



**Hochschule  
Bonn-Rhein-Sieg**  
University of Applied Sciences

Fachbereich Informatik  
Department of Computer Science

Master Project  
Computer Science

# LitQEval: Measuring the Effectiveness of Litreature Search Queries

Computer Science

Mohammad Sakinini

Supervisor	Philipp Baaden Fraunhofer INT
Examiner 1	Prof. Dr. Jörn Hees
Examiner 2	Dr. Milos Jovanovic

Draft as of	2024-11-19 13:00:03+01:00 (For submission: set <b>final</b> option in thesis.tex!)
To be submitted on	2024-07-26

## Abstract

This work is based on a larger initiative known as the Search Query Writer (SQW), an internal tool developed at Fraunhofer INT to aid scientific researchers in creating comprehensive literature search queries. These queries are intended to provide researchers with a strong starting point in a topic area they may have limited knowledge about.

The current state of the SQW tool presents a key challenge: the absence of a mechanism to evaluate the quality of the generated queries. As a result, the evaluation has so far been conducted subjectively. This project aims to address this issue by introducing a dataset that contains publications deemed relevant to specific topics. Additionally, it introduces several metrics to account for different aspects of query evaluation, given the complexity of the task. **(Explain performed experiments after completing them)**

# Contents

<b>List of Figures</b>	<b>iv</b>
<b>1. Introduction</b>	<b>1</b>
1.1. Motivation . . . . .	1
1.2. Research Questions . . . . .	2
1.3. Structure of this Work . . . . .	2
<b>2. Foundations</b>	<b>3</b>
2.1. Search Query Writer . . . . .	3
2.2. Related Work . . . . .	5
<b>3. LitQEval</b>	<b>7</b>
3.1. Dataset . . . . .	7
3.1.1. Dataset Analysis . . . . .	9
3.2. Evaluation metrics . . . . .	10
<b>4. Evaluation</b>	<b>13</b>
4.1. Experimental Setup . . . . .	13
4.2. Results . . . . .	13
4.3. Discussion . . . . .	13
<b>5. Conclusion</b>	<b>14</b>
5.1. Summary and Contributions . . . . .	14
5.2. Outlook . . . . .	14
<b>A. Appendix</b>	<b>15</b>
A.1. Further Details on Something . . . . .	15
<b>Bibliography</b>	<b>19</b>
<b>Declaration</b>	<b>20</b>

# List of Figures

2.1. Search Query Writer . . . . .	5
3.1. Dataset overview of the research topics . . . . .	8
3.2. Core Publications Clustering . . . . .	9
3.3. Semantic Precision using Cosine Similarity . . . . .	11
3.4. Semantic Precision using MVEEE . . . . .	12
A.1. SQW Knowledge Enrichment . . . . .	16
A.2. Field citation ratio per topic . . . . .	17
A.3. Distribution of publication years per topic . . . . .	18

# 1. Introduction

The Fraunhofer Institute for Technological Trend Analysis (INT)<sup>1</sup> is a research institute that frequently undertakes new tasks and research questions across various fields. Often, these inquiries require systematic and robust scientific responses, even when the initial knowledge in the area may be limited at the time. Given the recurring nature of this challenge, a tool that can support researchers by providing a head start and entry point into unfamiliar fields is essential.

To address this, several internal tools have been developed to analyze large volumes of scientific data from sources like Dimensions.ai<sup>2</sup> and Web of Science<sup>3</sup>. The rise of large language models (LLMs) has further enhanced the appeal and accessibility of automation across numerous domains, including scientific research, spanning from idea generation and experimental iteration to paper composition[6].

In the realm of search queries, the main focus has been on text-to-SQL[2], where an LLM is prompted via natural language to generate a precise and valid SQL query. However, to our knowledge, there has been limited and non-diverse effort dedicated to the development of text-to-literature search queries. Thus this work introduces a pipeline and curates a dataset designed to address this gap, with a particular focus on enhancing the Fraunhofer Search Query Writer tool.

## 1.1. Motivation

The SQW tool is currently under development by the company and has generated significant interest among researchers. However, a primary challenge we face after testing earlier versions is evaluating the quality of the generated queries. Initially, we considered gathering human feedback from users by requesting them to rate the generated query on a scale of 0 to 5. While this approach could be useful for fine-tuning the underlying model, the quantity of feedback has so far been limited and remains subjective. This is especially problematic because the tool’s purpose is to generate queries for researchers who are new to a given topic. Consequently, if the query quality is poor, the researcher may not immediately recognize this.

Identifying suitable evaluation metrics and datasets to assess the quality of the generated queries is a complex task, which forms the basis of this master’s project. The project’s objective is to find a robust solution for assessing the quality of literature search queries, enabling the further development of the SQW tool to

---

<sup>1</sup><https://www.int.fraunhofer.de/>

<sup>2</sup><https://www.dimensions.ai/>

<sup>3</sup><https://clarivate.com/>

provide more accurate results and improve productivity through the integration of large language models.

## 1.2. Research Questions

There are three main research questions that we aim to address while curating the dataset and formulating the metrics to evaluate the generated queries. These questions are based on the following hypothesis: Given that we know which publications are the most important for a given field, which we will refer to as **Core Publications** (CP).

- **RQ1:** How many of the core publications can the generated search query recall?
- **RQ2:** How many of the non-core publications are relevant?
- **RQ3:** Which metric can we use to penalize the model for exploiting the ability to generate large queries to achieve high recall?

## 1.3. Structure of this Work

The remainder of this work is structured as follows:

After this introduction, we will first focus on the foundations in [Chapter 2](#), where we will briefly explain the SQW tool, primarily focusing on the format of the input and output. We will also cover the necessary basics and information about the main data source, Dimensions.ai, which will be used to curate the dataset ??, as well as common metrics ?? used for evaluation in the community, before discussing related works in [Section 2.2](#).

Next, in [Chapter 3](#), we will detail our approach and its components, which involve curating a dataset that contains core publications and establishing a pipeline to streamline the evaluation process, thereby accelerating the development of the SQW tool.

We will evaluate our approach in [Chapter 4](#), beginning with a description of our experimental setup, where we acquire the generated query via the SQW tool and explain the reasoning behind the selection of topics for which the queries are generated.

Finally, we will conclude this work with a summary of our main contributions and an outlook in [Chapter 5](#).

## 2. Foundations

In this section, we will begin by briefly introducing the SQW tool to provide a foundational understanding of how its settings may influence the overall results. This introduction will also establish the groundwork for designing an evaluation process that ensures a fair and accurate assessment.

Next, we will review prior works that tackle the challenge of using large language models (LLMs) to generate literature search queries, examining their potential in this domain.

Following this, we will describe the dataset curation process, including relevant statistics and exploratory data analysis. This analysis will help us evaluate the dataset’s quality and determine the topics suitable for evaluation.

Then, we will introduce new metrics designed to assess the quality of generated search queries. These metrics are intended not only to identify true positives from core publications but also to account for other potentially relevant publications. Once the evaluation pipeline is established, we will conduct a comparative analysis between a baseline and the queries generated by the SQW tool.

Finally, we will provide an outlook on the next steps and potential optimization options for the SQW tool. Additionally, we will discuss other tools that could be developed to build upon these advancements.

### 2.1. Search Query Writer

The Search Query Writer is a tool based on a large language model (LLM), specifically using GPT-4o, to systematically generate literature search queries. The only required input for this tool, which is the main focus of this work, is the **Topic**. Users are required to provide a topic for generating a search query, irrespective of the scientific field—for example, *Synthetic Biology*.

Several optional inputs are available to enhance the quality of the generated query, including:

- **Negative Keywords:** Terms that should be excluded to avoid unwanted results.
- **Description:** A description that serves as an alignment mechanism to clarify the task’s intent.
- **Modes:** Three selectable modes (Strict, Moderate, Creative) that control the temperature of the LLM to manage the level of randomness in responses.

- **Depth:** A parameter that specifies how comprehensively the topic should be analyzed.
- **Supporting Documents:** Users can upload a PDF, ideally a survey or overview document on the topic, which helps the tool acquire knowledge about the scientific field and better align with the research intent.

These additional inputs are intended to refine and tailor the search query to more closely match the user’s research goals, but will not be extensively tested in this work.

To generate a literature search query, we designed the SQW to take a human-like approach, divided into two main steps: **Knowledge Enrichment** and **Iterative Scientific Fine-Tuning**.

The objective of the Knowledge Enrichment step is to provide the LLM with contextual information about the research topic. This is achieved by first retrieving information from Wikipedia based on the given topic. Specifically, the first 4,000 characters from the top- $k$  pages are collected and summarized before being passed into the LLM’s memory. ArXiv is queried in a similar manner to gather relevant research content. Additionally, we perform an online search using DuckDuckGo<sup>1</sup>, aggregating results to offer a broader understanding of the topic.

Notably, each of the steps is conducted within a separate memory session, with results stored independently for future use. This setup allows the model to explore the topic using various sources, helping to mitigate any potential recency bias and ensure a well-rounded context.

The output of this first stage will be a list of keywords that are usually presented in a transfer-list as shown in Figure A.1, that contains two lists; specific and general keywords, along side some additional information such as the number of publications found per keyword. The goal is to let the user decide whether a keyword is too broad in which he is supposed to move it to the general list, and if it targets the specified topic quite well then it should stay in the specific list, and at the end they user should also provide an overarching topic for which the scope of the general keywords is limited to be more focused to words the research intent. The output of this step will be the queries that will be used for the evaluation at the end.

The iterative scientific fine-tuning on the other hand approaches more scientific sources, namely dimensions.ai, which is a literature database that offers quick access to publications across a wide range of journals. The query generated in the earlier stage is then used to prompt dimensions three times, once to retrieve the most cited 1k literature, a second time to retrieve the newest 1k literature, and one last 1k to retrieve the most relevant documents based on their altmetric rating. This leaves us with a total of 3k publications in which we extract the title and abstract for, and use an Openai’s embedding model, and apply a simple RAG pipeline to retrieve the most relevant keywords based on the extracted passages.

---

<sup>1</sup><https://duckduckgo.com/>





Figure 2.1.: A simplified overview of the Search Query Writer. The process begins with the Knowledge Enrichment stage, where the model receives input data and sends it to an LLM agent equipped with a suite of tools to gain insights into the topic. Based on this understanding, the model generates a well-structured search query, formatted and executed across multiple dimensions to retrieve a relevant selection of literature. Within this literature set, RAG is applied to identify the most pertinent keywords, which are compiled into an optimized search query. This query can be iteratively refined to enhance overall search quality.

## 2.2. Related Work

Systematic literature reviews are widely used across various fields, allowing researchers to conduct a comprehensive manual review of scientific topics and identify publications that answer a set of important research questions. However, one significant challenge with this approach has been the exponential growth in the number of publications, which makes conducting unbiased reviews increasingly difficult. In the age of technological advancements, we can now leverage these technologies to assist in investigating topics without the need to manually sift through extensive lists of publications. To address this issue, a series of works have been proposed within the Conference and Labs of the Evaluation Forum (CLEF) [3–5]. These works focus on the evaluation of empirical medical research, utilizing a dataset of medical literature. They introduce two primary tasks: Task 1, which involves identifying relevant studies from the PubMed medical database, and Task 2, which assesses the ranking of studies following title and abstract screening. Notably, the evaluation pipeline, along with the dataset and descriptions of these tasks, are publicly accessible on GitHub<sup>2</sup>.

Large Language Models (LLMs) have had a significant impact on modern technology, including in scientific research, where they have provided remarkable

<sup>2</sup><https://github.com/CLEF-TAR/tar/tree/master>

improvements in speed. While the processing speed of LLMs is unprecedented, the quality of their output in various domains is still being explored. The work by Wang [10] investigated the performance of ChatGPT in generating Boolean search queries for literature reviews. Specifically, it evaluated the effectiveness of ChatGPT in constructing queries for systematic literature reviews using different prompting techniques, including zero-shot, few-shot and iterative guided approaches. The evaluation used the CLEF datasets [3–5] and an additional medical dataset containing a collection of seeds [9]. Although the results highlight the limitations of ChatGPT’s performance, this work underscores the potential of LLMs to aid literature review, especially when supported by examples or more advanced, structured pipelines.

A broader and more diverse evaluation of the quality of automatically generated literature search queries for systematic literature reviews was conducted by Badami [1]. In this work, they introduced a pipeline that generates literature search queries based on a given research question and abstracts from previously identified relevant publications. The evaluation was performed against a dataset they constructed, which contains the results of 10 systematic literature reviews, including candidate papers, queries used, and relevant papers identified in each review. For example, in the review  $SLR_1$ , a total of 7,002 candidate papers were retrieved using search query  $S$ , from which a subset of 59 relevant papers  $RP$  was identified. To assess their proposed approach, they compared the generated queries in various settings using recall and precision metrics, benchmarking them against the original search query  $S$ . The dataset is publicly available on Zenodo<sup>3</sup>.

---

<sup>3</sup><https://tinyurl.com/496zuar3>

### 3. LitQEval

Despite ongoing research on automatic literature query generation and related evaluations with medical datasets, such as CLEF [3–5] and the Collection of Seeds [9], the insights gained from these evaluation metrics are not particularly compelling for our use case. This limitation arises from two main factors.

First, the CLEF and Collection of Seeds datasets are exclusively focused on medical data. Although Badami’s work [1] offers a more diverse dataset, it lacks a suitable evaluation metric. Their evaluation primarily aims to maximize recall, with minimal consideration for precision, as literature search queries often yield far more results than necessary, making precision a less effective measure in this context.

A second limitation arises when recall is prioritized exclusively. For example, if we aim to train a model to generate queries that maximize recall, there is no penalty for generating overly broad queries, such as those that exploit wildcards, which could lead to an excessive number of irrelevant results.

To address these issues, we introduce a dataset structured similarly to that of Badami [1] but designed to be more comprehensive and covering a wider range of topics. Alongside this dataset, we propose new evaluation metrics that account for the inherently broad nature of literature search queries while penalizing excessively large queries. These metrics also emphasize the importance of accurately identifying core publications that are deemed highly relevant within the domain.

#### 3.1. Dataset

The dataset we aim to create has two primary goals: First, it should encompass a wide range of randomly selected scientific research fields. Second, for each selected field, it should contain a set of highly relevant publications to serve as anchors for evaluating additional publications found in these areas.

Selecting new research topics is straightforward; however, to avoid bias from ongoing research interests, we used ChatGPT to generate a list of scientific fields that are recent and not overly broad. For instance, a topic like *Artificial Intelligence* is vast, making it challenging to accurately and comprehensively identify core publications. Instead, we chose a more specific, problem-focused topics such as *Drones in Agriculture*. To search for the corresponding bibliometric analysis we used the following query: *<TERM> (“Bibliometric” OR “Scientometric” OR “Systematic literature” OR “Most Influential” OR “Most Cited” OR “Scientific Landscape” OR “Literature Landscape” OR “Core Literature”)*

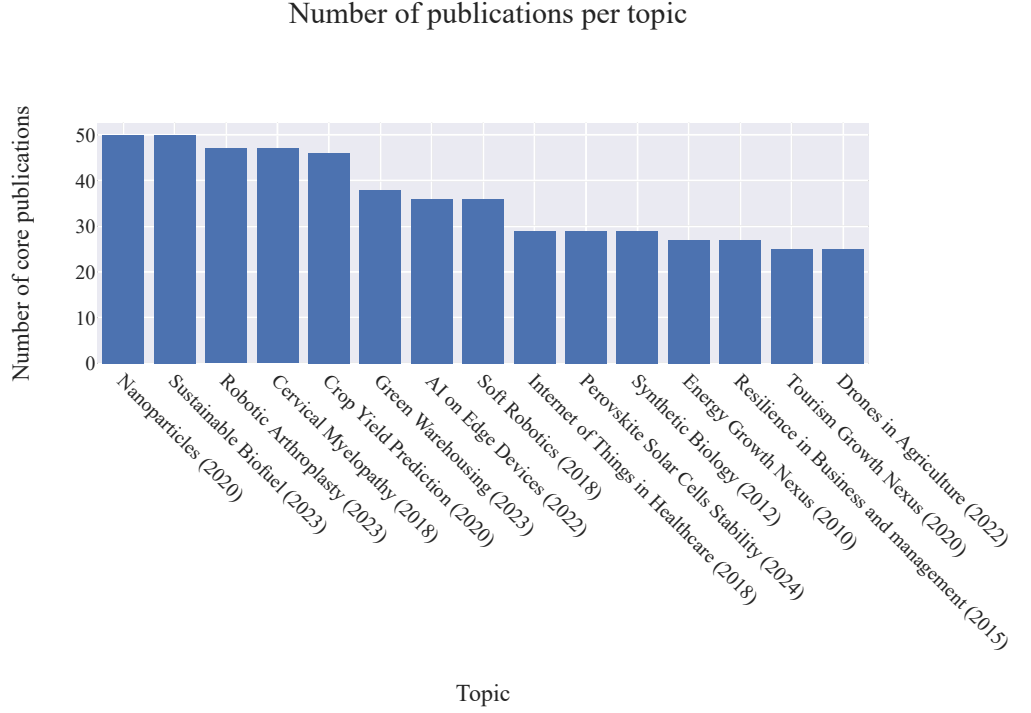


Figure 3.1.: An overview of the dataset and the selected 15 research fields with respective core publications identified through bibliometric analyses. The number in brackets following the field name on the x-axis represents the year of the analysis.

After identifying a sufficient number of diverse fields, 15 in our case, we sought to collect core publications for each field. Due to the difficulty of gathering core publications across a broad array of topics, we leveraged the bibliometrics community’s expertise. Specifically, we searched for bibliometric studies that identify the most relevant publications within each research area. For example, a bibliometric analysis of *Drones in Agriculture* [7] lists the most cited publications from 1990 to 2021. In this case, 40 core publications were identified, which we manually located on Dimensions.ai and added to our dataset, omitting any publications not found in Dimensions.

This process was repeated across all selected research fields, resulting in a dataset containing 15 topics, each with 25–50 core publications, as illustrated in Figure 3.1.

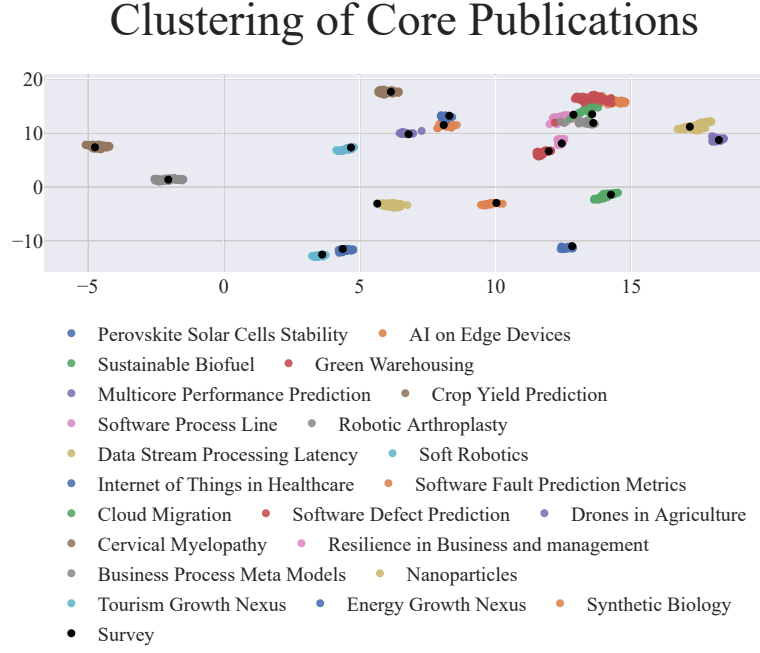


Figure 3.2.: This figure shows clusters of publication embeddings based on titles and abstracts. Embeddings were generated with OpenAI’s small model and reduced in dimensionality with t-SNE, then clustered using k-means with  $k = 15$  (indicating topic count). Clusters group core publications by semantic similarity, with overlaps in topics like *IoT in Healthcare* and *AI on Edge Devices*, as well as *Tourism Growth Nexus* and *Energy Growth Nexus*, due to the similarity in the research field.

#### 3.1.1. Dataset Analysis

We recognize that potential biases may exist in our dataset due to its complete reliance on the bibliometric community for identifying core publications. This often implies that publications with higher citation counts are deemed more relevant. To assess this, we analyzed the citation distribution per topic, as provided by Dimensions, shown in Figure A.2. Additionally, we examined the distribution of publication years per topic, illustrating the time span considered in the bibliometric analyses, as shown in Figure A.3. If we compare the distribution of publication years for the medical research field *Cervical Myelopathy* with that of *IoT in Healthcare*, both of which were published in 2018, we can observe distinct differences in the year distributions of their core publications. These variations may be attributed to factors such as the recency of the field, changes in terminology over time, or the nature of the research area, where one field may prioritize more established works while the other focuses on recent advancements.

For the evaluation pipeline that we will introduce, the embeddings of the core documents are essential for effectively assessing the search query, as detailed in

[Section 3.2](#). To validate this approach, we examine the clustered embeddings of the titles and abstracts for each core topic, as well as the bibliometric analyses in which these documents were initially referenced. This enables us to assess whether core publications within each field exhibit semantic similarity while also demonstrating some degree of dissimilarity from publications in other topics. The resulting clusters, shown in [Figure 3.2](#), were generated using k-means clustering, where  $k$  is set to the number of topics. For this, we use OpenAI’s small embeddings model alongside t-SNE[8] to reduce the dimensionality to 2D.

### 3.2. Evaluation metrics

The standard evaluation metrics for query evaluation are recall and precision. We argue that while recall is of high importance, particularly within the community, precision in this context becomes less feasible. Specifically, retrieving only the exact core publications via a search query would be impractical without explicitly using DOIs to target them directly, which renders this metric largely obsolete and likely to be consistently low. However, we still aim to account for the number of matched publications when executing a search query to prevent models from exploiting overly large queries. To address this, we introduce the concept of *Semantic Precision*.

The idea behind Semantic Precision is to evaluate the relevance of retrieved publications in comparison to the core publication set. If the retrieved publications are sufficiently similar to those in the core set, they are deemed to hold some relevance rather than being entirely unrelated. To achieve this, we propose that the core publications, encompass sufficient semantic breadth to gauge the quality of literature relevant to a specific field. We calculate Semantic Precision in two ways.

The first approach involves averaging the embeddings of the core publications. We then set an acceptance threshold based on the cosine similarity to the least similar core publication, given by. This means that if the embedding of a retrieved publication is more similar to the center than the least similar core publication, we consider it a relevant publication, as shown in [Figure 3.3](#). We define:

- $CPs$  as the set of core publications.
- $\mathbf{c}_i$  as the embedding vector of the  $i$ -th core publication.
- $\mathbf{p}$  as the embedding vector of a retrieved publication.
- $\cos(\mathbf{a}, \mathbf{b})$  as the cosine similarity between two vectors  $\mathbf{a}$  and  $\mathbf{b}$ .

First, compute the centroid of the core publication embeddings:

$$\mathbf{c}_{\text{centroid}} = \frac{1}{|CPs|} \sum_{\mathbf{c}_i \in CPs} \mathbf{c}_i$$

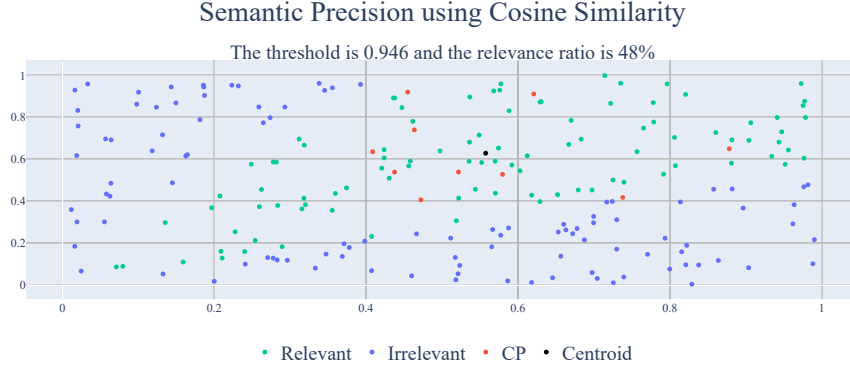


Figure 3.3.: This illustration demonstrates the effect of cosine similarity on randomly generated data in a 2D space. The core publications (CP) are shown in red, positioned between 0.4 and 0.9 on both the x- and y-axes. When we set the threshold to 0.946, based on the cosine similarity of the least similar core publication from the centroid, many retrieved publications on the opposite side of the spectrum are still assigned as relevant. This effect occurs because cosine similarity considers only the angle between vectors, ignoring their magnitude. In this case, this results in 48% of the retrieved publications being considered relevant.

Then, let the threshold similarity,  $\theta$ , be the cosine similarity of the least similar core publication to the centroid:

$$\theta = \min_{\mathbf{c}_i \in CP} \cos(\mathbf{c}_{\text{centroid}}, \mathbf{c}_i)$$

Finally, Semantic Precision using cosine similarity ( $SP_{cos}$ ) is defined as [Equation 3.1](#), where  $\mathbb{I}$  is an indicator function that equals 1 if the retrieved publication  $\mathbf{r}$  meets the similarity criterion and 0 otherwise:

$$SP_{cos} = \frac{\sum_{\mathbf{p} \in \text{pubs}} \mathbb{I}(\cos(\mathbf{c}_{\text{centroid}}, \mathbf{p}) \geq \theta)}{|\text{retrieved}|} \quad (3.1)$$

For the second approach we omit the averaging of the embeddings and use Minimum Volume Enclosing Ellipsoid (MMVE), which creates the smallest ellipsoid that includes the our CP, which we then use to determine which of the retrieved publications are relevant by checking whether they are within MMVE or not, as illustrated in [Figure 3.4](#). This approach allows us to take into account all the dimensions by not only considering the angle but also the magnitude

The Minimum Volume Enclosing Ellipsoid (MVEE) for the core publication set  $CP$  is centered at  $\delta$  with shape matrix  $A$ . To determine whether a retrieved publication  $\mathbf{r}$  is relevant, we check if it lies within the ellipsoid by testing the following condition:

$$(\mathbf{p} - \delta)^T A (\mathbf{p} - \delta) \leq 1$$

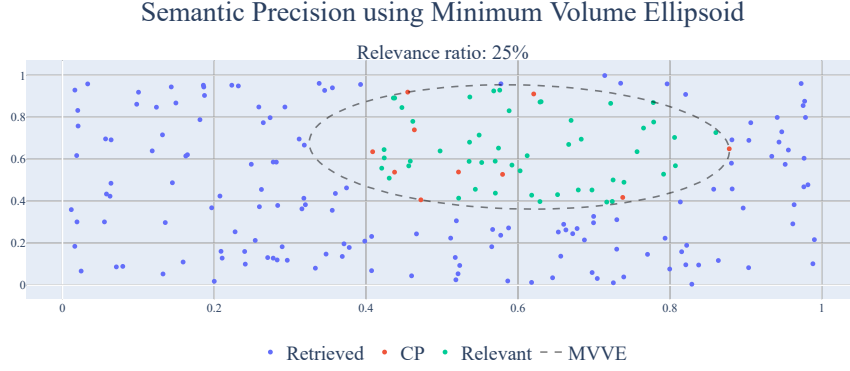


Figure 3.4.: This illustration demonstrates the effect of using the Minimum Volume Enclosing Ellipsoid (MVVE) on randomly generated data in a 2D space. The core publications (CP) are shown in red, positioned between 0.4 and 0.9 on both the x- and y-axes. An ellipsoid is generated using MVVE to define the scope of relevant publications, ensuring that only those within the maximal angles and magnitudes of the core publications are considered relevant. In this case, this approach results in only 25% of the retrieved publications being classified as relevant.

Semantic Precision (SP) for this approach is then:

$$SP_{MVVE} = \frac{\sum_{\mathbf{p} \in \text{pubs}} \mathbb{I}((\mathbf{p} - \delta)^T A(\mathbf{p} - \delta) \leq 1)}{|\text{pubs}|} \quad (3.2)$$

Now that we have a metric that can be used to punish the model in case of generating a too broad of a query, we can use it as a factor to calculate the F-Score, the goal of the standard F-score is to balance out between the recall and precision, but in our case we use  $F - \beta$  instead, whereby the  $\beta$  is the weighting factor of the recall, meaning the higher it is the more important the recall will be, in our case we set it to be 2, meaning that the recall is twice as important as the precision.

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{Precision} \cdot \text{Recall}}{(\beta^2 \cdot \text{Precision}) + \text{Recall}} \quad (3.3)$$

For our specific case where  $\beta = 2$ , emphasizing the importance of recall, it is:

$$F_2 = 5 \cdot \frac{\text{Precision} \cdot \text{Recall}}{(4 \cdot \text{Precision}) + \text{Recall}}$$



## 4. Evaluation

### 4.1. Experimental Setup

### 4.2. Results

### 4.3. Discussion

## 5. Conclusion

### 5.1. Summary and Contributions

### 5.2. Outlook

## A. Appendix

### A.1. Further Details on Something

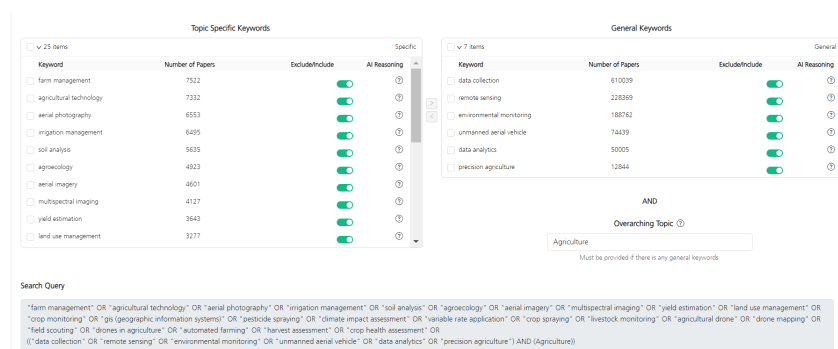


Figure A.1.: A screenshot of the SQW UI after completing the Knowledge Enrichment stage. On the left, a list of keywords is displayed alongside the number of publications associated with each keyword when used as a search term. The keywords on the right-hand side were manually categorized as general and can be roughly assessed by the number of associated publications. To narrow the scope of general keywords, we selected “agriculture” as the overarching topic. The final generated query is displayed and updated interactively as values in the transfer lists are adjusted.

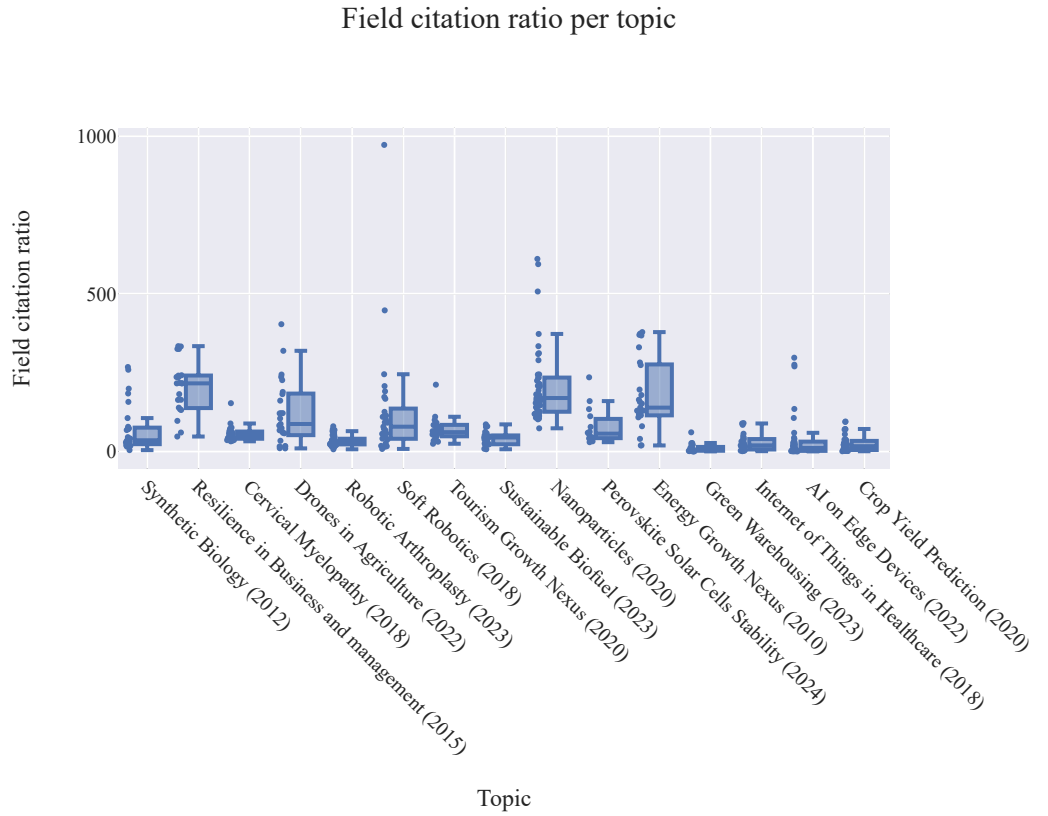


Figure A.2.: The citation ratio per topic, showing the relative citation counts of core publications compared to the average citation frequency within their respective research fields. This illustrates how the prominence of each publication compares to typical citation levels in its field.

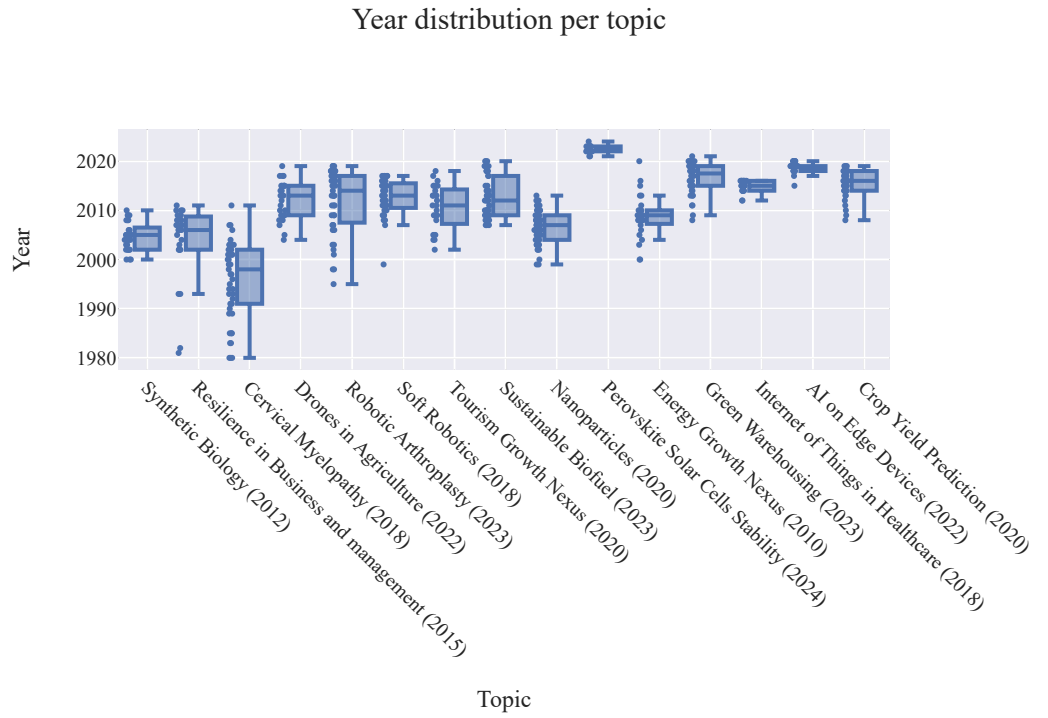


Figure A.3.: The distribution of publication years for core publications across various research topics, highlighting the historical range of studies considered in the bibliometric analyses for each field. Notably, for *Cervical Myelopathy*, the lower bound of publication years was set to 1980 for improved readability, although the actual range goes back to 1953.

# Bibliography

- [1] M. Badami, B. Benatallah, and M. Baez. “Adaptive search query generation and refinement in systematic literature review”. In: *Information Systems* 117 (2023), p. 102231 (cit. on pp. 6, 7).
- [2] X. Dong, C. Zhang, Y. Ge, Y. Mao, Y. Gao, J. Lin, D. Lou, et al. “C3: Zero-shot text-to-sql with chatgpt”. In: *arXiv preprint arXiv:2307.07306* (2023) (cit. on p. 1).
- [3] E. Kanoulas, D. Li, L. Azzopardi, and R. Spijker. “CLEF 2017 Technologically Assisted Reviews in Empirical Medicine Overview”. In: Jan. 2017 (cit. on pp. 5–7).
- [4] E. Kanoulas, D. Li, L. Azzopardi, and R. Spijker. “CLEF 2018 technologically assisted reviews in empirical medicine overview”. In: *CEUR workshop proceedings*. Vol. 2125. 2018 (cit. on pp. 5–7).
- [5] E. Kanoulas, D. Li, L. Azzopardi, and R. Spijker. “CLEF 2019 technology assisted reviews in empirical medicine overview”. In: *CEUR workshop proceedings*. Vol. 2380. 2019, p. 250 (cit. on pp. 5–7).
- [6] C. Lu, C. Lu, R. T. Lange, J. Foerster, J. Clune, and D. Ha. *The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery*. 2024. arXiv: 2408.06292 [cs.AI]. URL: <https://arxiv.org/abs/2408.06292> (cit. on p. 1).
- [7] A. Rejeb, A. Abdollahi, K. Rejeb, and H. Treiblmaier. “Drones in agriculture: A review and bibliometric analysis”. In: *Computers and Electronics in Agriculture* 198 (2022). <https://doi.org/10.1016/j.compag.2022.107017>, p. 107017. DOI: 10.1016/j.compag.2022.107017. URL: <https://app.dimensions.ai/details/publication/pub.1147958699> (cit. on p. 8).
- [8] L. Van der Maaten and G. Hinton. “Visualizing data using t-SNE.” In: *Journal of machine learning research* 9.11 (2008) (cit. on p. 10).
- [9] S. Wang, H. Scells, J. Clark, B. Koopman, and G. Zuccon. “From Little Things Big Things Grow: A Collection with Seed Studies for Medical Systematic Review Literature Search”. In: *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR '22. ACM, July 2022, pp. 3176–3186. DOI: 10.1145/3477495.3531748. URL: <http://dx.doi.org/10.1145/3477495.3531748> (cit. on pp. 6, 7).

- [10] S. Wang, H. Scells, B. Koopman, and G. Zuccon. “Can ChatGPT write a good boolean query for systematic review literature search?” In: *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2023, pp. 1426–1436 (cit. on p. [6](#)).



# Declaration

I declare that I have written this work by myself. I have identified as such all passages taken verbatim or in meaning from published or unpublished works by third parties. All sources and aids that I have used for the work are indicated.

(Example formulations follow, which you must adapt to your work for the sake of transparency. Of course, you should have discussed about the acceptability of such aids with your supervisor in advance.) In particular, the following AI systems were also used to create this work:

- ChatGPT in version ... was used for the initial text drafting based on bullet points given by me in the chapters ... / of the entire work.
- ChatGPT was consulted on the following topics: ... / was used to generate ideas regarding ... / for the structuring of ... / for the conception of the system ... .

The wording of the dialogs and the version used were documented in the appendix of this work. Passages used are marked as such in the text.

- ChatGPT was used to create source code for ... . The wording of the dialogs and the version used were documented in the appendix of this work. The use is indicated in the header of the respective source file / class / method / parts.
- Copilot in version ... was used to create source code / auto-complete for ... . The use is documented in the header of the respective source file / class / method / parts.

I am aware that content generated by AI systems is no substitute for careful scientific work, which is why all such generated content has been critically reviewed and finalized by me.

This work has neither been submitted with the same content nor in essential parts to any other examination authority.

*Your City, 2024-07-26*

---

Mohammad Sakinini