

# Exploring efficient zero-shot synthetic dataset generation for Information Retrieval

Anonymous ACL submission

## Abstract

The broad integration of neural retrieval models into Information Retrieval (IR) systems is significantly impeded by the high cost and laborious process associated with the manual labelling of training data. Similarly, synthetic training data generation, a potential workaround, often requires expensive computational resources due to the reliance on large language models. This work explored the potential of small language models for efficiently creating high-quality synthetic datasets to train neural retrieval models. We aim to identify an optimal method to generate synthetic datasets, enabling training neural reranking models in document collections where annotated data is unavailable. We introduce a novel methodology, grounded in the principles of information theory, to select the most appropriate documents to be used as context for question generation. Then, we employ a small language model for zero-shot conditional question generation, supplemented by a filtering mechanism to ensure the quality of generated questions. Extensive evaluation on five datasets unveils the potential of our approach, outperforming unsupervised retrieval methods such as BM25 and pretrained monoT5. Our findings indicate that an efficiently generated "silver-standard" dataset allows effective training of neural rerankers in unlabeled scenarios. To ensure reproducibility and facilitate wider application, we will release a code repository featuring an accessible API for zero-shot synthetic question generation.

## 1 Introduction

Deep Learning is at the heart of many current breakthroughs in AI in a wide range of fields. Typically, such progress is attributed to better computational capabilities, superior algorithms, and a larger corpus of high-quality training data. Particularly in the Information Retrieval (IR) field, significant gains against traditional baselines are obtained when a

large amount of labelled data is available. However, manual data labelling is expensive and labor-intensive, highlighting the urgency to devise methods that can automatically produce higher quality training data to unlock the potential of neural retrieval models for unlabelled data collections.



Figure 1: Overview of the process of generating synthetic questions with LM for information retrieval.

Recent strides in large language models offer a new avenue of generating synthetic training data to train neural retrieval models. Present strategies largely fall into two categories, finetune-based and prompt-based. The former necessitates annotated data to train a language model to craft questions given a document text and, optionally, a correct answer. In contrast, the prompt-based method capitalizes on expensive language models to generate a question in a zero-shot fashion, using a document as context. Although both techniques are effective, they still have some drawbacks.

The finetune-based approach is a supervised method, thus requiring the acquisition of labelled data. Moreover, even though publicly available models can be adopted, these inevitably bear inherent biases from their training dataset, which can be a limiting factor in adapting to the target domain. On the other hand, the prompt-based approach, often linked to large models, comes with steeper costs, be it for model execution or through paid APIs. This particularly restricts its applicability in low-resource environments. Another overlooked problem that is rooted in both approaches is that in IR the target document collection for which synthetic questions are being generated usually contains millions of documents. It is therefore common to randomly select some documents as

seeds to generate the synthetic dataset. However, some documents can be bad examples, leading the generator to produce unuseful or invalid questions, wasting computation resources.

In this work, we explore the limits of prompt-based small language models in generating high-quality synthetic training data. Specifically, we hypothesize that these models can efficiently and quickly create a synthetic dataset, which can then empower neural retrieval models to outperform traditional unsupervised techniques such as BM25. Our approach starts with an innovative filtering technique rooted in information theory measures to identify and exclude non-representative documents. We then investigate various small language models and generation strategies across diverse document collections, gauging their capacity for producing relevant questions. To further improve the quality of the generated dataset, we also explore filtering techniques to remove less suitable questions. Lastly, we assess the performance of simple neural retrieval models trained with the best synthetic datasets.

Our contributions can be summarized as follows: (1) an innovative method grounded in information theory principles for discovering outliers within a document collection; (2) the development and validation of techniques to estimate the quality of synthetic generated questions; (3) an extensive benchmark of the quality of synthetic datasets for document retrieval, derived from several small language models and generation strategies, totalling 150 unique configurations; (4) publicly available off-the-shelf software tool for creating synthetic datasets for a given document collection.

## 2 Related Work

The field of synthetic data generation has seen significant advances with the advent of deep learning, mostly thanks to the transformer-based large language models capability of generating coherent text (Brown et al., 2020a; Chowdhery et al., 2022). Following the same trend, generating synthetic training data for Information Retrieval became a viable option to replace the labour-intensive data annotation process.

On the one hand, we have the finetune-based approaches initially popularized by Nogueira et al. 2019a,b as the Doc2Query technique, where the main idea was to train a sequence-to-sequence model to generate a question given a document

as input. However, its purpose was not to build a synthetic dataset, but rather to add the generated questions to the document to aid lexical models. Then, Nogueira and Lin 2019 improved the initial approach by adopting T5 as the generator model. More recently, Gospodinov et al. 2023 showed that sequence-to-sequence models are prone to “hallucination”, suggesting the incorporation of pretrained relevance models to weed out inaccurate questions. Meanwhile, Ma et al. 2021; Thakur et al. 2021; Wang et al. 2022 adopted a similar methodology, but with the primary objective to construct a synthetic dataset for training neural retrieval models in unlabelled document collections.

Opposed to the previous trend, zero-shot question generation, also known as prompt-based, has recently emerged as a promising alternative that involves generating questions without training a generation model specifically for that task. Large language models (LLMs) are typically used in zero-shot question generation, given their capability of generating coherent text and being easily conditioned to produce the desired output without needing extra training. For instance, Bonifacio et al. 2022 and Dai et al. 2023 obtained promising results in the creation of zero-shot synthetic datasets for information retrieval by using LLMs, namely GPT-3 (Brown et al., 2020a). Nevertheless, the deployment of LLMs on a larger scale remains challenging due to their extensive computational resource requirements.

Our work resonates most with the approach presented by Bonifacio et al. 2022, given the shared focus on zero-shot question generation utilizing language models for IR. Notwithstanding, in this work, we focused on only exploring small language models (from 70M to 1.3B parameters) while entirely concentrating on the problem of effectively and efficiently producing a synthetic dataset for information retrieval. As such, contrary to previous works, herein we explore the limits of zero-shot question generation with small language models by evaluating the impact of different language models and generation strategies, as well as a mechanism for document outlier detection.

## 3 Methods

This section details all the individual components that we explored in order to generate a synthetic dataset for document retrieval, followed by the evaluation methodology.

### 3.1 Document sampling method

In real-world retrieval scenarios with document collections spanning millions of documents, it is impractical to generate questions for every single document. As a result, a common approach has been to randomly select a subset of documents. However, this carries the issue of potentially selecting unrepresentative documents (i.e., documents that are considerably different from the rest of the collection or contain errors), leading to questions with poor quality.

To mitigate this, we propose to estimate the information content of each document and contrast it with the collection’s average. This facilitates the identification of outlier documents, which would be documents that substantially diverge from the average. By excluding these documents from the sampling process, we enhance the likelihood of choosing good documents. We leverage the information theory framework, which states that the amount of information of an event,  $x$ , can be computed as the negative log-likelihood of that event, as shown in Equation 1. For clarity, in our information estimation we adopt a notation akin to Lesne 2014.

$$I(x) = -\log(P(x)). \quad (1)$$

In our context, we consider that the event,  $x$ , represents the sequence of tokens that compose each document,  $x = \{w_1, w_2, \dots, w_N\}$ , where  $w_i$  represents the  $i$ -th token and  $N$  as the total number of tokens in the document. Then, the associated probability of that document text can be estimated by any language model through  $P(x) = \prod_{i=1}^N P(w_i|w_1, \dots, w_{i-1})$ . When plugging this into the previous equation, we obtain a formula to estimate each document’s information, as shown in Equation 2.

$$I(x) = -\sum_{i=1}^N \log(P(w_i|w_1, \dots, w_{i-1})). \quad (2)$$

One challenge with the above measure is its dependence on document length, potentially causing discrepancies when comparing diverse documents. Namely, lengthier documents might seem more informative solely due to their increased token count. To rectify this, we normalize the measure by the information estimated from a uniform model, resulting in the Normalized Information (NI) measure

defined in Equation 3. This type of normalization is not new and is commonly adopted in genetics in the context of complexity and compression, and is known as Normalized Compression (Pinho et al., 2010).

$$NI(x) = \frac{-\sum_{i=1}^N \log(P(w_i|w_1, \dots, w_{i-1}))}{|x| \times \log(|V|)}. \quad (3)$$

Here,  $V$  represents the vocabulary set comprising all valid tokens and  $|\cdot|$  is the length operator. While NI’s lower-bound is zero, its maximum is theoretically unbounded. However, a good probabilistic model would typically yield NI values that are bounded between  $[0, 1]$ . Intuitively, higher values of NI would represent documents that are close to randomness, while lower values should correspond to documents that are highly repetitive.

To estimate NI, we propose to adopt small transformer open-domain language models and finite-context-models trained directly on the corpus. In Appendix A we address the differences between both approaches.

### 3.2 Question generation with small LM

To synthesize questions for a given document, we use an engineered prompt that conditions a language model to produce a question based on the information contained within the document. More formally, we construct the prompt, denoted as  $p$ , that maximises the likelihood of the language model generating a question, denoted as  $y$ . This process is conducted according to Equation 4, where  $y_1$  represents a question initiator as discussed later,

$$\hat{y} \sim P(y|p_1, \dots, p_M, y_1). \quad (4)$$

Although prompt engineering is a relatively recent topic, there is already a vast literature on the topic, ranging from simple zero-shot to few-shot (Brown et al., 2020b), chain-of-thought (Wei et al., 2022a) and ReAct (Yao et al., 2023) techniques. The central idea behind these techniques is to gradually increase the prompt complexity with actual task-related examples, such that the generated text would be better aligned with the desired output. However, while these techniques have shown promising results in large language models, the same cannot be said for small language models (Wei et al., 2022b). Coupled with the observation that the memory requirements of transformer-based

models grows quadratically with input size, we opted for a simple zero-shot prompting technique in our experiments.

To steer the model towards question generation, we infused the prompt with question-initiating phrases. By doing so, the model is more inclined to proceed with contextually appropriate wording rooted in the starting phrase. Common initiators include: {What, How, When, Is, Does}. Prompt 1 showcases our approach for questions commencing with "What." To further refine outputs, only questions culminating in a question mark were deemed valid.

```
Article: {selected_article}
Question: What
```

Prompt 1: Zero-shot prompt for generation questions that start with the word "What".

As previously mentioned, we explored several language models and generation strategies. Specifically, we investigated beam search (Freitag and Al-Onaizan, 2017), contrastive search (Su et al., 2022), and random sampling (Fan et al., 2018) as potential methods for question generation. Random sampling, while preferred for larger models owing to its efficiency and adeptness at harnessing their robust probabilistic knowledge, may fall short with smaller models (Su et al., 2022). Consequently, we seek to ascertain if deterministic algorithms like beam and contrastive search can strike a more optimal balance between efficiency and output quality than random sampling.

### 3.3 Accessing the question quality

Although we enforce the model to generate questions, there is still a need to ensure the quality of these questions, specially considering that language models are prone to produce erroneous or unrelated outputs, a phenomenon referred to as "hallucination". Numerous studies have focused on preventing or filtering out these wrong synthetic samples. With special interest for question generation, (Lu et al., 2022; Alberti et al., 2019; Dai et al., 2023; Gospodinov et al., 2023) have suggested solutions based on retrieval methods and probability-based methods. The former employs neural relevance models to estimate the relevance of the question-document pairs, discarding those with lower relevance. Meanwhile, the latter ranks each generated question by its conditional probability, eliminating those that fall below a pre-defined K-cut-off region.

In this work, we propose two primary criteria that a good synthetic question must meet:

- **Relevance to the Article:** Each generated question should pertain directly to the content of the article provided in the prompt.
- **Suitability for Retrieval:** Each generated question must be suitable for retrieval, i.e., must look for information within the collection.

The first criterion ensures that the generated question-article pairs serve as training examples, given that the article contains the answer to the question. The second criterion prevents overly generic questions, such as "What is this document about?", which are non-representative of genuine retrieval scenarios. In practice, we adopted unsupervised retrieval methods to fulfill both criteria. Although probability-based methods may remove questions unrelated to the article, they would struggle to filter out questions unsuitable for retrieval, as these methods do not incorporate any retrieval concept. Hence, we defined a binary function  $f_k(x; m)$ , in Equation 5, that based on the model,  $m$ , and the threshold,  $k$ , evaluates if the question-document pair,  $x = (q, d)$ , has quality (1) or not (0).

$$f_k(x; m) = \begin{cases} 1, & \text{if (type}(m)=\text{prob and } m(x) \geq k \\ & \text{or (type}(m)=\text{rank and } m(x) \leq k) \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

During our experiments, we utilized both BM25 (Robertson and Zaragoza, 2009) and monoT5 (Nogueira et al., 2020) as potential models, represented by  $m$ . It is noteworthy that while monoT5 functions as a relevance model, BM25 is a retrieval-based model. As such, the threshold  $k$  for monoT5 is defined in terms of probability, whereas for BM25, it pertains to ranking position.

### 3.4 Evaluation procedure

Our main goal is to explore several small language models, generation strategies and quality assessment mechanism to discover the most cost-efficient configuration for creating a synthetic dataset for document retrieval. To accomplish this, we first propose a two-step benchmarking process. In the first step, we benchmark all configurations based on the number of good questions that are generated (1). This initial evaluation will give us insight



362 into which configuration performs best. Then, as  
363 a second step, we aim to evaluate a more realistic  
364 scenario by benchmarking the best configurations  
365 in a downstream reranking evaluation task (2).

### 366 3.4.1 Question quality benchmark

367 Before delving into the first benchmark, let us  
368 define a synthetically generated dataset contain-  
369 ing a set of positive question-document pairs as  
370  $\mathbb{D}_s = (q_0, d_0), \dots, (q_N, d_N)$ . Likewise, let us repre-  
371 sent  $f_k(x; m)$  as a function capable of estimating  
372 question quality, as introduced in Section 3.3.

373 To assess the synthetic datasets quality, we pro-  
374 pose a hits-ratio-based evaluation metric, defined  
375 in Equation 6. This metric quantifies the proportion  
376 of valid question-document pairs.

$$377 \text{hitsR}_k(\mathbb{D}_s) = \frac{\sum_{x \in \mathbb{D}_s} f_k(x; m)}{|\mathbb{D}_s|}. \quad (6)$$

378 Additionally, to account for each configuration’s  
379 runtime, we propose using a hits-per-second vari-  
380 ant, defined in Equation 7. This metric incorporates  
381 the elapsed time,  $\Delta t$ , of each configuration, giving  
382 us the estimated number of good questions per sec-  
383 ond that each configuration produced. We chose  
384 to rely on elapsed time rather than counting the  
385 floating-point operations, as all experiments were  
386 conducted on the same hardware, described in Ap-  
387 pendix B.3. Furthermore, elapsed time provides a  
388 more intuitive value for readers to comprehend.

$$389 \text{hits-per-sec}_k(\mathbb{D}_s) = \frac{\sum_{x \in \mathbb{D}_s} f_k(x; m)}{\Delta t}. \quad (7)$$

390 It’s worth noting that this preliminary bench-  
391 mark, while insightful, carries inherent subjectivity.  
392 This subjectivity stems from our defined metrics  
393 of quality, which rely on other retrieval models.  
394 Nevertheless, its primary aim remains exploratory,  
395 since benchmarking all the configuration directly  
396 on the downstream task would be time-consuming.  
397 Moreover, Section 4.2.2 details experiments gaug-  
398 ing our question quality assessment method’s effec-  
399 tiveness. These experiments offer further evidence  
400 of the reliability of this approach.

### 401 3.4.2 Downstream reranking benchmark

402 To obtain a more realistic assessment of the ex-  
403 pected quality of the generated synthetic dataset,  
404  $\mathbb{D}_s$ , we use it to train a BERT-base (Devlin et al.,  
405 2019) top-100 reranker model for each document

406 collection. Subsequently, we compare the perfor-  
407 mance of the trained model against the BM25 base-  
408 line and other state-of-the-art works. We evaluate  
409 the results in terms of NDCG@10 metric.

410 We adopt the standard BERT base checkpoint  
411 when training to keep the experiment simple and  
412 accessible. Furthermore, we also adopt a simple  
413 random negative sampling strategy for selecting  
414 negative documents for each question. We consider  
415 this setup reasonable given that our objective is not  
416 to achieve state-of-art results, but rather to show  
417 that it is possible to train neural reranker models  
418 in unlabelled collections with cheaply obtainable  
419 synthetic datasets.

## 420 4 Experiments and Results

421 This section outlines the performed experiments  
422 and their outcomes. We first introduce the doc-  
423 ument collections used for the benchmarks. Fol-  
424 lowing this, we present experiments that validate  
425 our assumptions: the use of information theory for  
426 outlier document elimination and the employment  
427 of retrieval models for question quality assessment.  
428 Lastly, we disclose the results of the benchmarks  
429 themselves.

### 430 4.1 Data

431 During our experiments, we considered five  
432 datasets, namely, BioASQ (Tsatsaronis et al.,  
433 2015), MSMARCO (Bajaj et al., 2016), NQ  
434 (Kwiatkowski et al., 2019), SciDocs (Cohan et al.,  
435 2020) and HotpotQA (Yang et al., 2018), that rep-  
436 resent various data domains. See Appendix B.1 for  
437 more information regarding the datasets and the  
438 selection criteria.

### 439 4.2 Validation experiments

440 We present now experiments that allow us to vali-  
441 date our framework for discovering document out-  
442 liers and our mechanism for assessing question  
443 quality based on retrieval models.

#### 444 4.2.1 Validating document outlier detection

445 Regarding document outlier detection, we follow  
446 the methodology presented in Section 3.1, in which  
447 we compute the normalized information (NI) mea-  
448 sure using a transformer language model (gpt-neo-  
449 125M (Gao et al., 2020)) and an FCM. To validate  
450 the effectiveness of this approach, we contrasted  
451 the NI distribution of documents in each collection  
452 against the distribution of the gold standard docu-  
453 ments, which comprises documents acknowledged

as relevant. This comparison can be visualized in Figure 5 for each dataset. The objective was to analyse the overlap of both distributions, where a complete overlap would imply that the documents in both extremities of the collection distribution are less likely to be relevant according to the gold-standard distribution.

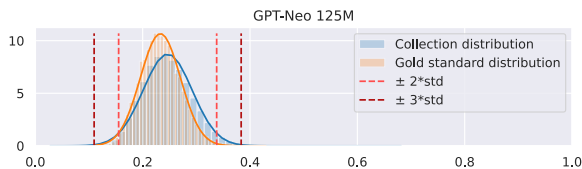


Figure 2: NI distribution of the BioASQ dataset using GPT-Neo 125M.

As an example, Figure 2 shows the distributions for the BioASQ dataset obtained with the gpt-neo-125M model. As observable, there is a clear overlap between the collection distribution and the gold standard distribution, meaning that removing documents at the extremities effectively eliminates potentially non-relevant documents. Based on this observation, we consider removing outliers that are at  $k$ -standard deviation away of the mean, denoted by the vertical lines on the Figure. Regarding the adopted language models, pretrained transformer LM is preferable due to their ability to produce better dataset distributions and the advantage of direct use, whereas FCMs require prior training. See Appendix C for a follow-up discussion regarding the remaining datasets and FCM model.

#### 4.2.2 Validating question quality method

To validate the efficacy of Equation 5 as a means of estimating the quality of questions, we propose to directly use the gold standard data of each dataset. By leveraging these already established question-document pairs, we examined how accurately Equation 5 identifies authentic questions for different values of the threshold  $k$ . Another way to interpret this experiment is to imagine that a language model synthetically generated the gold questions, and, therefore, we can estimate their quality because we have manually annotated data. Additionally, it is crucial to mine for strong negative questions, since the gold standard data typically only includes positive question-document pairs. To address this, we employ semantic search among the gold questions to identify questions with linguistic similarities but different positive document associations. We argue that these questions serve

as strong negative examples, as they share many common words while being distinct questions. We adopted SimCSE (Gao et al., 2021) to find semantic similar questions that do not share gold standard answer documents. See Appendix D for examples of negative questions.

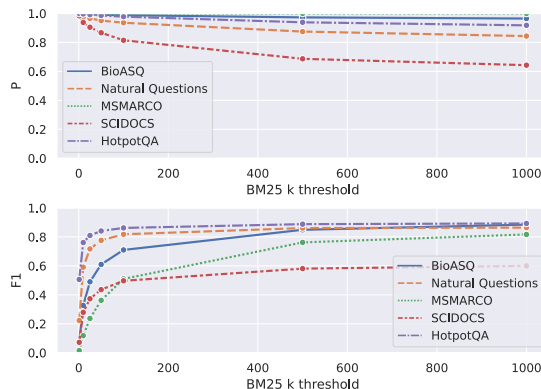


Figure 3: F1-score and precision ( $p$ ) values for varying threshold  $k$  with BM25 as our model.

Figure 3 depicts precision and F1-score values as functions of the threshold  $k$  when adopting BM25 as our model  $m$ . For the rest of the paper, we opt for BM25 due to its CPU efficiency and reusability for mining for negative documents in the downstream reranking benchmark. However, a comparison of BM25 and the monoT5 model for question quality estimation is presented in Appendix E. As observed in Figure 3, aside from the SciDocs dataset, the method can effectively distinguish correct questions from incorrect ones for thresholds exceeding 100. Notably, this approach favours higher precision values, enhancing our confidence in this method for question quality assessment.

#### 4.3 Benchmarking experiments

Here, we present two performed benchmarks: the first concerns a comprehensive analysis targeting all configurations for question generation, and the second assesses the best configurations within a reranking scenario where the synthetic questions are used as training data.

##### 4.3.1 Question quality benchmark

As previously mentioned, we adopted the hitsR and hits-per-sec as the main metrics to order our benchmark. We mainly adopted well-known publicly available small language models that range from 70M to 1.3B parameters, namely pythia-70M/160M/410M (Biderman et al., 2023), gpt-neo-350M/1.3B (Gao et al., 2020),

opt-125M/350M/1.3B (Zhang et al., 2022) and bloom(z)-560M (Muennighoff et al., 2022) totalling 10 models from 4 families. We selected 16K representative documents from each dataset, according to Section 4.2.1, and generated 5 questions for each document, conditioned on the starting words What, How, Where, Is, Why, totalling 80K expected questions from each model. Additionally, we also studied the impact of the generation method by considering three different strategies, Random Sampling (RS), Contrastive Search (CS)<sup>1</sup> and Beam Search (BS)<sup>2</sup>.

Figure 4 represents a parallel plot for all the 150 benchmarked runs that summarizes the impact of each model and generation strategy, see Appendix F for a comparison between datasets. Regarding the hits-per-sec measurement, it is clear that, independently of the model, the RS strategy largely outperforms the other generation methods, being almost 5x more efficient on average than BS and almost 6x than CS. On the other hand, when looking at hitsR, with  $k = 100$ , the best-performing generation strategy was BS reaching an average ratio of 0.68, against 0.48 and 0.47 for RS and CS, respectively. Another interesting observation is that, for all strategies, the amount of good synthetic questions seems to increase with model size, except for the opt family, where the results were similar independently of model size. The results regarding the CS strategy were surprising, since we expected them to be on par with BS. However, this could be related to less optimal hyperparameters.

### 4.3.2 Downstream reranking benchmark

Following the results obtained in the previous section, we proceeded to evaluate the synthetic datasets produced by gpt-neo-1.3B with BS and pythia-70m with RS in a downstream retrieval task, see Appendix G for additional combinations and further discussion. We believe that these two combinations cover the spectrum of configurations tested, namely, gpt-neo-1.3B with BS was the best configuration in terms of hitsR but one of the worst at hits-per-sec, while pythia-70m with RS showed the opposite behaviour.

Table 1 summarizes the results and compares them with relevant approaches from the literature<sup>3</sup>. Our approach consistently improves over the BM25 baseline, supporting our main hypoth-

<sup>1</sup>We choose topK of 4 and topP of 0.6.

<sup>2</sup>We adopted a beam-width of 5.

<sup>3</sup>Results for BM25+monoT5 were obtained by us.

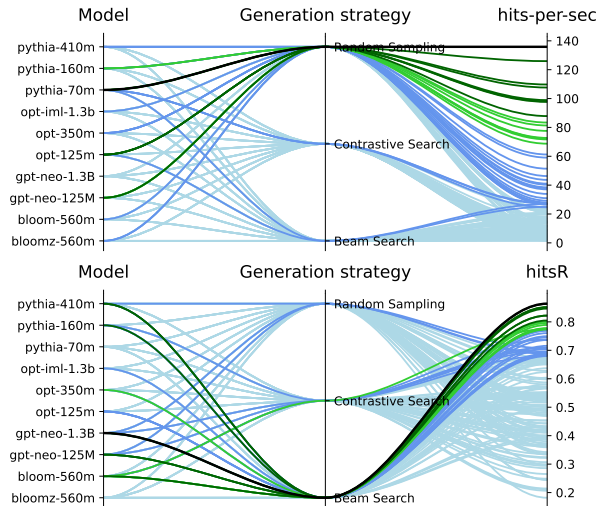


Figure 4: Parallel plot of benchmarked run impacts. Colors: black (best), dark green (top 5%), green (top 10%), blue (top 25%), light blue (rest).

Table 1: IR downstream task results.

Models	BioASQ nDCG@10	MSMARCO nDCG@10	NQ nDCG@10	HotpotQA nDCG@10	SciDocs nDCG@10
<b>Baseline (Unsupervised)</b>					
BM25	0.353	0.230	0.281	0.585	0.157
<b>Ours: BM25+BERT-base trained with following synthetic dataset</b>					
BS gpt-neo-1.3B	0.436	<b>0.336</b>	0.416	0.681	<b>0.228</b>
RS pythia-70m	0.438	0.307	0.407	0.730	0.187
<b>Supervised on synthetic data</b>					
GenQ (TAS-B)	-	-	0.358	0.534	0.143
InPars (220M)	-	0.259	0.335	-	-
InPars (3B)	-	0.297	0.513	-	-
<b>Supervised on MSMARCO</b>					
ANCE	-	-	0.446	0.456	0.122
<b>BM25 + Supervised on MSMARCO</b>					
BM25+MiniLM	-	-	0.533	0.707	0.166
BM25+monoT5	<b>0.444</b>	-	<b>0.639</b>	<b>0.7645</b>	0.183

esis that cheaply generated datasets can be used to train neural retrieval models. Remarkably, even when compared to InPars, which uses GPT-3 for synthetic generation, we achieved better results when considering a similarly sized reranker model (monoT5 220M vs. BERT-base 110M). Additionally, we achieved better results than the GenQ method, which employs a trained T5 model for synthetic generation and TAS-B as dense retrieval model. Lastly, we compared our approach to out-of-domain reranker models trained on MSMARCO, achieving competitive results. Importantly, these competitive results were obtained without extensively optimizing the training of our models and expensive architectures. Concretely, we trained the vanilla BERT-base checkpoint on the synthetic dataset using the huggingface trainer with default hyperparameters.

As a final discussion, we believe this work complements the findings of InPairs (Bonifacio et al., 2022), where they demonstrate that larger models produce better synthetic dataset. However, in this work, we show that by applying a robust question quality filter, smaller and more efficient models can be harnessed to generate synthetic datasets that rival the ones produced by larger models.

## 5 Ablation studies

In this section, we present an ablation study designed to understand the impact of each proposed method on the overall pipeline.

### 5.1 Document outliers

Central to our approach for document outlier detection is the assumption that documents located at the tails of the distribution of NI values in a collection may not be truly representative. To validate this, we conducted the experiment outlined in Table 2. Here, we deliberately generate questions for documents possessing the highest and lowest NI values across each collection. Subsequently, we computed HitsR ( $k = 100$ ) for these documents and compare it against our synthetic datasets that avoid such documents.

Table 2: Comparison of HitsR for questions from extreme NI documents vs the synthetic dataset (Synth DS).

Models	BioASQ HitsR	MSMARCO HitsR	NQ HitsR	HotpotQA HitsR	SciDocs HitsR
<b>Gpt-neo-1.3B BS</b>					
Lowest NI	0.625	0.371	0.535	0.838	0.879
Highest NI	0.568	0.447	0.343	0.718	0.845
Synth DS	<b>0.894</b>	<b>0.714</b>	<b>0.880</b>	<b>0.881</b>	<b>0.905</b>
<b>Pythia-70m RS</b>					
Lowest NI	0.358	0.101	0.034	<b>0.285</b>	<b>0.707</b>
Highest NI	0.058	0.064	0.027	0.120	0.439
Synth DS	<b>0.391</b>	<b>0.196</b>	<b>0.672</b>	0.267	0.641

The table clearly shows that the synthetic dataset (Synth DS) consistently achieves a higher HitsR than questions from both the lowest and highest NI documents. This disparity is pronounced in larger collections like BioASQ, MSMARCO, and NQ, which are more affected by irregular documents. Notably, for HotpotQA and SciDocs, the models yielded comparable rate of good questions for lower NI documents and the synthetic dataset, suggesting a cleaner dataset for these collections. Moreover, it is also observable that the models find it more challenging to generate useful questions

from documents with elevated NI values than those with lower NIs.

### 5.2 Question quality

Lastly, as a form to understand the impact of our question quality filtering, we trained the reranker model in two additional scenarios: using only the rejected questions (Only rejected) and without any filtering (All questions). The performance is then compared against the previously trained model (Only accepted).

Table 3: Comparison of reranker models across question subsets.

Questions	BioASQ nDCG@10	MSMARCO nDCG@10	NQ nDCG@10	HotpotQA nDCG@10	SciDocs nDCG@10
<b>Gpt-neo-1.3B BS</b>					
Only rejected	0.331	0.277	0.358	0.612	0.154
Only accepted	<b>0.436</b>	0.336	<b>0.416</b>	<b>0.681</b>	<b>0.228</b>
All questions	0.433	<b>0.340</b>	0.381	0.658	0.176
<b>Pythia-70m RS</b>					
Only rejected	0.105	0.223	0.313	0.237	0.160
Only accepted	<b>0.438</b>	<b>0.307</b>	<b>0.407</b>	<b>0.730</b>	<b>0.187</b>
All questions	0.373	0.276	0.406	0.507	0.185

In summary, Table 3 shows the importance of our question quality filtering mechanism. This approach not only contributes to a better performance of the reranker model, but this is also achieved more cheaply by avoiding the noise and inconsistencies present in the rejected questions. In other words, the overall positive differences in performance between ‘Only accepted’ and ‘All questions’ shows that the filtering mechanism was capable of removing questions that did not contribute to the overall results, at the same time improving performance and accelerating the training.

## 6 Conclusion and Future work

This work demonstrated that smaller language models can efficiently generate high-quality synthetic datasets for neural retrieval model training. Our approach shows that utilizing information theory principles for document selection and a small language model for zero-shot question generation can outperform methods like BM25 and pretrained monoT5 in certain scenarios.

Future work could focus on refining the downstream benchmark by also leveraging dense retrieval models and adopting stronger reranker models. Our findings bring us closer to broader neural retrieval model integration, mitigating data labelling and computational resource challenges.



## 670 **Limitations**

671 Although our study shows meaningful progress to-  
672 wards efficient synthetic dataset creation for neural  
673 retrieval models, it presents some limitations that  
674 should be considered for completeness and to guide  
675 future research directions.

676 Firstly, our method has not been applied to dense  
677 retrieval models. Owing to the substantial compu-  
678 tational resources required for encoding the col-  
679 lections, the decision was made to exclude dense  
680 retrieval from the scope of our research. Evaluating  
681 the performance on downstream tasks with dense  
682 retrieval models could further bridge the gap in the  
683 direction of adopting neural retrieval models as the  
684 default solution for information retrieval.

685 Secondly, we have not pursued the path of care-  
686 fully optimizing every hyperparameter for metric  
687 maximization, therefore, the presented results are  
688 obtained with default parameters. For instance, we  
689 did not fine-tune the BM25 component of our sys-  
690 tem. While BM25 serves as a key baseline in our  
691 evaluations, performance may be further optimized  
692 through additional fine-tuning. Additionally, we  
693 also did not fine-tune the prompt for question gen-  
694 eration. The design of prompts is a crucial aspect  
695 in many language model tasks, potentially influenc-  
696 ing the quality of generated questions. Therefore,  
697 our method’s effectiveness could depend on the  
698 prompt’s quality.

699 Thirdly, we have not explored the applicability  
700 of our approach within a Doc2Query-like scenario.  
701 In contrast to our goal of creating synthetic datasets,  
702 Doc2Query generates questions from a document  
703 and appends them to aid index-based retrieval mod-  
704 els like BM25.

705 Lastly, despite using small language models, the  
706 current setup may still require the usage of a GPU  
707 with at least 8GB of VRAM. This might also affect  
708 the scalability to longer texts, as the computational  
709 burden will increase with the length of the text.

## 710 **Ethics Statement**

711 This study presents a methodology to efficiently  
712 generate synthetic datasets for training neural re-  
713 trieval models, particularly beneficial for document  
714 collections lacking annotated data. Its broader im-  
715 pact lies in enabling effective neural information re-  
716 trieval adoption in retrieval scenarios that lack label  
717 data. It is essential to acknowledge the possibility  
718 of the model to generate inappropriate or harm-  
719 ful questions, leading to harmful retrieval training

720 data that can be learnt by models. To mitigate this  
721 problem, we used a filtering mechanism to ensure  
722 question quality. However, it is still important to be  
723 aware of the propagation of harmful information.  
724 Furthermore, we aimed to contribute to sustainable  
725 AI practices using small language models requiring  
726 fewer computational resources. Towards that goal,  
727 we will release a code repository for zero-shot syn-  
728 thetic question generation, promoting transparency  
729 and reproducibility. While we have strived to ad-  
730 dress the ethical implications, users should conduct  
731 a specific risk assessment based on their use-case  
732 scenarios to minimize potential harm and enhance  
733 filtering mechanisms if needed.

## 734 **References**

- 735 Chris Alberti, Daniel Andor, Emily Pitler, Jacob Devlin,  
736 and Michael Collins. 2019. *Synthetic QA corpora  
737 generation with roundtrip consistency*. In *Proceed-  
738 ings of the 57th Annual Meeting of the Association for  
739 Computational Linguistics*, pages 6168–6173, Flo-  
740 rence, Italy. Association for Computational Linguis-  
741 tics.
- 742 Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng,  
743 Jianfeng Gao, Xiaodong Liu, Rangan Majumder,  
744 Andrew McNamara, Bhaskar Mitra, Tri Nguyen,  
745 et al. 2016. Ms marco: A human generated ma-  
746 chine reading comprehension dataset. *arXiv preprint  
747 arXiv:1611.09268*.
- 748 Stella Biderman, Hailey Schoelkopf, Quentin Anthony,  
749 Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mo-  
750 hammad Aflah Khan, Shivanshu Purohit, USVSN Sai  
751 Prashanth, Edward Raff, Aviya Skowron, Lintang  
752 Sutawika, and Oskar van der Wal. 2023. *Pythia:  
753 A suite for analyzing large language models across  
754 training and scaling*.
- 755 Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, and  
756 Rodrigo Nogueira. 2022. Inpars: Data augmentation  
757 for information retrieval using large language models.  
758 *arXiv preprint arXiv:2202.05144*.
- 759 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
760 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
761 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
762 Askell, et al. 2020a. Language models are few-shot  
763 learners. *Advances in neural information processing  
764 systems*, 33:1877–1901.
- 765 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
766 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
767 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
768 Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
769 Gretchen Krueger, Tom Henighan, Rewon Child,  
770 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,  
771 Clemens Winter, Christopher Hesse, Mark Chen,  
772 Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin







probability of the next outcome depends on a finite number of recent past outcomes, known as the context (Pinho et al., 2010). One difference to the previous transformer-based LM is that we need to estimate the parameters for the FCM.

The primary benefit of Finite Context Models (FCM) lies in their capability to consider the whole document collection when estimating probabilities for individual documents, as the parameters of the FCM are derived from a comprehensive traversal of the entire collection. However, for either small or excessively diverse collections, FCMs might yield sub-optimal probability estimates.

The process of building an FCM model consists in iterating through the target collection and building a co-occurrence table,  $MT$ , between the current token,  $w_i$ , and the previous  $k$ -tokens, denoted as  $c = \{w_{i-1-k}, \dots, w_{i-1}\}$  (context). The probability estimation is given by Equation 8, where Laplace smoothing,  $\alpha$ , assigns small probability values to unseen co-occurrences. In  $MT$ , the rows correspond to the context tokens  $c$ , while the columns are associated with the current token  $w_i$ . Each entry within the  $MT$  specifies the frequency of instances where the context  $c$  is succeeded by the token  $w_i$ .

$$P(w_i|k) = \frac{MT(k, w_i) + \alpha}{\sum_{j=1}^{|V|} MT(k, w_j) + \alpha|V|}. \quad (8)$$

## B Experimental details

### B.1 Dataset details

Regarding the dataset selection, we mainly rely on the pool of datasets offer by BEIR (Thakur et al., 2021) benchmark. Then, to build our pool of datasets, we decided to only include datasets used in the evaluation of models that retrieve information to answer questions. Furthermore, we would also like to have varied datasets in terms of domain and number of documents.

Several datasets were excluded based on these criteria. For instance, Quora and CQADupStack, centred around retrieving similar questions, which did not fit our purpose. The Robust dataset, although important, dates back to 2004 and its questions are not framed in natural language. Practical constraints, like time and computational resources, also limited our choices.

Ultimately, we settled on five datasets: BioASQ, MSMARCO, NQ, HotPotQA, and Scidocs. It’s

worth noting that while BEIR offers a version of the BioASQ dataset, we opted for the official 2022 BioASQ dataset. This comprehensive version comprises 33M documents (tripling the BEIR variant) and includes 38k question-document pairs. Below is a more detailed breakdown:

- BioASQ: An annual challenge focused on biomedical document retrieval and question answering. We make use of the dataset from the 10th edition of the BioASQ, which contains 38,933 question-document pairs and uses the 33 million document 2022 PubMed baseline as the document collection (Tsatsaronis et al., 2015).
- MSMARCO: A well-known dataset for benchmarking deep learning neural reranking models in open-domain scenarios. It includes 4,102 question-document pairs and a document collection of over 8 million documents (Bajaj et al., 2016).
- NQ (Natural Questions): An open-domain dataset aimed at benchmarking question answering systems. It consists of 4,201 question-document pairs and a document collection of over 2 million documents (Kwiatkowski et al., 2019).
- Scidocs: A dataset primarily focused on scientific documents. It contains 4,928 question-document pairs, with a document collection of approximately 25,000 documents (Cohan et al., 2020).
- Hotpotqa: A challenging question answering dataset designed to test models capabilities for multi-hop reasoning and answering complex questions. It includes 14,810 question-document pairs, with a document collection of over 5 million documents (Yang et al., 2018).

### B.2 Software

Here we present the main packages used during the development of our work. For BM25 we adopted pyterrier (Macdonald and Tonellotto, 2020), a python wrapper of the Terrier (Macdonald et al., 2012) search engine. Regarding the training, inference and generation with neural models, we mainly rely on HuggingFace package (Wolf et al., 2020). More precisely, the BERT-base model that we trained corresponds to the “bert-base-uncased” checkpoint, while for monoT5 we used



1085 the “castorini/monot5-base-msmarco-10k” check- 1134  
1086 point. Regarding the generative models, we also 1135  
1087 used the checkpoints that were publicly available 1136  
1088 on the HuggingFace hub. 1137

### 1089 **B.3 Hardware** 1138

1090 All of our experiments run on the following 1139  
1091 desktop, Intel(R) Core(TM) i9-9900K CPU @ 1140  
1092 3.60GHz, 2x NVIDIA GeForce RTX 2070 8GB 1141  
1093 VRAM and 32GB of RAM. Although the machine 1142  
1094 is equipped with two RTX 2070, during our exper- 1143  
1095 iments we did not take advantage of a multiGPU 1144  
1096 setup. Therefore, all the experiments presented in 1145  
1097 this paper would run on a single GPU. For pro- 1146  
1098 ducing the results for both ablation studies, we 1147  
1099 relied on a DGX A100 system to streamline the 1148  
1100 experiences in parallel. However, the code and the 1149  
1101 parameters were the same as the ones used in our 1150  
1102 previous machine to keep the experiments compa- 1151  
1103 rable. 1152

### 1104 **C Document outlier detection for each** 1153 1105 **dataset** 1154

1106 Figure 5, similarly to Figure 2, shows the distribu- 1155  
1107 tion of NI values for each individual dataset. More 1156  
1108 precisely, each row corresponds to a dataset, the 1157  
1109 left column panels correspond to the NI estima- 1158  
1110 tive produced by the gpt-neo-125M model, and 1159  
1111 right column panels correspond to the NI estima- 1160  
1112 tive from the FCM model. 1161

1113 Starting by analysing the distributions produced 1162  
1114 by the gpt-neo-125M model, it is evident that each 1163  
1115 dataset exhibits a bell-shaped distribution with a 1164  
1116 high degree of alignment compared to the gold stan- 1165  
1117 dard distribution. Notably, the NQ dataset shows 1166  
1118 the most significant deviation in terms of an align- 1167  
1119 ment. Inclusively, it is observable that the gold 1168  
1120 standard data tends to favour lower NI values com- 1169  
1121 pared to the dataset distribution. This may be in- 1170  
1122 dicative that the documents in the gold standard 1171  
1123 are potentially more easily discoverable than the 1172  
1124 average ones from the entire collection. However, 1173  
1125 more experiments would be required to examine 1174  
1126 this. 1175

1127 Moving on to the FCM, it produced distributions 1176  
1128 that deviate slightly from a bell curve, specially, in 1177  
1129 the case of the MSMARCO dataset. We attribute 1178  
1130 this deviation to the dataset’s high diversity, which 1179  
1131 encompasses multiple sources from different do- 1180  
1132 mains, making it challenging to obtain accurate 1181  
1133 estimates when building the FCM. 1182

1134 Nevertheless, the alignment between dataset dis- 1135  
1136 tribution and the gold standard distribution is still 1137  
1138 present. This further supports the notion that we 1139  
1139 can exclude the trailing documents from the distri- 1140  
1140 bution, as they are less likely to be considered as 1141  
1141 gold documents. 1142

### 1140 **D Similarity between questions for** 1140 1141 **negative mining** 1141

1142 Table 4 show some examples of different gold stan- 1143  
1143 dard questions that are similar but do not share any 1144  
1144 positive document. As previously described, the 1145  
1145 fundamental assumption is that the set of positively 1146  
1146 labeled gold standard documents for one question 1147  
1147 should serve as a robust set of negatively labeled 1148  
1148 documents for a similar question. To illustrate, let 1149  
1149 us consider the first example in Table 4 from the 1150  
1150 NQ dataset. We can observe that both questions per- 1151  
1151 tain to movies from the Planet of the Apes trilogy, 1152  
1152 where the question on the left relates to the 2017 1153  
1153 film, while the question on the right pertains to the 1154  
1154 2011 film. Consequently, the positive documents 1155  
1155 for the first question should be regarded as strong 1156  
1156 negative documents for the second question, and 1157  
1157 vice versa, given that both documents address the 1158  
1158 same topic but do not contain the correct answer. 1159

1159 Moreover, it becomes evident that this negative 1160  
1160 mining technique is most effective when applied 1161  
1161 to a gold standard with a deep set of relevance 1162  
1162 per question, If the gold standard has a shallow 1163  
1163 set of relevance the probability of finding similar 1164  
1164 questions that share positive documents which are 1165  
1165 not annotated in the dataset would be too high. 1166  
1166 Lastly, due to the limited number of questions in 1167  
1167 the gold set for MSMARCO (only 43 questions), 1168  
1168 we were unable to mine strong negatives, as the 1169  
1169 number of questions was insufficient to find any 1170  
1170 match. 1171

### 1171 **E Comparison between BM25 and** 1171 1172 **monoT5 for estimating question quality** 1172

1173 Firstly, it is important to make a distinction in terms 1174  
1174 of both models. More precisely, BM25 is a retrieval 1175  
1175 model that provides a ranked order of documents 1176  
1176 for each question, while monoT5 predicts the rele- 1177  
1177 vance between question-document pairs. Therefore, 1178  
1178 based on our definition of question quality, BM25 1179  
1179 appears to be the more suitable model. It directly 1180  
1180 encodes the notion of retrieval, while monoT5 is 1181  
1181 trained solely to differentiate between relevant and 1182  
1182 irrelevant question-document pairs. For instance,

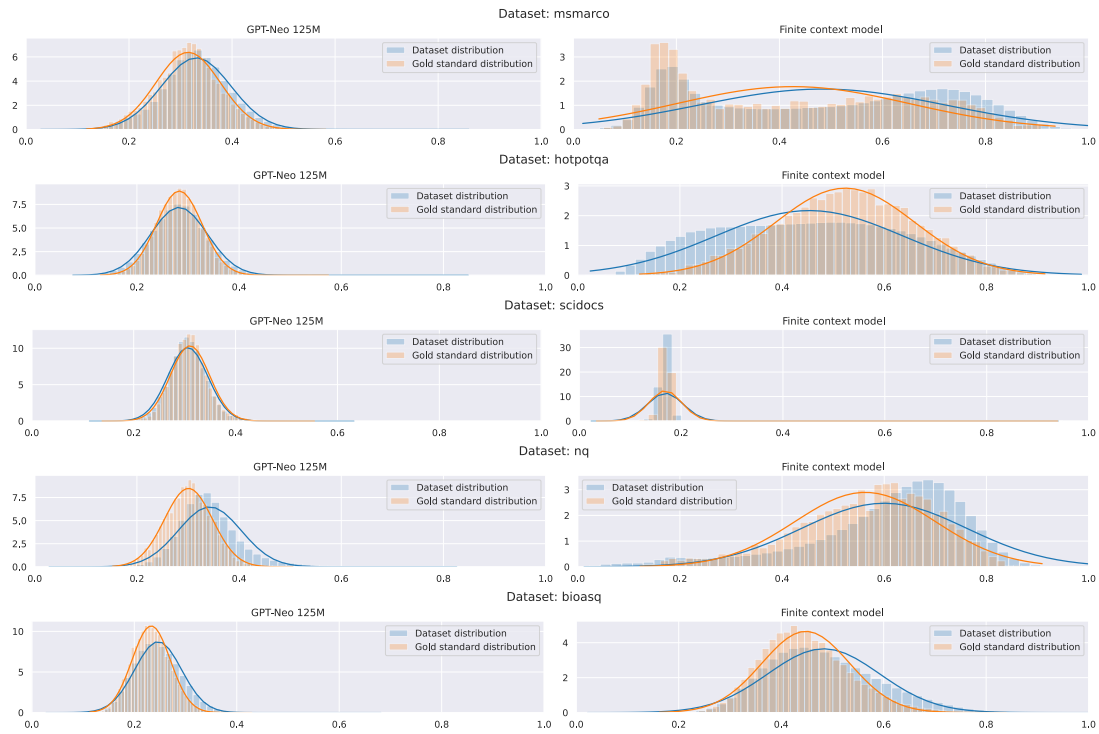


Figure 5: NI distribution for every dataset using the GPT-Neo 125M and a finite context model.

Table 4: Examples of different gold questions that are similar from the NQ, BioASQ, HotpotQA and SciDocs datasets.

Gold question	Similar gold question	SimCSE
<b>Natural Questions (NQ) dataset</b>		
where was the war of the planet of apes filmed	where was the rise of the planet of the apes filmed	0.905
when did world war 2 end in the pacific	who did us fight in world war 1	0.703
<b>BioASQ dataset</b>		
What is the mechanism of action of Fremanezumab?	What is mechanism of action of Benralizumab?	0.930
Which mutations of alpha-myosin heavy chain gene are implicated in hypertrophic cardiomyopathy?	which mutations of phospholamban gene have been found to cause hypertrophic cardiomyopathy?	0.910
<b>HotpotQA dataset</b>		
Which genus has more species, Xanthoceras or Ehretia?	Which Genus has more species Eucryphia or Lepidozamia ?	0.924
Between Greyia and Calibanus, which genus contains more species?	Which has more species, Clianthus or Callicoma?	0.866
<b>SciDocs dataset</b>		
Wideband millimeter-wave SIW cavity backed patch antenna fed by substrate integrated coaxial line	Broadband millimetre-wave passive spatial combiner based on coaxial waveguide	0.845
Reinforcement Learning for Coreference Resolution	Deep Reinforcement Learning for Dialogue Generation	0.776
<b>MSMARCO dataset</b>		
types of dysarthria from cerebral palsy	causes of left ventricular hypertrophy	0.608
when was the salvation army founded	who formed the commonwealth of independent	0.562

let’s consider an article that is a literature review discussing information retrieval (IR), and the question is “What is the main subject of this literature review?”. Since monoT5 is a relevance model, it would likely predict this as relevant, violating the second criterion in our definition. Nonetheless, monoT5 is trained using retrieval data, which might compel the model to capture a weak notion of retrieval. Therefore, we decided to make a judgment analysis against the BM25.

Secondly, it is equally important to consider the computation complexity of both solutions, since we aim to benchmark multiple configuration and therefore a high-performing method is preferable. BM25 is a CPU-bounded algorithm that can be easily scalable by the number of available CPU(s), while monoT5 is a GPU-bounded algorithm, that can be also easily scalable by the number of GPU(s). In a general point of view, we consider BM25 as the method with the lower computation cost, given that CPU-time is more easily accessible than GPU-time.

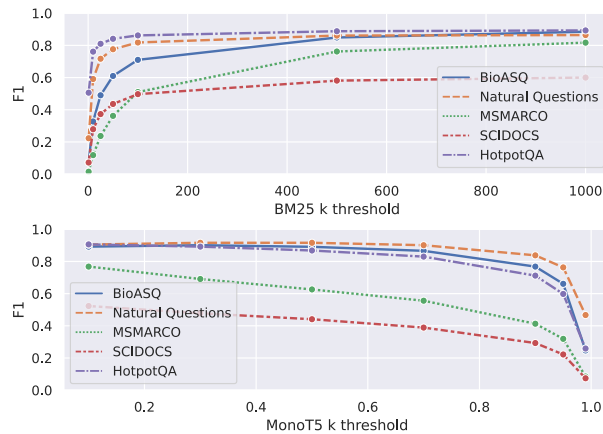


Figure 6: F1-score for varying threshold  $k$  for BM25 and monoT5.

Figure 6 and 7 present a comparison of both models, following the same methodology outlined in Section 4.2.2, in terms of F1 and precision, respectively, across varying thresholds. Overall, it appears that monoT5 performs comparably to BM25 for the different thresholds. However, considering the aforementioned points, we have decided to proceed with BM25 for the remainder of our experiments.

Furthermore, another advantage of BM25 is that when used as a quality filter, we also store all the retrieved documents during that process. This allows us to reuse these list of previously retrieved documents for subsequent negative document sampling

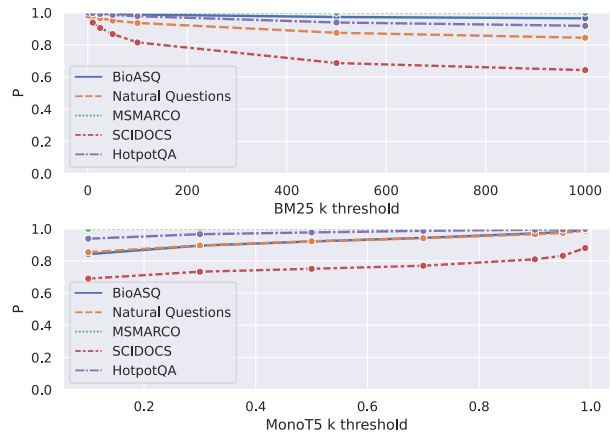


Figure 7: Precision ( $p$ ) for varying threshold  $k$  for BM25 and monoT5.

during the training of neural retrieval models.

## F QuestionQA quality benchmark per dataset

Figure 8, presents a more complete visualization of our benchmark metrics over each individual dataset. In general, the conclusions previously mentioned in Section 4.3.1 remain consistent. However, a more detailed analysis per dataset reveals that the models faced the most difficulty in generating questions for the MSMARCO dataset, as indicated by the relatively lower values of hitsR. The SciDocs dataset also posed challenges for the models. On the other hand, the dataset with the highest overall question generation success rate was BioASQ, meaning it was easier for the models to generate questions.

One possible explanation for this difference in performance may be the nature of the BioASQ dataset, which uses abstracts from biomedical scientific articles. These abstracts condense a large amount of diverse information, providing the models with a broader range of valid questions to generate.

Another interesting observation is that the difficulty in generating questions seems to be aligned with the average NI value of each dataset. For instance, recalling Figure 5, the dataset with the lowest average NI value was also the BioASQ dataset, while the MSMARCO was the dataset with the highest NI value. This suggests a possible relationship between the NI value and the difficulty of question generation by the models.

This relationship could be attributed to the model’s ability to comprehend the documents used as context for question generation, which should be

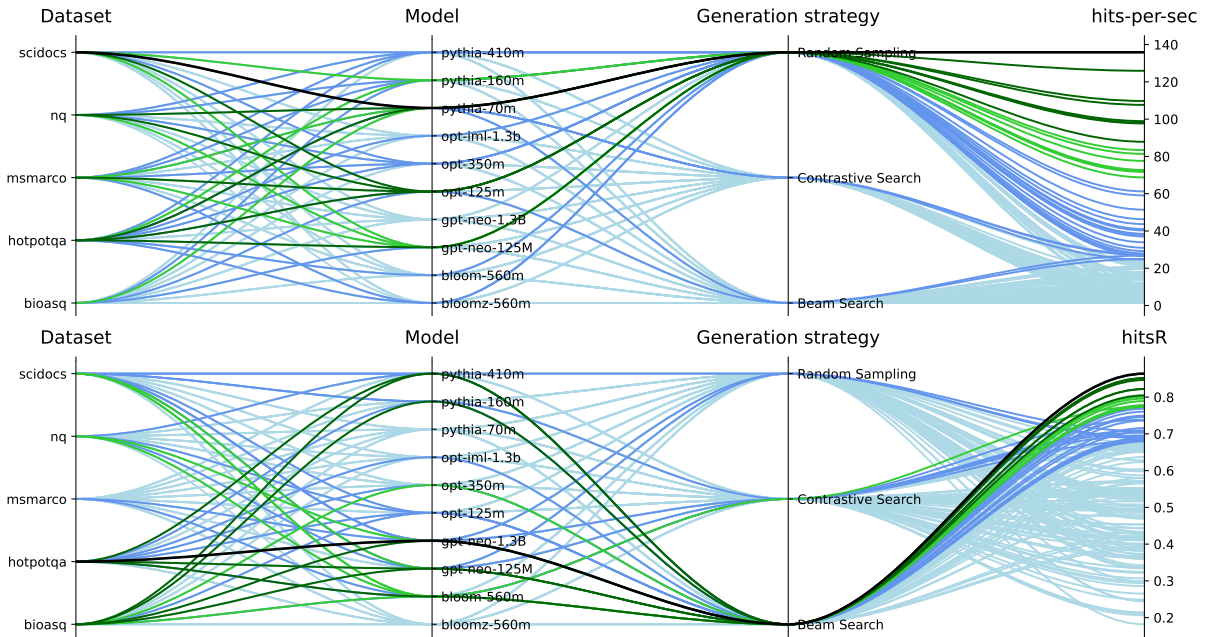


Figure 8: Parallel plot summarizing impacts across benchmarked runs. Color coding: black (best value), dark green (top 5%), green (top 10%), blue (top 25%), and light blue (remaining).

1253 captured by the NI measurement. In other words, 1281  
 1254 a lower NI value may be indicative that a document 1282  
 1255 is more easily interpreted by the language 1283  
 1256 model, because the language model itself was able 1284  
 1257 to produce better probability estimation for that 1285  
 1258 document. However, further experiments are nec- 1286  
 1259 essary to draw any definitive conclusions. 1287

## 1260 G Additional results on the downstream 1288 1261 IR task 1289

1262 Table 5 presents two additional results for the same 1290  
 1263 synthetic generative models, but with different 1291  
 1264 generation strategies, RS for gpt-neo-1.3B and BS for 1292  
 1265 pythia-70m. Upon comparing these strategies, it 1293  
 1266 appears that RS achieves slightly better results, ex- 1294  
 1267 cept for the SciDocs dataset. This unexpected out- 1295  
 1268 come raises an interesting point that the synthetic 1296  
 1269 dataset obtained with RS may exhibit better qual- 1297  
 1270 ity than that of BS. Initially, we believed that the 1298  
 1271 BS generation strategy would produce more co- 1299  
 1272 herent questions, therefore, resulting in a stronger 1300  
 1273 dataset. However, we hypothesize that this observa-  
 1274 tion could be explained by dataset diversity. When  
 1275 employing the BS strategy, the model generates  
 1276 5 questions for each document based on different  
 1277 starting words. Consequently, there is a higher like-  
 1278 lihood of generating semantically similar questions  
 1279 for different starting words. On the other hand,  
 1280 the stochastic nature of RS avoids such repetition.

To further investigate this, we propose analyzing  
 the diversity of each synthetic generated dataset.  
 Furthermore, we also believe that would be benefi-  
 cial to conducting a downstream evaluation under  
 a time budget constraint. By doing so, we may  
 gain additional insights into the performance of the  
 different methods, since when recalling Figure 4,  
 we observe significant variations in the number of  
 questions generated per second across the different  
 generation methods.



Table 5: IR downstream task results with both generation strategies for gpt-neo-1.3B and RS-pythia-70m.

Models	BioASQ nDCG@10	MSMARCO nDCG@10	NQ nDCG@10	HotpotQA nDCG@10	SciDocs nDCG@10
<b>Baseline (Unsupervised)</b>					
BM25	0.353	0.230	0.281	0.585	0.157
<b>Ours: BM25+BERT-base</b> trained with following synthetic dataset					
BS gpt-neo-1.3B	0.436	<b>0.336</b>	0.416	0.681	<b>0.228</b>
RS gpt-neo-1.3B	<b>0.451</b>	-	0.448	0.727	0.194
BS pythia-70m	0.418	-	0.379	0.691	0.181
RS pythia-70m	0.438	0.307	0.407	0.730	0.187
<b>Supervised on synthetic data</b>					
GenQ (TAS-B) <sup>a</sup>	-	-	0.358	0.534	0.143
InPars (220M) <sup>b</sup>	-	-	0.335	-	-
InPars (3B) <sup>b</sup>	-	-	0.513	-	-
<b>Supervised on MSMARCO</b>					
ANCE <sup>a</sup>	-	-	0.446	0.456	0.122
<b>Supervised on MSMARCO + BM25 Reranking</b>					
BM25+MiniLM <sup>a</sup>	-	-	0.533	0.707	0.166
BM25+monoT5 <sup>c</sup>	0.444	-	<b>0.639</b>	<b>0.7645</b>	0.183

<sup>a</sup> These results are from [Thakur et al., 2021](#)

<sup>b</sup> These results belong to [Bonifacio et al., 2022](#)

<sup>c</sup> This result was obtained by us.