

GenDec: A Generative Question-decomposition method for Multi-hop Question-answering

Anonymous ACL submission

Abstract

Multi-hop QA involves step-by-step reasoning to answer complex questions and find multiple relevant supporting facts. Previous question-decomposition research on multi-hop QA has shown that performance can be boosted by first decomposing questions into simpler, single-hop sub-questions (QD), and then answering them one by one in a specific order. However, such decomposition often leads to error propagation during QA: 1) incorrect QD leads to wrong QA results; 2) wrong answers to a previous sub-question compromise the next sub-question. In this work, we propose GenDec, a generative QD-based model for multi-hop QA from the perspective of explainable QA by generating independent and complete sub-questions based on incorporating additional extracted evidence. This approach first introduces sub-questions in retrieving relevant passages for each hop and fuses features of sub-questions into QA reasoning, which enables it to provide an explainable reasoning process for its answers. We evaluate GenDec by comparing it with existing QD-based and other strong QA models and the results show GenDec outperforms all QD-based multi-hop QA models for answer spans on the HotpotQA, 2WikiHopMultiHopQA and MuSiQue datasets. We also conduct experiments with the large language models (LLMs) ChatGPT and LLaMA to illustrate the impact of QD on QA tasks in the LLM era.

1 Introduction

Multi-hop QA (MQA) is a task that requires multiple reasoning steps over multiple information sources (e.g., text paragraphs). While explicit question decomposition (QD), which involves breaking down complex questions into simpler and more straightforward sub-questions, has long been an approach in developing robust and interpretable question-answering (QA) models and systems, most MQA models, e.g., DFGN (Qiu et al., 2019), Decomprc (Min et al., 2019a), CogQA (Ding

et al., 2019), HGN (Fang et al., 2019b), C2F Reader (Shao et al., 2020a), and BFR-Graph (Huang and Yang, 2021) illustrate how demonstrating the reasoning ability of a model in multi-hop questions remains a challenge. Tang et al. (2020b) proposes a human-verified sub-question dataset derived from HotpotQA (Yang et al., 2018a) and conducts experiments on sub-question reasoning. Their results indicate that DFGN, Decomprc, and CogQA performed badly on answering sub-questions, even when they can correctly answer the final multi-hop question, as it is common for models to bypass the correct reasoning process and fail to reason the intermediate answers to sub-questions.

Thus, understanding and potentially decomposing multi-hop questions into finer-grained sub-questions is a key desired step in QA. To accurately answer a multi-hop question, traditionally QD + QA methods start by decomposing the given multi-hop question into simpler sub-questions, attempting to answer them in a specific order, and then finally aggregating the information obtained from all sub-questions.

Through a preliminary investigation, we find that QD remains a major bottleneck in MQA. Previous QD methods Min et al. (2019b); Perez et al. (2020a) first decompose multi-hop questions into **dependent** sub-questions, e.g., in figure 1, the original question is decomposed into "Who is the record holder for Argentine PGA Championship tournaments?" and "How many tournaments did [Ans of Sub Q1] win?" and QA models need to correctly answer sub-question 1 and fill it into sub-question 2 and then answer it to get the final answer. Such QD+QA method suffers from error propagation, where incorrectly answering any of the sub-questions may lead to a wrong final answer. GenDec mitigates this error-propagation problem during reasoning since the decomposed sub-questions are independent and complete, thus not requiring answers in a specific order as was the case in previ-

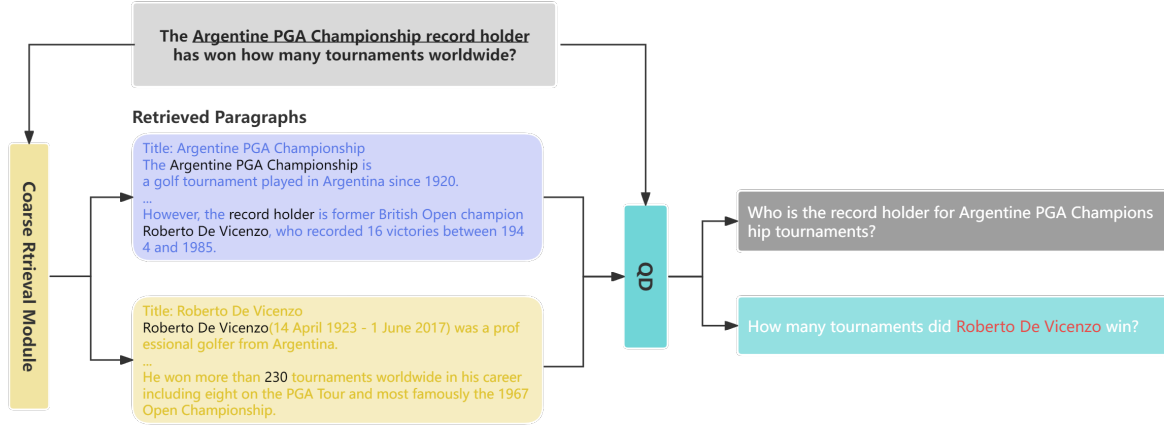


Figure 1: Example of multi-hop and decomposed sub-questions from the HotpotQA dataset. The original question is shown in light grey and the decomposed ones are in deep gray and cyan. "Roberto de Vincenzo" in the retrieved paragraph is the answer to sub-question Q1 and also part of the sub-question Q2. The literal "230" is the answer to sub-question Q2. Since the paragraphs are too long, we here only list the sentences that contain supporting facts.

ous models. We fuse the sub-questions into the QA model to provide the appropriate reasoning chain.

We propose **GenDec**, a generative-based QD method that incorporates retrieved paragraphs including evidence for decomposing independent sub-questions that do not require answers in order. After QD, GenDec combines the features of sub-questions into relevant paragraph retrieval, supporting facts prediction, and QA modules. Figure 1 shows the decomposition results of GenDec over the HotpotQA dataset. The original multi-hop question "The Argentine PGA Championship record holder has won how many tournaments worldwide?" is decomposed into independent sub-questions: "Who is the record holder for Argentine PGA Championship tournaments?" and "How many tournaments did Roberto De Vincenzo win?".

GenDec shows more robustness in answering sub-questions than other QA models as it only needs retrieved paragraphs as extra decomposing information and neither need to consider hop relations or the order of sub-questions. We further evaluate the effectiveness of our system in multi-hop QA to illustrate that QD still plays a vital role in QA in the large language model (LLM) era.

Our contributions are as follows:

- We develop a generative QD-based model that can directly generate natural language sub-questions by incorporating retrieved paragraphs that hide the reasoning chain.
- Detailed experimental results show that incorporating the generated sub-questions into

paragraph retrieval and QA modules allow GenDec to outperform all QD-based QA models and most other strong baselines.

- We explore the potential usage of LLMs (e.g., LLaMA or ChatGPT) and demonstrate QD still plays a vital role in QA in the LLM era.

2 Related Work

2.1 Multi-hop Question-answering

Multi-hop QA requires more than one reasoning step in multiple paragraphs to answer a question. For example, multi-hop QA in DROP (Dua et al., 2019) requires numerical reasoning such as addition and subtraction. Yang et al. (2018b) proposed the HotpotQA dataset that contains 113K multi-hop QA pairs collected from Wikipedia articles by crowd-sourcing. Ho et al. (2020a) presented 2WikiMultiHopQA, which uses structured and unstructured data and introduces the evidence information containing a reasoning path for multi-hop questions.

2.2 Question Decomposition

Several studies conducted QD in complex QA tasks by using different methods. Wolfson et al. (2020a) and Talmor and Berant (2018), inspired by SQL and SPARQL query, proposed rule-based methods. However, they failed to generalize into different types of questions because of the limited rules. Min et al. (2019b) proposed a supervised QD method with human-labeling data to predict the text span of sub-questions. ONUS (Perez et al.,

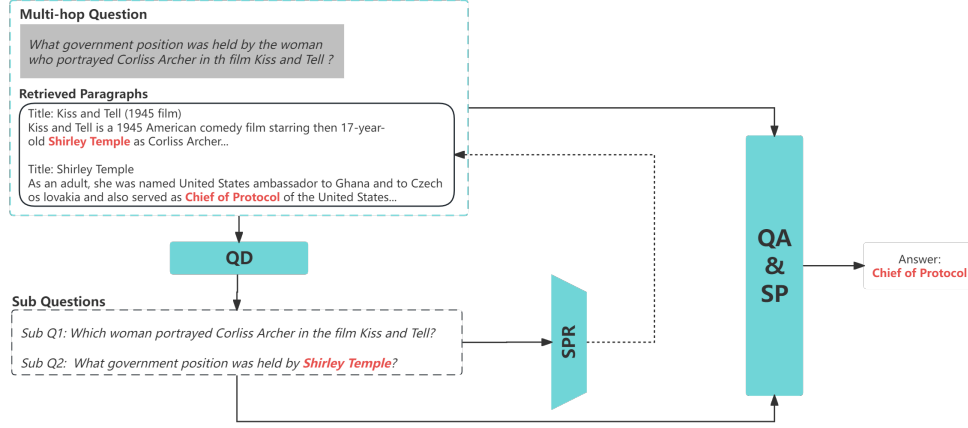


Figure 2: Pipeline of GenDec. From top to bottom. We first carry out Question Decomposition (QD) to decompose a multi-hop question into its sub-questions and then train a Sub-question-enhanced Paragraph Retrieval module (SPR). We then input multi-hop questions, sub-questions, as well as retrieved paragraphs, into the sub-question-enhanced QA module to extract the final answers.

2020a) is a one-to-N unsupervised sequence transduction method that uses supervision information of pseudo-decompositions from Common Crawl to map complex questions into simpler questions and recompose intermediate answers of sub-questions for reasoning final answers. These supervised and unsupervised QD methods decompose complex questions into two sub-questions but are not applicable to real scenarios. Deng et al. (2022b) proposed an Abstract Meaning Representation (AMR)-based QD method that trains an AMR-to-text generation model on the QDMR (Wolfson et al., 2020b) dataset. The entity description graph (EDG)-based QD method (Hu et al., 2021b) represents the structure of complex questions to solve the question-understanding and component-linking problems of knowledge base QA tasks. Zhou et al. (2022) pre-trained Decomp-T5 on human-collected parallel news to improve the ability of semantic understanding for QD. Instead of answering sub-questions one by one, Guo et al. (2022) directly concatenated sub-questions with the original question and context to leverage the reading-comprehension model to predict the answer. Wang et al. (2022) propose a step-by-step sub-question generation that generates sub-questions at each intermediate step. However, such stepwise reasoning and generation methods suffered from error propagation, while ours can directly generate the sub-questions and reasoning at the same time.

2.3 Large Language Models on Complex Reasoning

LLMs have shown reasoning abilities over several tasks, such as multi-hop QA (Bang et al., 2023), commonsense reasoning (Liu et al., 2022), and table QA (Chen, 2022). Chain-of-thought (CoT) (Wei et al., 2022) leverages a series of intermediate reasoning steps, achieving better reasoning performance on complex tasks. Jin and Lu (2023) proposed a framework called Tabular Chain of Thought (Tab-CoT) that can perform step-by-step reasoning on complex tableQA tasks by creating a table without fine-tuning by combining the table header with related column names as a prompt. Khot et al. (2022) proposed an approach called Decomposed Prompting to solve complex tasks by decomposing them into simple sub-tasks that can be delegated to a shared library of prompting-based LLMs dedicated to these sub-tasks.

However, these studies only decomposed questions into sub-questions and the latter sub-questions always rely on previous sub-questions. When the previous sub-questions are incorrectly answered, the latter sub-questions are also prone to be incorrectly answered.

3 GenDec

As discussed in the preceding section, previous QD-based QA methods fail to solve the error-propagation problem during the answer reasoning process as they decompose questions into sub-questions. GenDec’s approach consists of three

main components: (1) a generative QD module, to generate independent sub-questions with paragraphs that contain evidence; (2) a sub-question-enhanced paragraph-filtering module, that serves both the QD and QA modules; and (3) a sub-question enhanced QA module, which fuses features of sub-questions for QA and supporting-facts prediction. Figure 2 shows the overall framework of GenDec.

3.1 Question Decomposition Module

We explore different model architectures for the QD module, i.e., generative language models (e.g., BART, T5), LLMs, and traditional syntactic-parsing models. We use BART-large (Lewis et al., 2019) and T5-large (Raffel et al., 2020) as the generative language models in GenDec. Considering the computing resources and model availability, we also use LLaMA-7B (Touvron et al., 2023) with the Low-Rank Adaptation (LoRA) technique (Hu et al., 2021a) for training an LLM-based QD, as a design alternative for evaluation.

We leverage the coarse retrieval module proposed by Yin et al. (2023) to retrieve relevant paragraphs to serve the Question Decomposition module. In the coarse retrieval module, each question Q is typically combined by a set of N paragraphs P_1, P_2, \dots, P_N , but only a small number of paragraphs (e.g., two in HotpotQA) are labeled as relevant to the question Q . We model the paragraph retrieval as a binary classification task. Specifically, for each question-paragraph pair, we concatenate it as “[CLS], question, [SEP], paragraph, [SEP]” in sequence.

Finally, we make use of syntactic parsing, including constituency parsing and dependency parsing, to directly break multi-hop questions into sub-questions to compare the impact of not incorporating retrieved paragraphs with other generative QD-based QA models.

3.1.1 Generative Question Decomposition

To ensure the sub-questions are answerable by the QA module, we train a text-to-text generation model on the sub-question dataset from HotpotQA Khot et al. (2021).

We use BART-large and T5-large models as backend models and fine-tune them on the sub-question dataset to generate sub-questions. We use the retrieved paragraphs that contain supporting facts p and question q as input to train a question-generator model $G : (p, q) \Rightarrow sub_qs$, where

sub_qs is the generated sub-question set. Such a generator, G , produces the two sub-questions in the example in Figure 1. The details of finetuning T5-large and BART-large are given in Appendix A.

3.1.2 Large Language Models in Question Decomposition

Differently from typical QD-based QA models, we also explore leveraging powerful LLMs with few-shot prompting as a plugin for GenDec to decompose complex multi-hop questions. Despite the remarkable advancements brought about by LLMs, commercial models come with certain limitations that hinder transparent and open research. Therefore, we fine-tune LLaMA-7B (Touvron et al., 2023) with LoRA (Hu et al., 2021a) under low resource conditions as our LLM of use¹. The details of finetuning LLaMA are presented in Appendix A.

3.2 Sub-question-enhanced Paragraph Retrieval

Multi-hop question answering takes textual context into account and usually, MQA datasets include multiple paragraphs as question context (e.g., HotpotQA and 2WikiMultiHopQA datasets include 10 paragraphs per question). However, including all such paragraphs is not ideal due to noise and size (length). Therefore, paragraph retrieval plays a vital role in both QA and QD modules, since GenDec utilizes information from sub-questions and can thus focus on the more relevant data.

We propose sub-question-enhanced paragraph retrieval (SPR) for refined retrieval, which utilizes an encoder and a classification head to compute scores for each paragraph to help supporting facts prediction. Given a k -hop question Q , generated k sub-questions q_1, \dots, q_k , and a candidate set with n passages as $\mathcal{P} = \{p_1, p_2, \dots, p_n\}$, SPR aims to retrieve a relevant paragraph set $(\hat{p}_1, \hat{p}_2, \dots, \hat{p}_k)$ that relates to the k sub-questions and the k -hop question Q . While most previous work formulates it as a one- or two-step sequence labeling task, classifying every passage $p_i \in \mathcal{P}$ as relevant or not.

A passage $p_i \in \mathcal{P}$ corresponds to the question Q and j -th sub-question $q_j \in \mathcal{S}$. Consequently, we also denote the output score of SPR as $S(\hat{p}_i | Q, q_j)$, given the concatenated sequence of question, sub-question, and passages identified so far, (Q, q_j, \hat{p}_i) .

We use the DeBERTa model (He et al., 2021) as

¹<https://huggingface.co/decapoda-research/llama-7b-hf>

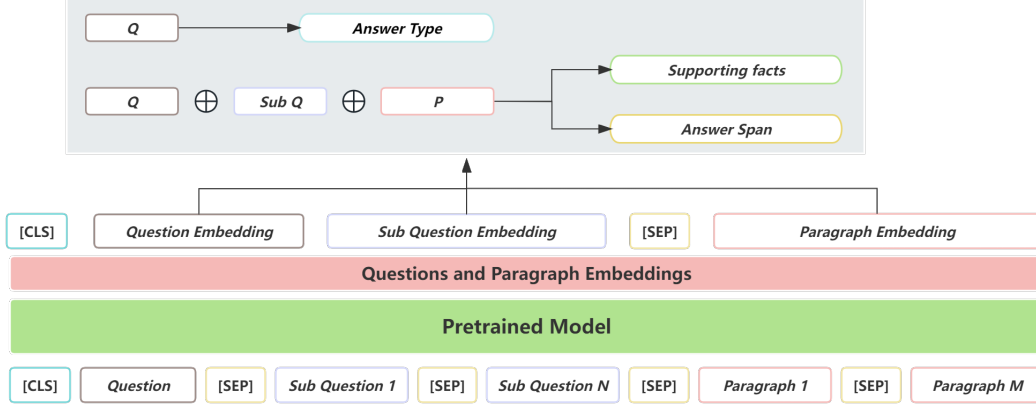


Figure 3: Architecture of QA module.

an encoder to derive embeddings for the concatenated sequence (Q, q_j, \hat{p}_i) and the output $\hat{o}_i \in \mathbb{R}^n$. Subsequently, a fully connected layer is added after DeBERTa to project the final dimension of the “[CLS]” representations of these embeddings into a 2-dimensional space, representing “irrelevant” and “relevant” respectively. The logit in the “relevant” side serves as the score for each paragraph. This scoring process is denoted by a function $S(\hat{p}_i|Q, q_j)$. In SPR, we optimize the classification of each combination of question, sub-question, and paragraph using Cross-Entropy loss.

$$\mathcal{L}_j = - \sum_{q_i \in \mathcal{S}} \sum_{\hat{p}_i \in \mathcal{P}} l_{j,p} \log S(\hat{p}_i|Q, q_j) + (1 - l_{j,p}) \log(1 - S(\hat{p}_i|Q, q_j)) \quad (1)$$

where $l_{j,p}$ is the label of \hat{p}_i and $S(\hat{p}_i|Q, q_j)$ is the score function predicted by the model.

Thus, we train a paragraph retrieval model based on DeBERTa (He et al., 2021) to execute binary classification and rank the scores of paragraphs containing the gold supporting facts.

3.3 Sub-question-enhanced QA module

In the QA module, we use multi-task learning to simultaneously predict supporting facts, and extract answer spans by incorporating sub-questions. In order to better evaluate the role of sub-question incorporation, we do not include other additional modules in our model. Instead, we focus on the effects of sub-question incorporation on the performance of the QA module. Additionally, as both HotpotQA and 2WikiMultiHopQA datasets also contain questions with yes/no answers, a common

scenario, we include an answer type task. The architecture of our QA module is illustrated in figure 3.

The QA module obtains an initial representation by first combining all retrieved paragraphs into context C , which is concatenated with question Q and sub-questions $\{Sub_Qs\}$ and fed into DeBERTa. We denote the encoded question and sub-question representations as $\mathbf{Q} = \{\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_{Q-1}\} \in \mathbf{R}^{m \times d}$ and the encoded context representation as $\mathbf{C} = \{\mathbf{c}_0, \mathbf{c}_1, \dots, \mathbf{c}_{C-1}\} \in \mathbf{R}^{C \times d}$, where Q is the length of the question. Each \mathbf{q}_i and $\mathbf{c}_j \in \mathbf{R}^d$.

$$\mathbf{P}^i = \text{DeBERTa} \left(S^{(i)}[d:] \right) \quad (2)$$

$$\text{sub_}\mathbf{q}^i = \text{DeBERTa} \left(\text{Sub_}Q^{(i)}[d:] \right)$$

$$\mathbf{q} = \text{DeBERTa}(\mathbf{Q}),$$

where $P^{(i)} \in \mathbf{R}^d$, $\text{Sub_}Q^{(i)} \in \mathbf{R}^d$, $\mathbf{Q} \in \mathbf{R}^d$ respectively denote the i -th paragraph, sub-question, and question representations.

To extract answer spans, we use a linear prediction layer on the contextual representation to identify the start and end positions of answers and employ cross-entropy as the loss function. The corresponding loss terms are denoted as \mathcal{L}_{start} and \mathcal{L}_{end} , respectively.

The classification loss for the supporting facts is denoted as \mathcal{L}_{sup} , and we jointly optimize all of these objectives in our model.

We also introduce an answer-type classification module trained with cross-entropy loss function.

$$\mathcal{L}_{type} = \mathbb{E} \left[- \sum_{i=1}^3 y_i^{type} \log(\hat{y}_i^{type}) \right] \quad (3)$$

where \hat{y}_i^{fine} denotes the predicted probability of question types classified by our model, and y_i^{fine} represents the corresponding one-hot encoded ground-truth distribution. y_i^{type} has three values: 0 denotes a negative answer, 1 denotes a positive answer, and 2 denotes the answer is a span.

The multi-task prediction model’s total loss is:

$$\mathcal{L}_{reading} = \lambda_1 \mathcal{L}_{type} + \lambda_2 (\mathcal{L}_{start} + \mathcal{L}_{end}) + \lambda_3 \mathcal{L}_{sup} \quad (4)$$

Similarly, we set λ_1 , λ_2 , and λ_3 all to 1, giving equal importance to each module for multitask learning. The implementation details of the Sub-question-enhanced QA module are described in Appendix A.

4 Experiments and Analysis

This section describes the different utilized datasets to analyze the different characteristics of the problem and our experimental setup.

4.1 Datasets

Question Answering (QA) We evaluate GenDec on the MuSiQue (Trivedi et al., 2022), 2WikiMultiHopQA (Ho et al., 2020b) and HotpotQA (Yang et al., 2018a) datasets, which contain 20K, 160K and 90K training instances. These three multi-hop QA datasets consist of questions, answers, supporting facts, and a collection of 10 paragraphs as context per question.

Question Decomposition (QD) To train and evaluate GenDec’s QD module, we use the sub-questions and answers data processed from the multi-hop HotpotQA dataset Khot et al. (2021) - here named SQA for clarity. These sub-questions are relatively high quality, in that we are able to use them to train a sub-question generator that achieves high task performance on multi-hop QD.

Sub-question Reasoning To evaluate the reasoning ability of GenDec, we also utilize a human-verified sub-question test dataset derived from HotpotQA Tang et al. (2020a) - here named HVSQA for clarity; which provides a strong benchmark to evaluate QA models in answering complex questions via sub-question reasoning.

4.2 Experiment Results

4.3 Quantitative Analysis

We use Exact Match (EM) and F1 scores as evaluation metrics for answer span prediction and support-

ing facts prediction on the HotpotQA, 2WikiMultiHopQA, and MuSiQue-Ans datasets to compare the performance of GenDec with that of QD-based, GNN-based, and other SOTA QA models.

As shown in Table 1, GenDec outperforms most models in both metrics, including the strong baseline consisting of our Question Decomposition method combined with HGN-large (Fang et al., 2019b) (itself a strong GNN-based QA model), and performs very competitively to the latest SOTA on the HotpotQA dataset.

The middle and bottom sections of the table also shows that GenDec significantly outperforms most previous work on the 2WikiMultiHopQA and MuSiQue-Ans (Trivedi et al., 2022) datasets. GenDec’s performance is only lower than the contemporary Beam Retrieval Zhang et al. (2023), which takes a retrieval approach that can be complementary to GenDec itself.

Table 2 shows the SOTA paragraph retrieval performance of GenDec’s SPR method against previous strong paragraph retrieval model baselines. SPR reaches very competitive results against Beam Retrieval (slightly higher in F1 vs slightly lower EM). Moreover, combining our QD approach with Beam Retrieval further improves performance and showcases the efficacy of leveraging sub-questions.

Table 4 shows the performance of GenDec and baseline models on the HVSQA dataset (human-verified sub-questions). GenDec achieves SOTA performance compared with the other QA models. Moreover, it is important to note that GenDec also outperforms all other models on sub-question reasoning (1 and 2), which highlights the benefits of our approach in reasoning chains. Lastly, with the help of our QD module, relative F1 scores are boosted by +6.82% and EM by +5.45% compared with ONUS (Perez et al., 2020b), which is also a QD-based model. We further verify the effectiveness of GenDec’s QD module in an ablation study discussed in the next section.

To analyze whether existing multi-hop QA models can demonstrate the right reasoning process, we compare the percentage of correct final answers and intermediate answers obtained by DFGN (Xiao et al., 2019), DecompRC (DecRC) (Min et al., 2019b), HGN (Fang et al., 2019a), PCL (Deng et al., 2022a), Beam Retrieval (BR) (Zhang et al., 2023), and our GenDec. Table 3 summarizes the percentage of correct answers for intermediary sub-questions and final multi-hop questions on the

Model	Ans		Sup		Joint	
	EM	F1	EM	F1	EM	F1
HotpotQA test set						
QD-based QA Models						
DecompRC (Min et al., 2019b)	55.20	69.63	-	-	-	-
ONUS (Perez et al., 2020a)	66.33	79.34	-	-	-	-
GNN-based Models						
DFGN (Xiao et al., 2019)	56.31	69.69	51.50	81.62	33.62	59.82
SAE-large (Tu et al., 2020)	66.92	79.62	61.53	86.86	45.36	71.45
C2F Reader(Shao et al., 2020b)	67.98	81.24	60.81	87.63	44.67	72.73
HGN-large (Fang et al., 2019a)	69.22	82.19	62.76	88.47	47.11	74.21
BRF-graph (Huang and Yang, 2021)	70.06	82.20	61.33	88.41	45.92	74.13
AMGN+ (Li et al., 2021)	70.53	83.37	63.57	88.83	47.77	75.24
Other SOTA Models						
FE2H on ALBERT (Li et al., 2022b)	71.89	84.44	64.98	89.14	50.04	76.54
PCL (Deng et al., 2022a)	71.76	84.39	64.61	89.20	49.27	76.56
Smoothing R3 (Yin et al., 2023)	72.07	84.34	65.44	89.55	49.73	76.69
Beam Retrieval (Zhang et al., 2023)	72.69	85.04	66.25	90.09	50.53	77.54
QD + HGN-large	71.73	84.23	64.32	89.46	49.22	75.63
GenDec (DeBERTa-large)	<u>72.39</u>	<u>84.69</u>	<u>65.88</u>	90.31	<u>50.34</u>	<u>77.48</u>
2WikiMultiHotpotQA test set						
CRERC (Fu et al., 2021)	69.58	72.33	82.86	90.68	49.80	58.99
NA-Reviewer (Fu et al., 2022)	76.73	81.91	89.61	94.31	52.75	65.23
BigBird-base model (Ho et al., 2023)	74.05	79.68	77.14	92.13	39.30	63.24
Beam Retrieval (Zhang et al., 2023)	88.47	90.87	95.87	98.15	-	-
GenDec (DeBERTa-large)	<u>86.47</u>	<u>88.15</u>	<u>93.28</u>	<u>96.45</u>	<u>56.87</u>	<u>68.38</u>
MuSiQue-Ans test set						
Beam Retrieval (beam size 2) (Zhang et al., 2023)	-	69.20	-	91.40	-	-
Beam Retrieval (beam size 1) (Zhang et al., 2023)	-	66.90	-	90.00	-	-
Ex(SA) (Trivedi et al., 2022)	-	49.00	-	78.10	-	-
Ex(Ee) (Trivedi et al., 2022)	-	46.40	-	80.60	-	-
GenDec (DeBERTa-large)	-	<u>65.40</u>	-	<u>87.90</u>	-	-

Table 1: Performance of different QA models on test distractor settings of HotpotQA, 2WikiMultihopQA and MuSiQue Answerable datasets. GenDec outperforms all QD-based and other GNN-based QA models.

Model	EM	F1
SAE _{large} (Tu et al., 2020)	91.98	95.76
S2G _{large} (Wu et al., 2021)	95.77	97.82
FE2H _{large} (Li et al., 2022a)	96.32	98.02
C2FM _{large} (Yin et al., 2023)	96.85	98.32
Beam Retrieval _{large} (Zhang et al., 2023)	97.52	98.68
GenDec_SPR (ours)	97.13	98.78
QD + Beam Retrieval	98.02	99.17

Table 2: Comparison of our sub-question enhanced paragraph retriever with previous baselines on HotpotQA dev set.

HVSQA (Tang et al., 2020a) dataset. We observe that, although Beam Retrieval (BR) achieves a relatively high percentage of correct final answers (88.3% vs GenDec’s 84.3%), 36.2% of these are

obtained from incorrect intermediate answers and reasoning chain. While GenDec’s performance is less affected by this case, at 31.4%. Comparing fully correct answer chains, Beam Retrieval and GenDec reach 52.2%, and 52.9% respectively; a small advantage to GenDec.

4.4 Ablation Studies

To evaluate the impact of GenDec’s QD module, we conduct an ablation study testing the performance of answering all sub-questions and original questions, with and without the QD module. The results, shown in Table 4, indicate that the QD module shows consistent and significantly improved results; improving the F1 score and EM by 3.36 and 2.16, respectively, in the original QA. In answering

q	q_{sub1}	q_{sub2}	DFGN	DecRC	HGN	PCL	BR	GenDec
c	c	c	23	31.3	39.5	43.6	52.2	52.9
c	c	w	9.7	7.2	5.1	6.8	7.8	6.4
c	w	c	17.9	19.1	19.6	21.3	21.7	18.8
c	w	w	7.5	5.5	3.8	2.1	6.7	6.2
w	c	c	4.9	3.0	2.8	1.7	1.2	1.1
w	c	w	17	18.6	16.7	16.3	7.3	10.3
w	w	c	3.5	3.4	2.6	1.1	0.7	0.9
w	w	w	16.5	11.9	9.9	7.1	2.4	3.4

Table 3: Categorical EM statistics (%) of sub-question evaluation for six multi-hop QA models over HVSQA (Tang et al., 2020a). c/w denotes questions answered correctly/wrongly. For example, the fourth row shows the percentage of multi-hop questions that can be correctly answered while sub-questions cannot.

intermediate answers to sub-questions GenDec w/ QD also improves over w/o QD (improving the F1 score and EM by 2.07 and 3.78, and 4.49 and 4.45 on sub-questions 1 and 2 respectively). The results indicate that the QD module plays an important role in GenDec in not only its QA ability but also in intermediate answer reasoning to support answering the final question. We also evaluate the impact of different backend models in our QD module and compare the performances of T5-large, BART-large, SynDec, and LLaMA-7B on the dev distractor setting of HotpotQA. LLaMA-7B achieves the best overall performance on both answer span prediction and supporting facts prediction since it had the best QD performance, with BART-large (even being a much smaller model) presenting a very competitive performance, as shown in Table 7.

4.5 Qualitative Analysis

We compare the QD performance of different LMs (Table 7 in the Appendix) and their impact on QA performance (Table 5). LLaMA-7B achieves SOTA performance in F1 score and EM (6.81 and 9.47 higher, respectively). We also compare F measure, ROUGE-1, ROUGE-L, and BLEU scores of generated sub-questions and LLaMA-7B significantly improves the quality of sub-questions reaching 80.57, 69.48, and 31.32, respectively. Likely due to the different max input lengths of T5-large (512) and BART-large (1024), BART outperforms it since some inputs contain many sentences (including both the multi-hop questions themselves and their supporting facts). GenDec with LLaMA-7B also improves QA performance on the distractor setting dev set, as shown in Table 5, but not substantially. We also evaluate the impact of sub-question

on LLM reasoning in table 6. Further analysis of ChatGPT is discussed in Appendix B.

We further compare the decomposition quality of GenDec and previous QD models in Table 8, which illustrates how GenDec can produce natural language questions with higher quality, more fluent and complete sentences.

4.6 Error Analysis

We conduct an error analysis of GenDec’s performance by selecting 20 samples from the dev set for evaluation, with 10 correct and 10 incorrect answers to analyze the impact of supporting facts prediction on QD and QA. We find a total of 12 correct supporting facts predictions and 8 incorrect supporting facts predictions among these 20 samples. For the 12 correct Supporting fact Predictions (SPs), we obtain 10 correct and 2 incorrect QD results. For the 10 correct QD results, we finally obtained 8 correct answers and 2 incorrect answers. And for the 8 incorrect SPs, we obtain 7 incorrect QD results and 6 incorrect answers. We then also select 20 samples from the dev set, with 10 correct and 10 incorrect QD results. For the 10 correct QD results, we obtained 8 correct answers, while for the 10 incorrect QD results, we obtained 5 correct answers. We list 6 samples of our selection in Table 9 in the Appendix, showing that questions be well answered based on high-quality QD.

5 Conclusion

We proposed GenDec, a generative-based QD method that generates independent sub-questions based on incorporating supporting facts. Intuitively, the supporting facts inform the reasoning chain of multi-hop questions. To explore this intuition, we train a sub-question-enhanced paragraph retrieval and QA module that incorporates sub-questions and shows that it significantly improves QA. We also explore the possible role of LLMs in QD and QA tasks. Lastly, while GenDec reaches new SOTA results in multi-hop QA, it can still face errors due to incorrect supporting fact predictions influencing the model to incorrectly predict both sub-questions and final answers.

6 Limitations

In this paper, we focus on the impact of QD in multi-hop QA, where the answers to most questions can be decomposed into several independent sub-questions via the fusion of supporting facts.

Although GenDec performs very well on QD and QA, one of its limitations is that it is still sensitive to errors in paragraph filtering. The QD results would be affected when given incorrect paragraphs are selected. For future work, we plan to focus on tackling this problem.

References

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-task, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.

Wenhu Chen. 2022. Large language models are few(1)-shot table reasoners. In *Findings*.

Zhenyun Deng, Yonghua Zhu, Yang Chen, Qianqian Qi, Michael Witbrock, and Patricia Riddle. 2022a. Prompt-based conservation learning for multi-hop question answering. *arXiv preprint arXiv:2209.06923*.

Zhenyun Deng, Yonghua Zhu, Yang Chen, M. Witbrock, and Patricia J. Riddle. 2022b. Interpretable amr-based question decomposition for multi-hop question answering. In *International Joint Conference on Artificial Intelligence*.

Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, and Jie Tang. 2019. Cognitive graph for multi-hop reading comprehension at scale. *arXiv preprint arXiv:1905.05460*.

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In *North American Chapter of the Association for Computational Linguistics*.

Yuwei Fang, S. Sun, Zhe Gan, Rohit Radhakrishna Pillai, Shuohang Wang, and Jingjing Liu. 2019a. Hierarchical graph network for multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing*.

Yuwei Fang, Siqu Sun, Zhe Gan, Rohit Pillai, Shuohang Wang, and Jingjing Liu. 2019b. Hierarchical graph network for multi-hop question answering. *arXiv preprint arXiv:1911.03631*.

Ruiliu Fu, Han Wang, Xuejun Zhang, Jun Zhou, and Yonghong Yan. 2021. *Decomposing complex questions makes multi-hop QA easier and more interpretable*. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 169–180, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ruiliu Fu, Han Wang, Jun Zhou, and Xuejun Zhang. 2022. *Na-reviewer: Reviewing the context to improve the error accumulation issue for multi-hop qa*. *Electronics Letters*, 58(6):237–239.

Xiao-Yu Guo, Yuan-Fang Li, and Gholamreza Haffari. 2022. Complex reading comprehension through question decomposition. *ArXiv*, abs/2211.03277.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. DeBERTa: Decoding-enhanced BERT with disentangled attention. In *ICLR 2021*.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020a. *Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps*. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020b. *Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps*. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2023. *Analyzing the effectiveness of the underlying reasoning tasks in multi-hop question answering*. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1163–1180, Dubrovnik, Croatia. Association for Computational Linguistics.

J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021a. Lora: Low-rank adaptation of large language models. *ArXiv*, abs/2106.09685.

Xixin Hu, Yiheng Shu, Xiang Huang, and Yuzhong Qu. 2021b. Edg-based question decomposition for complex question answering over knowledge bases. In *International Workshop on the Semantic Web*.

Yongjie Huang and Meng Yang. 2021. Breadth first reasoning graph for multi-hop question answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5810–5821.

Ziqi Jin and Wei Lu. 2023. Tab-cot: Zero-shot tabular chain of thought. *arXiv preprint arXiv:2305.17812*.

Tushar Khot, Daniel Khoshabi, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2021. Text modular networks: Learning to decompose tasks in the language of existing models. *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, page 1264–1279.

670	Tushar Khot, H. Trivedi, Matthew Finlayson, Yao Fu,	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 . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	724
671	Kyle Richardson, Peter Clark, and Ashish Sabharwal.		725
672	2022. Decomposed prompting: A modular approach		726
673	for solving complex tasks. <i>ArXiv</i> , abs/2210.02406.		727
674	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan		728
675	Ghazvininejad, Abdel rahman Mohamed, Omer Levy,		729
676	Veselin Stoyanov, and Luke Zettlemoyer. 2019. Bart:		
677	Denoising sequence-to-sequence pre-training for nat-	Nan Shao, Yiming Cui, Ting Liu, Shijin Wang, and	730
678	ural language generation, translation, and compre-	Guoping Hu. 2020a. Is graph structure necessary	731
679	hension. In <i>Annual Meeting of the Association for</i>	for multi-hop question answering? <i>arXiv preprint</i>	732
680	<i>Computational Linguistics</i> .	<i>arXiv:2004.03096</i> .	733
681	Ronghan Li, Lifang Wang, Shengli Wang, and Zejun		
682	Jiang. 2021. Asynchronous multi-grained graph net-	Nan Shao, Yiming Cui, Ting Liu, Shijin Wang, and	734
683	work for interpretable multi-hop reading comprehen-	Guoping Hu. 2020b. Is graph structure necessary for	735
684	sion . In <i>Proceedings of the Thirtieth International</i>	multi-hop reasoning? <i>ArXiv</i> , abs/2004.03096.	736
685	<i>Joint Conference on Artificial Intelligence, IJCAI-21</i> ,		
686	pages 3857–3863. International Joint Conferences on	A. Talmor and J. Berant. 2018. The web as a knowledge-	737
687	Artificial Intelligence Organization. Main Track.	base for answering complex questions. In <i>North</i>	738
		<i>American Association for Computational Linguistics</i>	739
		(<i>NAACL</i>).	740
688	Xin-Yi Li, Wei-Jun Lei, and Yu-Bin Yang. 2022a.		
689	From easy to hard: Two-stage selector and reader	Yixuan Tang, Hwee Tou Ng, and Anthony K. H. Tung.	741
690	for multi-hop question answering . <i>ArXiv preprint</i> ,	2020a. Do multi-hop question answering systems	742
691	abs/2205.11729.	know how to answer the single-hop sub-questions?	743
		In <i>Conference of the European Chapter of the Asso-</i>	744
692	Xin-Yi Li, Weixian Lei, and Yubin Yang. 2022b. From	<i>ciation for Computational Linguistics</i> .	745
693	easy to hard: Two-stage selector and reader for multi-		
694	hop question answering . <i>ArXiv</i> , abs/2205.11729.	Yixuan Tang, Hwee Tou Ng, and Anthony KH Tung.	746
		2020b. Do multi-hop question answering systems	747
695	Jiacheng Liu, Skyler Hallinan, Ximing Lu, Pengfei He,	know how to answer the single-hop sub-questions?	748
696	Sean Welleck, Hannaneh Hajishirzi, and Yejin Choi.	<i>arXiv preprint arXiv:2002.09919</i> .	749
697	2022. Rainier: Reinforced knowledge introspec-		
698	tor for commonsense question answering. <i>ArXiv</i> ,	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	750
699	abs/2210.03078.	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	751
		Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	752
700	Sewon Min, Victor Zhong, Luke Zettlemoyer, and Han-	Azhar, Aur'elien Rodriguez, Armand Joulin, Edouard	753
701	naneh Hajishirzi. 2019a. Multi-hop reading compre-	Grave, and Guillaume Lample. 2023. Llama: Open	754
702	hension through question decomposition and rescoring.	and efficient foundation language models. <i>ArXiv</i> ,	755
703	In <i>ACL</i> .	abs/2302.13971.	756
704	Sewon Min, Victor Zhong, Luke Zettlemoyer, and Han-		
705	naneh Hajishirzi. 2019b. Multi-hop reading compre-	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	757
706	hension through question decomposition and rescoring.	and Ashish Sabharwal. 2022. MuSiQue: Multi-	758
707	<i>ArXiv</i> , abs/1906.02916.	hop questions via single-hop question composition.	759
		<i>Transactions of the Association for Computational</i>	760
		<i>Linguistics</i> .	761
708	Ethan Perez, Patrick Lewis, Wen tau Yih, Kyunghyun		
709	Cho, and Douwe Kiela. 2020a. Unsupervised ques-	Ming Tu, Kevin Huang, Guangtao Wang, Jing Huang,	762
710	tion decomposition for question answering. In <i>Con-</i>	Xiaodong He, and Bowen Zhou. 2020. Select, an-	763
711	<i>ference on Empirical Methods in Natural Language</i>	swer and explain: Interpretable multi-hop reading	764
712	<i>Processing</i> .	comprehension over multiple documents. In <i>Proceed-</i>	765
		<i>ings of the AAAI conference on artificial intelligence</i> ,	766
713	Ethan Perez, Patrick Lewis, Wen tau Yih, Kyunghyun	volume 34, pages 9073–9080.	767
714	Cho, and Douwe Kiela. 2020b. Unsupervised ques-		
715	tion decomposition for question answering . In	Siyuan Wang, Zhongyu Wei, Zhihao Fan, Qi Zhang,	768
716	<i>EMNLP</i> .	and Xuanjing Huang. 2022. Locate then ask: Inter-	769
		pretable stepwise reasoning for multi-hop question	770
717	Lin Qiu, Yunxuan Xiao, Yanru Qu, Hao Zhou, Lei Li,	answering. <i>arXiv preprint arXiv:2208.10297</i> .	771
718	Weinan Zhang, and Yong Yu. 2019. Dynamically		
719	fused graph network for multi-hop reasoning . In	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	772
720	<i>Proceedings of the 57th Annual Meeting of the Asso-</i>	Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and	773
721	<i>ciation for Computational Linguistics</i> , pages 6140–	Denny Zhou. 2022. Chain of thought prompting	774
722	6150, Florence, Italy. Association for Computational	elicits reasoning in large language models. <i>ArXiv</i> ,	775
723	Linguistics.	abs/2201.11903.	776

Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan Berant. 2020a. Break it down: A question understanding benchmark. *Transactions of the Association for Computational Linguistics*.

Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan Berant. 2020b. Break it down: A question understanding benchmark. *Transactions of the Association for Computational Linguistics*, 8:183–198.

Bohong Wu, Zhuosheng Zhang, and Hai Zhao. 2021. [Graph-free multi-hop reading comprehension: A select-to-guide strategy](#). *ArXiv preprint, abs/2107.11823*.

Yunxuan Xiao, Yanru Qu, Lin Qiu, Hao Zhou, Lei Li, Weinan Zhang, and Yong Yu. 2019. Dynamically fused graph network for multi-hop reasoning. *ArXiv, abs/1905.06933*.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018a. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018b. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Zhangyue Yin, Yuxin Wang, Xiannian Hu, Yiguang Wu, Hang Yan, Xinyu Zhang, Zhao Cao, Xuanjing Huang, and Xipeng Qiu. 2023. Rethinking label smoothing on multi-hop question answering. In *China National Conference on Chinese Computational Linguistics*, pages 72–87. Springer.

Jiahao Zhang, Haiyang Zhang, Dongmei Zhang, Yong Liu, and Shen Huang. 2023. Beam retrieval: General end-to-end retrieval for multi-hop question answering. *arXiv preprint arXiv:2308.08973*.

Ben Zhou, Kyle Richardson, Xiaodong Yu, and Dan Roth. 2022. Learning to decompose: Hypothetical question decomposition based on comparable texts. *ArXiv, abs/2210.16865*.

A Implementation Details

Question Decomposition We use the pre-trained T5-large and BART-large models with `max_input_length` $L = 512$, and `max_output_length` $O = 64$. During training, we use the Adam optimizer in the QD modules and set batch size to 32 and learning rate to $5e-5$. All experiments utilized two TITAN RTX GPUs and 5 hours in total.

Question Answering We choose DeBERTa-v2-large as backend model and set number of epochs to 12 and batch size to 4. We use BERTAdam with learning rate of $5e-6$ for the optimization and set max position embeddings to 1024.

Fine-tuning LLaMA To fine-tune LLaMA, considering computing resources, we select LLaMA-7B as backbone, batch size of 4, number of epochs is 3, learning rate is $3e-4$, LoRA alpha of 16, and LoRA dropout of 0.05.

B ChatGPT on Multi-hop QA

We also evaluated the performance of ChatGPT with and without QD on 1000 samples of dev distractor settings. Figure 4 shows the used with QD and without QD prompt settings. We selected the 1-shot setting in which ChatGPT is given one example from the training set with two prompts, one is reasoning over sub-questions and the other is directly reasoning answers. As shown in Table 6, ChatGPT with additional sub-question information performs better than without sub-questions. ChatGPT with QD prompting achieves higher answer span extraction on the F1 score (76.28) and EM (56.24). However, both ChatGPT with QD prompting and ChatGPT without QD prompting are still lower than current QA models.

Model	Q_ori		Q_sub1		Q_sub2	
	F1	EM	F1	EM	F1	EM
CogQA	67.82	53.2	69.65	58.6	68.49	54
DFGN	71.96	58.1	68.54	54.6	60.83	49.3
DecompRC	77.61	63.1	75.21	61	70.77	56.8
ONUS	79.25	67.43	77.56	63.89	72.21	57.62
PCL	73.8	87.15	68.4	83.62	68.5	81.07
GenDec w/o QD	82.81	70.72	87.45	72.65	80.12	70.38
GenDec w QD	86.17	72.88	90.52	76.43	84.61	74.83

Table 4: Performance comparison between GenDec (with and without the QD module) and other QA models on HVSQA (Tang et al., 2020a), a human-verified sub-question test dataset from HotpotQA.

Model	Ans		Sup		Joint	
	EM	F1	EM	F1	EM	F1
GenDec (BART-large)	70.13	84.47	63.51	89.47	46.12	75.52
GenDec (T5-large)	69.94	84.11	63.32	89.35	46.02	75.69
GenDec (LLaMA-7B)	70.23	84.76	63.41	89.78	46.28	76.05

Table 5: Performance of QD module with different generative LMs on SQA Khot et al. (2021), distractor dev set of sub-questions processed from HotpotQA.

Model	Ans		Sup		Joint	
	EM	F1	EM	F1	EM	F1
ChatGPT w/o QD	51.08	74.53	60.61	87.96	30.95	65.55
ChatGPT w QD	56.24	76.28	60.74	87.85	34.16	67.01

Table 6: Performance of ChatGPT (with and without QD) on 1000 samples from HotpotQA’s dev set distractor setting data.

Models	Metric			
	F Measure	Rouge1	Rouge-L	BLEU
BART-LARGE	74.41	73.85	62.68	26.94
T5-LARGE	72.85	71.27	60.12	24.37
LLAMA-7B	81.32	80.57	69.48	31.22

Table 7: Generative QD performance of different generative LMs on test instances of HOTPOTQA sub-questions. Results are averaged on 1549 test instances.

GENDEC (OURS)
Sub-question 1: Which South Korean boy group had their debut album in 2014?
Sub-question 2: WINNER was formed by who?
MODULARQA (Khot et al., 2021)
Sub-question 1: What is the name of the South Korean group that had their debut album in 2014?
Sub-question 2: What was WINNER formed by?
DECOMPRC (Min et al., 2019b)
Sub-question 1: 2014 S/S is the debut album of which South Korean boy group?
Sub-question 2: which formed by who ?

Table 8: QD examples produced by {GENDEC, MODULARQA, DECOMPRC} for question “2014 S/S is the debut album of a South Korean boy group that was formed by who?”.

Original Question	Sub-questions	Intermediate Answers	Answer
Were Scott Derrickson and Ed Wood of the same nationality?	What was Scott Derrickson’s nationality? What was Ed Wood’s nationality? ✓	American ✓	Yes ✓
What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell?	Who portrayed Corliss Archer in Kiss and Tell? What position was held by Shirley Temple? ✓	Shirley Temple ✓	Chief of Protocol ✓
The director of the romantic comedy Big Stone Gap is based in what New York City neighborhood?	Who is the director of the romantic comedy Big Stone Gap? In what New York City neighborhood is Adriana Trigiani based? ✓	Adriana Trigiani ✓	Greenwich Village ✓
Are Random House Tower and 888 7th Avenue both used for real estate?	The Random House Tower used as real estate? What is 888 7th Avenue used also for? ✗	Used ✗	No ✗
What is the name of the executive producer of the film that has a score composed by Jerry Goldsmith?	What is the name of the film of which Jerry Goldsmith composed the score? Which co-writer of Alien was also an executive producer? ✓	Alien ✓	Francis Coppola ✗ Ford
Alvaro Mexia had a diplomatic mission with which tribe of indigenous people?	Who was given a diplomatic mission to the native populations living south of St. Augustine and in the Cape Canaveral area? What is the name of the indigenous tribe of Florida? ✗	Alvaro Mexia ✗	Indigenous peoples of Florida ✗

Table 9: Examples of 3 correct samples and 3 incorrect samples from dev set of HotpotQA

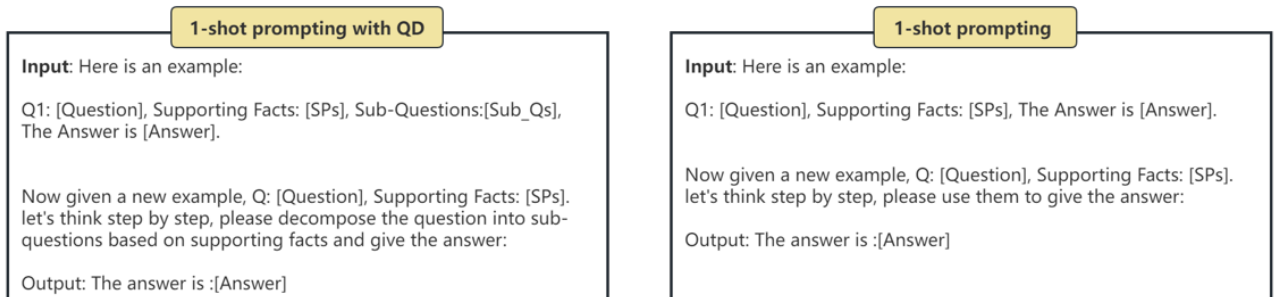


Figure 4: Prompting examples of different settings.