

# 605.744: Information Retrieval

## Efficient Phrase Querying with an Auxiliary Index

### Summary

Sabbir Ahmed

October 3, 2022

## Goals and Motivations

Phrase querying refers to querying search engines or other information retrieval systems with a sequence of more than one word. The most frequent number of words in a phrase query is two, and a significant number of phrase queries with more than two words contain stopwords. Since a significant portion of user queries include phrases, the ever-growing search engines need to evaluate such queries extremely fast.

Inverted indexes are typically helpful in querying single words. The use of nextword indexes has been shown to improve support for querying phrases, albeit at the expense of larger storage requirements. This paper proposes a method to optimize querying for phrases by using a combination of inverted indexes and nextword indexes.

## Background

Inverted indexes are the standard method for supporting querying on large corpuses. Typically, they are represented as B-trees of indexes mapped to their postings of triples that contain the document identifier, the term frequency in that document, and the positions of the term in those documents.

Nextword indexes consist of data structures that are relatively more expansive, where each term, or firstword, maps to data structures of its corresponding nextword. Each of the nextwords then maps to their posting lists of where the firstword-nextword pairs occur. The sizes of such structures increase much faster than conventional inverted indexes, easily surpassing twice the storage requirements.

With the added complexity of nextword indexes, it may appear that they may not be necessary to accomplish what inverted indexes can achieve using less space. Phrase querying in inverted indexes does involve additional steps of sorting the lists from the rarest to the most common terms and then merging those lists. However, inverted index files reward stopwords removal for smaller storage requirements and faster retrievals. Stopword removals contribute to a loss of accuracy for phrase querying, where phrases such as “end of days” or “the Who” get removed entirely.

Nextword indexes perform much better with stopwords present in phrases. Their postings lists are typically short since most firstword-nextword pairs only occur infrequently. Using the example from before, the pair “the”-“who” is significantly less frequent than those terms occurring in an inverted index.

Naturally, it appears that there needs to be an implementation where both the space efficiency of inverted indexes and the accuracy and speed of nextword indexes with stopwords can be achieved,

## **Approach**

The authors suggest a combination of inverted indexes with an auxiliary nextword for terms that are considered stopwords. They also propose a constraint on the nextword terms, where only stopwords can be used as firstwords in the pairs. This approach forces the engine to switch its querying method based on individual terms in the phrase; if a term is infrequent enough to not be considered a stopword, then it will be searched in the inverted index for its postings list. On the other hand, if the system encounters a stopword in the phrase, then it will be considered the firstword and its nextword will be searched. Once a match is found for the stopword and the subsequent term, its postings list will be retrieved.

## **Experiments and Results**

## **Discussion**