

Dynamic Few-Shot Learning for Knowledge Graph Question Answering

Anonymous ACL submission

Abstract

Large language models present opportunities for innovative Question Answering over Knowledge Graphs (KGQA). However, they are not inherently designed for query generation. To bridge this gap, solutions have been proposed that rely on fine-tuning or ad-hoc architectures, achieving good results but limited out-of-domain distribution generalization. In this study, we introduce a novel approach called Dynamic Few-Shot Learning (DFSL). DFSL integrates the efficiency of in-context learning and semantic similarity and provides a generally applicable solution for KGQA with state-of-the-art performance. We run an extensive evaluation across multiple benchmark datasets and architecture configurations.

1 Introduction

The growth of the Semantic Web has led to the creation and storage of vast amounts of structured knowledge (Hitzler, 2021; Shadbolt et al., 2006), organized into massive Knowledge Graphs (KGs) such as Wikidata (Pellissier Tanon et al., 2016), DBpedia (Lehmann et al., 2014), and FreeBase (Bollacker et al., 2008). The scale of these KGs, with over 109 million items in Wikidata alone,¹ has made extracting relevant information from them increasingly challenging. This led to the emergence of Knowledge Graph Question Answering (KGQA), whose goal is to answer natural language questions posed over KGs.

A typical KGQA system consists of three main components: Entity Linking (EL), Relation Linking (RL), and Query Generation (QG). Starting from a natural language question q , EL and RL return a set of entities \mathcal{E}_q and relations \mathcal{R}_q therein. The QG module, crucially, takes q , \mathcal{E}_q and \mathcal{R}_q and generates a SPARQL query that produces the answer.

¹<https://www.wikidata.org/wiki/Wikidata:Statistics>

This paper focuses on the QG component. State-of-the-art approaches to SPARQL query generation are based on fine-tuning language models like T5 (Qi et al., 2024), or ad-hoc architectures leveraging LLMs and dependency trees (Rony et al., 2022). Despite their success, such approaches have limited flexibility and scalability. Fine-tuning in particular may be computationally expensive and struggle with out-of-domain distributions.

This paper proposes a novel approach to KGQA, leveraging in-context learning with Large Language Models (LLMs). The main intuition is that *a significant number of errors could be addressed by making better use of the examples in the training set*. Our methodology, termed Dynamic Few-Shot Learning (DFSL), leverages semantic search to retrieve similar questions from the training set and enrich the prompt accordingly.

To evaluate the performance and robustness of DFSL, we run experiments on two widely-used Knowledge Bases, DBpedia and Wikidata, using four publicly available datasets: QALD-9, based on DBpedia, and QALD-9 plus, QALD-10 and LC-QuAD 2.0, based on Wikidata. As backbones, we use three state-of-the-art LLMs: Mixtral 8x7B, Llama-3 70B, and CodeLlama 70B. Our experimental results demonstrate that our model achieves new state-of-the-art results, with significant advantages in terms of speed and efficiency. We also run ablation studies to gauge the effectiveness of the approach without gold information from the EL and RL modules.

Our main contributions are: (1) a novel approach to KGQA, called DFSL, that leverages semantic search for dynamic few-shot learning; (2) state-of-the-art results by a significant margin in most benchmarks; (3) an extensive evaluation and ablation study to investigate quantitatively and qualitatively the impact of hyperparameters, backbones, embedding methods, answer selection strategies and gold entity/relation information.

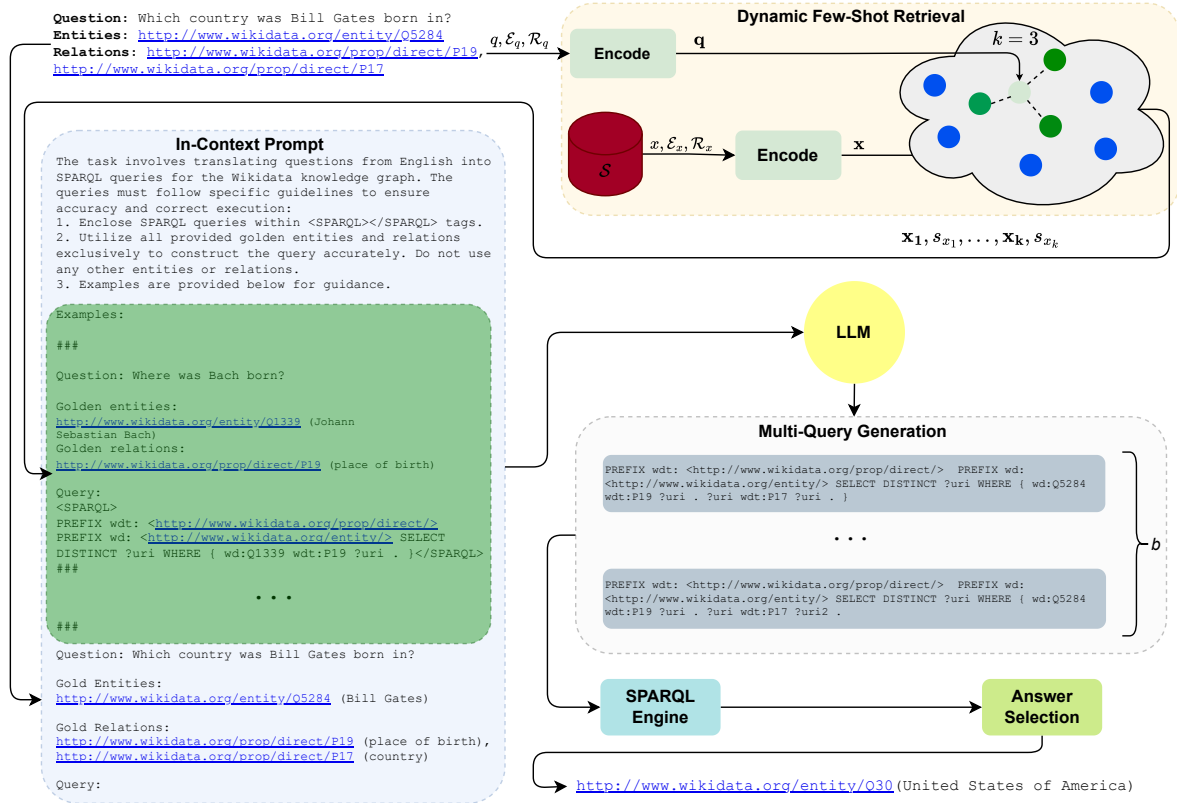


Figure 1: Sketch of DFSL. Given a question, its entities and its relations, k -most similar examples are retrieved from a text-to-SPARQL collection \mathcal{S} and injected into the in-context prompt. Then, the LLM generates one or more queries that are all executed by a SPARQL engine. An answer selection strategy identifies which response to pick.

2 Related work

Early research in KG query generation was rule-based (Guo et al., 2005; Owens et al., 2008), template-based (Zenz et al., 2009; Unger et al., 2012) or search-based. For example, Görlitz et al. (2012) developed a query generation heuristic to predict the final SPARQL representations by exhaustively checking all possible combinations of query patterns. However, manual or semi-manual approaches hit scalability issues with KGs like WikiData and DBpedia. More recent approaches belong to two main streams: information-retrieval based methods and Text-to-SPARQL approaches.

Information Retrieval KGQA. This family of methods involves the identification of sub-graphs relevant to q . Approaches include divide-and-conquer (Kim et al., 2023), fact retrieval based on linked entities (Baek et al., 2023), more complex methods involving hops, relation predictions, and triple sampling (Wu et al., 2023), or Evidence Pattern Retrieval (EPR) through structural dependency modeling (Ding et al., 2024).

Text-to-SPARQL. With the recent wave of decoder-based LLMs such as GPT (Brown et al., 2020), Mixtral (Jiang et al., 2024), and LLaMA (Touvron et al., 2023), generative AI was also used to translate q into SPARQL queries. Notably, Zou et al. (2021) introduced a text-to-SPARQL model that leverages a relation-aware attention decoder and a pointer network encoder, incorporating three separate scaled dot-product attention mechanisms to generate SPARQL queries that capture entity, relation, and keyword representations. Banerjee et al. (2022) experimented with various models, including T5 (Raffel et al., 2020), BART (Lewis et al., 2019), and Pointer Generation Networks (See et al., 2017), to explore their efficacy in KGQA tasks. Rony et al. (2022)’s SGPT employs a stack of transformer encoders to extract linguistic features from q and GPT-2 as a decoder. However, this architecture is limited by its inability to capture connections among entities and relations in the underlying knowledge graph, leading to errors in generating triple sequences in the final SPARQL queries. Pliukhin et al. (2023) presented a one-shot generative approach, where the prompt is aug-

mented with a KG fragment required to construct the query and a question-subgraph query example.

Despite promising results, these architectures are prone to systematic errors. One such error, the so-called “triple-flip”, refers to the reversal of subject and object positions in the generated SPARQL triples, yielding wrong, often empty answers. Qi et al. (2024) addressed this issue by developing TSET, a fine-tuned T5 model with a pretraining stage called Triplet Structure Correction. This approach aims to deepen the model’s understanding of triple order, establishing state-of-the-art performance on major KGQA datasets.

Example Selection in Few-Shot Learning. In-context learning (ICL) is a paradigm that leverages reasoning through analogies. A task description, question, and demonstration context are usually concatenated to create a prompt, which is then input into an LLM for prediction. Unlike fine-tuning, ICL performs predictions without gradient updates (Dong et al., 2023). Few-Shot Learning is a type of ICL where the demonstration context includes a few examples. Owing to the effectiveness of ICL and the obvious advantage of building systems that don’t need domain-specific training, a great deal of research and engineering efforts have been devoted to designing suitable prompts. ICL has been successfully applied to many NLP problems, including QA (Chada and Natarajan, 2021; Chen et al., 2023) and KGQA (Li et al., 2023). Some studies have also focused on the selection of in-context examples. In particular, Liu et al. (2022) developed KATE, an unsupervised retriever that utilizes k-nearest neighbors and distance metrics (e.g., L2 distance and cosine similarity) to select in-context examples for tasks such as sentiment analysis, table-to-text generation, and question answering. Levy et al. (2023) explored the incorporation of diverse demonstrations into prompts for compositional semantic parsing task, demonstrating that such diversity leads to better structural coverage in target utterances. Kim et al. (2022) leveraged the generative capabilities of pre-trained language models to generate demonstrations for each class in downstream tasks, conditioned on test inputs and class information. Gonen et al. (2022) found that selecting examples based on perplexity, in particular lower perplexity, is an effective strategy. However, to the best of our knowledge, example selection has not yet been applied to KGQA.

Text-to-SQL. A cognate domain, text-to-SQL, aims at the translation of natural language questions to SQL queries. There, Rajkumar et al. (2022) demonstrated a zero-shot and few-shot approach using simple prompts, achieving lower results compared to fine-tuned approaches with models such as GPT-3 (Brown et al., 2020) and CODEX (Chen et al., 2021). Nan et al. (2023) introduced various strategies for selecting examples based on similarities/dissimilarities, selecting similar questions with the same difficulty level and dissimilar questions by using k-means clustering to obtain k diverse examples close to each centroid. More recently, Zhang et al. (2023) proposed an automatic chain-of-thought (Wei et al., 2023) approach, where question slices are matched with all possible table and column names to identify the most relevant ones for a given question, using models such as GPT 3.5 and GPT 4. In spite of the similarities between text-to-SQL and text-to-SPARQL, the methods developed so far for the former are not applicable in the latter, where instead of a data model with a relatively small-sized set of tables and columns, the domain is modeled by a large-scale, semi-structured KG.

3 Method

Given a collection of natural language questions \mathcal{Q} and a knowledge graph $\mathcal{G} := (\mathcal{E}, \mathcal{R}, \mathcal{F})$, where \mathcal{E} are *entities*, \mathcal{R} are *relations*, and $\mathcal{F} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ are *facts*, KGQA is the problem of answering questions in \mathcal{Q} based on \mathcal{G} . KGQA can be framed as a **text-to-SPARQL** task, where a question q must be translated into a SPARQL query s_q to be executed on \mathcal{G} by a SPARQL engine, to return a (possibly empty) answer a . The entities and relations in q , denoted as \mathcal{E}_q and \mathcal{R}_q , may be, and usually are, extracted from q before generating s_q . Hence, query generation can be tackled as a conditional text generation problem given q, \mathcal{E}_q and \mathcal{R}_q . Within the scope of ICL, P_θ is a pre-trained LLM and the conditional input $\mathcal{E}_q, \mathcal{R}_q, q$ is combined with other contextual information C , such as additional instructions, guidelines, constraints and demonstrations, all expressed via natural language text. Accordingly, the generated query is:

$$s_q = \arg \max_s P_\theta(s|C, \mathcal{E}_q, \mathcal{R}_q, q). \quad (1)$$

3.1 Dynamic Few-Shot Retrieval

In few-shot ICL, the choice of demonstrations to inject in the prompt can significantly affect perfor-

Approach	Backbone	QALD-9 Plus	QALD-10	LC-QUAD 2.0	QALD-9 DB
Zero-shot Learning		49.90	33.76	40.66	65.73
Few-shot Learning	Mixtral 7x8	54.80 (+4.90)	50.26 (+16.50)	61.04 (+20.38)	63.86 (-1.87)
DFSL		71.75 (+21.85)	49.90 (+16.14)	81.81 (+41.15)	72.74 (+7.01)
Zero-shot Learning		63.01	58.31	54.21	70.49
Few-shot Learning	Llama-3 70B	67.69 (+4.68)	51.28 (-7.03)	68.52 (+14.31)	68.84 (-1.65)
DFSL		73.60 (+10.59)	56.59 (-1.72)	81.93 (+27.72)	72.66 (+2.17)
Zero-shot Learning		45.94	33.36	38.40	66.43
Few-shot Learning	CodeLlama 70B	64.49 (+18.55)	57.38 (+24.02)	64.46 (+26.06)	72.67 (+6.24)
DFSL		76.59 (+30.65)	57.69 (+24.33)	85.45 (+47.05)	75.14 (+8.71)

Table 1: Comparison between zero-shot, few-shot and DFSL with different backbones. Absolute F1 gains with respect to the naive zero-shot approach are reported between parenthesis.

Approach	QALD-9 Plus	QALD-10	LC-QUAD 2.0	QALD-9 DB
DFSL	76.59	57.69	85.45	75.14
DFSL-MQP _{LS}	73.67	58.85	85.06	73.25
DFSL-MQP _{FS}	74.40	58.34	85.38	73.92
DFSL-MQ _{LS}	83.21	60.48	85.54	72.06
DFSL-MQ _{FS}	84.45 (+7.86)	62.20 (+4.51)	89.10 (+3.65)	77.89 (+2.75)

Table 2: Multi-query Generation: comparing DFSL-MQ with DFSL and Multi-query prompting baselines. Absolute F1 gains with respect to DFSL are reported for the best performing configuration.

mance. Usually, few-shot examples are predetermined representative instances of the task, hand-picked during prompt design. Conversely, we aim to retrieve good examples dynamically, based on their relevance to the input question. Inspired by Liu et al. (2022), we adopt a retrieval approach based on the similarity between a question q and a set of previously answered text-to-SPARQL examples collected in a storage \mathcal{S} (see Figure 1), where each example is a tuple including a question x , its entities \mathcal{E}_x and relations \mathcal{R}_x , and the associated SPARQL query s_x . The question, its entities and relations $\langle q, \mathcal{E}_q, \mathcal{R}_q \rangle$ are mapped onto a vector representation $e_q \in \mathbb{R}^d$ using a sentence encoder. To properly feed such information to an encoder-only LM, we concatenate question, entities and relations in a single input sequence $\mathbf{q} := [q, \mathcal{E}_q, \mathcal{R}_q]$. Likewise, we encode each example $x \in \mathcal{S}$ into a vector $e_x \in \mathbb{R}^d$ and then compute the similarity between the target question and the storage:

$$\text{score}(\mathbf{q}, \mathbf{x}) = \text{sim}(e_q, e_x), \forall x \in \mathcal{S}, \quad (2)$$

where the sim is a similarity function. Based on such a scoring, we retrieve the k -most similar examples \mathcal{S} and include them as demonstrations in the in-context prompt.

3.2 In-Context Prompt

The in-context prompt has three parts. The first is the task description, instructing the LLM with a numbered list of guidelines on the output format and on the available information. The second, highlighted in Figure 1 with a green block, contains the k retrieved demonstrations. Each demonstration consists of a question, its entities and relations, denoted as *gold* entities/relations, all paired with their SPARQL query delimited by `<SPARQL></SPARQL>` tags. The `###` symbol delimits each single example. The final part is the input question, associated with its gold entities and relations. The answer returned by the LLM prompted as such is then parsed to extract the generated text enclosed in `<SPARQL></SPARQL>` tags. The resulting query s_q is executed by a SPARQL engine on \mathcal{G} to yield the answer to q . We call our approach Dynamic Few-Shot Learning (DFSL).

3.3 Multi-Query Generation

A typical challenge faced by LLMs in SPARQL query generation is the understanding of what is the subject and what is the object of a relation, an information the model does not have. This problem is called triple-flip error (Qi et al., 2024). LLMs often end up in swapping the subject with the object

in the query, almost choosing one way or the other randomly. Thanks to DFSL, this issue may be alleviated whenever there are similar demonstrations in the in-context prompt that clarify the subject-object roles. To further reduce triple-flip errors, we propose the generation of multiple SPARQL queries by retaining all the final hypotheses generated during beam search. The model uncertainty in placing subject and object is likely to be reflected in the beam search exploration. Intuitively, both triple-ordering hypotheses are considered plausible by the model. Thus, instead of just returning the most probable sequence s according to Equation 1, we keep the whole b queries $\{s_{q,1}, \dots, s_{q,b}\}$ formulated by beam search. We use DFSL-MQ to denote such a multi-query extension of DFSL.

Answer Selection. Executing multiple queries inevitably leads to multiple possible answers. Therefore, we must define an answer selection criterion. We designed two heuristics: Largest Set (LS) and First Set (FS). LS executes all the b queries, obtaining with each query $s_{q,j}$ a (possibly empty) answer set \mathcal{A}_j . LS then selects, among $\{\mathcal{A}_1, \dots, \mathcal{A}_b\}$, the largest one², i.e:

$$\mathcal{A} = \arg \max_{\mathcal{A}_i} (|\mathcal{A}_1|, \dots, |\mathcal{A}_b|),$$

the rationale being that incorrect candidates will likely have empty results. However, LS can be misled into selecting answers from under-constrained queries that return many irrelevant instances. FS adheres to the natural beams ordering by selecting the first query that yields a non-empty answer set.

4 Experiments

In this section, we aim to study the effects of each component involved in our DFSL approach. We evaluate DFSL and its extension DFSL-MQ on four KGQA datasets. In our investigation, we consider different backbones and we compare with multiple baselines and state-of-the-art solutions.

4.1 Datasets

To assess the flexibility and robustness of our approach, we evaluate it on four heterogeneous KGQA benchmarks based on two different Knowledge Graphs (Wikidata, DBpedia).

QALD-9 DB. QALD-9 (Ngomo, 2018) is a dataset from the Question Answering over Linked

Data (QALD) challenge series. It comprises 408 training questions and 150 test questions. Unlike the other KGQA benchmarks, the SPARQL queries are meant for a DBpedia Knowledge Graph. We refer to it as QALD-9 DB to emphasize that.

QALD-9 plus. QALD-9 plus extends QALD-9 on new languages and transfers SPARQL queries from DBpedia to Wikidata. Although some queries were not portable to Wikidata due to the absence of corresponding information, it still comprises 371 training questions and 136 test questions. In our experiments, we only consider English questions.

QALD-10. QALD-10 (Usbeck et al., 2023) is the latest dataset in the QALD series, designed to increase the complexity of gold SPARQL queries. It consists of 412 training questions extracted from QALD-9 plus Wikidata. The test set was created from scratch, comprising 394 test questions that express real-world information needs. Test questions significantly differ from those in training.

LC-QuAD 2.0. LC-QuAD 2.0 (Dubey et al., 2019) is a large-scale dataset grounded on Wikidata. It consists of 30,226 simple and complex questions: 24,180 in training, and 6,046 in test. Questions are diverse. They include single- and multi-fact, boolean, count, and other query types. LC-QuAD 2.0 allows us to gauge the DFSL performance against a large text-to-SPARQL storage.

4.2 Backbones

Mixtral 8x7B. Based on the Sparse Mixture of Experts (SMoE) architecture (Fedus et al., 2022), Mixtral 8x7B (Jiang et al., 2024) is a 46.7B parameters model. Among the backbones adopted in this paper, Mixtral is the smallest. Moreover, thanks to the characteristics of its SMoE architecture, less than 13B are active at each inference step, making Mixtral particularly efficient.

Llama-3 70B. Built upon the Llama architecture (Touvron et al., 2023), Llama-3 70B has been trained on 15T tokens, a 650% increase from its predecessor, Llama 2. At the moment we are writing, Llama-3 70B is one of the best-performing open-weights LLMs available.

CodeLlama 70B. Initialized from Llama2 70B, CodeLlama (Rozière et al., 2024) is a specialized version fine-tuned on 1T tokens of code-heavy data. Therefore, we expect CodeLlama to be particularly suitable for SPARQL query generation.

²In case of ties, we take the first largest set.

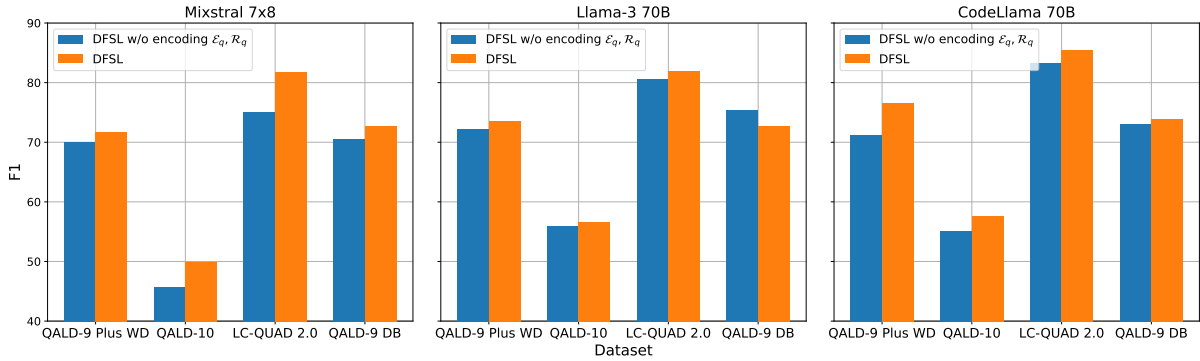


Figure 2: Comparison of Embeddings: DFSL (in orange) encoding that incorporates question, entities and relations versus an embedding solely based on the question q (in blue).

4.3 Baselines

Plain Question. This is a naive baseline where we feed an LLM only with the task description and the question q . Without in-context examples nor any entity or relation associated with q , the LLM can only rely on its parameter memory.

Zero-Shot Learning. Here we do not provide any demonstrative example in the prompt. However, unlike the plain question baseline, we do inject golden entities and relations into the prompt. With reference to Figure 1, the In-Context prompt remains the same but without the green-like block containing the demonstrations.

Few-Shot Learning. The prompt is filled with a single set of k manually selected examples, used for all the questions in the test set. The examples were chosen to maximize diversity and cover different kinds of queries³.

Multi Query Prompting (DFSL-MQP). As an alternative to our proposed multi-query generation (DFSL-MQ), this baseline consists in a naive multi-query prompting strategy. Essentially, we ask the model to generate more queries to answer the question. To ease the creation of inverted subject-object queries that can solve triple-flip errors, we extend the prompt to explicitly ask the model to produce this kind of SPARQL queries. Answer selection uses LS and FS heuristics, like with DFSL-MQ.

4.4 Experimental Setup

Implementation. In our experiments, the training set of each dataset serves as storage for the retrieval of the k most similar examples (see the next

paragraph for details on k tuning) with DFSL. Examples are encoded with a sentence transformer⁴, all-mpnet-base-v2⁵, and sim is defined as the cosine similarity. Inference is performed via beam search in both DFSL, where b is set to 3, and DFSL-MQ, with b set to 10. All the experiments were run on a cluster of 4 NVIDIA A100 GPUs.

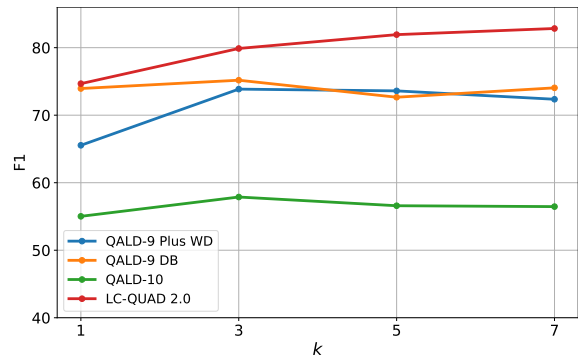


Figure 3: Impact of the number of in-context examples on the four benchmarks.

Number of Few-shot Examples. We first analyzed how the number of few-shot examples k retrieved by DFSL affects the performance. We chose among $k = \{1, 3, 5, 7\}$ and evaluated DFSL with Llama 3 70B backbone on the four datasets. The results shown in Figure 3 suggest that values of k greater than one perform comparably well on smaller benchmarks, while on LC-QUAD 2.0, where there are about 25 thousands examples as storage, increasing k seems to be beneficial. This may be due to the increased likelihood of finding similar examples in larger datasets as k grows. We set $k = 5$ for all the forthcoming experiments, which is a good trade-off across all the datasets.

⁴<https://www.sbert.net/index.html>

⁵<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

³The chosen examples and more details are provided in Appendix B.

Approach	QALD-9 Plus	QALD-10	LC-QUAD 2.0	QALD-9 DB
Plain Question	0.08	0.02	12.00	16.42
BART (Banerjee et al., 2022)	-	-	64.00	-
PGN-BERT-BERT (Banerjee et al., 2022)	-	-	86.00	-
SGPT (Rony et al., 2022)	-	-	89.04	67.82
TSET-small (Qi et al., 2024)	72.86	47.15	94.00	-
TSET-base (Qi et al., 2024)	75.85	51.37	95.00	-
Zero-shot Learning	45.94	33.36	38.40	66.43
Few-shot Learning	64.49	57.38	64.46	72.67
DFSL	76.59	57.69	85.45	75.14
DFSL-MQ beam FS	84.45 (+8.60)	62.20 (+10.83)	89.10 (-5.90)	77.89 (+10.07)

Table 3: DFSL and ICL approaches vs state-of-the-art fine-tuned models.

Prompt. The prompt illustrated in Figure 1 constitutes the default template in our experimentations. However, slight variations are required in certain cases. For example, when running experiments on DBpedia knowledge graph, we replace the Wikidata reference with DBpedia in the first text segment. When we study the absence of gold information instead, we remove all the references to gold entities/relations (according to the ablation) from the entire prompt. There are no differences in the prompts layout when running few-shot-learning baseline experiments. In zero-shot learning, only the in-context examples any reference to them are removed, all else being equal.

Evaluation metric. We follow a standard F1 score evaluation in KGQA benchmarks. The F1 is computed between the answer set returned by the target SPARQL query and the predicted one. When both the queries return an empty set, we assign an F1 score of 1. The F1 scores of all the examples are then averaged.

4.5 Results

Impact of Dynamic Examples. To measure the importance of retrieving few-shot examples dynamically, we compare DFSL on different backbones against Zero-Shot and Few-Shot Learning baselines. Results are outlined in Table 1.

In terms of backbones, Llama 3 consistently outperforms both Mixtral and CodeLlama in zero-shot learning scenario, whereas in few-shot, results are generally comparable between Llama-3 and CodeLlama. Such a strong Llama 3 zero-shot performance may be caused by some sort of data contamination, however we leave such an investigation for future works.

Both few-shot learning and DFSL generally yield substantial gains with respect to zero-shot baseline on all the backbones and datasets. An

exception occurs in QALD-10 with Llama-3. Notably, when comparing DFSL and Few-shot Learning baseline, we can see our approach improving F1 scores by a large margin in LC-QUAD 2.0, QALD-9 Plus and QALD-9 DB, with F1 increasing up to 21 absolute points⁶. In QALD-10 instead, where the test set has a different distribution from its training, there are no significant differences between DFSL and the standard few-shot learning approach. Indeed, an approach like DFSL brings little benefits when the storage only contains unrelated examples.

Overall, DFSL with CodeLlama3 achieved the greatest performance with respect to all the other configurations. Therefore, we adopt CodeLlama as our backbone in the following DFSL experiments.

Impact of Multi-Query Generation. Here we investigate DFSL-MQ, the multi-query approach extending DFSL. We evaluate both answer selection strategies, LS and FS, and compare them against the plain DFSL and the multi-query prompting baseline described in Section 4.3. All the results are outlined in Table 2.

Having multiple queries is not necessarily beneficial. Indeed, the multi-query prompting baseline under-performs in three datasets out of four with respect to (single query) DFSL, regardless of the answer selection method adopted. DFSL-MQ instead proves to be generally beneficial. Both Largest Set and First Set heuristics are effective when the hypotheses come from the beams. Furthermore, FS consistently outperforms LS, even by substantial margins in QALD-9 DB.

In-context Learning vs Fine-tuning. Up to this point, we have assessed DFSL in the scope of In-Context Learning approaches. In Table 3 instead, we compare our approach against state-of-

⁶Some qualitative examples illustrate the benefits of DFSL over few-shot learning in Appendix A (see Table 6).

the-art models trained and/or fine-tuned for specific downstream KGQA datasets. Without any training, DFSL-MQ outperforms current state-of-the-art approaches in three out of four benchmarks, namely QALD-9 Plus, QALD-10 and QALD-9 DB, even with the single query DFSL setup. DFSL-MQ does not obtain state-of-the-art results in LC-QUAD 2.0, the dataset most affected by triple-flip errors. This means that multi-query generation only alleviates the issue, but the problem still remains.

4.6 Ablation studies

Different Example Encoding. As described in Section 3.1, to compute the embeddings we concatenated the textual input made of the question and its list of entities and relations. Here, we gauge the impact of this additional information on DFSL performance. In Figure 2 we compare DFSL, with a variant where we only embed the natural language question q , without any additional data concatenated. The evaluation carried out in all the benchmarks and with all the backbones, demonstrates that such information improves the quality of the generated queries.

Approach	QALD-9 DB
Plain Question	16.42
DFSL	75.14
DFSL w/o \mathcal{R}_q	56.62 (-18.56)
DFSL w/o \mathcal{E}_q	60.92 (-14.22)
DFSL w/o $\mathcal{E}_q, \mathcal{R}_q$	49.59 (-25.55)

Table 4: DFSL in the absence of entities and/or relations.

Absence of gold information. In KGQA, text-to-SPARQL generation usually relies not only on the question itself, but also on entities and relations associated to it. Here we assess DFSL when either the entities \mathcal{E}_q or the relations \mathcal{R}_q , or both are missing. The information is removed throughout the entire process. For example, when removing entities, we discard them from both the storage and the prompt. Even the embeddings for the retrieval are computed by encoding an input without any entity concatenated in q , i.e. becoming $q = [q, \mathcal{R}_q]$. We report this on QALD-9 DB dataset. By observing the results outlined in Table 4, it is clear that, without full knowledge of the entities and the relations required for generating the query, the LLM perfor-

mance drops significantly. Nonetheless, even in the case where no information is given (DFSL w/o $\mathcal{E}_q, \mathcal{R}_q$), the presence of dynamic demonstrations is essential, yielding a 33+ absolute F1 increase compared to plain question baseline.

5 Conclusion

In this paper, we introduced DFSL, a novel approach to Knowledge Graph Question Answering. This method leverages semantic search to dynamically retrieve relevant examples from the training set, enriching the prompt for LLMs to improve the generation of SPARQL queries. We conducted comprehensive experiments on four publicly available datasets based on two widely-used KBs, DBpedia and Wikidata. By employing three different state-of-the-art LLMs as backbones, we demonstrated that DFSL achieves superior performance compared to both standard in-context learning techniques and state-of-the-art models fine-tuned on the downstream task. We further conducted an extensive evaluation of DFSL through ablation studies to measure the impact of hyper-parameters, different backbones, embedding methods, answer selection strategies, and the inclusion or exclusion of entities and relations information associated to a question. The code will be released publicly upon acceptance of the paper. In the future, we plan to study the effectiveness of DFSL in cognate domains like text-to-SQL.

Limitations

We recognize some limitations in our work. Our experiments are all on English-based datasets, where notoriously LLMs are better performing. Moreover, the massive pre-training of those LLMs on a vast portion of the Web, may expose those models to unintended data contamination. Experiments only focused on LLMs with large number of parameters, without investigating the behaviour of smaller models. To encode examples, we limited the investigation to what kind of text to encode (just the question, or the question and its entities and relations), without exploring different embedding models, similarity criteria or other input concatenation strategies. We leave these investigations to future work.

References

Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. Knowledge-augmented language model prompting

580		for zero-shot knowledge graph question answering. Preprint, arXiv:2306.04136.	
581			
582	Debayan Banerjee, Pranav Ajit Nair, Jivat Neet Kaur,		
583	Ricardo Usbeck, and Chris Biemann. 2022. Modern		
584	baselines for sparql semantic parsing . In Proceedings		
585	of the 45th International ACM SIGIR Conference on		
586	Research and Development in Information Retrieval,		
587	SIGIR '22 . ACM.		
588	Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim		
589	Sturge, and Jamie Taylor. 2008. Freebase: a col-		
590	laboratively created graph database for structuring		
591	human knowledge . In Proceedings of the 2008 ACM		
592	SIGMOD International Conference on Management		
593	of Data, SIGMOD '08 , page 1247–1250, New York,		
594	NY, USA. Association for Computing Machinery.		
595	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie		
596	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind		
597	Neelakantan, Pranav Shyam, Girish Sastry, Amanda		
598	Askell, Sandhini Agarwal, Ariel Herbert-Voss,		
599	Gretchen Krueger, Tom Henighan, Rewon Child,		
600	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,		
601	Clemens Winter, Christopher Hesse, Mark Chen,		
602	Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin		
603	Chess, Jack Clark, Christopher Berner, Sam Mc-		
604	Candlish, Alec Radford, Ilya Sutskever, and Dario		
605	Amodei. 2020. Language models are few-shot learn-		
606	ers . CoRR , abs/2005.14165.		
607	Rakesh Chada and Pradeep Natarajan. 2021. Fewshotqa:		
608	A simple framework for few-shot learning of ques-		
609	tion answering tasks using pre-trained text-to-text		
610	models. In Proceedings of the 2021 Conference on		
611	Empirical Methods in Natural Language Processing,		
612	pages 6081–6090.		
613	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming		
614	Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka-		
615	plan, Harri Edwards, Yuri Burda, Nicholas Joseph,		
616	Greg Brockman, Alex Ray, Raul Puri, Gretchen		
617	Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-		
618	try, Pamela Mishkin, Brooke Chan, Scott Gray,		
619	Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz		
620	Kaiser, Mohammad Bavarian, Clemens Winter,		
621	Philippe Tillet, Felipe Petroski Such, Dave Cum-		
622	ings, Matthias Plappert, Fotios Chantzis, Eliza-		
623	beth Barnes, Ariel Herbert-Voss, William Hebgan		
624	Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie		
625	Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,		
626	William Saunders, Christopher Hesse, Andrew N.		
627	Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan		
628	Morikawa, Alec Radford, Matthew Knight, Miles		
629	Brundage, Mira Murati, Katie Mayer, Peter Welinder,		
630	Bob McGrew, Dario Amodei, Sam McCandlish, Ilya		
631	Sutskever, and Wojciech Zaremba. 2021. Evaluat-		
632	ing large language models trained on code . Preprint,		
633	arXiv:2107.03374 .		
634	Xiusi Chen, Yu Zhang, Jinliang Deng, Jyun-Yu Jiang,		
635	and Wei Wang. 2023. Gotta: generative few-shot		
636	question answering by prompt-based cloze data		
637	augmentation. In Proceedings of the 2023 SIAM		
638	International Conference on Data Mining (SDM),		
639	pages 909–917. SIAM.		
	Wentao Ding, Jinmao Li, Liangchuan Luo, and Yuzhong		640
	Qu. 2024. Enhancing complex question answering		641
	over knowledge graphs through evidence pattern re-		642
	trieval . Preprint, arXiv:2402.02175 .		643
	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong		644
	Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and		645
	Zhifang Sui. 2023. A survey on in-context learning .		646
	Preprint, arXiv:2301.00234 .		647
	Mohnish Dubey, Debayan Banerjee, Abdelrahman Ab-		648
	delkawi, and Jens Lehmann. 2019. Lc-quad 2.0: A		649
	large dataset for complex question answering over		650
	wikidata and dbpedia . page 69–78, Berlin, Heidel-		651
	berg. Springer-Verlag.		652
	William Fedus, Jeff Dean, and Barret Zoph. 2022. A		653
	review of sparse expert models in deep learning .		654
	Preprint, arXiv:2209.01667 .		655
	Hila Gonen, Srini Iyer, Terra Blevins, Noah A.		656
	Smith, and Luke Zettlemoyer. 2022. Demystifying		657
	prompts in language models via perplexity estima-		658
	tion . Preprint, arXiv:2212.04037 .		659
	Olaf Görlitz, Matthias Thimm, and Steffen Staab. 2012.		660
	Splodge: Systematic generation of sparql benchmark		661
	queries for linked open data. In The Semantic Web		662
	– ISWC 2012 , pages 116–132, Berlin, Heidelberg.		663
	Springer Berlin Heidelberg.		664
	Yuanbo Guo, Zhengxiang Pan, and Jeff Heflin. 2005.		665
	Lubm: a benchmark for owl knowledge base systems .		666
	Web Semantics: Science, Services and Agents on		667
	the World Wide Web , 3:158–182.		668
	Pascal Hitzler. 2021. A review of the semantic web		669
	field . Commun. ACM , 64(2):76–83.		670
	Albert Q. Jiang, Alexandre Sablayrolles, Antoine		671
	Roux, Arthur Mensch, Blanche Savary, Chris		672
	Bamford, Devendra Singh Chaplot, Diego de las		673
	Casas, Emma Bou Hanna, Florian Bressand, Gi-		674
	anna Lengyel, Guillaume Bour, Guillaume Lam-		675
	ple, L�lio Renard Lavaud, Lucile Saulnier, Marie-		676
	Anne Lachaux, Pierre Stock, Sandeep Subramanian,		677
	Sophia Yang, Szymon Antoniak, Teven Le Scao,		678
	Th�ophile Gervet, Thibaut Lavril, Thomas Wang,		679
	Timothe�e Lacroix, and William El Sayed. 2024. Mix-		680
	tral of experts . Preprint, arXiv:2401.04088 .		681
	Hyuhng Joon Kim, Hyunsoo Cho, Junyeob Kim, Taek		682
	Kim, Kang Min Yoo, and Sang goo Lee. 2022.		683
	Self-generated in-context learning: Leveraging auto-		684
	regressive language models as a demonstration gen-		685
	erator . Preprint, arXiv:2206.08082 .		686
	Jiho Kim, Yeonsu Kwon, Yohan Jo, and Edward Choi.		687
	2023. Kg-gpt: A general framework for reasoning		688
	on knowledge graphs using large language models .		689
	Preprint, arXiv:2310.11220 .		690
	Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch,		691
	Dimitris Kontokostas, Pablo Mendes, Sebastian Hell-		692
	mann, Mohamed Morsey, Patrick Van Kleef, S�ren		693

694	Auer, and Christian Bizer. 2014. Dbpedia - a large-scale, multilingual knowledge base extracted from wikipedia . Semantic Web Journal , 6.	
695		
696		
697	Itay Levy, Ben Bogin, and Jonathan Berant. 2023. Diverse demonstrations improve in-context compositional generalization . Preprint , arXiv:2212.06800.	
698		
699		
700	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2019. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . CoRR , abs/1910.13461.	
701		
702		
703		
704		
705		
706	Tianle Li, Xueguang Ma, Alex Zhuang, Yu Gu, Yu Su, and Wenhui Chen. 2023. Few-shot in-context learning on knowledge base question answering . In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) , pages 6966–6980.	
707		
708		
709		
710		
711		
712	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. What makes good in-context examples for GPT-3? In Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures , pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.	
713		
714		
715		
716		
717		
718		
719		
720	Linyong Nan, Yilun Zhao, Weijin Zou, Narutatsu Ri, Jaesung Tae, Ellen Zhang, Arman Cohan, and Dragomir Radev. 2023. Enhancing few-shot text-to-sql capabilities of large language models: A study on prompt design strategies . Preprint , arXiv:2305.12586.	
721		
722		
723		
724		
725		
726	Ngonga Ngomo. 2018. 9th challenge on question answering over linked data (qald-9). language , 7(1):58–64.	
727		
728		
729	Alisdair Owens, Nick Gibbins, and m.c Schraefel. 2008. Effective benchmarking for rdf stores using synthetic data .	
730		
731		
732	Thomas Pellissier Tanon, Denny Vrandečić, Sebastian Schaffert, Thomas Steiner, and Lydia Pintscher. 2016. From freebase to wikidata: The great migration . In Proceedings of the 25th international conference on world wide web , pages 1419–1428.	
733		
734		
735		
736		
737	Dmitrii Pliukhin, Daniil Radyush, Liubov Kovrigina, and Dmitry Mouromtsev. 2023. Improving subgraph extraction algorithms for one-shot sparql query generation with large language models . In Scholarly-QALD-23: Scholarly QALD Challenge at The 22nd International Semantic Web Conference (ISWC 2023) (Athens, Greece, volume 3592, pages 1–10.	
738		
739		
740		
741		
742		
743		
744		
745	Jiexing Qi, Chang Su, Zhixin Guo, Lyuwen Wu, Zanwei Shen, Luoyi Fu, Xinbing Wang, and Chenghu Zhou. 2024. Enhancing sparql query generation for knowledge base question answering systems by learning to correct triplets . Applied Sciences , 14(4).	
746		
747		
748		
749		
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer . Journal of Machine Learning Research , 21(140):1–67.	750
		751
		752
		753
		754
		755
	Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. 2022. Evaluating the text-to-sql capabilities of large language models . Preprint , arXiv:2204.00498.	756
		757
		758
		759
	Md Rashad Al Hasan Rony, Uttam Kumar, Roman Teucher, Liubov Kovrigina, and Jens Lehmann. 2022. Sgpt: A generative approach for sparql query generation from natural language questions . IEEE Access , 10:70712–70723.	760
		761
		762
		763
		764
	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code llama: Open foundation models for code . Preprint , arXiv:2308.12950.	765
		766
		767
		768
		769
		770
		771
		772
		773
		774
	Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks . Preprint , arXiv:1704.04368.	775
		776
		777
	N. Shadbolt, T. Berners-Lee, and W. Hall. 2006. The semantic web revisited . IEEE Intelligent Systems , 21(3):96–101.	778
		779
		780
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models . Preprint , arXiv:2302.13971.	781
		782
		783
		784
		785
		786
		787
	Christina Unger, Lorenz Bühmann, Jens Lehmann, Axel-Cyrille Ngonga Ngomo, Daniel Gerber, and Philipp Cimiano. 2012. Template-based question answering over rdf data . WWW’12 - Proceedings of the 21st Annual Conference on World Wide Web .	788
		789
		790
		791
		792
	Ricardo Usbeck, Xi Yan, Aleksandr Perevalov, Longquan Jiang, Julius Schulz, Angelie Kraft, Cedric Möller, Junbo Huang, Jan Reineke, Axel-Cyrille Ngonga Ngomo, et al. 2023. Qald-10—the 10th challenge on question answering over linked data . Semantic Web , (Preprint):1–15.	793
		794
		795
		796
		797
		798
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models . Preprint , arXiv:2201.11903.	799
		800
		801
		802
		803

804 Yike Wu, Nan Hu, Sheng Bi, Guilin Qi, Jie Ren, An-
805 huan Xie, and Wei Song. 2023. [Retrieve-rewrite-
806 answer: A kg-to-text enhanced llms framework
807 for knowledge graph question answering](#). [Preprint](#),
808 arXiv:2309.11206.

809 Gideon Zenz, Xuan Zhou, Enrico Minack, Wolf Siber-
810 ski, and Wolfgang Nejdl. 2009. [From keywords to
811 semantic queries—incremental query construction
812 on the semantic web](#). [Journal of Web Semantics](#),
813 7(3):166–176. The Web of Data.

814 Hanchong Zhang, Ruisheng Cao, Lu Chen, Hongshen
815 Xu, and Kai Yu. 2023. [ACT-SQL: In-context learn-
816 ing for text-to-SQL with automatically-generated
817 chain-of-thought](#). In [Findings of the Association
818 for Computational Linguistics: EMNLP 2023](#), pages
819 3501–3532, Singapore. Association for Computa-
820 tional Linguistics.

821 Jianyun Zou, Min Yang, Lichao Zhang, Yechen Xu,
822 Qifan Pan, Fengqing Jiang, Ran Qin, Shushu Wang,
823 Yifan He, Songfang Huang, and Zhou Zhao. 2021.
824 [A chinese multi-type complex questions answering
825 dataset over wikidata](#). [Preprint](#), arXiv:2111.06086.

826
827
828
829
830
831
832
833

A Qualitative Analysis

In this appendix we provide some qualitative analyses of DFSL and DFSL-MQ. First of all, we report some examples in Table 6 that highlight the benefits from introducing similar examples with DFSL with respect to standard few-shot learning approach. Then, we show some examples in Table 7 where the multi-query approach solves triple-flip errors. In Table 5 instead, we showcase errors caused by employing LS answer selection heuristic. Notably, by choosing larger sets, LS sometimes selects queries that are often relegated to latter positions in the beam hypotheses, which tend to be more general, thus more prone to returning a greater number of results.

Question	Target Query	Answer Selection	Predicted Query	Beam
What is manufactured NEC PC-9800 series whose sector is electronics?	SELECT ?answer WHERE { wd:Q183505 wdt:P176 ?answer . ?answer wdt:P452 wd:Q11650 }	LS	SELECT ?answer WHERE { wd:Q183505 wdt:P176 ?X . ?X wdt:P452 wd:Q11650 . ?X wdt:P31 ?answer }	7
		FS	select distinct ?obj where { wd:Q183505 wdt:P176 ?obj . ?obj wdt:P452 wd:Q11650 }	4
What is the enthalpy of vaporization for potassium hydroxide?	select distinct ?answer where { wd:Q132298 wdt:P2116 ?answer }	LS	ASK WHERE { wd:Q14982 wdt:P2116 ?obj filter(?obj = 37400) }	3
		FS	select distinct ?answer where { wd:Q132298 wdt:P2116 ?answer }	1
Which is the rural city of Fukushima?	SELECT ?answer WHERE { wd:Q161176 wdt:P131 ?answer . ?answer wdt:P150 wd:Q1347240 }	LS	SELECT ?answer WHERE { wd:Q161176 wdt:P131 ?X . ?X wdt:P150 ?answer }	3
		FS	SELECT ?answer WHERE { wd:Q161176 wdt:P131 ?answer . ?answer wdt:P150 wd:Q1347240 }	1

Table 5: Qualitative comparison between different answer selection strategies in DFSL-MQ.

834
835
836
837

B Few-shot Learning Examples

We report in Figure 4 the examples selected for the Few-shot learning baseline prompt. The five examples were chosen to be the most representative of the training set, including queries of different kind and structure, such as ASK, COUNT and SELECT.

Question	Target Query	Approach	Predicted Query	Similar In-context Example
Who is the daughter of Robert Kennedy married to?	SELECT DISTINCT ?uri WHERE { wd:Q25310 wdt:P40 ?daughter . ?daughter wdt:P21 wd:Q6581072 . ?daughter wdt:P26 ?uri . }	Few-Shot	SELECT DISTINCT ?uri WHERE { ?uri wdt:P40 wd:Q25310 ; wdt:P21 wd:Q6581072 ; wdt:P26 ?spouse . }	-
		DFSL	SELECT DISTINCT ?uri WHERE { wd:Q25310 wdt:P40 ?child . ?child wdt:P21 wd:Q6581072 . ?child wdt:P26 ?uri . }	SELECT DISTINCT ?uri WHERE { wd:Q43247 wdt:P40 ?child . ?child wdt:P26 ?uri . }
Which countries are connected by the Rhine?	SELECT DISTINCT ?uri WHERE { wd:Q584 wdt:P17 ?uri . ?uri wdt:P31 wd:Q6256 . }	Few-Shot	SELECT DISTINCT ?uri WHERE { ?uri wdt:P31 wd:Q6256 ; wdt:P17 wd:Q584 . }	-
		DFSL	SELECT DISTINCT ?uri WHERE { wd:Q584 wdt:P17 ?uri . ?uri wdt:P31 wd:Q6256 . }	SELECT DISTINCT ?res WHERE { wd:Q3392 wdt:P885/wdt:P17 ?res . }
Give me the capitals of all countries in Africa.	SELECT DISTINCT ?uri WHERE { ?country wdt:P31 wd:Q6256 . ?country wdt:P30 wd:Q15 . ?country wdt:P36 ?uri . }	Few-Shot	SELECT DISTINCT ?uri WHERE { ?uri wdt:P31 wd:Q6256 ; wdt:P30 wd:Q15 ; wdt:P36 ?capital . }	-
		DFSL	SELECT DISTINCT ?uri WHERE { ?country wdt:P31 wd:Q6256 . ?country wdt:P30 wd:Q15 . ?country wdt:P36 ?uri . }	SELECT DISTINCT ?uri WHERE { wd:Q5451 wdt:P17 ?country . ?country wdt:P36 ?uri . }

Table 6: A qualitative comparison between DFSL and Few-shot Learning. Thanks to the similar in-context examples retrieved with DFSL, the LLMs generates the correct query.

Question	Target Query	Approach	Predicted Query
Who is the enclave within of Montreal?	select distinct ?subj where { ?subj wdt:P501 wd:Q340 . ?subj wdt:P31 wd:Q171441 }	DFSL	select distinct ?obj where { wd:Q340 wdt:P501 ?obj . ?obj wdt:P31 wd:Q171441 }
		DFSL-MQ	select distinct ?subj where { ?subj wdt:P501 wd:Q340 . ?subj wdt:P31 wd:Q171441 }
The trachea is of what anatomical branch?	select distinct ?answer where { ?answer wdt:P3261 wd:Q175449 }	DFSL	select distinct ?answer where { wd:Q175449 wdt:P3261 ?answer }
		DFSL-MQ	select distinct ?answer where { ?answer wdt:P3261 wd:Q175449 }
What revolution caused the destruction of the Russian Empire?	select distinct ?obj where { wd:Q34266 wdt:P770 ?obj . ?obj wdt:P31 wd:Q10931 }	DFSL	select distinct ?subj where { ?subj wdt:P770 wd:Q34266 . ?subj wdt:P31 wd:Q10931 }
		DFSL-MQ	select distinct ?obj where { wd:Q34266 wdt:P770 ?obj . ?obj wdt:P31 wd:Q10931 }

Table 7: Some triple-flip errors that are solved by DFSL-MQ.

```

Examples:

Question: Give me all companies in Munich.

Entities:
http://www.wikidata.org/entity/q4830453 (business), http://www.wikidata.org/entity/q1726 (Munich)

Relations:
http://www.wikidata.org/prop/direct/p279 (subclass of), http://www.wikidata.org/prop/direct/p31 (instance of),
http://www.wikidata.org/prop/direct/p159 (headquarters location)

Query:
<SPARQL>
PREFIX wdt: <http://www.wikidata.org/prop/direct/> PREFIX wd: <http://www.wikidata.org/entity/> SELECT DISTINCT ?uri WHERE { ?type wdt:P279*
wd:Q4830453 . ?uri wdt:P31 ?type ; wdt:P159 wd:Q1726 . }
</SPARQL>
###

Question: Was Marc Chagall a jew?

Entities:
http://www.wikidata.org/entity/q93284 (Marc Chagall), http://www.wikidata.org/entity/q7325 (Jewish people)

Relations:
http://www.wikidata.org/prop/direct/p172 (ethnic group)

Query:
<SPARQL>
PREFIX wdt: <http://www.wikidata.org/prop/direct/> PREFIX wd: <http://www.wikidata.org/entity/> ASK WHERE { wd:Q93284 wdt:P172 wd:Q7325 . }
</SPARQL>
###

Question: How many films did Leonardo DiCaprio star in?

Entities:
http://www.wikidata.org/entity/q11424 (film), http://www.wikidata.org/entity/q38111 (Leonardo DiCaprio)

Relations:
http://www.wikidata.org/prop/direct/p31 (instance of), http://www.wikidata.org/prop/direct/p161 (cast member)

Query:
<SPARQL>
PREFIX wdt: <http://www.wikidata.org/prop/direct/> PREFIX wd: <http://www.wikidata.org/entity/> SELECT (COUNT(DISTINCT ?uri) AS ?c) WHERE { ?uri
wdt:P31 wd:Q11424 ; wdt:P161 wd:Q38111 . }
</SPARQL>
###

Question: Give me all libraries established earlier than 1400.

Entities:
http://www.wikidata.org/entity/q7075 (library)

Relations:
http://www.wikidata.org/prop/direct/p31 (instance of), http://www.wikidata.org/prop/direct/p571 (inception)

Query:
<SPARQL>
PREFIX wdt: <http://www.wikidata.org/prop/direct/> PREFIX wd: <http://www.wikidata.org/entity/> SELECT DISTINCT ?uri WHERE { ?uri wdt:P31 wd:Q7075
; wdt:P571 ?date . FILTER (YEAR(?date) < 1400 ) }
</SPARQL>
###

Question: Is Christian Bale starring in Batman Begins?

Entities:
http://www.wikidata.org/entity/q166262 (Batman Begins), http://www.wikidata.org/entity/q45772 (Christian Bale)

Relations:
http://www.wikidata.org/prop/direct/p161 (cast member)

Query:
<SPARQL>
PREFIX wdt: <http://www.wikidata.org/prop/direct/> PREFIX wd: <http://www.wikidata.org/entity/> ASK WHERE { wd:Q166262 wdt:P161 wd:Q45772 }
</SPARQL>

```

Figure 4: Examples injected in the Few-shot-learning baseline prompt.