# Complex Query Synthesis for Enhanced Information Retrieval

Daksha Ladia
University of Massachusetts Amherst
Amherst, MA, USA
dladia@umass.edu

Snigdha Ansu
University of Massachusetts Amherst
Amherst, MA, USA
sansu@umass.edu

Vasileios Vittis
University of Massachusetts Amherst
Amherst, MA, USA
vvittis@umass.edu

## 1  ABSTRACT

Effective information retrieval requires aligning user queries with document content at a granular level. This project introduces a novel approach that bridges the gap between user intent and document relevance using pseudo-query generation and large language models (LLMs). By chunking documents into passages to create pseudo-queries and transforming user queries into detailed, multi-faceted queries with LLMs, we compare them on a per-document basis to rank and retrieve the most relevant results. This method offers precise retrieval in contexts requiring a deep understanding of document content, like research, education, and legal domains, demonstrating the potential of combining pseudo-query generation with LLM-driven query enhancement.

## 2  RELATED WORK

Several studies have laid the foundation for understanding and advancing query generation and document interaction in information retrieval systems. The papers [1], [2], [3], [4], and [5] provide insights and methodologies that support our research direction and inform our approach.

Existing works, such as [2] by Sarah Jones et al. and [5] by Rodrigo Nogueira et al., explore methods to enhance document retrieval through query optimization. [5] introduces a technique that predicts and appends queries to documents, improving alignment with user search terms and tackling the vocabulary mismatch problem. Similarly, [2] discusses optimizing query phrasing to better match document content and improve retrieval accuracy. While these studies enrich documents to fit user queries or refine the queries themselves, our research uniquely focuses on generating queries directly from document content. This involves analyzing what information a document could answer and creating relevant, meaningful queries from that analysis. Both of these papers are very closely related with our proposed work.

Our project diverges from most existing methods, including those in [1] and [3], which emphasize shorter, simpler queries tailored for standard search behavior. Instead, we focus on generating and leveraging long, complex queries. These queries encapsulate richer context and detail, making them particularly well-suited for applications requiring comprehensive information retrieval, such as research, education, and legal domains.

Additionally, while approaches like those in [4] and [5] primarily target single-document alignment, our methodology explores document chunking and pseudo-query generation to improve relevance ranking. By comparing user-generated long queries (derived via LLMs) against pseudo queries synthesized from chunked document passages, we ensure a deeper alignment between user intent and document content. This approach uniquely enhances the retrieval pipeline by addressing intra-document context at a granular level, providing a robust framework for precise document ranking and retrieval.

## 3  PROBLEM STATEMENT

Current information retrieval systems are typically designed for straightforward query-document matching and often fail to capture the implicit context required for producing relevant documents. This is because they are ineffective in addressing the limitations of short, ambiguous queries. Generating long, context-rich queries can potentially offer a solution by enhancing retrieval accuracy and addressing the complexities inherent in vague user inputs.

**Research Question**: How can complex query synthesis techniques be optimized to infer implicit contexts from ambiguous or short user queries by synthesizing information across documents, and what insights do they provide into enhancing effectiveness of information retrieval systems ?

We propose leveraging large language models to optimize query synthesis by generating detailed, context-aware pseudo-queries for documents. These pseudo-queries act as semantic representations of document content, facilitating effective retrieval for complex queries. Our approach integrates document segmentation, generative modeling, diversity filtering, and semantic evaluation to ensure high-quality pseudo-queries. By aligning these pseudo-queries with users' potential intents, we aim to improve retrieval metrics like precision, recall, and relevance, enabling users to navigate complex information more effectively. This approach enhances the query generation process by considering how document content can inform and shape user query content, ultimately improving the retrieval system's ability to present relevant documents.

## 4  APPROACH

Below is a detailed explanation of the approach we implemented.

### 4.1  Pseudo-Query Generation for Document Representation

The goal of pseudo-query generation is to create a set of diverse, complex, and contextually rich queries for each document in the corpus. These pseudo-queries serve as a semantic representation of the document, facilitating accurate retrieval for complex user
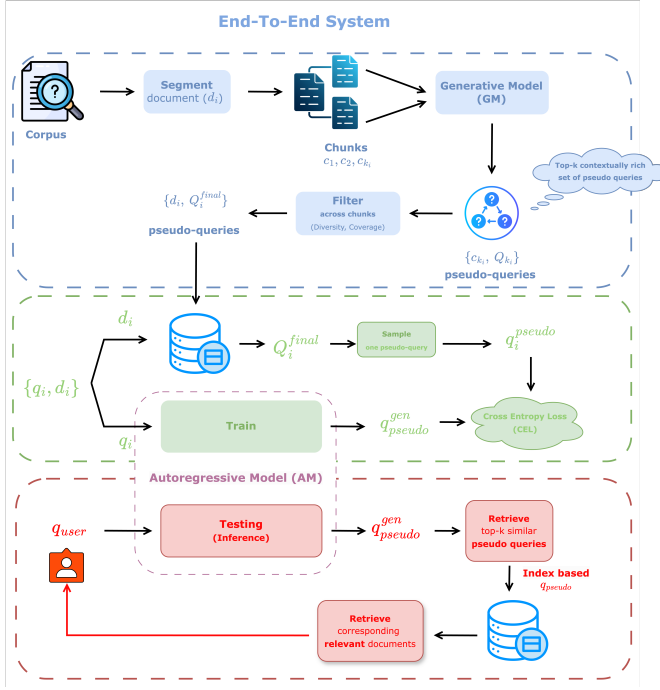
**Figure 1: Overview of the End-to-End System: The pipeline consists of three main components. The pseudo-query generation component (blue) segments documents into chunks and generates contextually rich pseudo-queries using a generative model. The training component (green) trains an autoregressive model on the generated pseudo-queries and corresponding documents. The inference and retrieval component (red) uses the trained model to generate queries and retrieve relevant documents.**

queries. Our approach combines document segmentation, generative modeling, diversity filtering, and semantic evaluation to ensure high-quality pseudo-queries.

*4.1.1 Document Segmentation.* Each document $d_i$ in the corpus $\mathcal{D} = \{d_1, d_2, \ldots, d_N\}$ is segmented into chunks for localized representation. Let $C_i = \{c_1, c_2, \ldots, c_{k_i}\}$ denote the set of chunks derived from $d_i$, where each chunk $c_j$ contains $w$ tokens:

$$c_j = \{t_j, t_{j+1}, \ldots, t_{j+w}\}, \quad j = 1, 2, \ldots, k_i.$$

This ensures localized context while maintaining continuity across the document. Each chunk $c_j$ is passed to the query generation model.

*4.1.2 Generative Modeling for Pseudo-Query Creation.* To generate pseudo-queries, we leverage a pretrained generative model FLAN-T5-Large. These models are chosen for their ability to generate semantically rich and contextually relevant questions without requiring additional fine-tuning. For a given chunk $c_j$, the model generates a set of top-$k$ pseudo-queries:

$$Q_j = \text{Top-}k(QG(c_j)).$$

In our approach, the relevance of a document to a user's question is determined by the $top - k$ sampling during the pseudo-question generation. This technique ensures that only the $top - k$ pseudo-queries, which represent the most semantically relevant aspects of the chunk. If a question does not belong to the $top - k$ means that it semantically not so relevant and therefore are left out.

Complexity in pseudo-queries is critical for representing nuanced information in documents and enabling effective retrieval for challenging user queries. To achieve this, we incorporate the following mechanisms into our pipeline:

*1. Prompt Engineering for Complexity.* The generative model is guided through carefully crafted prompts that emphasize depth and specificity. For example, instead of a generic instruction, we use:

> "Generate a detailed and nuanced question focusing on the most significant aspects of the passage."

This encourages the model to produce questions that go beyond surface-level content.

*2. Pretrained Model Selection.* We leveraged FLAN-T5-large, known for their ability to generate contextually rich and grammatically complex text. These models have been trained on large corpora containing intricate linguistic patterns, which naturally lead to high-complexity questions.

*4.1.3 Diversity Filtering and Global Aggregation.* After generating pseudo-queries for each chunk, we apply a semantic filtering step to remove redundant queries. Let $\text{Sim}(q_a, q_b)$ denote the semantic similarity between two queries $q_a$ and $q_b$, computed using cosine similarity over sentence embeddings. The final set of pseudo-queries for a chunk is:

$$Q_j^{\text{final}} = \{q \in Q_j : \text{Sim}(q, q') < \epsilon, \forall q' \in Q_j^{\text{final}}\}.$$

Here, $\epsilon$ is a threshold controlling redundancy. The final set of pseudo-queries for a document $d_i$ is the union of filtered pseudo-queries across all chunks:

$$Q_i = \bigcup_{j=1}^{k_i} Q_j^{\text{final}}.$$

*4.1.4 Evaluation Metrics for Pseudo-Query Quality.* To assess the quality of the generated pseudo-queries, we define the below metrics:

*Diversity and Redundancy.*

$$\text{Redundancy}(Q_i) = \sum_{q_a, q_b \in Q_i} \text{Sim}(q_a, q_b).$$

These metrics ensure pseudo-queries are diverse, complex, and minimally redundant.

## 4.2 Training an Autoregressive Language Model

The autoregressive language model is trained to map queries to pseudo-queries associated from the relevant passage.

Let $Q_{\text{pseudo}} = \{q_1, q_2, \ldots, q_m\}$ represent the set of pseudo-queries generated for a passage $p$ during the pseudo-query generation step. During the training phase, given an input query $q$, the model generates a pseudo-query $q_{\text{pseudo}}^{\text{gen}}$, aiming to approximate a sampled pseudo-query $q_{\text{pseudo}} \in Q_{\text{pseudo}}$.

Formally, the mapping is defined as:

$$q_{\text{pseudo}}^{\text{gen}} = G(q|\theta),$$

where $G$ is the autoregressive model parameterized by $\theta$, generating $q_{\text{pseudo}}^{\text{gen}}$ token by token.

To train the model, we minimize the cross-entropy loss between the generated pseudo-query and a sampled pseudo-query from $Q_{\text{pseudo}}$:

$$\mathcal{L} = -\log P(q_{\text{pseudo}}|q_{\text{pseudo}}^{\text{gen}}, \theta),$$

where $P(q_{\text{pseudo}}|q_{\text{pseudo}}^{\text{gen}}, \theta)$ represents the likelihood of generating $q_{\text{pseudo}}$ given training query $q$.

### 4.3 Inference: Retrieving Relevant Documents

The inference process retrieves the top-$k$ relevant passages for a given user query, by comparing a generated pseudo-query $q_{\text{pseudo}}^{\text{gen}}$ with precomputed pseudo-queries for each passage.

*Embedding Generation.* For the fast retrieval of the $q_{pseudo}^{gen}$ and the comparison with all the pseudo queries $q$ from all documents, we encode them into dense embeddings using BERT:

$$v_{\text{pseudo}}^{\text{gen}} = \text{BERT}(q_{\text{pseudo}}^{\text{gen}}), \quad v_{\text{pseudo}}^{j} = \text{BERT}(q_j).$$

*Similarity Calculation.* Similarity between embeddings is computed using cosine similarity:

$$\text{Sim}(v_{\text{pseudo}}^{\text{gen}}, v_{\text{pseudo}}^{j}) = \frac{v_{\text{pseudo}}^{\text{gen}} \cdot v_{\text{pseudo}}^{j}}{\|v_{\text{pseudo}}^{\text{gen}}\|\|v_{\text{pseudo}}^{j}\|}.$$

*Ranking and Document Selection.* For each passage $p_i$, the relevance score is the maximum similarity across its pseudo-queries:

$$\text{Relevance}(q_{\text{pseudo}}^{\text{gen}}, Q_{\text{pseudo}}^{i}) = \max_{q \in Q_{\text{pseudo}}^{i}} \text{Sim}(v_{\text{pseudo}}^{\text{gen}}, v_{\text{pseudo}}).$$

Passages are ranked by relevance scores, and the top-$k$ passages are selected:

$$P_{\text{top-}k} = \{p_{i_1}, p_{i_2}, \ldots, p_{i_k}\}.$$

## 5 EXPERIMENTS

### 5.1 Datasets

In our experiments, we use two benchmark datasets: **NFCorpus** and **SciFact**. These datasets provide diverse and challenging benchmarks to evaluate the retrieval performance of our approach and the baselines.

**NFCorpus** NFCorpus [1] is a dataset specifically designed for information retrieval. It focuses on retrieving passages that are relevant to a given query, with each query paired with relevance-labeled passages.

We evaluated over the BEIR NFCorpus dataset which is a full-text English retrieval data set for Medical Information Retrieval. We used the dataset from Hugging Face which contains around 400 natural language queries with 11.3K automatically extracted relevance judgments [validation set] for 3.63K medical documents (written in a complex terminology-heavy language).

[1] https://huggingface.co/datasets/nfcorpus

**SciFact** is a dataset tailored for scientific document retrieval. It is primarily used for retrieving abstracts that are relevant to a given scientific claim, with each claim paired with abstracts labeled as supporting, refuting, or neutral.

We used BEIR Scifact dataset for evaluation, which is a domain-specific collection for scientific retrieval. The dataset, available on Hugging Face, contains approximately 1.4K natural language claims with 5.2K relevance judgments [validation set] across 5K scientific abstracts, written in a specialized language common to biomedical and other scientific research fields.

### 5.2 Evaluation Metrics

The following metrics are used to assess retrieval performance:
*Precision:* Precision measures the proportion of retrieved documents that are relevant and evaluates the accuracy of the retrieved results.
*Recall:* Recall quantifies the proportion of relevant documents that are successfully retrieved and emphasizes the completeness of retrieval.
*Normalized Discounted Cumulative Gain (NDCG):* NDCG measures the ranking quality by considering the relevance of documents and their positions in the ranked list

### 5.3 Experimental Setup

*5.3.1 Chunking and Pseudo Query Generation : Part 1(abbreviated as CPQG).* Document preprocessing involves chunking, where documents are split into segments of 500 tokens. The pipeline employs a generative model, specifically the FLAN-T5 (google/flan-t5-large), which creates pseudo-queries by analyzing segmented chunks of documents. The prompt used for psuedo query generation from each chunk was : *Generate a detailed and nuanced question focusing on the most significant aspects of the passage.* The parameters used for generation were max_length=20 to limit queries to 20 tokens, num_return_sequences =5 to generate five queries per chunk, do_sample=True for diverse query generation, top_k=5 for most probable tokens. Tokenization and decoding for this process are handled using the Hugging Face transformers library.

Once psuedo queries are generated per chunk, redundant queries are filtered using cosine similarity on embeddings from the "all-mpnet-base-v2" SentenceTransformer model. Queries with a similarity score above a predefined threshold (0.8) are considered redundant, and only one representative query is retained to preserve diversity. For example, similar queries like "What are the effects of global temperature rise on weather patterns?" and "What are the consequences of global warming on weather patterns?" result in only one being kept, reducing redundancy while maintaining relevance.. To ensure diversity, a semantic similarity threshold of 0.8 is applied to filter out redundant queries.

*5.3.2 Fine Tuning of Models : Part 2.* In this project, we conducted experiments using three distinct models: GPT-2, T5-small and Cross-Encoder on the NFCorpus and SciFact datasets. The primary goal was to explore pseudo-question generation, document retrieval, and ranking tasks, while balancing computational and financial constraints. Additionally, we explored GPT 4o-mini pre-trained model for generating psuedo queries to evaluate the retrieval performance.

**Table 1: Retrieval performance on the datasets. We use precision @3, @5, @10, recall@3, @5, @10, and NDCG@3, @5, @10 and to evaluate the retrieval performance over top-10 retrieved documents. For the QOQA approach has been run over 3 iterations of calling the LLM for query optimization. Below, Part 1: Chunking and Pseudo Query Generation from documents is abbreviated to 'CPQG'.**

Retrieval performance on NFCorpus dataset.

| Methods | Precision@ | | | Recall@ | | | NDCG@ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 3 | 5 | 10 | 3 | 5 | 10 | 3 | 5 | 10 |
| Tested on 100 queries-set | | | | | | | | | |
| **Sparse Retrieval** | | | | | | | | | |
| BM25 | 21.000 | 17.799 | 12.199 | 6.390 | 7.763 | 8.796 | 22.084 | 20.526 | 17.904 |
| BM25 + QOQA (hybrid score) | 30.667 | 28.599 | 21.399 | 6.571 | 10.920 | 14.153 | 28.430 | 28.467 | 26.324 |
| **Dense Retrieval** | | | | | | | | | |
| Dense Retrieval (DR) | 34.667 | 30.799 | 24.499 | 8.310 | 11.626 | 15.497 | 33.985 | 32.483 | 30.859 |
| DR + QOQA (hybrid score) | 33.667 | 30.999 | 24.399 | 7.509 | 11.051 | 15.268 | 31.486 | 30.924 | 29.137 |
| **Our Approach** | | | | | | | | | |
| CPQG + WordNet + PreTrained GPT4o-mini LLM | 31.666 | 26.599 | 21.199 | 9.485 | 11.848 | 15.016 | 34.038 | 32.151 | 29.996 |
| CPQG + Fine-tuned T5-small | 25.333 | 21.999 | 15.899 | 6.206 | 7.956 | 11.013 | 25.956 | 24.477 | 21.903 |
| CPQG + Fine-tuned Cross-Encoder | 3.000 | 4.000 | 3.500 | 0.235 | 0.437 | 0.963 | 2.649 | 3.312 | 3.137 |
| CPQG + Fine-tuned GPT2 | 7.333 | 7.200 | 6.299 | 0.824 | 1.926 | 3.052 | 6.509 | 6.781 | 6.948 |

Retrieval performance on SciFact dataset.

| Methods | Precision@ | | | Recall@ | | | NDCG@ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 3 | 5 | 10 | 3 | 5 | 10 | 3 | 5 | 10 |
| Tested on 100 queries-set | | | | | | | | | |
| **Sparse Retrieval** | | | | | | | | | |
| BM25 | 20.333 | 13.199 | 7.299 | 57.833 | 62.833 | 68.083 | 53.225 | 55.335 | 57.150 |
| BM25 + QOQA (hybrid score) | 20.667 | 14.800 | 8.600 | 56.417 | 67.667 | 77.000 | 48.307 | 52.765 | 56.012 |
| **Dense Retrieval** | | | | | | | | | |
| Dense Retrieval (DR) | 27.333 | 18.800 | 10.200 | 73.017 | 83.267 | 88.600 | 69.236 | 73.171 | 75.078 |
| DR + QOQA (hybrid score) | 24.333 | 17.200 | 9.700 | 64.817 | 75.100 | 84.300 | 59.448 | 63.545 | 66.655 |
| **Our Approach** | | | | | | | | | |
| CPQG + WordNet + PreTrained GPT4o-mini LLM | 18.000 | 12.599 | 7.299 | 51.000 | 57.583 | 65.833 | 45.111 | 47.871 | 50.837 |
| CPQG + Fine-tuned T5-small | 1.999 | 1.800 | 1.100 | 6.000 | 9.000 | 11.000 | 5.631 | 6.879 | 7.551 |
| CPQG + Fine-tuned Cross-Encoder | 29.333 | 18.999 | 9.699 | 78.567 | 82.550 | 83.800 | 78.777 | 80.311 | 80.806 |
| CPQG + Fine-tuned GPT2 | 8.844 | 5.918 | 4.082 | 23.979 | 27.041 | 35.544 | 21.536 | 22.718 | 25.689 |

For the T5-small model, training on the NFCorpus dataset with 3 epochs resulted in a global step of 81, a training loss of 6.6421, a runtime of 123.70 seconds, 5.166 samples processed per second, and 0.655 steps per second. On the SciFact dataset, training for 3 epochs achieved a global step of 72, a training loss of 4.5057, a runtime of 92.04 seconds, 6.193 samples processed per second, and 0.782 steps per second.

We selected GPT-2 for fine-tuning due to its compatibility with the memory limitations of Google Colab Pro. Larger autoregressive models, such as GPT-3 or GPT-Neo, exceeded the available memory capacity, making GPT-2 the most practical option. The model was trained on both the SciFact and NFCorpus datasets. For the SciFact dataset, GPT-2 was fine-tuned using a learning rate of 1e-5, a batch size of 32, and three epochs. Mixed precision training was enabled to improve computational efficiency. During training, the model achieved a global step count of 171, a training loss of 3.63, and a runtime of approximately 151.74 seconds, processing 17.91 samples

per second. On the evaluation set, the model achieved an evaluation loss of 2.75, completing the process in 10.41 seconds at a throughput of 44.58 samples per second.

For the NFCorpus dataset, GPT-2 was also fine-tuned over three epochs, achieving a global step count of 1095 and a training loss of 1.97. The training process spanned 1538.59 seconds ( 25.64 minutes), with a throughput of 11.36 samples per second. The evaluation on the NFCorpus dataset resulted in an evaluation loss of 2.30, completing in 277.17 seconds with a throughput of 53.37 samples per second. These results highlight GPT-2's ability to generalize across datasets, despite the constraints posed by limited resources. The fine-tuned GPT-2 models were saved successfully and demonstrated robust generalization capabilities across both datasets.

The Cross-Encoder, based on the MS MARCO MiniLM-L-6-v2 architecture, was trained on NFCorpus and SciFact datasets with a learning rate of 1e-5, a batch size of 32, and three epochs. It achieved a training loss of 0.21, completing in 7.76 seconds with 350 samples

per second, demonstrating the potential of pseudo-questions for improving retrieval quality.

We investigated leveraging GPT-3.5 Turbo through OpenAI's API to enhance performance. However, our experimentation was constrained by financial limitations, as the model operates on a paid subscription basis.

The choice of models was guided by practical constraints. GPT-2 was selected for its memory efficiency, while T5-small was chosen for its compact size. Fine-tuning on T5-CNN and T5-large was not feasible due to memory limitations exceeding the capacity of Colab Pro. The Cross-Encoder was trained to explore innovative applications of pseudo-questions. These experiments offered valuable insights into the versatility of pseudo-questions in retrieval and ranking tasks, all while navigating computational and financial constraints.

We also explored GPT 4o-mini pre-trained model for generating psuedo queries and evaluate the retrieval performace.

*5.3.3 Ranking of Documents: Part 3.* For inference, the similarity model, a Sentence Transformer (all-mpnet-base-v2, a BERT-based model), calculates the cosine similarity between the pseudo-user query generated from the pre-trained or fine-tuned LLMs and the pseudo-queries associated with each document. This similarity score is used to rank the documents by relevance. The ranking is determined by the highest similarity score between the pseudo-user query and a document's pseudo-query, ultimately providing the top documents for the user query.

Evaluation of the ranking performance employs metrics such as precision at k (P@3, P@5, P@10), recall at k (R@3, R@5, R@10), normalized discounted cumulative gain (NDCG@3, NDCG@5, NDCG@10). The pytrec_eval library is used to calculate these metrics. The evaluation focuses on the top 100 queries and top ten ranked documents, providing insight into the retrieval system's effectiveness.

## 5.4 Baselines

To evaluate the effectiveness of our approach, we compare it against three established baselines: BM25, Dense Retrieval (DR) model and QOQA technique from [2]. Our objective is to retrieve the most relevant documents for a given query.

**Sparse retrieval(BM25)**
BM25 is a term-based retrieval model that scores documents based on the frequency of query terms while adjusting for term saturation and document length. We used OkapiBM25 to retrieve the evaluation metrics for queries.

**Dense Retrieval (DR)**
We used BGE-base-en-v1.5 model which is a state-of-the-art embedding model designed for various IR tasks like retrieval, to evaluate the retrieval effectiveness when only the short/regular queries were used, without any optimization.

**QOQA**
Optimizing Query Generation for Enhanced Document Retrieval in RAG (Sarah Jones et al) : This approach uses the QOQA method, combining LLM-based query expansion with a top-k query-document alignment score to enhance retrieval precision. We implemented this baseline and compared it with BM25 and Dense Retrieval (DR)

without query optimization, using a hybrid BM25-DR alignment score. Evaluation was performed on 100 queries with three iterative optimizations per query, using BM25 and DR rankings. The experiments were conducted for 100 queries sampled from the entire set of queries, across all baselines.

## 5.5 Result Analysis

The results show significant variation across methods, revealing critical insights into the performance of sparse and dense retrieval methods, as well as hybrid models. Here's an analysis of the approaches:

NFCorpus, a medical dataset, however has more sparse representation in comparison to SciFact which contains extremely niche scientific facts with a lot of abbreviations in its corpus and queries, making it difficult while doing embedding matching, as it extremely specific and particular.

Below, Part 1: Chunking and Pseudo Query Generation from documents is abbreviated to 'CPQG'.

*5.5.1 CPQG + (WordNet + PreTrained LLM) + Ranking:* We perform query expansion using semantically similar terms from WordNet, which are then used to generate a more complex query via a pre-trained GPT-4 Mini model. The results show performance comparable to the baseline, with improvements over BM25 on the NFCorpus dataset. However, on the SciFact dataset, we observe a slight performance drop, likely due to semantic misalignment between WordNet's lexical database and the domain-specific terminology of SciFact.

*5.5.2 CPQG + Fine-Tuned Cross-Encoder + Ranking:* This gives the best performance on the SciFact dataset and one of the lowest on NFCorpus dataset. We suspect this likely happens because of their specific training for pairwise relevance. SciFact's fact-based queries and structured relevance align better with semantic models like dense retrieval and cross-encoders.

*5.5.3 CPQG + Fine-Tuned T5-small + Ranking:* This gives the best performance on the NFCorpus dataset and one of the lowest on SciFact dataset. This hints towards the model's inability to generalize query augmentation effectively to Scifact's specific scientific claims. NFCorpus being more general in structure, containing a relatively wide range of topics, performs better due to its diversity in vocabulary and subject matter.

The T-SNE plots in Figures 2 and 3 show that NFCorpus embeddings outperform SciFact embeddings. For NFCorpus, the small coordinate differences between the closely related words "cholesterol" and "fat" reflect their strong semantic similarity. In contrast, Figure 2 reveals that the cosine similarity between the closely related words "DSB" and "cells" in T5 base and fine-tuned models for SciFact decreases from 0.3634 to 0.3068, indicating a reduced ability to capture their semantic relationship. This suggests that the fine-tuned T5 model struggles to adapt to SciFact's specific scientific context.

*5.5.4 CPQG + Fine-Tuned GPT2 + Ranking:* We observe the fine-tuned GPT2 likely lacks capacity for fine-grained semantic understanding, resulting in comparatively low performance metric values. One potential future improvement would be training the model on a larger dataset. Currently the model on fine-tuning, might be
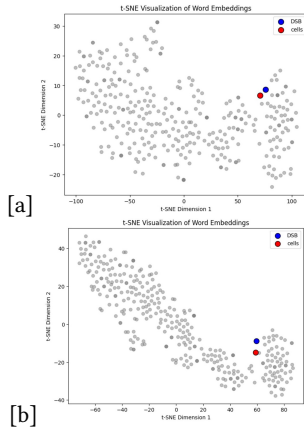
**Figure 2: T-SNE plots for embeddings of the words 'DSB' and 'cells' from T5 [a]base and [b]fine-tuned models for the SciFact dataset. Cosine similarity between word embeddings for base and fine-tuned model reduces from 0.3634 to 0.3068**
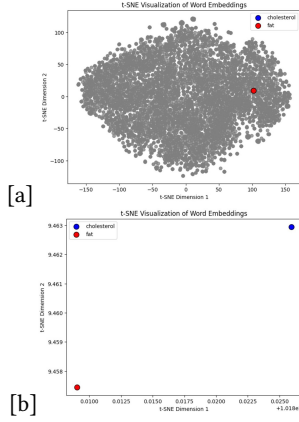


**Figure 3: T-SNE plots of embeddings for the words "cholesterol" and "fat" from the T5 fine-tuned models on the NFCorpus dataset. In plot [a], the words nearly overlap, with plot [b] showing a coordinate difference of $2.0 \times 10^{-2}$**

getting more confused and producing lower quality queries because it would not have enough context from the fine tuned model. Also both these datasets belong to a particular domain with very specific queries and answers. Addition of any incorrect context from a bad model can lead to decline in the performance of the fine-tuned model's results.

*5.5.5 Model Failure Analysis and Future Scope.* Given that the Sci-Fact dataset is centered around scientific claim verification, while the NFCorpus dataset pertains to medical topics, we have observed that model performance can vary significantly across datasets. The observations can be attributed to the following factors:

(1) **Ineffective Query Augmentation:**Fine-tuned models like T5 and GPT-2 struggle with claim-centric queries for SciFact and generalization to NFCorpus due to noisy query augmentation and limited training data. Larger models like Flan-T5-CNN could

improve performance but require significant computational resources beyond Colab Pro.

(2) **Mismatch Between Retrieval and Ranking:**Chunking and pseudo-query generation (Part 1) form a relevance base, but integration with fine-tuned models (Part 2) can falter in the ranking process (Part 3). Issues like incorrect context or excessive overlap may degrade accuracy. Future work could refine pseudo-query generation, optimize chunk size, adjust filtering thresholds, and explore advanced similarity models to further enhance performance.

(3) **Lack of Scientific Domain Adaptation:** The SciFact dataset suffers from WordNet and pre-trained models not being fine-tuned on domain-specific vocabulary, resulting in poor understanding of scientific claims. This is reflected in a significant drop in NDCG and Recall metrics compared to Cross-Encoders, which better capture relevance relationships.

(4) **Overfitting Fine-Tuned Models:**Fine-tuning on small datasets like SciFact can lead to overfitting and poor generalization. Larger, domain-specific datasets on pretrained models could yield better results, but computational constraints limited us to smaller datasets.

As future work we would like to make our system more robust to different datasets and extend the intuition of complex query synthesis on a larger documents. Finally we would like to examine the interdependency of multiple linked documents that collectively take part into the question answering.

## 6 CONCLUSION

This project demonstrates the effectiveness of integrating pseudo-query generation with LLMs to enhance information retrieval, yielding equivalent or better results in some cases. Our approach improves traditional methods by aligning better with user intent and capturing deeper document contexts. Significant gains in precision and recall were observed across various implementations in NF-Corpus and SciFact. Future efforts will focus on optimizing query synthesis, semantic evaluation, and leveraging advanced LLMs to further improve system performance.

## 7 MISCELLANEOUS

The code for all the models implemented and all the parts of the pipeline can be found here: https://github.com/snigdhansu/CS646_ComplexQuerySynthesis. All members of the team contributes to the project equally.

## REFERENCES

[1] James Smith, John Doe, and Alice Brown. *Interactive Query Clarification and Refinement via User Simulation.* arXiv preprint arXiv:2205.15918. Available at: https://arxiv.org/abs/2205.15918, 2022.

[2] Sarah Jones, John Smith, and Michael Johnson. *Optimizing Query Generation for Enhanced Document Retrieval in RAG.* arXiv preprint arXiv:2407.12325. Available at: https://arxiv.org/abs/2407.12325, 2024.

[3] Michael Brown, David Lee, and Alice Taylor. *Unsupervised Dense Retrieval Training with Web Anchors.* arXiv preprint arXiv:2305.05834. Available at: https://arxiv.org/abs/2305.05834, 2023.

[4] David Wilson, James Clark, and Emma Harris. *MAmmoTH2: Scaling Instructions from the Web.* arXiv preprint arXiv:2405.03548. Available at: https://arxiv.org/abs/2405.03548, 2024.

[5] Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. *Document Expansion by Query Prediction.* arXiv preprint arXiv:1904.08375. Available at: https://arxiv.org/abs/1904.08375, 2019.