# User-Centric Evaluation of Query Suggestions

MSc UX Engineering
Academic Project 2023
Anyu Wang

## 1/Abstract

As a search assistant, query suggestions have the potential to improve search efficiency, enhance user experience, etc. Academic search is a complex and difficult task and query suggestions can address some of its challenges. The purpose of this study is to investigate the extent to which query suggestions can enhance academic searches conducted by students. Using a user-centric evaluation method, the performance of four sources of query suggestions was measured in terms of both the number of valid terms and the selection rate. The results show that query suggestions from Wordvectors outperform the others on both quantity and quality, and perform excellently under multiple topics.

## 2/Background

Query suggestions can improve the efficiency of academic searches, but there is a lack of user-centric evaluation of query suggestions from different sources.

### Query Suggestions

Query suggestions, commonly seen in search engines, help users by offering related search terms. This makes searching faster and less stressful. However, not all suggestions improve search results. Some techniques, like adding related terms to a search, can help, but the best methods for improving searches are still debated.

### Academic Search

Academic searching is a complex behaviour and varies depending on the background of the searcher. Most people will overestimate their level of academic searching. Too few keywords in a query can lead to poor search results. The high level query suggestions provided by query expansion techniques can eliminate the above issues.

### System-centric VS User-centric

The common evaluation methods and corresponding metrics are accuracy, precision, recall, MRR, nDCG, etc. These are all system-centric evaluation methods. The sample sizes are very large, with the largest having data collection spanning several years. In terms of evaluating applications of query suggestions, many user-centric approaches aim to evaluate interaction and interface design. I have learned from and refined some of the user-centric terms evaluation methodology.

## 5/Conclusion

After comprehensive analyses, Wordvetors was found to be outstanding in terms of the quantity and quality of query suggestions, as well as applicability, and performed well under keywords across a wide range of topics.There were significant differences between Wordvectors' performance and the other three sources. But the remaining three, MeSH, PubMed, and WebIsA, did not show significant differences between each other. MeSH performs well for terms in specific domains, but generally does not work well for academic searches across disciplines. In the future, it is possible to deeply explore the reasons for the variability of the different sources and how they can be combined with actual searches to achieve the best results.
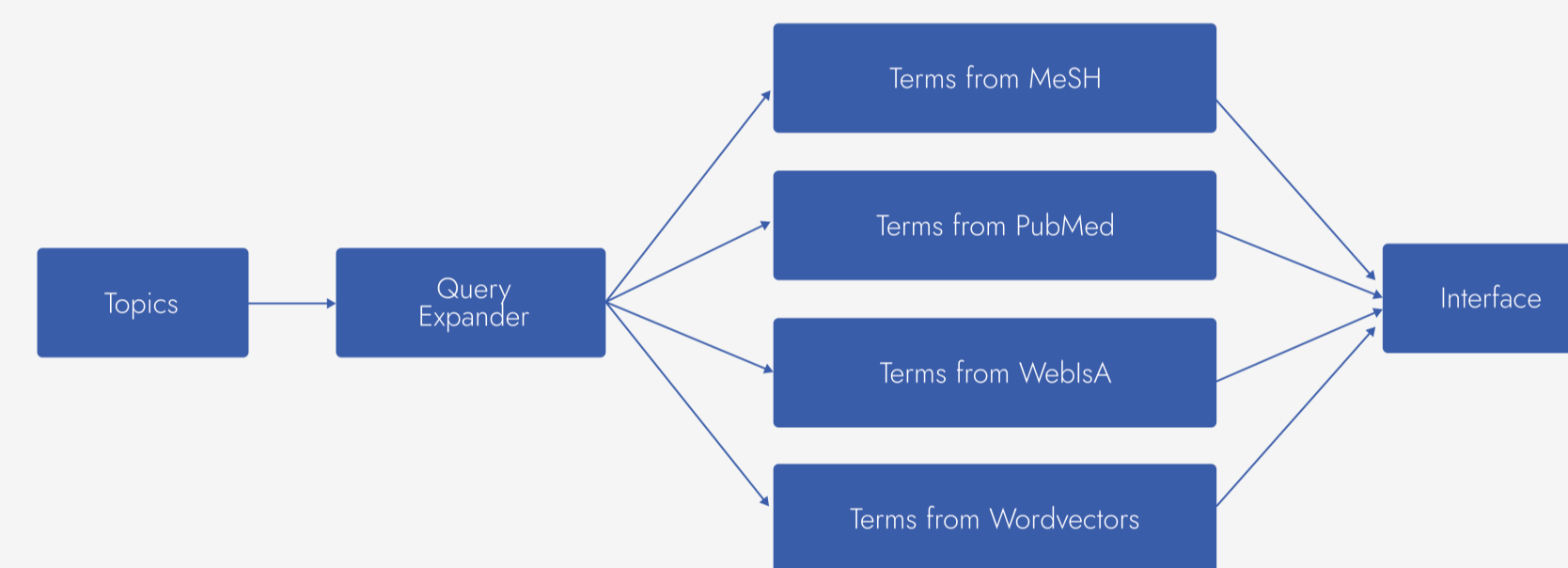
## 3/Methodology

Participants were asked to complete 20 academic search tasks on different topics, and they were free to search with the help of query suggestions until they were satisfied with the results. The performance of the sources was evaluated by the number of query suggestions generated and the selection rate of the participants.

### Sources



Source1 **MeSH**  Source2 **PubMed**  Source3 **WebIsA**  Source4 **Wordvectors**
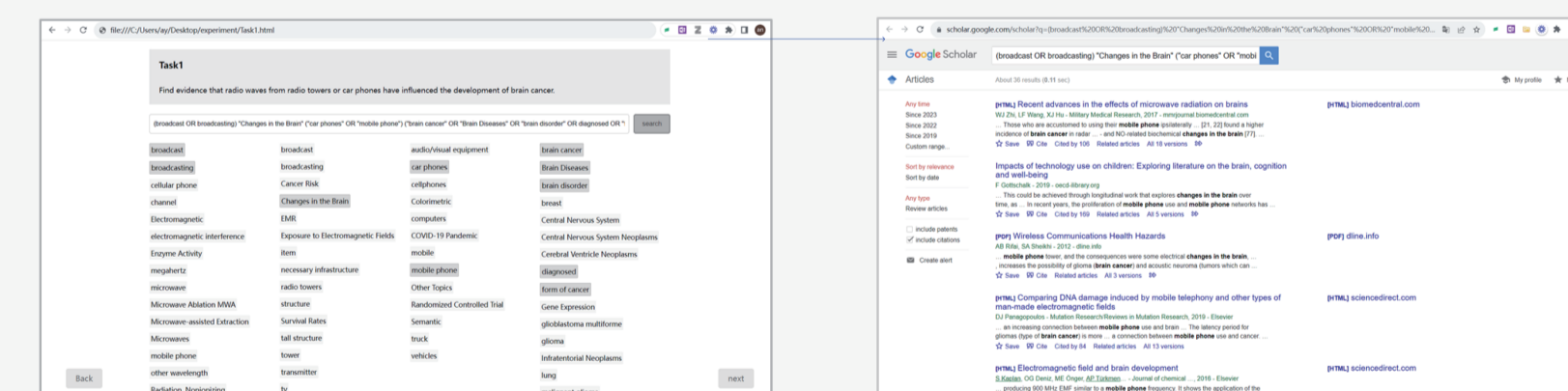
### Terms Generation

I generated query suggestions from different sources based on keywords on the same topic and put them on the interface where the sources were not visible to the participants. Participants conducted academic searches with the help of query suggestions provided by the study. The performance of the sources was evaluated by the number of query suggestions generated and the selection rate of the participants.
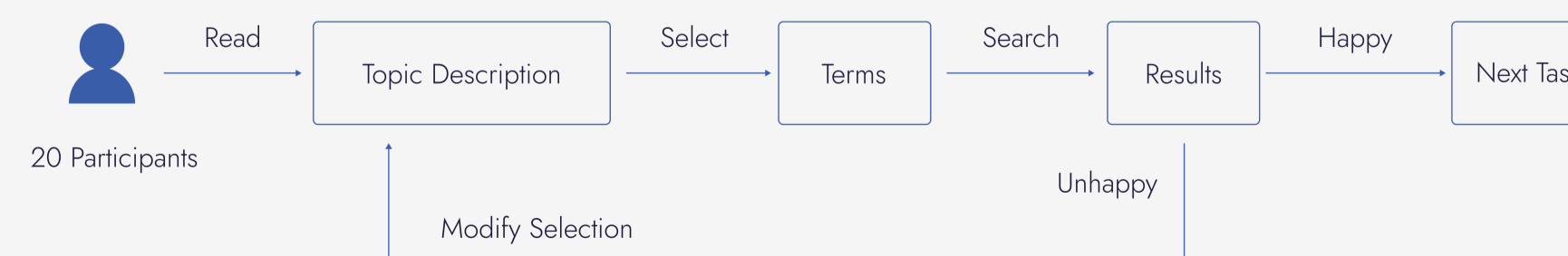


### Interface

The experiment has 22 pages. Participants performed the 20 main tasks after passing the example task on the first page. Each task was on a single page. Participants clicked to select terms that were combined by Boolean operators and displayed in the search bar. Through the Google Scholar API, the content of the search bar is searched with the Google Scholar search engine.



### Operations

Academic searches are usually multiple queries to obtain a single piece of expected information. Therefore this experiment allows for multiple searches to obtain more accurate results.
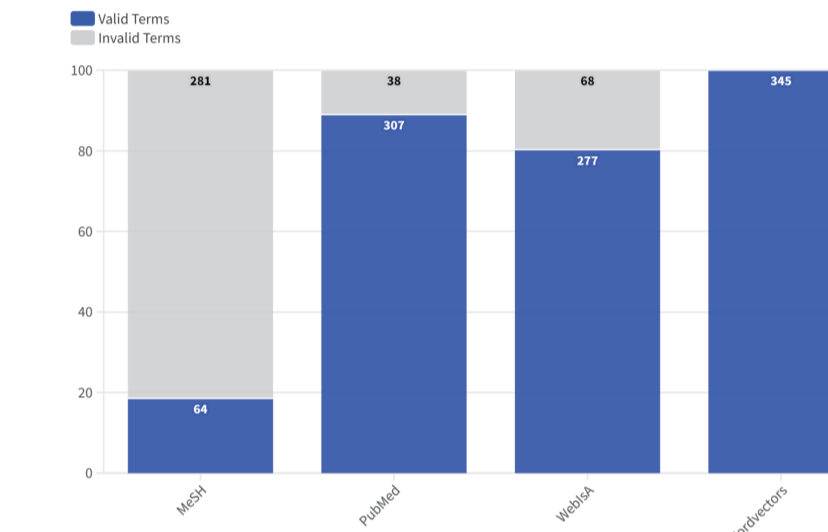


## 4/Results

The number of valid terms and the selection rate show the performance of the different sources from the aspect of quantity and quality of query suggestions. The more valid terms and the higher the selection rate, the better the performance.
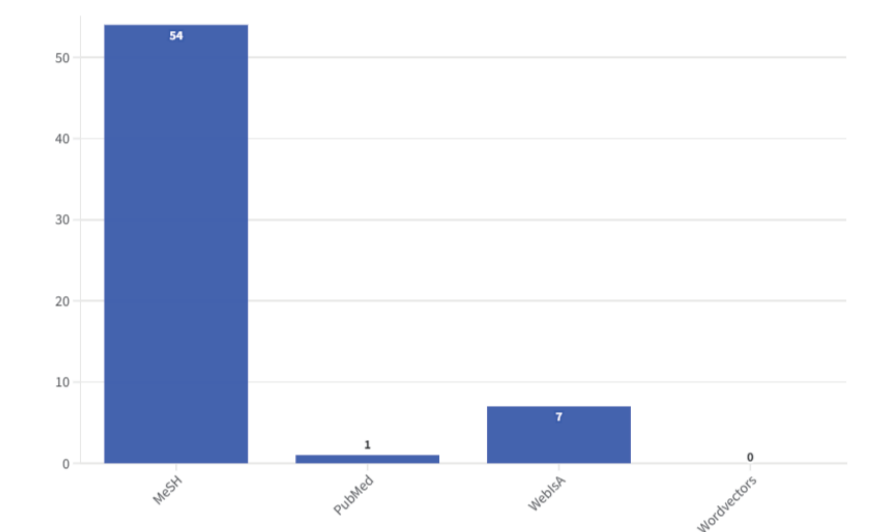
### Valid Terms

The number of valid terms objectively reflects the performance of the source.Wordvectors performed best, generating a sufficiently high number of query suggestions for each keyword. Out of a total of 69 keywords, 54 were not able to get query suggestions through MeSH.The number of valid terms for PubMed and WebIsA was close and in an acceptable range.
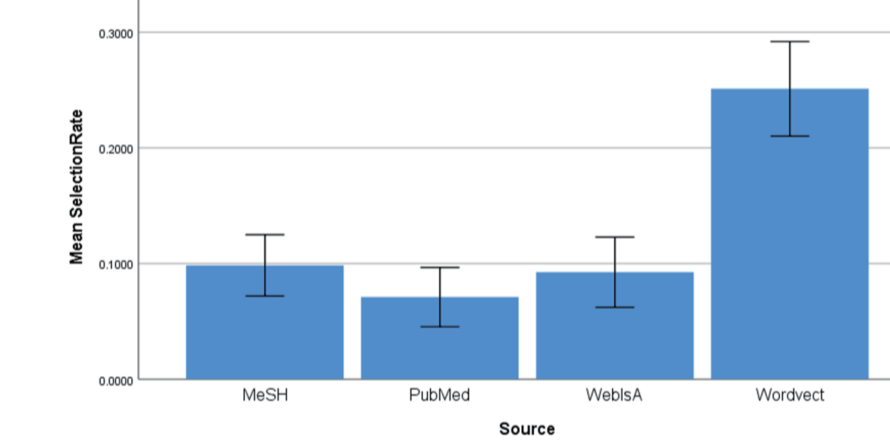


### Selection Rate

The selection rate reflects the quality of valid terms. A higher selection rate indicates a higher quality of terms from the source. As can be seen from the figure, the selection rate of Wordvectors is significantly different from that of all three other sources, and the average selection rate is much higher. Although MeSH generates fewer terms, it has the second highest selection rate following Wordvectors. But there is no significant difference between PubMed, MeSH and WebIsA.



Wordvectors had high selection rates in general and could be used in a wide range of topics. PubMed and WebIsA also had similar and average overall performance. MeSH had very high selection rates in a few tasks and zero selection rates in nearly half of the tasks. Furthermore, MeSH performed well in specific specialised terms, but poorly in generic ones. A few participants chose more terms and spent more time on them, but overall almost all participants tended to select terms from Wordvectors.