

# Corpus-informed Retrieval Augmented Generation of Clarifying Questions

Antonios Minas Krasakis  
a.m.krasakis@uva.nl  
University of Amsterdam  
The Netherlands

Andrew Yates  
a.c.yates@uva.nl  
University of Amsterdam  
The Netherlands

Evangelos Kanoulas  
e.kanoulas@uva.nl  
University of Amsterdam  
The Netherlands

## ABSTRACT

This study aims to develop models that generate corpus informed clarifying questions for web search, in a way that ensures the questions align with the available information in the retrieval corpus. We demonstrate the effectiveness of Retrieval Augmented Language Models (RAG) in this process, emphasising their ability to (i) jointly model the user query and retrieval corpus to pinpoint the uncertainty and ask for clarifications end-to-end and (ii) model more evidence documents, which can be used towards increasing the breadth of the questions asked. However, we observe that in current datasets search intents are largely unsupported by the corpus, which is problematic both for training and evaluation. This causes question generation models to “hallucinate”, ie. suggest intents that are not in the corpus, which can have detrimental effects in performance. To address this, we propose dataset augmentation methods that align the ground truth clarifications with the retrieval corpus. Additionally, we explore techniques to enhance the relevance of the evidence pool during inference, but find that identifying ground truth intents within the corpus remains challenging. Our analysis suggests that this challenge is partly due to the bias of current datasets towards clarification taxonomies and calls for data that can support generating corpus-informed clarifications.

## 1 INTRODUCTION

Open-domain search and Question Answering (QA) systems make the best effort to respond to any user’s question or query. Recent work attempts to quantify the uncertainty in the question-answering models in order to defer from answering a question they are uncertain about [6, 38]. In a parallel line of research, limited work in open-domain conversational and ad-hoc search systems has investigated enabling them to ask Clarifying Questions (CQs) [2, 19, 41]. The majority of this work employs static models of clarifying question generation, i.e. models that generate a clarifying question independent of the ability of the system to locate the right answer in the underlying corpus, or the potential different answers present therein [2, 9, 10, 40, 47].

In this work, we emphasize the importance of generating *corpus-informed* clarifying questions, that are *dynamic* with respect to the collection. We argue that it is critical to defer the ambiguity of the user query by modeling the aspects<sup>1</sup> of it in the document collection. Since the primary goal of search is to retrieve relevant information, clarifying questions should be generated as a function of the corpus and the relevant information therein. Failing to be dynamic and corpus-informed poses a risk to user experience, (a) due to the disruption caused by asking generic questions that offer no relevant information, and more severely (b) due to “hallucinations” while

generating clarifying questions, ie. presenting users with options or facets that do not exist in the collection.

Some research explicitly models the underlying corpus by relying on pipeline methods that (i) extract keywords or features from the document collection and (ii) generate questions based on these features [30, 32, 36, 45–47]. However, this separation prevents the joint modeling of the query and its ambiguity, the information in the corpus and the various aspects present therein, and the clarification questions. In this work we advocate for Retrieval Augmented Generation of Clarifying Questions, in a way that jointly models queries and the retrieval corpus to generate questions. Another line of work achieves state-of-the-art performance using text generation models without intermediate feature extraction [21, 29], but is limited by the size of the retrieval pool. This is a crucial factor, as it determines the breadth of the possible clarifications.

Since the generation model cannot read the entire corpus, Retrieval Augmentation helps to inform the generator of the possible information needs present in the collection. Hence, along with the generator, a retriever is responsible for obtaining a representative sample that covers the uncertainty of the user’s query in the collection. To apply Retrieval Augmentation effectively, it is crucial to (a) select a sample of the corpus that is representative of users’ potential information needs, and (b) account for as many relevant documents as possible to cover those. For this reason, we propose using *Fusion-in-Decoders (FiD)*, a family of models [12] that are computationally efficient in modelling multiple evidence documents. We demonstrate their effectiveness in simultaneously modeling the collection (user queries and retrieved documents) and generating questions, in a way that eliminates the need for an intermediate step of facet or keyword extraction and increases the ability to model larger parts of the corpus when generating questions.

Furthermore, to ensure dynamic and adaptive nature of clarifying question generation and avoid “hallucinations”, we enhance the training setup of question generation models. We find that current datasets [40] suffer from a disconnect between document evidence gathered from search engine results pages and the ground truth clarification questions derived from user reformulations, negatively impacting the question generator. In light of this, we further explore the relationship between question generation and evidence documents. We find that aligning evidence and generation during training is critical towards preventing “hallucinations” [15] and ensuring the generated question remains faithful to evidence documents, ie. the retrieval corpus.

Having trained an effective and grounded question generator, our focus shifts to retrieving evidence during inference. To ensure evidence documents encompass all existing information facets of

<sup>1</sup>Intents, facets and aspects are used interchangeably in this paper.

the original queries, we experiment with inducing novelty in evidence documents to improve facet generation. Yet, we find that capturing the ground truth facets present in current datasets with such methods remains a challenge.

Our main research question is *how can we train an end-to-end system that explicitly models the retrieval corpus and generates Corpus-informed clarifying questions?* Specifically, we aim to answer the following research questions:

- RQ1** Can we generate Clarifying Questions end-to-end with Retrieval Augmented Generation models?
- RQ2** What is the optimal evidence set for training and evaluating that support *Corpus-informed* Retrieval-Augmented Generators of clarifying questions?
- RQ3** How can we enhance the evidence set at test time to assist clarifying question generation?

In particular, we make the following contributions to the research of asking clarifying questions: (a) we define the desired properties of clarifying questions to be tightly dependent on both user queries, as well as documents retrieved for those queries; (b) to this end we develop a dataset that allows for training faithful and adaptive clarifying question models; (c) we propose using Retrieval Augmented Generation for generating clarifying questions end-to-end and show that Fusion-in-Decoder is an effective approach on this task, and (d) we investigate how to enhance search results for the purpose of better informing the generation of clarifying questions.

## 2 RELATED WORK

We first discuss datasets used in clarifying question generation and evaluation (Section 2.1), followed by methods for generating clarifying questions with a focus on web search (Section 2.2).

### 2.1 Clarifying question datasets

Recent work on clarifying questions differs on several aspects, including the type of questions asked, the way presented to the user, and the evaluation methodology for automatic question generation. Below we discuss several types of clarifying questions datasets.

Clarifying Questions for Web Search were first introduced in Qulac [2], an open-domain information-seeking conversational search dataset. This work builds upon faceted or ambiguous queries from the TREC Web Track dataset [7, 8]. The authors crowd-sourced clarifying questions to be asked from a search system to the user, as well as answers to these questions from users with different intents. The quality of automatically generated clarifying question was evaluated by the relevance of the retrieved documents. ClariQ [1] extended Qulac with additional topics and built synthetic multi-turn conversations from single-turn clarifications.

MIMICS introduces a clarification pane on the search engine result page (SERP) [40, 41]. This pane includes a template-based question (eg. “Who are you shopping for”) with multiple candidate answers or search intents (eg. “men, women, kids”). Intents are extracted from query reformulation logs and a defined taxonomy. The offline evaluation framework compares ground truth intents to the generated ones [9, 29, 35], while later work also introduces an online evaluation setup based on user interactions [33, 34, 42].

Another line of research investigates clarification in community forums like StackExchange [16, 24, 25, 27, 34]. In these datasets,

clarifying questions are written by expert users and the task is often defined as clarifying question retrieval from a pool of questions. Last, there are efforts to create clarifying questions on the product search domain, either towards assisting users in product search [44, 48], but also to clarify ambiguities originating from product descriptions [28, 45]. These works differ from our line of work, since they often revolve around structured product metadata.

For a more comprehensive overview on clarifying questions datasets we refer readers to Rahmani et al. [26].

### 2.2 Clarifying questions generation for web search

A number of approaches have been proposed for generating clarifying questions for web search results. Broadly, these works use a combination of rule-based systems, keyword extraction or topic modelling approaches, and Large Language Models (LLMs). Hashemi et al. [9] proposed NMIR, a transformer architecture that learns multiple intent representations for web queries by matching different document clusters to query intents. Later works found that a transformer encoder-decoder approach based on the BART model can outperform NMIR [29]. Other research tried combining LLMs with facet extraction methods. Sekulić et al. [30] constructed the ClariQ-FKW (Facet KeyWords) dataset and used it to guide GPT-2 in generating clarifications. They find that using facet keywords to guide the LLM helped with grounding and generating useful questions. In follow-up work, Sekulić et al. [32] try to improve upon the facet/keyword extraction part by using part-of-speech tags, entities and LDA topics. They combine those approaches by ranking their output using entropy-based methods and generate a template-based question using the top ranked keyword. In contrast to our work, (a) these methods only generate questions addressing one facet and (b) keep the facet extraction part disconnected from question generation. Samarinas et al. [29] combines various facet extraction methods, such as autoregressive generation, sequence tagging, extreme multi-labelling classification and LLM prompting in an ensemble, concluding that these methods are often complementary.

Another line of work tries to inform the CQ generation using descriptions of queries and lists of attributes from retrieved web pages [47]. Those are extracted with heuristics and filtering, and ranked using learning-to-rank. The top ranked are given to a seq2seq model (QLM [40]) to generate the final clarifying question. Similarly, Zhao et al. [46] focuses on complementing web-search results with relational information (eg. “Windows” → Operating system) extracted from web search results to inform the generation. Wang et al. [36] introduces another approach that uses LLMs and extracted keywords to generate clarifications. Instead of prompting, they use Neurologic decoding, a constrained generation approach that biases the generation towards these keywords. Multiple clarifications are generated and then ranked with a CQ ranking system. In contrast to those our work focuses on modelling the collection and generating a question in an end-to-end way, without the need for intermediate steps. Finally, another line of work focuses on ranking clarifying questions from a large pool of existing questions [1, 22]. However, such question pools do not exist in an open-domain setting and we do not discuss those further. Those approaches include

many different components (eg. keyword extraction, generation) that work in a disconnected way from each other.

Last there is work that focuses on enhancing the training of autoregressive models for generating clarification facets, which are by nature unordered. Next token prediction objectives could unfairly penalize the models for predicting facets in a different ordering. To tackle this, Hashemi et al. [10] generate all possible permutations of facets and backpropagate the minimum loss during training. Ni et al. [21] does a thorough comparison of training objectives and conclude that training on one-facet-at-a-time basis increases performance but hurts facet diversity. This work is orthogonal to ours and can also be applied to our models.

### 3 METHODOLOGY

*Problem definition.* In this work we investigate the problem of clarifying question generation for Web Search, where clarifying questions are presented through a clarification pane in the Search Engine Result Page (SERP). In this setting, the user has a search intent  $g_i \in G$  and searches for a set of relevant documents  $R_i$  in a corpus of documents  $C$ . To express his information need  $g_i$ , the user issues an ambiguous or faceted query  $q$  to the search engine. The SERP presents the top ranking results along with a clarification pane that tries to clarify the ambiguity of  $q$ . The clarification pane consists of a clarifying question  $CQ$  and some potential answers  $F$  to this clarifying question  $CQ$ . Those answers are basically possible search intents, facets or aspects of query  $q$ .

The goal of asking a good clarification question becomes directly related to finding the possible search intents (facets)  $F$  that query  $q$  engulfs. To find such search intents  $F$ , it is important to take into account the document collection  $C$  the user searches through. In failing to do so, we risk presenting search engine users with facets that are not found in  $C$ .

*RAG models.* Retrieval Augmented Generation models have become the standard choice for Knowledge-Intensive tasks where generation is required. Their strength lies on combining LLMs, that are proficient in text generation, with retrieval, that allows them to access information from a large non-parametric memory like a document collection [3]. That also motivates their use for our task, end-to-end and Corpus-informed clarifying question generation. Specifically, we use a Fusion-in-Decoder (*FiD*) model [12] as a question generator, due to its efficiency and effectiveness. *FiD* is an encoder-decoder model, but its encoder models input documents independently (no cross-attentions), producing individual embeddings. The decoder fuses the information from those embedding, generating an output answer. Due to the lack of cross-attentions between documents, *FiD* models can model longer context, ie. multiple retrieved documents with the same GPU memory requirements.

*Retrieval.* We experiment with various types of retrieval techniques, such as lexical (*BM25*), semantic (*Contriever*), as well as the documents that originated from the BING SERP provided with the MIMICS dataset (*BING*). For *BM25* search, we use the Pyserini toolkit [17] using the default parameters. For semantic retrieval, we use the *Contriever* architecture that has been jointly pretrained with the *Atlas* checkpoints, but observe in our preliminary experiments that retrieval performance of the retriever jointly pretrained with

*Atlas* models is sub-optimal. To this end, we initialize the retriever from the unsupervised pretraining of *Contriever* (*Contriever*) and the checkpoint finetuned on MSMarco[20] (*Contriever - FT*)<sup>2</sup>. Additionally, we consider two retrieval variants, namely non-aligned ( $*|Q$ ) or facet-aligned ( $*|Q, F$ ). For the former, we use only the original user query to retrieve, while in case of the latter we do multiple retrieval rounds using the query and each facet and interleave the retrieved documents (without replacing duplicates).

When experimenting with novelty methods (Section 5.3), we train the retriever with knowledge distillation from the question generator [11]. In practice, this method uses attention scores of document embeddings in the decoder as a proxy for document relevance. In Question Answering, this signal corresponds to detecting relevant passages for a question, while in our setting it promotes novel documents that contain the ground truth facets.

### 4 EXPERIMENTAL SETUP

In this section we outline our experimental setup.

#### 4.1 Datasets and Evaluation

As discussed in Section 2.1, a number of Clarifying Question datasets and setups exist. In this work, we focus on Web Search Clarifying Questions and perform our experiments on the MIMICS dataset. Our choice is based on a number of factors, namely: (a) the presentation of clarifying questions that happens via a clarification pane on a SERP. In our view, this presentation is less intrusive than a conversational system that interrupts the user’s journey without presenting results to ask a question, and hence less likely to harm user experience [4], and (b) a more straightforward and robust evaluation method that directly evaluates generated facets [14]. In contrast, the QuLaC [2] evaluation framework involves a user answering the clarifying question and ranking documents to judge question quality. However, evaluation can be particularly noisy, since both of these parts are challenging. User answers largely depend on their cooperativeness [31], while ranking documents with clarification-based queries is particularly challenging [14].

Following prior work on *MIMICS*, we train on *MIMICS - Click* and test on *MIMICS - Manual* [9, 10, 29].

*Evaluation metrics.* Following previous research [9, 10, 29, 46], we focus on a number of lexical and semantic metrics that measure overlap of generated facets to the ground truth ones. Given a sequence of ground truth facets,  $G = g_1, g_2, g_n$ , and a sequence of generated facets,  $F = f_1, f_2, \dots, f_m$ , we assess aspect generation quality using:

- Term Overlap ( $P, R, F1$ ): measures lexical overlap at the word level, i.e. overlap between words in  $G$  and  $F$ .
- Exact Match ( $P, R, F1$ ): measures exact lexical match at the facet level, i.e. whether  $f_i = g_j$ .
- Set-BERT [9]: measures semantic overlap at the facet level based on BERT-score [43]
- Set-BLEU-ngram[9, 23]: measures n-gram overlap at the facet level

<sup>2</sup><https://huggingface.co/facebook/contriever>  
<https://huggingface.co/facebook/contriever-msmarco>

| Model  | Co-Gen | Term Overlap                 |                            |                            | Exact Match                |               |                            | Set-BERT                   |                            |                            | Set-BLEU                   |                            |                            |                            |
|--|--------|------------------------------|----------------------------|----------------------------|----------------------------|---------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|
|  |        | Prec.                        | Rec.                       | F1                         | Prec.                      | Rec.          | F1                         | Prec.                      | Rec.                       | F1                         | BLEU1                      | BLEU2                      | BLEU3                      | BLEU4                      |
| Our Models (FiD)                             |        |                              |                            |                            |                            |               |                            |                            |                            |                            |                            |                            |                            |                            |
| FiD-CQGen                                    | ✓      | <b>0.3459</b>                | 0.3041                     | <b>0.3095</b>              | 0.0781                     | 0.0593        | 0.0645                     | 0.4232                     | 0.4507                     | 0.4323                     | 0.3815                     | 0.3265                     | 0.2969                     | 0.2777                     |
| FiD-RevGen <sup>a</sup>                      | ✓      | <u>0.3418</u> <sup>b,c</sup> | 0.3061                     | <u>0.3093</u> <sup>c</sup> | <b>0.0854</b> <sup>c</sup> | 0.0677        | <b>0.0723</b> <sup>c</sup> | <u>0.4322</u> <sup>c</sup> | <u>0.4571</u> <sup>b</sup> | <u>0.4400</u> <sup>b</sup> | <u>0.3858</u> <sup>c</sup> | <u>0.3310</u> <sup>c</sup> | <b>0.3020</b> <sup>c</sup> | <b>0.2826</b> <sup>c</sup> |
| FiD-AspGen                                   | ✗      | 0.3416                       | 0.2998                     | 0.3061                     | <u>0.0829</u>              | 0.0613        | 0.0671                     | 0.4217                     | 0.4466                     | 0.4296                     | 0.3762                     | 0.3224                     | 0.2937                     | 0.2754                     |
| Baselines                                    |        |                              |                            |                            |                            |               |                            |                            |                            |                            |                            |                            |                            |                            |
| <i>BART</i> (query + docs) [29] <sup>b</sup> | ✗      | 0.3003                       | <u>0.3199</u> <sup>a</sup> | 0.2984                     | 0.0770                     | <u>0.0694</u> | <u>0.0711</u>              | <b>0.4452</b> <sup>a</sup> | <b>0.4673</b>              | <b>0.4516</b>              | <b>0.3994</b> <sup>a</sup> | <b>0.3355</b>              | <u>0.3000</u>              | <u>0.2778</u>              |
| <i>Faspects</i> ensemble [29] <sup>c</sup>   | ✗      | 0.2252                       | <b>0.3296</b> <sup>a</sup> | 0.2588                     | 0.0499                     | <b>0.0795</b> | 0.0596                     | 0.3558                     | 0.3475                     | 0.3473                     | 0.2625                     | 0.2076                     | 0.1826                     | 0.1681                     |
| Oracles/Ablations                            |        |                              |                            |                            |                            |               |                            |                            |                            |                            |                            |                            |                            |                            |
| <i>BART</i> – closedbook                     | ✗      | 0.0779                       | 0.0530                     | 0.0591                     | 0.0294                     | 0.0226        | 0.0245                     | 0.3706                     | 0.3071                     | 0.3312                     | 0.1280                     | 0.0632                     | 0.0396                     | 0.0333                     |
| FiD-AspGen-closedbook*                       | ✗      | 0.2973                       | 0.2499                     | 0.2616                     | 0.0293                     | 0.0202        | 0.0227                     | 0.3956                     | 0.4188                     | 0.4034                     | 0.3448                     | 0.2927                     | 0.2648                     | 0.2480                     |
| FiD-AspGen-oracle – compressed               | ✗      | 0.9968                       | 0.9138                     | 0.9480                     | 0.9867                     | 0.8916        | 0.9296                     | 0.9001                     | 0.8997                     | 0.8999                     | 0.8991                     | 0.8988                     | 0.8965                     | 0.8825                     |
| FiD-AspGen-oracle                            | ✗      | 0.9999                       | 0.9999                     | 0.9999                     | 0.9999                     | 0.9999        | 0.9999                     | 1.0000                     | 1.0000                     | 1.0000                     | 1.0000                     | 1.0000                     | 0.9974                     | 0.9810                     |

**Table 1: Aspect generation performance when using *bing*-snippets [41] as evidence documents. Co-Gen stands for co-generating the template-based questions and target facets. Bold designates the top system while underline the second top excluding the oracle systems. Superscripts indicate statistically significant improvements wrt. FiD-RevGen <sup>a</sup>, BART <sup>b</sup> and the *Faspects* ensemble <sup>c</sup> (paired T-test,  $p$ -value < 0.05 and Bonferroni correction for multiple hypothesis testing).**

For metrics computed on the facet level (Set-BERT and Set-BLEU), we follow prior work and create the best matching pairs between *F* and *G* using the BLUE-1 score [9].

## 4.2 Baselines

Due to variations in (a) metric implementation, (b) BERT-Scores’ variability to package versions and (c) the use of different test sets <sup>3</sup>, reported results for a method often differ across papers, indicating a serious reproducibility issue [9, 29, 46]. To circumvent this issue and ensure a fair comparison, we only compare with publicly available baselines and use the evaluation code of *Faspect*[29]<sup>4</sup>.

We compare with *BART*, a SOTA encoder-decoder aspect extraction model that uses queries and evidence documents to generate facets from [29]. In comparison with our *FiD* models, *BART* is more heavy computationally, since its encoder computes cross-attentions between the evidence documents. Additionally, we compare with *Faspect* [29], a strong Recall-optimised ensemble model that relies on multiple diverse models for generating facets. Those models include the *BART* baseline we used here, as well as other Sequence Labelling, Extreme Facet Classification or Unsupervised Facet extraction methods. Ensembling is done using the Round-Robin algorithm [29], that practically interleaves results coming from different models until generating the maximum amount of facets (5). We do not compare with baselines that focus on optimising the training objectives for generating set-based predictions, since those are orthogonal to our work and beyond the scope of this paper [21].

## 4.3 Implementation details

We initialize the question generator (*FiD*) model from an *Atlas*-base checkpoint, which is pretrained with an unsupervised language modelling objective and exhibits good few-shot abilities [13]<sup>5</sup>. Unless stated differently, we use a training batch size of 32 to maximize our GPU memory usage, and a maximum generation length of 64 to fit the longest generation output. The maximum length of the encoded documents is also set to 64, since our collection consists of

small *bing* snippets and *MSMarco*-passage. We do early stopping while optimizing for Exact Match F1 (chosen due to reliability and computational efficiency) on a held out validation set, taken from *MIMICS*-Click. We perform our experiments on one Nvidia RTX A6000 GPU.

## 5 RESULTS

In this Section, we discuss our experimental results and the answers to our Research Questions.

### 5.1 Can we generate end-to-end Clarifications with Retrieval Augmented Generation?

In this section, we try to answer **RQ1**, that is whether we can generate Clarifying Questions end-to-end with Retrieval Augmented Generation models. In this part, we use as Retrieval Augmentation (evidence) the *BING* snippets provided along with the *MIMICS* dataset [40]. We compare *Fusion-in-Decoder* (*FiD*) models with other RAG and closed book models and present results on Table 1. Following prior work, we only report facet generation performance so results are comparable across different methods.

We experiment with three variants of *FiD*, depending on whether the model only predicts facets (*FiD-AspGen*) or jointly generates facets and the template-based questions of *MIMICS* (*FiD-CQGen*, *FiD-RevGen*). *FiD-CQGen* predicts question and facets in the original order of *MIMICS* (question followed by the aspects), while *FiD-RevGen* first outputs the aspects. We find that all variants are competitive, with *FiD-RevGen* performing best across Exact Match metrics while being able to co-generate questions and aspects. For this reason, in the rest of our experiments we use this variant, unless stated differently.

Next we compare *FiD* models with *BART*, the SOTA seq2seq baseline. We observe that results are mixed, with significant differences occurring for both models on different metrics. *FiD* significantly outperforms in Term Overlap Precision metrics, but loses in Term Overlap Recall and Set-BERT Precision, Overall, *BART* seems to be better in terms of recall, and also when compared using the Set-BERT metric, which computes facet-to-facet semantic

<sup>3</sup>See footnote 6 of Samarinas et al. [29]

<sup>4</sup><https://github.com/algoprog/Faspect/>

<sup>5</sup><https://github.com/facebookresearch/atlas/>

similarities. The Set-BERT metric, although used in relevant literature [9, 10, 29, 43], has significant shortcomings, as it favors method that generate a large number of facets. This is attributed to the fact that our methods generate less facets on average compared to the baselines. The FiD method is also on par with BART when using Set-BLUE, outperforming BART for larger n-grams. In general, BART achieves higher Recall while FiD is stronger in terms of Precision. This comparison between *FiD* and *BART* is important, because *FiD* fuses document embeddings in the decoder and therefore can model lengthier or more evidence documents than *BART*, with the same GPU memory footprint. The benefits of this efficiency are not visible here and are explored in Section 5.2.4, due to the evidence snippets being few and very short in length.

Comparing *FiD* with *Faspects*, a Recall-optimized ensemble that contains BART, we observe that this trend is magnified further. *Faspects* produces more facets which at times leads to better Recall, at the cost of much lower metrics in terms of Precision as well as Set-BERT and Set-BLEU measures.

Furthermore, we assess whether FiD can successfully compose answers that depend on multiple evidence documents. It is important to assess how the information bottleneck of FiD affects aspect generation performance, because FiD models are widely used on open-domain QA, a task that in contrast to ours does not necessarily require compositionality across evidence documents. To do so, we run an oracle experiment where we provide the target facets directly as evidence documents, either in the form of a single evidence document (*oracle*) or in multiple and independently encoded evidence documents (*oracle – compressed*). We observe only a slight drop in performance, and hence we conclude that FiD models can successfully utilise all evidence documents towards the generation. Lastly, we perform closed-book generation (without using evidence passages), we observe a big drop across performance metrics. By manual inspection, we observe that closed-book models are able to generate facet words for some queries, while they often resort in capturing reoccurring patterns from the training set (eg. on shopping related queries: “Are you shopping for: *men, women, kids*?”).

Our results in this Section suggest that RAG models and in particular *FiD* can be used to effectively and efficiently model the retrieval corpus and generate clarifying questions end-to-end.

Overall, we answer **RQ1** positively and conclude that Retrieval-Augmented-Generation models such as *FiD* are effective in jointly modelling the query, collection and generating clarifying questions in an end-to-end way. We find that those *FiD* is competitive with more computationally expensive baselines such as *BART*, suggesting that it can be effectively used to model larger parts of the collection and therefore generate more informed clarifying questions. Further, we verify the finding of previous works [29, 32] that evidence documents are crucial for generating good clarifying questions and strengthen our motivation towards performing *Corpus-informed Clarifying Question Generation*.

| Collection         | Retriever                   | Term Overlap<br>Recall | Exact Match<br>Recall | Evidence set type |
|--------------------|-----------------------------|------------------------|-----------------------|-------------------|
| Bing snippets [40] | <i>Bing</i>                 | 0.518                  | 0.185                 | diversified       |
| MSMarco-passage    | <i>BM25 Q</i>               | 0.448                  | 0.151                 | non-diversified   |
| MSMarco-passage    | <i>Contriever – FT Q</i>    | 0.434                  | 0.124                 | semantic          |
| MSMarco-passage    | <i>BM25 Q, F</i>            | 0.813                  | 0.380                 | aligned-lexical   |
| MSMarco-passage    | <i>Contriever – FT Q, F</i> | 0.761                  | 0.345                 | aligned-semantic  |

**Table 2: Alignment statistics between evidence documents and target facets in clarifying question**

## 5.2 Building evidence sets to support Corpus-informed Retrieval-Augmented Generation of Clarifying Questions

In this section, we focus on the role of retrieval and evidence documents, as those are fundamental to generating good clarifications. Our hypothesis here is that if evidence documents and ground truth facets are not aligned during training, that is if ground truth facets do not appear in the evidence documents, then models learn to produce these facets out of their internal model knowledge and eventually stop relying on the retrieved evidence [15]. This means that the model becomes corpus-agnostic and static, with a high risk to “hallucinate”, ie. present facets to the user that are either irrelevant to her query or do not exist in the retrieval corpus. Either of these cases would have *detrimental effects for user experience*. This would happen if users are presented with irrelevant facets, or facets that are irretrievable on this search collection *C*, causing the ultimate goal of search (retrieving a relevant document) to fail. Hence, we emphasize the importance of generating Corpus-informed Clarifying Questions. Overall, we aim to answer **RQ2**, that is how to build evidence sets that can support *Corpus-informed Clarifying Question* generation models. Specifically we look into the following questions: (a) does the corpus contain information on all ground truth facets, and can retrieval algorithms bring them up (Section 5.2.1), (b) what is the effect of missing or irretrievable facets in the clarifying question generation (Section 5.2.2), (c) do clarifying generation model remain faithful to given evidence documents and how does training affect that (Section 5.2.3), and (d) does increasing the document pool help in recovering missing facets 5.2.4)?

**5.2.1 Measuring alignment between target facets and evidence documents.** First, we investigate to what extent the evidence set (*BING* snippets) we used in Section 5.1, and is commonly used in the literature [9, 10, 29], is aligned with the target ground truth facets we try to generate. As a proxy for relevance, we use *Term-Overlap-Recall* that measures how many of all facet words appear in the evidence pool, and *Exact-Match-Recall* that measures whether the entire facet appears verbatim.

In Table 2, we measure the alignment of the provided *Bing-snippets* as well as evidence pools retrieved using lexical (*BM25*) and semantic (*Contriever*) methods. For *Bing snippets*, we observe that only 50% of facet words appear in the evidence pool. When retrieving documents from the *MSMarco* passage collection [20] using the query, alignment is even lower with 43% and 45% of the facet words appearing in the evidence set. Last when we use both the original query and corresponding ground truth facets to construct the evidence pool (by expanding the original query with a single facet and interleaving the top-K results of all facets into a

| Model | evidence (train&test)            | Term Overlap                 |                              |                              | Exact Match                  |                              |                              | Set-BERT                     |                              |                              | Set-BLEU                     |                              |                              |                              |
|-------|----------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|------------------------------|
|       |                                  | Prec.                        | Rec.                         | F1                           | Prec.                        | Rec.                         | F1                           | Prec.                        | Rec.                         | F1                           | BLEU1                        | BLEU2                        | BLEU3                        | BLEU4                        |
| BART  | [Bing]                           | 0.3003                       | <b>0.3199</b>                | 0.2984                       | 0.0770                       | <b>0.0694</b>                | 0.0711                       | <b>0.4452</b>                | <b>0.4673</b>                | <b>0.4516</b>                | <b>0.3994</b>                | <b>0.3355</b>                | 0.3000                       | 0.2778                       |
| FiD   | [Bing] <sup>a</sup>              | <b>0.3418</b> <sup>c</sup>   | 0.3061 <sup>c</sup>          | <b>0.3093</b> <sup>c</sup>   | <b>0.0854</b>                | 0.0677                       | <b>0.0723</b>                | 0.4322                       | 0.4571                       | 0.4400                       | 0.3858                       | 0.3310 <sup>c</sup>          | <b>0.3020</b> <sup>c</sup>   | <b>0.2826</b> <sup>c</sup>   |
| FiD   | [BM25 Q,F]                       | 0.4782                       | 0.4104                       | 0.4220                       | 0.1605                       | 0.1154                       | 0.1290                       | 0.4599                       | 0.4920                       | 0.4700                       | 0.4114                       | 0.3726                       | 0.3504                       | 0.3329                       |
| FiD   | [Contriever Q,F]                 | 0.4714                       | 0.4338                       | 0.4326                       | 0.1674                       | 0.1296                       | 0.1413                       | 0.4847                       | 0.5197                       | 0.4964                       | 0.4354                       | 0.3926                       | 0.3694                       | 0.3507                       |
| FiD   | [Contriever-FT Q,F] <sup>b</sup> | <b>0.5022</b> <sup>a,c</sup> | <b>0.4517</b> <sup>a,c</sup> | <b>0.4556</b> <sup>a,c</sup> | <b>0.1886</b> <sup>a,c</sup> | <b>0.1497</b> <sup>a,c</sup> | <b>0.1609</b> <sup>a,c</sup> | <b>0.4879</b> <sup>a,c</sup> | <b>0.5204</b> <sup>a,c</sup> | <b>0.4980</b> <sup>a,c</sup> | <b>0.4376</b> <sup>a,c</sup> | <b>0.3990</b> <sup>a,c</sup> | <b>0.3775</b> <sup>a,c</sup> | <b>0.3597</b> <sup>a,c</sup> |
| FiD   | [BM25 Q]                         | <b>0.3216</b>                | <b>0.2912</b>                | <b>0.2928</b>                | 0.0659                       | 0.0540                       | 0.0568                       | 0.4244                       | <b>0.4485</b>                | <b>0.4320</b>                | 0.3754                       | <b>0.3184</b>                | <b>0.2883</b>                | <b>0.2698</b>                |
| FiD   | [Contriever Q]                   | 0.3095                       | 0.2772                       | 0.2799                       | 0.0544                       | 0.0433                       | 0.0464                       | 0.4218                       | 0.4434                       | 0.4281                       | 0.3688                       | 0.3108                       | 0.2805                       | 0.2625                       |
| FiD   | [Contriever-FT Q] <sup>c</sup>   | 0.3202                       | 0.2905                       | 0.2917                       | <b>0.0742</b>                | <b>0.0592</b>                | <b>0.0630</b>                | <b>0.4245</b>                | 0.4462                       | 0.4309                       | <b>0.3754</b>                | 0.3175                       | 0.2868                       | 0.2679                       |

**Table 3: Effect of different evidence pools in facet generation performance (evidence used during training and inference). Superscripts <sup>a,b,c</sup> indicate significant improvements wrt. best model from each category (paired T-test,  $p$ -value < 0.05 and Bonferroni correction for multiple hypothesis testing. Best results across group are boldfaced.**

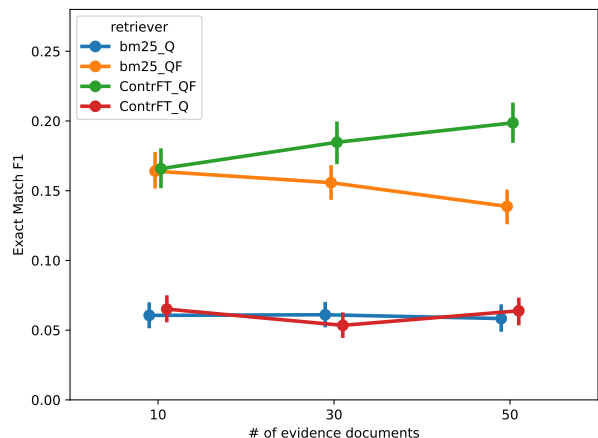
single ranking), the alignment statistics almost double. However, even in this case and when performing lexical retrieval ( $BM25|Q, F$ ), a large number of facet words ( $\sim 20\%$ ) are still absent from the evidence pool. This shows that many facets in this dataset might be irretrievable or not even existing, creating a shortcoming in terms of model training (possible hallucinations) and also in terms of evaluation and result reliability.

**5.2.2 Effect of retriever in generation performance.** In Table 3 we use these different evidence pools to explore how they affect downstream performance. It is evident that aligning evidence and generation ( $*|Q, F$ ) during training and inference brings a big boost and significant performance improvements across metrics. Most notably the strict Exact Match metric that increases from  $\sim 6$ -7% to up to  $\sim 16\%$  when our best retriever (*Contriever* – *FT*) is used. Overall, facet extraction scores follow the trend of the alignment statistics of Table 2. Using only the query ( $*|Q$ ) to retrieve evidence has slightly smaller facet alignment than the Bing snippets and this same trend is reflected in facet extraction performance.

We observe that using better retrievers (ie. the finetuned *Contriever*) only results in gains in the ( $*|Q, F$ ) setting. This shows that retrievers have to be biased towards the ground truth facets to find them, and suggests the presence of other prevalent facets within the collection. We explore this issue further in Sections 5.3.2 and 6. It is also noteworthy that the generator benefits significantly more from high-quality semantic retrieval (*Contriever* – *FT*) than lexical (*BM25*), even when measuring lexical metrics. This is important because the semantic pool of documents has less lexical overlap with the facets (0.345 vs 0.380, Table 2), yet *FiD* benefits more from this retriever pool, rather than the lexical one.

Therefore, we conclude that *FiD* can successfully extract and paraphrase facets from evidence pools without the need to find them verbatim in the evidence, and verify that retrieval quality is of great importance when generating clarifications.

**5.2.3 Faithfulness of Question Generators towards evidence documents.** Next, we test the faithfulness of question generators to the evidence pool. We emphasize the importance of this property to ensure clarifications are grounded in the retrieval corpus, where search will take place. To probe model faithfulness, we design a leave-one-out (*LOO*) evaluation, where documents corresponding to a facet are included or dropped from the evidence. We start from an evidence pool containing all facets (similar to  $*|Q, F$ ) and measure *Recall* of a random ground truth facet before or after removing its documents. We probe the generators trained on the Bing



**Figure 1: Effect of increasing the number of evidence documents in performance**

snippets evidence set, the non-diversified  $BM25|Q$  and the aligned generator  $BM25|Q, F$ .

As we can see on Table 4, aligned Fusion-in-Decoders are much more likely to capture the facet if the evidence documents contain it, in contrast to the non aligned ones (*TermOverlap* and *ExactMatch* – *Recall*). *Recall* – *LOO* performance is roughly the same across models, with facet recall of aligned models being slightly lower.<sup>6</sup>

At the same time, the drop in Recall of the facet that is removed from the evidence set ( $\Delta Recall$ ) is much larger in aligned models, demonstrating substantially more faithful clarifications towards the evidence documents. This hints that previous question generators lack sufficient evidence grounding, which is essential for creating *Corpus-informed Clarifications*.

**5.2.4 Increasing number of evidence documents.** One of the main advantages of Fusion-in-Decoder models over other encoder-decoder models is that they can take into account lengthier context when generating a response. In practice *FiD* models can be viewed as compressed encoder-decoders, which can model more evidence documents with the same GPU memory footprint. Here, we explore whether *FiD* can benefit from modelling more evidence documents when generating questions. Due to computational restrictions, in this experiment we reduce the Generation Length to 32 tokens instead of 64 and skip the question prediction using *FiD* – *AspGen*.

<sup>6</sup>Note that since facets often include query words on mimics (eg. “*Leiden* zip code”) Term Overlap metrics remain relatively high.

| Model           | Train evidence set | Term Overlap |              |                    | Exact Match |              |                    |
|-----------------|--------------------|--------------|--------------|--------------------|-------------|--------------|--------------------|
|                 |                    | Recall       | Recall – LOO | $\Delta$ Recall(%) | Recall      | Recall – LOO | $\Delta$ Recall(%) |
| Non-Diversified | BM25 Q             | 38.138       | 35.987       | -5.64%             | 6.047       | 3.369        | -44.30%            |
| Diversified     | Bing snippets [40] | 40.233       | 36.700       | -8.78%             | 7.833       | 3.531        | -54.92%            |
| Aligned         | BM25 Q, F          | 45.804       | 35.247       | -23.05%            | 14.692      | 3.044        | -79.28%            |

**Table 4: Faithfulness of question models to input evidence documents, using a leave-one-out (LOO) evaluation.**

| Model   | train. evidence | Term Overlap  |               |               | Exact Match   |               |               | Set-BERT      |               |               | Set-BLEU      |               |               |               |
|---|-----------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|   |                 | Prec.         | Rec.          | F1            | Prec.         | Rec.          | F1            | Prec.         | Rec.          | F1            | BLEU1         | BLEU2         | BLEU3         | BLEU4         |
| Test-set evidence set: bing snippets                  |                 |               |               |               |               |               |               |               |               |               |               |               |               |               |
| FiD   | [BING]          | <b>0.3418</b> | <b>0.3061</b> | <b>0.3093</b> | <b>0.0854</b> | <b>0.0677</b> | <b>0.0723</b> | <b>0.4325</b> | <b>0.4574</b> | <b>0.4403</b> | <b>0.3857</b> | <b>0.3308</b> | <b>0.3018</b> | <b>0.2826</b> |
| FiD   | [BM25 Q]        | 0.3257        | 0.3005        | 0.2993        | 0.0774        | 0.0615        | 0.0652        | 0.4218        | 0.4465        | 0.4296        | 0.3746        | 0.3197        | 0.2904        | 0.2719        |
| FiD   | [BM25 Q,F]      | 0.3353        | 0.2951        | 0.3015        | 0.0753        | 0.0542        | 0.0595        | 0.4185        | 0.4451        | 0.4269        | 0.3721        | 0.3185        | 0.2897        | 0.2716        |
| FiD   | [Contr-FT Q,F]  | 0.3217        | 0.3012        | 0.2983        | 0.0743        | 0.0604        | 0.0633        | 0.4232        | 0.4494        | 0.4313        | 0.3783        | 0.3212        | 0.2911        | 0.2725        |
| FiD   | [Contr-FT Q]    | 0.3269        | 0.2998        | 0.2993        | 0.0817        | 0.0665        | 0.0701        | 0.4252        | 0.4483        | 0.4321        | 0.3753        | 0.3201        | 0.2903        | 0.2717        |
| Test-set evidence set: non-diversified [BM25 Q]       |                 |               |               |               |               |               |               |               |               |               |               |               |               |               |
| FiD   | [Bing]          | <b>0.3228</b> | 0.2798        | 0.2868        | 0.0571        | 0.0451        | 0.0482        | 0.4022        | 0.4299        | 0.4117        | 0.3634        | 0.3068        | 0.2770        | 0.2584        |
| FiD   | [BM25 Q]        | 0.3216        | <b>0.2912</b> | <b>0.2928</b> | <b>0.0659</b> | <b>0.0540</b> | <b>0.0568</b> | <b>0.4244</b> | <b>0.4485</b> | <b>0.4320</b> | <b>0.3751</b> | <b>0.3183</b> | <b>0.2883</b> | <b>0.2698</b> |
| FiD   | [BM25 Q,F]      | 0.3126        | 0.2661        | 0.2763        | 0.0517        | 0.0386        | 0.0424        | 0.4067        | 0.4314        | 0.4146        | 0.3640        | 0.3060        | 0.2752        | 0.2567        |
| Test-set evidence set: non-diversified [Contr – FT Q] |                 |               |               |               |               |               |               |               |               |               |               |               |               |               |
| FiD   | [Contr-FT Q]    | <b>0.3202</b> | <b>0.2905</b> | <b>0.2917</b> | <b>0.0742</b> | <b>0.0592</b> | <b>0.0630</b> | <b>0.4245</b> | <b>0.4462</b> | <b>0.4309</b> | <b>0.3754</b> | <b>0.3175</b> | <b>0.2868</b> | <b>0.2679</b> |
| FiD   | [Contr-FT Q,F]  | 0.3133        | 0.2778        | 0.2831        | 0.0594        | 0.0477        | 0.0507        | 0.4172        | 0.4402        | 0.4243        | 0.3744        | 0.3149        | 0.2834        | 0.2644        |

**Table 5: Effect of evidence grounding during training (for different evidence sets at inference time)**

| Model                    | Inference evidence pool | Term Overlap  |               |               | Exact Match   |               |               | Set-BERT      |               |               | Set-BLEU      |               |               |               |
|--------------------------|-------------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                          |                         | Prec.         | Rec.          | F1            | Prec.         | Rec.          | F1            | Prec.         | Rec.          | F1            | BLEU1         | BLEU2         | BLEU3         | BLEU4         |
| [Contr-FT Q]             | [Contr-FT Q]            | <b>0.3202</b> | <b>0.2905</b> | <b>0.2917</b> | <b>0.0742</b> | <b>0.0592</b> | <b>0.0630</b> | <b>0.4245</b> | <b>0.4462</b> | <b>0.4309</b> | 0.3754        | 0.3175        | <b>0.2868</b> | <b>0.2679</b> |
| [Contr-FT Q,F]           | MMR ( $\lambda = 0.5$ ) | 0.3052        | 0.2763        | 0.2790        | 0.0559        | 0.0458        | 0.0485        | 0.4216        | 0.4441        | 0.4284        | 0.3772        | <b>0.3179</b> | 0.2861        | 0.2670        |
| [Contr-FT Q,F]           | MMR ( $\lambda = 0.7$ ) | 0.3085        | 0.2767        | 0.2807        | 0.0575        | 0.0471        | 0.0496        | 0.4217        | 0.4444        | 0.4286        | 0.3776        | 0.3174        | 0.2852        | 0.2660        |
| [Contr-FT Q,F]           | MMR ( $\lambda = 0.9$ ) | 0.3132        | 0.2784        | 0.2835        | 0.0598        | 0.0484        | 0.0512        | 0.4187        | 0.4416        | 0.4257        | <b>0.3758</b> | 0.3164        | 0.2844        | 0.2653        |
| Contr-FT-Train           |                         | 0.3001        | 0.2715        | 0.2738        | 0.0560        | 0.0451        | 0.0482        | 0.4198        | 0.4416        | 0.4262        | 0.3655        | 0.3070        | 0.2762        | 0.2578        |
| Contr-FT-Train-WarmedUpQ |                         | <b>0.3150</b> | <b>0.2935</b> | <b>0.2909</b> | <b>0.0723</b> | <b>0.0639</b> | <b>0.0652</b> | <b>0.4340</b> | <b>0.4553</b> | <b>0.4402</b> | <b>0.3795</b> | <b>0.3209</b> | <b>0.2899</b> | <b>0.2710</b> |

**Table 6: Effect of evidence diversification methods in performance**

In Figure 1, we observe that increasing the number of context documents leads to performance improvements only for *Contriever*|Q, F. This shows that the question generator can be improved with more documents, as long as noise is minimised in the evidence pool. This requires not only that retrieval is bounded towards the facets to be generated, but also that the retriever is of high quality (Contriever finetuned on the ranking task). For instance, performance quickly deteriorates when adding more BM25|Q, F documents, probably due to adding irrelevant documents to the pool. All in all, these results strengthen our motivation to use *Fusion-in-Decoder* models for this task, since they can effectively and efficiently model larger parts of the corpus than other seq2seq counterparts.

**Section Conclusions.** Experimental results on this Section highlight a significant gap between ground truth clarifications and the retrieval corpus in current datasets. When this gap appears models do not remain faithful to the facets present in the collection, i.e. clarifications are not corpus-informed, which can severely harm search engine user experience. We show that evidence grounding can mitigate this issue and significantly boost performance. We

further verify the suitability of *FiD* models for this task by showing that (i) it can effectively extract facets/topics from evidence without the need of finding them verbatim in the text, and (ii) its efficiency benefits can be used towards modelling larger parts of the collection and leading to generating better clarifying questions.

### 5.3 Introducing novelty in the evidence pool to assist facet extraction at test time

In previous sections, we showed that *FiD* is effective for generating grounded Clarifying Questions end-to-end. So far, we bounded the evidence pool towards the target facets, which is not possible at inference time. In this Section, we explore how we can relax this condition and capture the ground truth facets in an open-domain setting. To that end, we first examine whether grounded generators perform better than ungrounded generators in this non-bounded setting (Section 5.3.1). Second, we use a grounded generator that, as shown in Section 5.2 is capable to extract the facets of the evidence

pool, and investigate whether diversifying the evidence pool without constraining it towards the ground-truth facets can improve results.

**5.3.1 Is evidence-facet alignment during training enough to improve facet extraction?** In Table 5 we use three different evidence sets at inference time and measure performance of models trained in different ways. We observe that differences in performance metrics are minimal within groups, while models trained and tested with evidence pools from the same evidence distribution as the test distribution is somewhat better. This comes in contrast to the large differences observed in Table 3, where alignment brought more than double increases in Exact Match metrics, as well as Table 4 that showed that aligned generators are much more effective in extracting the facets present in the evidence. Therefore, our hypothesis remains that aligned generators do not perform better in those evidence sets because the facets do not exist there, and seek to improve this in the following Section.

**5.3.2 Evidence diversification.** Our results in the previous subsection showed that grounded question generators cannot find the ground truth facets when the retrieval is not bounded towards those facets. In this section we explore whether introducing novelty in the evidence set can help.

To induce novelty in the inference evidence pool, we use two sets of methods. First, we use the Maximal Marginal Relevance (MMR) algorithm [5] to rerank documents from an initially retrieved evidence pool of  $n = 50$  documents. MMR reranks documents taking into account their relevance with respect to a query, but also how similar they are to the rest of the ranking list. In practice, this is achieved with a  $\lambda$  parameter that is bounded between 0 and 1, where higher values of  $\lambda$  give more importance to relevance than novelty. Second, we experiment with training the dense retriever (Contriever-FT) with a diversification objective that uses knowledge distillation from the question generator [11]. This method has been commonly used in Question Answering, allowing retrievers to be trained on question-answer pairs, rather than query-document relevance pairs. However, we investigate this in a different task, where we seek to optimize retrieval for novelty or diversity of retrieved passages. In practice, our retriever here is trained on a distillation signal from the reader, which shows which documents were attended to generate the ground truth facets. In this way, we try to bias the retriever towards retrieving the ground truth facets.

**Introducing novelty using MMR.** The first part of Table 6 shows that inducing diversity through MMR in the evidence pool does not make the ground truth facets more likely to be retrieved and therefore extracted. Specifically, we see that increasing diversification results in lower performance. This is consistent with the findings of Samarinas et al. [29], that show that diversifying a list of extracted facets with MMR does not lead to improvements in MIMICS. Interestingly, in the related area of Multi-Answer retrieval, other works also report that MMR fails to improve results, since it increases diversity but hurts relevance significantly [18].

**Introducing novelty using Retriever-Reader Knowledge Distillation.** In the last part of Table 6, we experiment with finetuning the retriever jointly with the reader, using knowledge distillation from the ground truth facets. As usual, the retriever is initialised from the checkpoint finetuned on MSMarco. For initialising the

reader, we have two different variants: *Contr – FT – Train* is initialised directly from the Atlas checkpoint, pre-trained with the unsupervised Masked Language Modelling objective, while *Contr – FT – Train – WarmedUp* is initialised from *[Contr – FT|Q]*. We observe that directly finetuning a reader and retriever from the Atlas checkpoint has a detrimental effect on performance, which may be because the reader has to first learn how to adapt the generation task, from the unsupervised Masked Language Modelling objective used in pretraining to the task of Clarifying Question generation. However, when we initialise the model from *[Contr – FT|Q]*, training the retriever for diversification has a slight positive impact on certain performance metrics, particularly the semantic *Set – BERT* and *Set – Blue* metrics. In the rest of the metrics, performance remains at worst competitive with *[Contr – FT|Q]*, ie. the model we initialised from.

We conclude that training the retriever with knowledge distillation from the reader can have a slight positive impact in performance, yet bringing up documents related to the ground truth facets remains a challenge in this dataset and task.

## 6 ANALYSIS AND DISCUSSION

**Implications of Taxonomy Bias.** Our preliminary qualitative analysis suggested that certain query-facet patterns are heavily repeated within the MIMICS dataset. This can be attributed to the dataset construction process, during which facets were extracted based on search engine query logs and a taxonomy defined by the authors. To detect such patterns, we extract the top-20 frequent facet words (excluding stop-words) and inspect a couple of queries containing them. Based on those, we report the most frequent detected patterns in Table 7 and try to quantify to what extent this taxonomy bias affects the dataset. Our conservative estimate suggests that at least 1/5-th of dataset queries are biased towards the predetermined taxonomy. We arrive to this conclusion given that roughly 20% of dataset queries contain facets from this incomplete taxonomy (defined just on the top-20 facet keywords).

The implications of this for evaluation are important. Our experiments and related work, show that model performance has an upper-bound of 5% – 15% in Exact Match metrics (depending on the type of Retrieval Augmentation). Yet, 20% of the dataset is (at least to some extent) described by a fixed taxonomy. In such a setting, it is reasonable to assume that models over-fitting to the patterns of Table 7 might outperform models that produce more diverse but equally reasonable query facets. Further, models trained on this dataset are very likely to be biased towards the underlying taxonomy rather than being *Corpus-informed*, raising questions regarding their ability to generalize to open-domain settings.

This, in combination with the finding that many of the facets are irretrievable (Table 2) suggests that current test collections and evaluation frameworks are not suitable for evaluating *Corpus-Informed* clarifying question generation.

**Qualitative analysis.** In Table 8 we analyse a handful of random queries, presenting their ground-truth and generated facets. We show facets generated by our aligned model (train. evidence: *Contr-FT|Q,F*), and tested in an open-domain setting without diversification (test evidence: *Contr-FT|Q*) (second-to-last row of Table



| Re-occurring patterns  | Example  |
|--|--|
| <i>query type</i> : [facets]                                       |  |
| <i>software-related</i> :<br>[operating system]                    | <i>media creation tool</i> :<br>windows [10,7,8]                       |
| <i>shopping query</i> : [review, manual, specs, for sale]          | <i>huawei phones</i> : [manual, specs, battery, review, accessories]   |
| <i>location</i> : [hotels, restaurants, population]                | <i>jersey city</i> : [restaurants, hotels, population, homes for sale] |
| <i>location</i> : [zip code, weather, things to do]                | <i>leiden</i> : [things to do, weather, zip code, what time is it]     |
| <i>movie or series</i> : [cast, trailer, quotes, review]           | <i>her alibi movie</i> : [cast, trailer, review, quotes]               |
| <i>medical condition</i> : [symptom, treatment, causes, diagnosis] | <i>adhd</i> : [symptom, treatment, causes, diagnosis, diet]            |
| <i>act</i> : [checklist, tips, hacks]                              | <i>getting a kitten</i> : [tips, checklist, life hacks, pros]          |
| <i>various</i> : [men, women, kids]                                | <i>short layered hairstyles</i> :<br>[women, men, teens, kids]         |
| Dataset percentage: 20.74%   |  |

Table 7: Reoccurring facet patterns in the MIMICS dataset

5). Firstly, we observe that even the model trained with alignment learns to reproduce common re-occurring facet patterns from the training set (eg. “windows 7,8,10”). In fact, this still happens when evidence documents are unrelated to those facet patterns, as only 2/10 evidence documents mention operating systems versions, while the topic of the documents is not specifically related to those. Another observation that stands out is that the ground-truth facets are not exhaustive and other equally good or better facets can be generated. Such examples include **episodes** of penny appearing in the big bang theory, or referring to the leiden **thrombophilia** disorder rather than aspects related to the city.

Given those static ground truths, evaluation metrics would be very low, despite the overall high quality. For the queries “network drive”, “leiden”, “sickle cell anemia” and “consumer price index”, both Precision and Recall would be zero, despite the quality and greater diversity of the produced facets. This highlights the the presence of **multiple ground truths**, which are in fact ignored in the current evaluation framework.

This analysis highlights an important gap in current datasets and evaluation protocols, that prohibits the generation of *Corpus-informed* clarifying questions: Ground truth facets are constructed from query reformulations and static taxonomies, but do not reflect the facets that are often present in search collections.

## 7 CONCLUSIONS AND FUTURE DIRECTIONS

In this paper, we investigated the task of generating *Corpus-Informed* Clarifying Questions end-to-end, based on Retrieval Augmented Generation models. We showed that *Fusion-in-Decoder* models are able to effectively and efficiently model queries and evidence documents when generating clarifying questions. This efficiency advantage allows them to model larger parts of the corpus when asking clarification questions, potentially improving question quality.

Further, we investigated the role of retrieval in this task, showing that retrieval quality is more important for the generator than finding exact lexical matches of the facets to be generated. We

| query                 | facets<br>(ground truth)  | facets<br>(generated)              |
|-----------------------|---|------------------------------------|
| network drive         | set up, remove  | windows 7, windows 8, windows 10   |
| leiden                | what time is it, zip code, things to do, weather                          | <b>leiden thrombophilia</b>        |
| penny big bang theory | cast, quotes  | <u>dvd</u> , cast, <b>episodes</b> |
| pemetrexed            | side effects, chemotherapy injection, mechanism of action, package insert | side effects, <b>cost</b>          |
| sickle cell anemia    | diagnosis, causes, diet, symptom, treatment                               | <b>in children, in adults</b>      |
| consumer price index  | U.S., india, japan  | <b>chart, calculator</b>           |

Table 8: Qualitative analysis. Generated facets with an aligned question generator (train. evidence: *Contr-FT|Q,F*) in an open domain setting (test evidence: *Contr-FT|Q*). Underlined indicate irrelevant facets, and **boldfaced** indicate high-quality facets outside of the ground truth.

showed that current datasets suffer from a lack of grounding between ground truth facets and evidence documents, which has serious implications. Specifically, it hinders their ability of models to ask *Corpus-informed* questions, making them prone to “hallucinations” (ie. ignoring evidence documents while generating facets), which can severely harm user experience.

Lastly, we attempt to construct an evidence pool that contains the ground truth facets by inducing novelty in the retriever. We experiment with two novelty methods, (i) *MMR* and (ii) training the retriever with knowledge distillation, but we conclude that doing so remains a challenge. Our analysis sheds light in those challenges and highlights the gap between ground truth facets and those present in existing ranking corpora, calling for the need of datasets that can support generating *Corpus-Informed* clarifications. Future work should focus in developing appropriate train and test collections that better reflect the objective of creating *Corpus-Informed* Clarifying Questions. Novelty and diversity of the evidence pool remains important for generating good questions. To this end, exploring probabilistic retrieval methods [37, 39], and combining them with retriever training techniques, such as reader distillation remains a promising direction.

## REFERENCES

- [1] Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeff Dalton, and Mikhail Burtsev. 2020. ConvAI3: Generating clarifying questions for open-domain dialogue systems (ClariQ). *arXiv preprint arXiv:2009.11352* (2020).
- [2] Mohammad Aliannejadi, Hamed Zamani, Fabio Crestani, and W Bruce Croft. 2019. Asking clarifying questions in open-domain information-seeking conversations. In *Proceedings of the 42nd international acm sigir conference on research and development in information retrieval*. 475–484.
- [3] Akari Asai, Sewon Min, Zexuan Zhong, and Danqi Chen. 2023. Retrieval-based Language Models and Applications. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 6: Tutorial Abstracts)*. 41–46.
- [4] Sandeep Avula, Bogeum Choi, and Jaime Arguello. 2022. The effects of system initiative during conversational collaborative search. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW1 (2022), 1–30.
- [5] Jaime Carbonell and Jade Goldstein. 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*. 335–336.
- [6] Jiahui Chen and Jonas Mueller. 2023. Quantifying uncertainty in answers from any language model via intrinsic and extrinsic confidence assessment. *arXiv preprint arXiv:2308.16175* (2023).
- [7] Charles LA Clarke, Nick Craswell, and Ian Soboroff. 2009. Overview of the TREC 2009 Web Track.. In *Trec*, Vol. 9. 20–29.
- [8] Kevyn Collins-Thompson, Craig Macdonald, Paul N Bennett, Fernando Diaz, and Ellen M Voorhees. 2014. TREC 2014 Web Track Overview.. In *TREC*, Vol. 13. 1–15.
- [9] Helia Hashemi, Hamed Zamani, and W Bruce Croft. 2021. Learning multiple intent representations for search queries. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 669–679.
- [10] Helia Hashemi, Hamed Zamani, and W Bruce Croft. 2022. Stochastic Optimization of Text Set Generation for Learning Multiple Query Intent Representations. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 4003–4008.
- [11] Gautier Izacard and Edouard Grave. 2020. Distilling knowledge from reader to retriever for question answering. *arXiv preprint arXiv:2012.04584* (2020).
- [12] Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282* (2020).
- [13] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299* (2022).
- [14] Antonios Minas Krasakis, Mohammad Aliannejadi, Nikos Voskarides, and Evangelos Kanoulas. 2020. Analysing the effect of clarifying questions on document ranking in conversational search. In *Proceedings of the 2020 ACM SIGIR on International Conference on Theory of Information Retrieval*. 129–132.
- [15] Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. Hurdles to progress in long-form question answering. *arXiv preprint arXiv:2103.06332* (2021).
- [16] Vaibhav Kumar et al. 2020. ClarQ: A large-scale and diverse dataset for clarification question generation. *arXiv preprint arXiv:2006.05986* (2020).
- [17] Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 2021. Pyserini: A Python toolkit for reproducible information retrieval research with sparse and dense representations. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2356–2362.
- [18] Sewon Min, Kenton Lee, Ming-Wei Chang, Kristina Toutanova, and Hannaneh Hajishirzi. 2021. Joint passage ranking for diverse multi-answer retrieval. *arXiv preprint arXiv:2104.08445* (2021).
- [19] Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. AmbigQA: Answering ambiguous open-domain questions. *arXiv preprint arXiv:2004.10645* (2020).
- [20] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A human generated machine reading comprehension dataset. In *CoCo@ NIPS*.
- [21] Shiyu Ni, Keping Bi, Jianfeng Guo, and Xueqi Cheng. 2023. A Comparative Study of Training Objectives for Clarification Facet Generation. In *Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region*. 1–10.
- [22] Wenjie Ou and Yue Lin. 2020. A clarifying question selection system from ntes along in convai3 challenge. *arXiv preprint arXiv:2010.14202* (2020).
- [23] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.
- [24] Gustavo Penha, Alexandru Balan, and Claudia Hauff. 2019. Introducing mantis: a novel multi-domain information seeking dialogues dataset. *arXiv preprint arXiv:1912.04639* (2019).
- [25] Chen Qu, Liu Yang, W Bruce Croft, Johanne R Trippas, Yongfeng Zhang, and Minghui Qiu. 2018. Analyzing and characterizing user intent in information-seeking conversations. In *The 41st international acm sigir conference on research & development in information retrieval*. 989–992.
- [26] Hossein A Rahmani, Xi Wang, Yue Feng, Qiang Zhang, Emine Yilmaz, and Aldo Lipani. 2023. A Survey on Asking Clarification Questions Datasets in Conversational Systems. *arXiv preprint arXiv:2305.15933* (2023).
- [27] Sudha Rao and Hal Daumé III. 2018. Learning to ask good questions: Ranking clarification questions using neural expected value of perfect information. *arXiv preprint arXiv:1805.04655* (2018).
- [28] Sudha Rao and Hal Daumé III. 2019. Answer-based adversarial training for generating clarification questions. *arXiv preprint arXiv:1904.02281* (2019).
- [29] Chris Samarinas, Arkin Dharawat, and Hamed Zamani. 2022. Revisiting Open Domain Query Facet Extraction and Generation. In *Proceedings of the 2022 ACM SIGIR International Conference on Theory of Information Retrieval*. 43–50.
- [30] Ivan Sekulić, Mohammad Aliannejadi, and Fabio Crestani. 2021. Towards facet-driven generation of clarifying questions for conversational search. In *Proceedings of the 2021 ACM SIGIR international conference on theory of information retrieval*. 167–175.
- [31] Ivan Sekulić, Mohammad Aliannejadi, and Fabio Crestani. 2022. Evaluating mixed-initiative conversational search systems via user simulation. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 888–896.
- [32] Ivan Sekulić, Mohammad Aliannejadi, and Fabio Crestani. 2022. Exploiting document-based features for clarification in conversational search. In *European Conference on Information Retrieval*. Springer, 413–427.
- [33] Leila Tavakoli, Johanne R Trippas, Hamed Zamani, Falk Scholer, and Mark Sanderson. 2022. MIMICS-Duo: Offline & Online Evaluation of Search Clarification. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 3198–3208.
- [34] Leila Tavakoli, Hamed Zamani, Falk Scholer, William Bruce Croft, and Mark Sanderson. 2022. Analyzing clarification in asynchronous information-seeking conversations. *Journal of the Association for Information Science and Technology* 73, 3 (2022), 449–471.
- [35] Jian Wang and Wenjie Li. 2021. Template-guided clarifying question generation for web search clarification. In *Proceedings of the 30th ACM international conference on information & knowledge management*. 3468–3472.
- [36] Zhenduo Wang, Yuancheng Tu, Corby Rosset, Nick Craswell, Ming Wu, and Qingyao Ai. 2023. Zero-shot Clarifying Question Generation for Conversational Search. In *Proceedings of the ACM Web Conference 2023*. 3288–3298.
- [37] Frederik Warburg, Marco Miani, Silas Brack, and Soren Hauberg. 2023. Bayesian metric learning for uncertainty quantification in image retrieval. *arXiv preprint arXiv:2302.01332* (2023).
- [38] Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. Do Large Language Models Know What They Don't Know?. In *Findings of the Association for Computational Linguistics: ACL 2023*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 8653–8665. <https://doi.org/10.18653/v1/2023.findings-acl.551>
- [39] Hamed Zamani and Michael Bendersky. 2023. Multivariate Representation Learning for Information Retrieval. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 163–173.
- [40] Hamed Zamani, Susan Dumais, Nick Craswell, Paul Bennett, and Gord Lueck. 2020. Generating clarifying questions for information retrieval. In *Proceedings of the web conference 2020*. 418–428.
- [41] Hamed Zamani, Gord Lueck, Everest Chen, Rodolfo Quispe, Flint Luu, and Nick Craswell. 2020. Mimics: A large-scale data collection for search clarification. In *Proceedings of the 29th ACM international conference on information & knowledge management*. 3189–3196.
- [42] Hamed Zamani, Bhaskar Mitra, Everest Chen, Gord Lueck, Fernando Diaz, Paul N Bennett, Nick Craswell, and Susan T Dumais. 2020. Analyzing and learning from user interactions for search clarification. In *Proceedings of the 43rd international acm sigir conference on research and development in information retrieval*. 1181–1190.
- [43] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675* (2019).
- [44] Yongfeng Zhang, Xu Chen, Qingyao Ai, Liu Yang, and W Bruce Croft. 2018. Towards conversational search and recommendation: System ask, user respond. In *Proceedings of the 27th acm international conference on information and knowledge management*. 177–186.
- [45] Zhiling Zhang and Kenny Zhu. 2021. Diverse and specific clarification question generation with keywords. In *Proceedings of the Web Conference 2021*. 3501–3511.
- [46] Ziliang Zhao, Zhicheng Dou, Yu Guo, Zhao Cao, and Xiaohua Cheng. 2023. Improving Search Clarification with Structured Information Extracted from Search Results. (2023).
- [47] Ziliang Zhao, Zhicheng Dou, Jiaxin Mao, and Ji-Rong Wen. 2022. Generating clarifying questions with web search results. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 234–244.

[48] Jie Zou, Yifan Chen, and Evangelos Kanoulas. 2020. Towards question-based recommender systems. In *Proceedings of the 43rd international ACM SIGIR conference*

*on research and development in information retrieval*. 881–890.