times d}$ , and the learning process is formulated as:

$$\mathbf{H}_i^G = \text{GNN}(G_i, \mathbf{X}_i^G), \quad (4)$$

where the learned node representations  $\mathbf{H}_i^G$  contain relational knowledge extracted from the KG as well as contextualized knowledge from the encoder. For the coherence of reading, the details of  $\text{GNN}(\cdot, \cdot)$  are described later in this subsection.

With the learned entity node representations  $\mathbf{H}_i^G$  containing knowledge from the fact relations, we leverage the special tokens to fuse them back into the reader. Specifically, we have  $\mathbf{H}_i^u = \mathbf{H}_i^b$  then

$$\mathbf{H}_i^u[I_t(\text{doc}_i)] = \mathbf{H}_i^b[I_t(\text{doc}_i)] + \mathbf{H}_i^G, \quad (5)$$

where  $[\cdot]$  is the indexing operation. The updated contextualized representations  $\mathbf{H}_i^u$  are then used as the input of the top part of the encoder to enable further information exchanges among regular tokens and the updated special tokens:

$$\mathbf{H}_i = \text{Enc}_{\text{top}}(\mathbf{H}_i^u). \quad (6)$$

Given the question  $q$ , GRAPE forwards all  $k$  retrieved documents through the above-described encoding process, and acquires the hidden states of all documents  $\{\mathbf{H}_i\}_{i=1}^k$ . These hidden states are then concatenated and sent to the decoder  $\text{Dec}(\cdot)$  for answer generation. Formally,

$$\text{answer} = \text{Dec}\left(\bigoplus_{i=0}^k \mathbf{H}_i\right). \quad (7)$$

To sum up, the workflow of our proposed GRAPE can be concluded as the following four steps: (i) get the initial contextualized representations via  $\text{Enc}_{\text{bot}}$  (Equation (2)) and the node attributes (Equation (3)), (ii) fuse fact relation by a relation-aware GNN (Equations (4) and (5)), (iii) exchange additional information via  $\text{Enc}_{\text{top}}$  (Equation (6)), and (iv) generate the answer by the decoder (Equation (7)).

**Relation-aware GNN** Here we elaborate the details of the aforementioned  $\text{GNN}(\cdot, \cdot)$ . Typically, each GNN layer (Hamilton et al., 2017; Kipf and Welling, 2017) can be formulated as

$$\mathbf{h}_v^{(n)} = \text{AGG}^{(n)}\left(\{\text{TRANS}^{(n)}(\mathbf{h}_u^{(n-1)}) : u \in \mathcal{N}(v)\}\right), \quad (8)$$

where  $\mathcal{N}(v)$  is the set of neighbors for node  $v$  including itself,  $n$  is the index of the current layer, and  $\mathbf{h}_v^n$  denotes the representation of node  $v$  at the  $n$ -th layer. The transform function  $\text{TRANS}(\cdot)$  projects node representations from the previous layer to a new vector space for message passing. The aggregation function  $\text{AGG}(\cdot)$  takes a set of node representations and aggregates them as a vector in a unified view (Kipf and Welling, 2017; Veličković et al., 2018; Zhang et al., 2019; Fan et al., 2022; Ju et al., 2022). Our proposed GRAPE uses a multi-layer perceptron as  $\text{TRANS}(\cdot)$  in each layer. That is,

$$\mathbf{A}^{(n)} = \sigma\left(\mathbf{H}^{(n-1)} \cdot \mathbf{W}_t^{(n)} + \mathbf{b}_t^{(n)}\right) \quad (9)$$

where  $\mathbf{A}^{(n)}$  is the intermediate embedding to be used by  $\text{AGG}(\cdot)$ ,  $\mathbf{H}^{(0)} = \mathbf{X}_i^G$  as aforementioned in Equation (3),  $\mathbf{W}_t^{(n)} \in \mathbb{R}^{d \times d}$  and  $\mathbf{b}_t^{(n)} \in \mathbb{R}^d$  denote the learnable parameters, and  $\sigma(\cdot)$  refers to the non-linear activation function.

For the aggregation function  $\text{AGG}(\cdot)$ , we explore a relation-aware attention mechanism. Different from GAT (Veličković et al., 2018) that considers only node representations for the edge attention weight, GRAPE also incorporates the relation representations between nodes. At layer  $n$ , for each node  $v$ , its representation  $\mathbf{h}_v^{(n)}$  is calculated by

$$e_{v,u}^{(n)} = \left(\mathbf{a}_v^{(n)} \bigoplus \text{avg}(\text{Enc}(r_{v,u})) \bigoplus \mathbf{a}_u^{(n)}\right) \cdot \mathbf{W}_e^{(n)},$$

$$\alpha_{v,u} = \frac{e_{v,u}^{(n)}}{\sum_{m \in \mathcal{N}(v)} e_{v,m}^{(n)}}, \quad \mathbf{h}_v^{(n+1)} = \sum_{u \in \mathcal{N}(v)} \alpha_{v,u} \cdot \mathbf{a}_u^{(n)}, \quad (10)$$

where  $\mathbf{a}_v^{(n)}$  is the node  $v$ ’s representation in  $\mathbf{A}^{(n)}$ ,  $\mathbf{W}_e^{(n)} \in \mathbb{R}^{3d \times 1}$  calculates the importance score of node  $u$  to node  $v$ , considering contextualized representations of the connected two nodes and language model’s understanding of their relationship (i.e.,  $\text{avg}(\text{Enc}(r_{u,v}))^2$ ). We further extend this

<sup>2</sup>In our implementation, we calculate and buffer  $\{\text{avg}(\text{Enc}(r)) : r \in R\}$  before every batched forward and simply query the representation of any relation as needed.

schema to the multi-head attention pipeline by having multiple operations as described in Equation (10) running in parallel. That is,

$$\mathbf{H}^{(n)} = \sum_{m=0}^M \mathbf{H}^{(n,m)}, \quad (11)$$

where  $M$  denotes the number of heads, and  $\mathbf{H}^{(n,m)}$  refers to the learned representations of the  $m$ -th head at  $n$ -th layer. Finally, the node representations  $\mathbf{H}_i^{(N)}$  are used as the output of  $\text{GNN}(\cdot, \cdot)$ , where  $N$  is the number of layers in the GNN.

In summary, our relation-aware GNN combines the current reader’s understanding of the factual relationships among nodes (i.e.,  $\text{avg}(\text{Enc}(r))$ ) with the intermediate hidden states  $\mathbf{X}_i^G$  from  $\text{Enc}_{\text{bot}}(\cdot)$ . Enriched by structured fact relations, entity node representations are then fused back into the reader’s encoder so that our GRAPE can comprehend facts between entities during the encoding process.

## 4 Experiments

In this section, we conduct comprehensive experiments on three community-acknowledged public open-domain QA benchmarks: Natural Questions (NQ) based on Google search queries, TriviaQA based on questions from trivia and quiz-league websites, and Web Questions (WebQ) based on questions from Google Suggest API (Kwiatkowski et al., 2019; Joshi et al., 2017; Berant et al., 2013). We explore the same train / dev / test splits and preprocessing techniques as used by (Izacard and Grave, 2021; Karpukhin et al., 2020).

### 4.1 Experimental Setup

**Retrieval Corpus** We followed the same process as used in (Karpukhin et al., 2020; Lewis et al., 2020) for preprocessing Wikipedia pages. We split each Wikipedia page into disjoint 100-word passages, resulting in 21 million passages in total. As for the knowledge graph used to construct our localized bipartite graphs, we used English Wikidata (Vrandečić and Krötzsch, 2014). The total number of aligned entities, relations, and triplets on Wikidata is 2.7M, 630, and 14M respectively<sup>3</sup>. We used ELQ (Li et al., 2020) to identify mentions in the question and retrieved passages, and link them to corresponding entities on Wikidata.

**Implementation Details** In GRAPE, involved hyper-parameters are the number of retrieved passages  $k$ , the number of GNN layers  $N$ , the number

<sup>3</sup>The Wikipedia and Wikidata were all collected in December of 2019. We only used the most visited top 1M entities.

of GNN head  $M$ , and the encoder layer index  $L$ , where  $\text{Enc}_{\text{top}}$  and  $\text{Enc}_{\text{bot}}$  are partitioned). We explore  $k = 100$ ,  $N = 2$ ,  $L = 3$  and  $M = 8$  as the default setup. The hidden dimension of the GNN in GRAPE are set to the dimension of its language model (i.e.,  $d = 768$  for the base configuration and  $d = 1024$  for the large configuration). Since  $N = 2$  and  $M = 8$  are the standard values for most GNNs with the attention mechanism, able to capture 2-hop neighbor information while being stable (Veličković et al., 2018), we simply follow the same principle.

For the encoder layer index  $L$ , we search the optimal value in the range of  $\{3, 4, 6, 8, 9\}$ . According to our experiment, we observe that different selections of  $L$  don’t have much impact on the performance, indicating that the infusion of structured knowledge doesn’t correlate with the contextualization of the entity embedding. However, we do observe a faster convergence rate for lower  $L$  values. So for faster convergence, we set  $L$  to 3. Other hyper-parameter selections related to training with best performance across all datasets are shown in Table 5 in §A.1. The software and hardware information can be found in §A.1 in the appendix.

**Evaluation Metrics** We use the standard evaluation metric for open-domain QA: exact match score (EM) (Rajpurkar et al., 2016; Zhu et al., 2021). An answer is considered correct if and only if its normalized form<sup>4</sup> has a match in the acceptable answer list. For all experiments, we conduct 3 runs with different random seeds and report the average.

### 4.2 Baseline Models

We compare GRAPE with four groups of baselines: (i) The first group includes closed-book models, where no Wikipedia document is provided during training and inference: T5-11B (Raffel et al., 2019) and GPT-3 (Brown et al., 2020). (ii) The second contains extractive models, which utilize passages extracted by enhanced retrievers and find the span of the answer: DPR (Karpukhin et al., 2020), RIDER (Mao et al., 2021), RECONSIDER (Iyer et al., 2021). (iii) The third includes approaches that utilize KG for retrieving: Graph-Retriever (Min et al., 2019), Path-Retriever (Asai et al., 2019), and KAQA (Zhou et al., 2020). (iv) The last group contains advanced generative readers: RAG (Lewis et al., 2020), REALM (Guu et al., 2020), Top-K

<sup>4</sup>We use the same normalization procedure as introduced in (Karpukhin et al., 2020).

Group	Model	#params	NQ	TriviaQA	WebQ
(i)	T5-11B	11B	32.6	42.3	37.2
	GPT-3 (64-shot)	175B	29.9	-	41.5
(ii)	DPR	110M	41.5	56.8	41.1*
	RIDER	626M	48.3	-	-
	RECONSIDER	670M	45.5	61.7	-
(iii)	Graph-Retriever	110M	34.7	55.8	36.4
	Path-Retriever	445M	31.7	-	-
	KAQA	110M	-	64.1	-
(iv)	REALM	330M	40.4	-	40.7
	RAG	626M	44.5	56.1	45.2*
	Joint Top-K	990M	48.1	59.6	-
	FiD (base)	440M	48.2	65.0	46.5
	FiD (large)	990M	51.4	67.6	50.5
(v)	GRAPE (base)	454M	48.7 (0.5 $\uparrow$ )	66.2 (1.2 $\uparrow$ )	48.1 (1.6 $\uparrow$ )
	GRAPE (large)	1.01B	<b>53.5</b> (2.1 $\uparrow$ )	<b>69.8</b> (2.2 $\uparrow$ )	<b>51.7</b> (1.2 $\uparrow$ )

Table 2: Exact match scores over the test sets of Natural Questions, TriviaQA and Web Questions. We put the training details (such as learning rate, batch size, dev performance, etc) corresponding to the performance of our GRAPE in Table 5. Numbers in parenthesis are improvements of GRAPE over the corresponding best-performing baseline. Note that  $\star$  means model is warmed with external training data from Natural Questions.

(Sachan et al., 2021) and FiD (Izacard and Grave, 2021). For FiD, we compare both its base and large versions, and for other baselines we compare their best-performing versions.

We note that our method is the first work using the knowledge graph to improve the reader performance for open-domain QA. Hence, this is orthogonal to existing works using the knowledge graph to improve passage retrieval (Liu et al., 2022) or re-ranking (Yu et al., 2021). Our experiments show that with the same retriever and the same set of retrieved passages, GRAPE can outperform the state-of-the-art reader FiD by a large margin.

### 4.3 Experimental Results

#### 4.3.1 Comparison with Baselines

Table 2 shows the model performance of 13 baselines as well as our GRAPE. We can observe that our proposed GRAPE can significantly outperform the best performing baselines across all datasets over both base and large configurations. Specifically, GRAPE improves FiD by 0.5, 1.2, and 1.6 EM score on the base model, and 2.1, 2.2, and 1.2 EM scores on the large model on NQ, TriviaQA and WebQ, respectively. Albeit being competitive on all datasets, GRAPE brings more improvements on TriviaQA and WebQ than NQ. We believe the reason is that on NQ, the percentage of questions

favorable from factual KG relations among all questions, as shown in Table 1 and Table 4, is relatively lower, compared to TriviaQA and WebQ. On the large configuration, even though some questions are not directly favored by structured facts, the additional information re-routing from GRAPE still benefits with additional learning capability, which outperforms FiD large for 2.1 EM score. We also note that the performance of RAG on WebQ is better than FiD base and close to GRAPE base, which is caused by a tremendous amount of additional training on other open-domain QA datasets.

#### 4.3.2 Ablation Study

We design two variants for our GRAPE. The first is GRAPE without considering relations between entity nodes, i.e.,  $\text{avg}(\text{Enc}(r_{u,v}))$  in Equation (10) is deleted. The goal is to validate the improvements brought by relational knowledge, which is denoted as w/o Rel in Table 3. The second is GRAPE without considering the relations as well as neighbor differences. The goal is to validate if GRAPE can differentiate the important neighbor without the attention mechanism, which is denoted as w/o Att. As shown in Table 3, we can observe that the performance drops when removing any of the two mechanisms, demonstrating their effectiveness and further validating the rich inductive bias brought

Datasets	FiD	GRAPE	w/o Rel	w/o Att
NQ	48.2	<b>48.7</b>	48.6	48.3
	51.4	<b>53.5</b>	53.4	53.1
TriviaQA	65.0	<b>66.2</b>	65.7	65.8
	67.6	<b>69.8</b>	69.5	69.6
WebQ	46.5	48.1	<b>48.4</b>	48.2
	50.5	<b>51.7</b>	51.5	51.0

Table 3: Ablation study of GRAPE without relation knowledge or attention mechanism. The first line refers to the performance on the base model; whereas the second line refers to the large model.

Dataset	Subset%	Model size	EM among subset	
			FiD	GRAPE
NQ	20.3%	base	68.1	<b>70.7</b> (2.6 $\uparrow$ )
		large	69.3	<b>71.9</b> (2.6 $\uparrow$ )
TriviaQA	33.0%	base	82.6	<b>86.4</b> (3.8 $\uparrow$ )
		large	85.6	<b>88.9</b> (3.3 $\uparrow$ )
WebQ	58.1%	base	57.5	<b>61.2</b> (3.7 $\uparrow$ )
		large	62.7	<b>65.1</b> (2.4 $\uparrow$ )

Table 4: Exact match score on the subset of questions that can be enhanced by factual triplets from the KG.

by the factual relations between entities. Besides, we also notice that the incorporation of relation is more important than the attention mechanism.

### 4.3.3 Improvement Analysis by KG Relation

Since GRAPE utilizes the factual relational knowledge from KG, to validate the legitimacy of our assumption, we analyze the performance gain on the subset of questions that can be directly solved by a factual triplet on KG (i.e., constructed graphs for these questions contain at least one edge that links the answer entity to entities in questions). From Table 4, we observe that GRAPE significantly improves performance on this subset that factual relations from KG naturally favors. For example, on the base model, GRAPE improves the overall performance by 0.5, 1.2, and 1.6 EM score respectively on these three datasets, and almost doubles the performance margin (i.e., 2.6, 3.8, and 3.7) on their subsets. This phenomenon demonstrates that GRAPE tends to utilize the inductive bias we introduce through graphs and the major performance gain can be rooted in the factual relational knowledge from KG, which further validates the legitimacy of GRAPE.

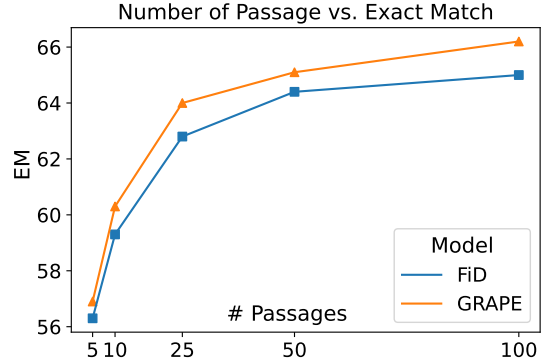


Figure 4: The performance on the test set of TriviaQA w.r.t. the number of passages.

### 4.3.4 Scaling with Number of Passages

We further evaluate the performance of GRAPE with respect to the different numbers of retrieved passages (i.e.,  $D = \{5, 10, 25, 50, 100\}$ ), as shown in Figure 4. We observe that given the same number of passages, GRAPE consistently outperforms FiD, with greater performance gains given more passages. Specifically, GRAPE performs on par with FiD with only the half amount of retrieved passages, starting from 25 retrieved passages. This phenomenon demonstrates that, when the answer is well presented in the retrieved passages, facts introduced by our curated graphs constructed from KG significantly help the reader answer questions.

### 4.3.5 Case Studies on KG Relations

To further validate the improvement gain induced by GRAPE, we analyze samples that are incorrectly answered by FiD but correctly answered by GRAPE, and visualize the constructed graphs for these samples, as shown in Figure 5. From these samples, we can observe the performance gain indeed comes from the strong enhancement brought by fact relations from their constructed graphs. For example, in the first example, with the fact relations, GRAPE understands that ‘‘Arges’’ is a member of ‘‘Cyclopes’’, which perfectly enhances answering for the given problem. In the third example, we can observe that FiD delivers an answer that is only partially correct. Whereas, enhanced by fact relation from KG, GRAPE correctly answers the question, because of the triplet (‘‘UK’’, ‘‘applies to jurisdiction’’, ‘‘Parliament of UK’’). This design is tractable yet effective. Because only entities highly correlated with useful facts will be included. Specifically, passage entities unrelated to question entities are very likely to be marginal and hence removed, and only the factual triplets helpful for answering the problems are kept. Relations



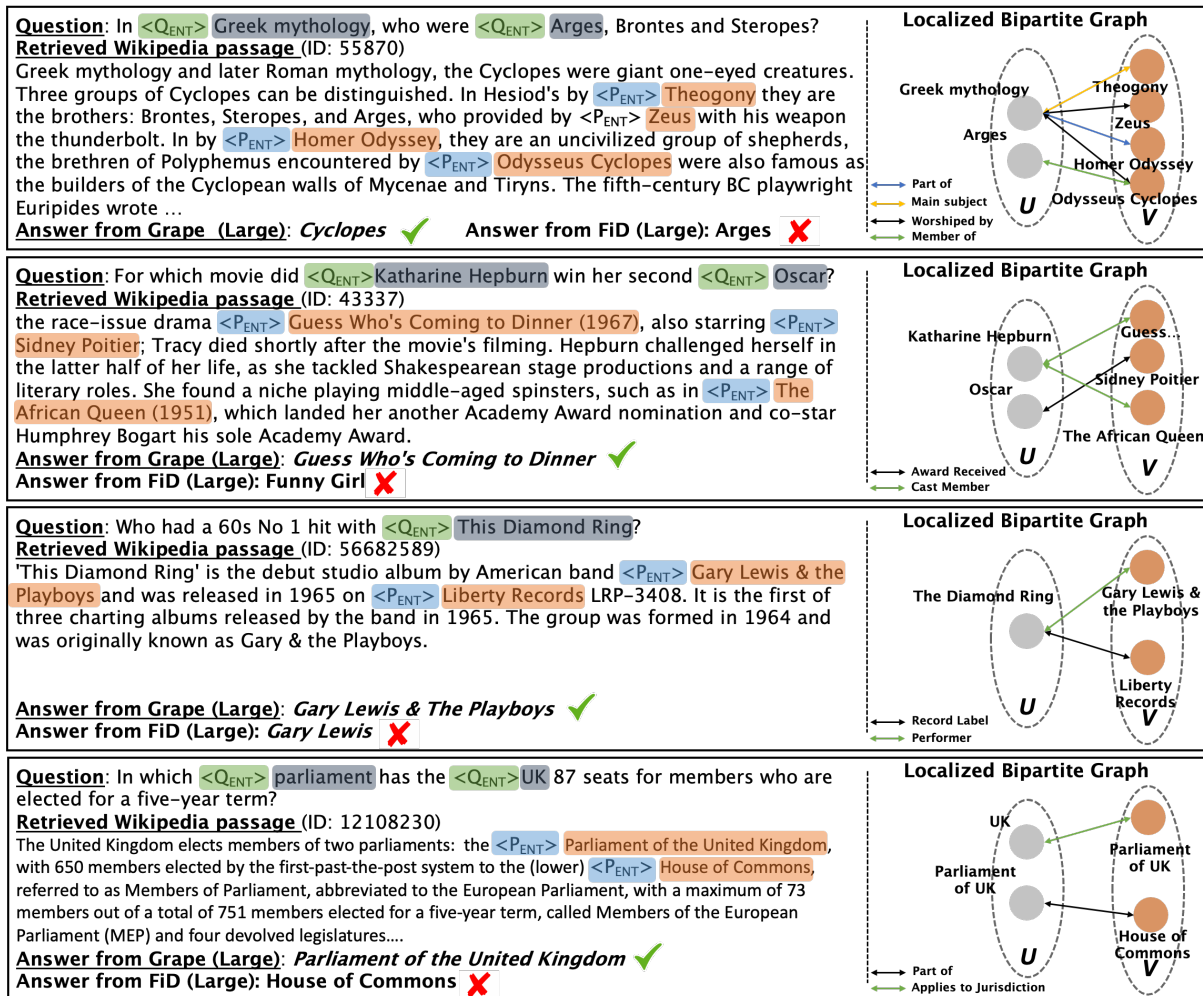


Figure 5: Case studies on samples that are incorrectly answered by FiD but correctly answered by GRAPE. Relations in green arrow indicate the factual relations from KG that enhance the question answering.

within passage entities are most likely peripheral and hence neglected in the bipartite graph.

## 5 Conclusion

In this work, we study the problem of open-domain QA. We discover that state-of-the-art readers fail to capture the complex relationships between entities appearing in questions and retrieved passages, resulting in produced answers that contradict the facts. To this end, we propose a novel knowledge Graph enhanced passage Reader (GRAPE) to improve the reader performance for open-domain QA. Specifically, for each pair of question and retrieved passage, we construct an informative localized bipartite graph and explore an expressive relation-aware GNN to learn entity representations that contain contextual knowledge from passages as well as fact relations from the KG. Experiments on three open-domain QA benchmarks show that GRAPE significantly outperforms state-of-the-art readers by up to 2.2 exact match score. In the future, we

plan to enrich the structured information contained in our graphs from other external resources.

## 6 Limitations

GRAPE only solves errors from fact-related examples. Besides, GRAPE explores fact relations from Wikidata, and hence we might omit fact relations from other sources such as Freebase.

## Acknowledgement

We appreciate Neil Shah and Yozen Liu from Snap Inc. and Zhihan Zhang from University of Notre Dame for valuable discussions and suggestions. This work is partially supported by the NSF under grants IIS-2209814, IIS-2203262, IIS-2214376, IIS-2217239, OAC-2218762, CNS-2203261, CNS-2122631, CMMI-2146076, and the NIJ 2018-75-CX-0032. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of any funding agencies.

## References

- Akari Asai, Kazuma Hashimoto, Hannaneh Hajishirzi, Richard Socher, and Caiming Xiong. 2019. Learning to retrieve reasoning paths over wikipedia graph for question answering. *arXiv preprint arXiv:1911.10470*.
- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic parsing on freebase from question-answer pairs. In *Procs. of EMNLP*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Procs. of NeurIPS*.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading wikipedia to answer open-domain questions. In *Procs. of ACL*.
- Danqi Chen and Wen-tau Yih. 2020. Open-domain question answering. In *ACL 2020: tutorial abstracts*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Procs. of NAACL*.
- Yujie Fan, Mingxuan Ju, Chuxu Zhang, and Yanfang Ye. 2022. Heterogeneous temporal graph neural network. In *Procs of SDM*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. Realm: Retrieval-augmented language model pre-training. *arXiv preprint arXiv:2002.08909*.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. *Procs. of NeurIPS*.
- Ziniu Hu, Yichong Xu, Wenhao Yu, Shuohang Wang, Ziyi Yang, Chenguang Zhu, Kai-Wei Chang, and Yizhou Sun. 2022. Empowering language models with knowledge graph reasoning for open-domain question answering. In *Procs. of EMNLP*.
- Srinivasan Iyer, Sewon Min, Yashar Mehdad, and Wen-tau Yih. 2021. Reconsider: Improved re-ranking using span-focused cross-attention for open domain question answering. In *Procs. of NAACL-HLT*.
- Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering. In *Procs. of EACL*.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Procs. of ACL*.
- Mingxuan Ju, Shifu Hou, Yujie Fan, Jianan Zhao, Yanfang Ye, and Liang Zhao. 2022. Adaptive kernel graph neural network. In *Procs of AAAI*.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *Procs. of EMNLP*.
- Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *Procs. of ICLR*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: A benchmark for question answering research. *TACL*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Procs. of NeurIPS*.
- Belinda Z. Li, Sewon Min, Srinivasan Iyer, Yashar Mehdad, and Wen-tau Yih. 2020. Efficient one-pass end-to-end entity linking for questions. In *Procs. of EMNLP*.
- Jiduan Liu, Jiahao Liu, Yang Yang, Jingang Wang, Wei Wu, Dongyan Zhao, and Rui Yan. 2022. Gnn-encoder: Learning a dual-encoder architecture via graph neural networks for passage retrieval. *arXiv preprint arXiv:2204.08241*.
- Yuning Mao, Pengcheng He, Xiaodong Liu, Yelong Shen, Jianfeng Gao, Jiawei Han, and Weizhu Chen. 2021. Reader-guided passage reranking for open-domain question answering. In *Findings of ACL-IJCNLP*.
- Sewon Min, Danqi Chen, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2019. Knowledge guided text retrieval and reading for open domain question answering. *arXiv preprint arXiv:1911.03868*.
- Barlas Oguz, Xilun Chen, Vladimir Karpukhin, Stan Peshterliev, Dmytro Okhonko, Michael Schlichtkrull, Sonal Gupta, Yashar Mehdad, and Scott Yih. 2022. Unik-qa: Unified representations of structured and unstructured knowledge for open-domain question answering. *Procs. of NAACL-HLT*.
- Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2021. Rocketqa: An optimized training approach to dense passage retrieval for open-domain question answering. In *Procs. of NAACL*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.

- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. In *Procs. of EMNLP*.
- Devendra Sachan, Mostofa Patwary, Mohammad Shoeybi, Neel Kant, Wei Ping, William L Hamilton, and Bryan Catanzaro. 2021. End-to-end training of neural retrievers for open-domain question answering. In *Procs. of ACL-IJCNLP*.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. In *Procs. of ICLR*.
- Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*.
- Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. 2019. End-to-end open-domain question answering with bertserini. In *NAACL 2019 (demo)*.
- Donghan Yu, Chenguang Zhu, Yuwei Fang, Wenhao Yu, Shuohang Wang, Yichong Xu, Xiang Ren, Yiming Yang, and Michael Zeng. 2021. Kg-fid: Infusing knowledge graph in fusion-in-decoder for open-domain question answering. *arXiv preprint arXiv:2110.04330*.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2022a. Generate rather than retrieve: Large language models are strong context generators. *arXiv preprint arXiv:2209.10063*.
- Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhiting Hu, Qingyun Wang, Heng Ji, and Meng Jiang. 2022b. A survey of knowledge-enhanced text generation. *ACM Computing Surveys (CSUR)*.
- Chuxu Zhang, Dongjin Song, Chao Huang, Ananthram Swami, and Nitesh V Chawla. 2019. Heterogeneous graph neural network. In *Procs. of KDD*.
- Mantong Zhou, Zhouxing Shi, Minlie Huang, and Xiaoyan Zhu. 2020. Knowledge-aided open-domain question answering. *arXiv preprint arXiv:2006.05244*.
- Fengbin Zhu, Wenqiang Lei, Chao Wang, Jianming Zheng, Soujanya Poria, and Tat-Seng Chua. 2021. Retrieving and reading: A comprehensive survey on open-domain question answering. *arXiv preprint arXiv:2101.00774*.

## A Supplementary Appendix

### A.1 Software and Hardware Environment

The software that GRAPE uses includes Huggingface Transformers 4.18.0 , Deep Graph Library (DGL) 0.8.2 and PyTorch 1.11.0 . And all our experiments are conducted on servers with 8 Tesla V100 32GB GPUs, with the RAM occupation less than 50GB. On average, updating GRAPE’s parameters costs 72 GPU hours for 10K training steps at the base configuration and 170 GPU hours at the large configuration, including intermediate validations, under the settings reported in Table 5.

### A.2 Hyper-parameter Sensitivity

In GRAPE, hyper-parameters related to our relation-aware GNN are the number of retrieved passages  $k$ , the number of GNN layers  $N$ , the number of GNN head  $M$ , and encoder layer index  $L$ , where  $Enc_{top}$  and  $Enc_{bot}$  are partitioned). We explore  $k = 100$ ,  $N = 2$ ,  $L = 3$  and  $M = 8$  as the default setup. The hidden dimension of the GNN in GRAPE are set to the dimension of its language model.  $N = 2$  and  $M = 8$  are the standard values for most GNNs with attention mechanism, able to capture 2-hop neighbor information while being stable (Veličković et al., 2018); and we simply follow the same principle. For the encoder layer index  $L$ , we search the optimal value in the range

of  $\{3, 4, 6, 8, 9\}$ . According to our experiment, we observe that different selection of  $L$  doesn’t have much impact on the performance, indicating that the infusion of structured knowledge doesn’t rely much on the contextualization of entity embedding. However, we do observe a faster convergence rate for lower  $L$  values. So for faster convergence, we set  $L$  to 3. Other hyper-parameter selections related to training with best performance across all datasets are shown in Table 5.

### A.3 Dataset and Graph Statistics

We explore the same training/dev/testing split as used by FiD and RAG (Izacard and Grave, 2021; Lewis et al., 2020). We further analyze the geometry of our constructed graphs for each pair of question and its corresponding retrieved passage, from the perspective of the number of available graphs, average node count and node degree, shown in Table 6. We notice that the number of constructed graphs per question is proportional to the percentage of questions favorable from factual KG triplets, and further proportional to the performance gain introduced by GRAPE. The statistics of our graphs align with the source of explored datasets: TriviaQA and WebQ are more entity-focused compared with NQ. Besides the number of graphs per question, we do not observe other difference in graphs among these three datasets.



Methods	GRAPE (base)			GRAPE (large)		
	NQ	TriviaQA	WebQ	NQ	TriviaQA	WebQ
Computing resources	8x32GB Nvidia Tesla V100			8x32GB Nvidia Tesla V100		
Peak memory cost	Around 30.08GB (94%)			Around 26.56GB (83%)		
Peak learning rate	1e-4	6e-5	1e-4	3e-5	3e-5	6e-5
learning optimizer	AdamW with 2,000 warmup			AdamW with 2,000 warmup		
Batch size (per device)	3	3	3	1	1	1
Total training steps	50K	30K	20K	50K	30K	20K
Best dev iterations	19,500	16,500	20,000	45,000	30,000	27,500
Best dev performance	48.17	65.67	51.00	50.73	69.46	57.00

Table 5: Best training hyper-parameters for results reported in Table 2.

Dataset	Split	# Graphs per Q	# Q Entity Node	# P Entity Node	# Nodes per Graph
NQ	Train	29.5 ± 25.9	1.2 ± 0.4	2.7 ± 2.8	3.9 ± 3.0
	Dev	29.0 ± 25.8	1.2 ± 0.4	2.7 ± 2.8	3.8 ± 2.9
	Test	28.4 ± 26.0	1.2 ± 0.5	2.7 ± 2.9	3.9 ± 3.0
TriviaQA	Train	34.8 ± 28.3	1.4 ± 0.7	2.7 ± 2.8	4.1 ± 3.1
	Dev	34.6 ± 28.7	1.3 ± 0.7	2.8 ± 2.8	4.2 ± 3.0
	Test	34.0 ± 28.3	1.4 ± 0.7	2.7 ± 2.8	4.2 ± 3.1
WebQ	Train	42.2 ± 27.8	1.1 ± 0.3	2.8 ± 2.7	3.9 ± 2.8
	Dev	43.9 ± 27.9	1.1 ± 0.4	2.9 ± 2.9	4.0 ± 3.0
	Test	41.2 ± 26.9	1.1 ± 0.3	2.7 ± 2.7	3.8 ± 2.8

Table 6: Statistics of our constructed graphs for all datasets. Mean and standard deviation are calculated for each attribute. 100 passages are retrieved for each questions.