

Effects of Stop Words Elimination for Arabic Information Retrieval: A Comparative Study

Ibrahim Abu El-Khair

Dept. of Library and Information Science
Faculty of Arts, Minia University-Egypt
iabuelkhair@gmail.com

Abstract: The effectiveness of three stop words lists for Arabic Information Retrieval---General Stoplist, Corpus-Based Stoplist, Combined Stoplist ---were investigated in this study. Three popular weighting schemes were examined: the inverse document frequency weight, probabilistic weighting, and statistical language modelling. The Idea is to combine the statistical approaches with linguistic approaches to reach an optimal performance, and compare their effect on retrieval. The LDC (Linguistic Data Consortium) Arabic Newswire data set was used with the Lemur Toolkit. The Best Match weighting scheme used in the Okapi retrieval system had the best overall performance of the three weighting algorithms used in the study, stoplists improved retrieval effectiveness especially when used with the BM25 weight. The overall performance of a general stoplist was better than the other two lists.

Keywords: Arabic Information Retrieval, Stoplists, Lemur Toolkit.

Received: December 30, 2005 | **Revised:** June 10, 2006 | **Accepted:** December 5, 2006

1. Introduction

Although most of the research in the field of information retrieval has focused on the English language, recently there has been a considerable amount of work and effort to develop information retrieval systems for languages other than English. Research and experimentation in the field of information retrieval in the Arabic language is relatively new and limited compared to the research that has been done in English, which has been dominant in the field of information retrieval for a long while. This is despite the fact that the Arabic Language is one of the five languages of the United Nations, the mother tongue of over 256 million people. In addition, because it is the language of the Qur'an, it is also the second language for many Moslems and Moslem countries around the world [4].

This study attempts to compare the use and effect of stop words for Arabic information retrieval. Using the Lemur Toolkit, a language modelling and information retrieval package (see Methodology for more details), multiple weighting schemes, and three stopword lists

are implemented in order to determine the effect of stop words elimination on an Arabic information retrieval system.

The weighting schemes to be used are the TF*IDF weight, the best match weight (BM25), and the statistical language modelling (KL). Three stop words lists will be created, a general list, a corpus-based list, and a combined list. Although Stemming is an important factor when dealing with Arabic information retrieval, it was not implemented in this study in order to isolate the effect of stop words from any other factor.

2. Related Studies

Stopwords are very common words that appear in the text that carry little meaning; they serve only a syntactic function but do not indicate subject matter. These stopwords have two different impacts on the information retrieval process. They can affect the retrieval effectiveness because they have a very high frequency and tend to diminish the impact of

frequency differences among less common words, affecting the weighting process. The removal of the stopwords also changes the document length and subsequently affects the weighting process. They can also affect efficiency due to their nature and the fact that they carry no meaning, which may result in a large amount of unproductive processing [9]. The removal of the stopwords can increase the efficiency of the indexing process as 30 to 50% of the tokens in a large text collection can represent stopwords [18].

Identifying a stopwords list or a *stoplist* that contain such words in order to eliminate them from text processing is essential to an information retrieval system. Stoplists can be divided into two categories; domain independent stoplists and domain dependent stoplists. They can be created using syntactic classes or using corpus statistics, which is a more domain dependent approach, used for well-defined fields. They can also be created using a combination of the syntactic classes and corpus statistics to obtain the benefits of both approaches.

Fox [6] was the first to create an English stoplist to be used for general text based on word usage in English. He generated a stoplist that consisted of 421 words which was used later with in the Okapi retrieval system. Fox's method in creating the list is the most frequently used method; it is a domain independent approach. The problem with this method is that there are several arbitrary decisions to be taken during the creation of the list, such as the cut-off point. The elimination of some words and addition of others is based on personal judgment, which requires a certain expertise with the language in hand.

There is no general standard stoplist to use in an IR experiment for the Arabic language. The stoplist used in the Lemur Toolkit is the one created by Khoja [8] when she was creating her Arabic stemmer and is relatively short (168 words). This list was used by Larkey and Connell [11] and Larkey, Ballesteros and Connell [10]. Chen and Gey [5] used a list they created by translating an English list and augmenting it with high frequency words from the corpus creating a rather large list, 1,131 words. They do not discuss the effect of the list.

Savoy and Rasolofo's stoplist¹ [17] is a domain dependent list which has three problems. First they use some words preceded by the letter waw "و" which means "and" in 17 words including 11 duplicates. This letter comes in its separate format in a large portion of words in the Arabic language and could precede all the words in the language with no exceptions. A more efficient way to do this is to remove it using a good stemming algorithm. Second, they remove several other single letters with the waw namely the hamza "ء", alef "أ", "إ", "آ", ba "ب", heh "ه". Due to the way the Arabic language is written, these letters can come separately

but they are still a part of the word and removing them changes the word meaning or leave it meaningless, e.g. the word "كُتَاب" which could mean book, writers, or a place for learning has the letter ba' as a single separate letter and when it is removed the word is meaningless. The third problem is that some of the words used in it are not stopwords even though they appeared frequently in the analysis of the corpus statistics, for example, "الولايات" States, "المتحدة" United, "القاهرة" Cairo, etc. In addition, it is a more domain dependent list so it may not be suitable for other collections.

3. Methodology

This study explores the use of stop words and their effect on Arabic information retrieval. It compares the use of three term weighting schemes, and three stoplists. These techniques are examined using a large corpus that was not available before the introduction of Arabic Cross-Language Retrieval at TREC 2001. The evaluation used the Lemur Toolkit with Arabic language capability.

The study evaluates these techniques using the standard recall and precision measures as the basis for comparison. It answers the following question:

- What is the effect of the stoplists on retrieval, i.e. how sensitive is retrieval to the use of stopwords; and which one of the lists, the general, the corpus based, or the combined list is superior to the other?

First, performance of term weighting schemes without elimination of stopwords was compared, and then combinations of weighting schemes, and stoplists were run. Using statistical analysis, the effectiveness of all techniques was evaluated to determine which combination achieves the optimal performance for Arabic language retrieval.

3.1. Data Set

This research used one Arabic test corpus, created in the Linguistic Data Consortium in Philadelphia, also used in the recent TREC experiments. The Arabic Newswire A corpus was created by David Graff and Kevin Walker at the Linguistic Data Consortium [13]. It is composed of articles from the Agence France Presse (AFP) Arabic Newswire. The source material was tagged using TIPSTER style SGML and was transcoded to Unicode (UTF-8). The corpus includes articles from 13 May 1994 to 20 December 2000. The data is in 2,337 compressed Arabic text data files. There are 209 Mbytes of compressed data (869 Mbytes uncompressed) with 383,872 documents containing 76 Million tokens over approximately 666,094 unique words.

3.2. Query Sets and Relevance Judgments

The query set associated with the LDC corpus was created for TREC 2001 and 2002 [12, 20, 21]. It

¹ The stoplists for all the languages are available at <http://www.unine.ch/info/clef>

consists of 75 queries, developed at the LDC by native Arabic speakers and translated to English and French. The relevance judgements for these queries were obtained using assessment pools from different runs at TREC 2001 and 2002, and using the top 70 documents from each run with an average size of 910 documents for each pool. For TREC 2001, the average number of relevant documents over the 25 queries was 164.9 with five topics having more than 300 relevant documents and another five with fewer than 25 relevant documents [22]. For TREC 2002, the average number of relevant documents over the 50 queries was 118.2 with eight topics having more than 300 relevant documents and 16 topics with less than 25 relevant documents.

3.3. Retrieval Engine

The Lemur Toolkit for Language Modelling and Information Retrieval was used. The results of the experiments were mapped against the relevance judgments that are available for the data set. Standard recall and precision measures were calculated using the *ireval.pl* routine in the Toolkit. The evaluation was based on the use of eleven levels of recall creating the recall/precision matrix.

The Lemur toolkit was chosen for several reasons. It supports the construction of basic text retrieval systems using language modelling methods, as well as traditional methods such as those based on the vector space model and Okapi. It is available on the Web as open source software written in C and C++, and runs on both UNIX and Windows (NT). It was developed by collaboration between the Computer Science Department at the University of Massachusetts and the School of Computer Science at Carnegie Mellon University.

Arabic language capability was recently added to the Toolkit by Leah Larkey. This addition has solved the problem of the availability of an Arabic retrieval system for research and experimentation. The Toolkit uses Windows CodePage 1256 encoding (CP1256).

The Toolkit comes equipped with *TextQueryRetMethods* that implement a basic TF*IDF vector space model, Okapi, and a language modeling method using the Kullback-Leibler similarity measure between document and query language models. Initially these algorithms were used without stemming or stoplists. Each of these algorithms requires several parameters to be set in order to function properly. The parameters that were used were the default parameters set in Lemur. Several unofficial runs were conducted in an attempt to tune these parameters to the collection in hand but this did not improve the results over the defaults in Lemur.

These parameters are as follows: TF*IDF parameters: $K1 = 1$, $B = 0.3$; BM25 Parameters: $K1 = 1.2$, $B = 0.75$, $K3 = 7$; KL-Divergence Model parameter: The KL

model needs only one parameter that is used for the smoothing algorithm applied in it. This study uses the simple KL model with the Dirichlet Prior smoothing algorithm and the Dirichlet Prior parameter was set to 2000 as a typical value for that parameter since currently there is no good way of estimating it [23].

3.4. Stoplists

A general stoplist was created, based on the Arabic language structure and characteristics without any additions. All possible words or articles that may be considered a stopword were collated from the different syntactic classes in Arabic in a systematic way to ensure the completeness of the list.

The word categories [1, 2] that were used are:

- Adverbs.
- Conditional Pronouns.
- Interrogative Pronouns.
- Prepositions.
- Pronouns.
- Referral Names/Determiners.
- Relative Pronouns
- Transformers (verbs, letters).
- Verbal Pronouns.
- Other.

Choosing a word from any of these categories was based on a personal judgment. Not all the words under these categories were used, as some of them were not considered stopwords.

The resulting list consisted of **1,377 words**. The list was checked against Khoja [8] and Alshehri's [3] lists, and two Standard English lists, the Okapi and SMART lists. The reason for having a large number of stopwords is due to the characteristics of the Arabic language. First, in Arabic several letters can be used as prefixes and may change the meaning of the word. These letters are ("أ", "ب", "ف", "ك", "ل"), and they were used on some of the words, not all of them, that they could be used with. Second, a considerable number of the original words from these categories could be joined together and used as suffixes or prefixes, especially the pronouns. Finally, the conjunction letter WAW meaning "and" could be used in the same way but was not because it could be used for all Arabic words with no exception and it would not be realistic to use it.

In order to test the effect of a corpus-based stoplist, a second list was created. A cut-off point determining a certain number of words at which the list will stop, was decided based on the corpus statistics. Words occurring more than 25,000 times were used to create this list. Preliminary examination of the corpus word-frequency statistics showed that this is a reasonable number to use as the cut-off point, even though it may appear to be an arbitrary decision.

Under this condition, 359 words occurred with a frequency of more than 25,000. Then doing a manual check to remove any content bearing word, which

may not be considered a conventional stopword, from this list. The removal of these words is another arbitrary decision based on personal judgment, and the reason for doing this is that there are no clear rules on stopwords list creation. The resulting list contained **235 words**.

A third stoplist created using both the general and corpus-based stoplist. Combining both lists resulted in a list of **1,529 words**. Of these 83 words were in common between the two lists.

3.5. Experimental Setup

The data and the query set for the experiments were processed as follows.

- The 383,872 files in the data set were converted from *UTF-8* format to Windows Code Page 1256 encoding (*CP1256*) for Lemur compatibility.
- The queries were converted from the *ASMO 708* encoding to *CP1256* encoding.
- Title and description for each of the 75 queries were extracted from the original query set.
- Several fatal spelling errors in the queries were corrected.
- The table Normalization function in the Light-10 stemmer in Lemur was implemented for all the runs regardless of the techniques that were used.
 - The letters (ı, !, ı) were replaced with (i).
 - The final (ç) was replaced with (c).
 - The Final (ö) was replaced with (o).
- In order to avoid confusion each technique was given a code to represent it throughout the experiments:
 - Term Frequency Weighting Scheme: *TFIDF*.
 - Okapi Weighting Scheme: *BM25*.
 - Language Modeling (Kullback-Leibler Divergence Model): *KL*.
 - General Stoplist: *GS*.
 - Corpus-Based Stoplist: *CBS*.
 - Combined Stoplist: *CS*.

Combinations of the techniques were coded starting with the weighting scheme, and then the stoplist. For instance, *BM25_CBS* represents a combination of the Okapi Weighting Scheme, and the Corpus-Based Stoplist.

3.6. Evaluation

The performance of each technique was evaluated using the standard measures of *Recall* and *Precision*, [9, 15]. A total of 27 runs were carried out where each run represents one or a combination of more than one of these techniques. The raw output results obtained from the *RetEval* application in Lemur were processed with the *ireval.pl* script to give recall and precision. The script does a TREC-style evaluation and the output includes:

- Total number of relevant documents.
- Total number of relevant documents retrieved.
- Average non-interpolated precision.
- Interpolated precision over the eleven levels of recall.
- Non-interpolated precision at document cut-off levels.
- Breakeven point (exact) precision.

3.7. Data Analysis

Retrieval results were analyzed by calculating the differences between the Recall and Precision scores, and plotting them in the R-P graph. The *Friedman Two-Way ANOVA test* and the *Wilcoxon Matched-Paired Signed-Rank test* were used for judging whether measured differences between different methods can be considered statistically significant or not.

Hull [7] has examined the validity of different statistical techniques that are used in comparing retrieval experimental results. He states that there are two non-parametric alternatives to the t-test that make no assumptions about the shapes of the distributions of the two variables, the *Wilcoxon Matched-Paired Signed-Rank test* and the sign test.

The sign test looks only at the sign of the difference, ignoring its magnitude. If one method performs better than the other far more frequently than would be expected on average, then this is strong evidence that it is superior. The Wilcoxon test replaces each difference with the rank of its absolute value. These ranks are then multiplied by the sign of the difference, and the sum of the ranks for each group is compared to its expected value under the assumption that the two groups are equal.

The reason for choosing the Wilcoxon signed rank test over the sign test is that it is a more powerful and indicative test as it considers the relative magnitude in addition to the direction of the differences considered [19]. It also assumes that as differences between pairs increase, significance also increases [16].

Hull also indicates that when comparing more than two retrieval methods, the *Friedman Two-Way ANOVA* is appropriate. It is a non-parametric equivalent to the One-Way ANOVA that does not require any assumptions. It is used to compare observations repeated on the same subjects, which is the case in hand. It uses the ranks of the data rather than their raw values to calculate the statistic [19].

In this study the *Friedman Two-Way ANOVA test* was used to indicate if there is a significant difference on multiple techniques, then it was followed by the *Wilcoxon Matched-Paired Signed-Rank test* which was used to test pair-wise differences.

4. Results and Data Analysis

This study compared alternative stop word list and their effect on retrieval effectiveness. Six different techniques with a total of 12 different combinations were examined. Results are compared using the Wilcoxon test (see Table 2), and recall and precision curves (figures 1-4).

The Friedman Two-way ANOVA test was used to determine if the differences are statistically significant (see table 1), the test statistic $\chi^2 = 70.471$ and the P-value* = 0.000 which indicates that the differences between techniques as a whole are statistically significant. This was not surprising considering the wide variety of techniques used but the test does not depict individual differences between two techniques.

Table 1: The Friedman Test for All Techniques

| Technique | Mean Precision | Mean Rank |
|--------------|----------------|-------------|
| BM25 | 0.2169 | 4.91 |
| KL_CBS | 0.2217 | 4.95 |
| KL_CS | 0.2243 | 5.76 |
| KL | 0.2264 | 6.04 |
| TFIDF_CBS | 0.2260 | 6.23 |
| KL_GS | 0.2248 | 6.29 |
| TFIDF | 0.2260 | 6.42 |
| TFIDF_GS | 0.2283 | 7.09 |
| BM25_CBS | 0.2433 | 7.11 |
| TFIDF_CS | 0.2285 | 7.3 |
| BM25_CS | 0.2474 | 7.73 |
| BM25_GS | 0.2496 | 8.45 |

In order to explore these differences, the techniques are grouped according to weighting scheme, and stoplists and combinations of them. For each group of retrieval techniques the significance testing starts with the Friedman Two-way ANOVA test to observe the differences between all techniques used are statistically significant. Then, a post-hoc test using the Wilcoxon Matched-Paired Signed-Rank test was performed to determine the difference between paired techniques. Due to the number of techniques used and the number of different combinations of them, it was impractical to do a pair-wise comparison on all of them. To set a baseline for comparison, the raw term weighting techniques were run without any additional techniques and the best of the three was used as the baseline.

4.1. Term Weighting

The results of the term weighting approach show that the three algorithms performed relatively well considering the difficulties of the Arabic language and the fact that no linguistic adaptation for it was implemented during retrieval.

Although the BM25 and the KL model are known to perform well compared to the TFIDF weight, in the current study, the overall performance of TFIDF weight was better than the performance of both the BM25 and KL model (see figure 1). The good performance of TFIDF is due to the way the term frequency portion of the weight is calculated in the Lemur Toolkit, using the TF function from the BM25 scheme, which improves its performance significantly.

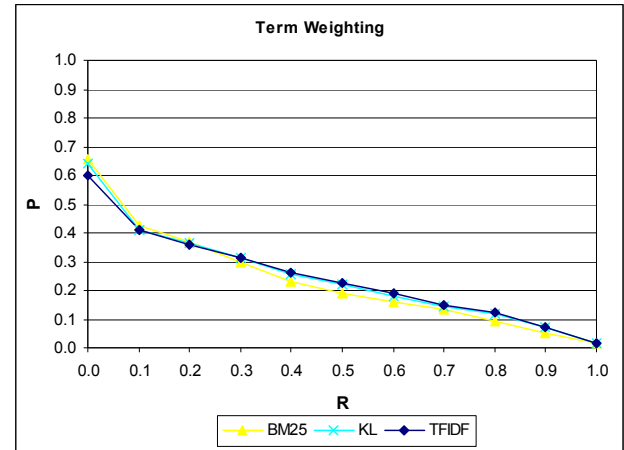


Figure 1: Recall-Precision Curves for Term Weighting

In the Friedman test, $\chi^2 = 5.946$; P-value = .051 and the P-value indicates that the differences between these three techniques are not statistically significant. The Wilcoxon test was used to determine if the pair-wise differences are statistically significant. Because the TFIDF had the best overall performance, it was used as the baseline for comparisons. The Wilcoxon test confirms the result of the Friedman test results, the differences between TFIDF and BM25, and TFIDF and KL were not statistically significant.

4.2. Term Weighting and Stoplists

This subsection presents the results obtained by using the three stoplists that were created for this study. The stoplists were created assuming that they would improve the retrieval efficiency when used with other techniques. The results illustrate how sensitive retrieval is to the use of stopwords. The stopwords will essentially affect the term weights used as they have a significant effect on the term frequency.

4.2.1. General Stoplist

After creating the list, it was initially used in combination with term weighting. The results obtained when using the general stoplist are presented in figure 2. The test statistic for the Friedman test $\chi^2 = 14.517$ and the P-value is .002. This indicates that the differences between these runs and the baseline precision are significant.

* The P-Value for both Wilcoxon and Friedman tests was set on the .05 level.

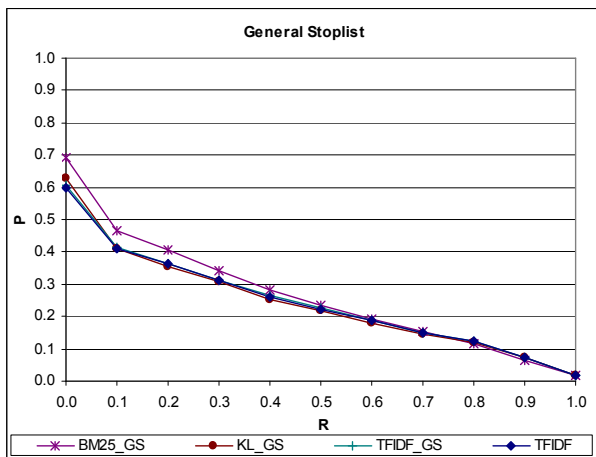


Figure 2: Recall-Precision Curves for the General Stoplist

The differences in the mean precision were minimal, but the increase in precision that was made by the list was apparent at low cut-off levels, especially for the BM25. The Wilcoxon test indicates that the differences between the KL_GS and the baseline precision were not statistically significant. There was a minimal improvement for only one query with the KL model. As for the BM25_GS, and TFIDF_GS there was a significant difference and change. After combining the GS with BM25 the results changed drastically going up from 30 queries better than the baseline precision to 47 queries, and the same thing happened with the TFIDF_GS run with 49 queries favoring it over the baseline precision. Both results indicate how sensitive these weights are, especially the BM25, to the use of stopwords, bearing in mind that the term frequency (TF) portion in the TFIDF is calculated using the BM25 term frequency function. Conversely, the KL model had poor performance with the stoplist as the results slightly deteriorated when it was combined with the stoplist.

4.2.2. Corpus-Based Stoplist

The list has improved the precision for the BM25 weight on the lower levels of recall. Comparing these results with the baseline precision, the test statistic for the Friedman test χ^2 is 12.957 and the P-value is .005. This indicates that the differences between these runs and the baseline precision are significant. The differences in the mean precision were minimal.

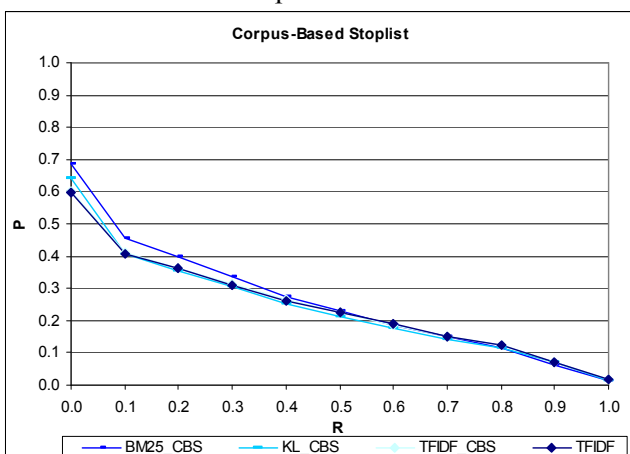


Figure 3: Recall-Precision curves for the Corpus-Based Stoplist

The Wilcoxon test indicates that the differences between the BM25_CBS, TFIDF_CBS and the baseline precision were not statistically significant, even though there was an apparent improvement with the BM25. Even though there was a slight improvement with the KL-Model (not more than 2.7 %), the corpus-based stoplist had a negative effect on the overall performance of the model as the results degraded from 31 queries in favor of the KL-Model to 25 when it was combined with the stoplist.

4.2.3. Combined Stoplist

Figure (4) presents the results obtained from using the combined stoplist. In this figure, the curves show that the results were also very close, as for the general stoplist, and that the list has improved the precision for the BM25 weight on the lower levels of recall. Comparing these results with the baseline precision, the test statistic for the Friedman test χ^2 is 13.327 and the P-value is .004. This indicates that the differences between these runs and the baseline precision are significant. The differences in the mean precision were minimal.

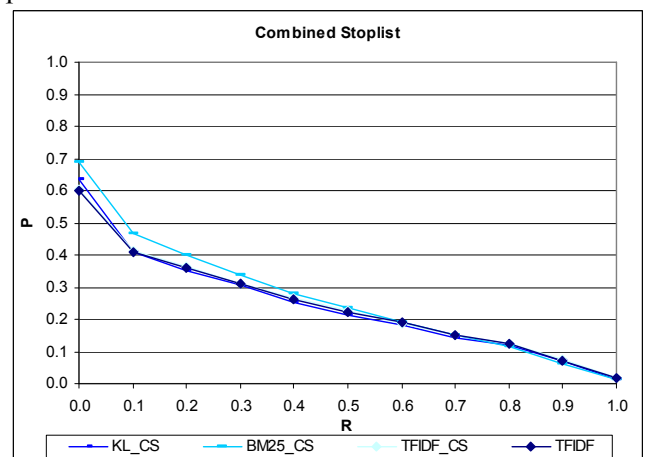


Figure 4: Recall-Precision curves for the combined Stoplist

The Wilcoxon test indicates that the differences between the BM25_CS, KL_CS, and the baseline precision were not statistically significant. However there was an improvement for the BM25 weight. Compared to the general stoplist, there were some differences but the results were almost identical to the general stoplist combinations despite the additions to it. Looking at the individual queries the differences in precision were in favor of the general stoplist, showing that the improvement in both was very much due to it. Even though there was an improvement in the KL-Model at low document cut-off levels, the overall performance of the model also had poor performance with the stoplist. The results deteriorated when it was combined with the stoplist. When combined with the TFIDF the combined list had a better overall performance than the baseline precision and the difference was significant.

5. Discussion

Using the standard recall and precision measures the above techniques were compared. Six techniques were used separately and combined, generating a total of 12 different indexing approaches.

Without any additional linguistic processing the three schemes, TF*IDF weighting, Okapi best match algorithm, and the Kullback-Leibler Divergence Model, had a good performance with Arabic which was not surprising considering their previous success with other languages as they depend only on the corpus and query statistics. The differences between the three weighting schemes were very minimal and not statistically significant.

The TF*IDF scheme is the best weighting scheme to be used with the Arabic language when used separately without stemming or stopwords removal. This contradicts some previous research indicating that the BM25 algorithm is better when used with Arabic [17].

One reason for this is that the term frequency portion of the TF*IDF scheme in Lemur is calculated using the term frequency portion in BM25 giving it the advantages of both of schemes. A second is that in Savoy and Rasolofo's experiment the BM25 was combined with a stemmer which particularly boosts the results with Arabic language.

Three stoplists were created for this study, a general stoplist, a corpus-based stoplist, and a combined stoplist. The assumption here was that stopwords affect the retrieval process but the extent of this effect was not known. The results showed that this effect varies from one weighting scheme to the other. Combining the stoplists with the TFIDF and KL model did not make a substantial difference, only 0.1%-0.2% increase in the former and decrease in the latter. When the stoplists were combined with the BM25 weight there was noticeable improvement of 7.67% with the corpus-based stoplist, 9.49% with the combined stoplist, and 10.44% with the general stoplist.

Table 2: The Wilcoxon Test for All Runs

| Technique | Mean Precision | QP > BP | QP < BP | QP = BP | P-Value |
|-----------|----------------|---------|---------|---------|---------|
| KL CBS | 0.2217 | 25 | 49 | 0 | 0.009 |
| KL CS | 0.2243 | 29 | 45 | 0 | 0.107 |
| BM25 | 0.2169 | 30 | 45 | 0 | 0.093 |
| KL | 0.2264 | 31 | 43 | 1 | 0.287 |
| KL_GS | 0.2248 | 32 | 42 | 0 | 0.229 |
| TFIDF CBS | 0.226 | 35 | 39 | 1 | 0.668 |
| TFIDF CS | 0.2285 | 45 | 29 | 1 | 0.037 |
| BM25 CBS | 0.2433 | 46 | 29 | 0 | 0.123 |
| BM25 CS | 0.2475 | 46 | 29 | 0 | 0.056 |
| BM25_GS | 0.2496 | 47 | 28 | 0 | 0.018 |
| TFIDF_GS | 0.2283 | 49 | 25 | 0 | 0.028 |

- QP: Query Precision.
- BP: Baseline Precision (TFIDF).

The results illustrate how sensitive the BM25 weight and KL model are to the use of stopwords. The use of stopwords has positively affected the BM25 weight while associating it with the KL model has affected the results negatively. The corpus-based list was the lower than the General list, which suggests that we should revisit the corpus-based list. Unfortunately there are no clear-cut rules on how to create a list like this and most of the decisions that were taken in creating this list were arbitrary.

Generally the overall performance of the general stoplist was better than the corpus-based stoplist and to some extent better than the combined stoplist. The list can be used as a standard list for Arabic retrieval regardless of the nature of the data used. The list will be added to the Lemur Toolkit making it available for research and further development as there are no publicly available stoplists for Arabic language.

As for the Lemur Toolkit, one of the main objectives of this study was to use it in experimenting with Arabic language retrieval. The addition of the Arabic language to the Lemur Toolkit will benefit the language as it facilitates the retrieval process whatever the approach that is followed. A major advantage of the toolkit is that it is open source software making it easy to add or modify applications. During this study few applications of the toolkit were used but it proved to be very efficient when used to work with the Arabic language in terms of time, the capability to handle the language, and the ease of use.

6. Conclusions and Further Research

Experimentation with Arabic language retrieval is still a relatively new area of research; it still requires exploring and more research. In this study several retrieval techniques and their potential in improving retrieval effectiveness were explored. The effects of term weighting, and stopwords on Arabic retrieval were examined and compared using the Lemur Toolkit.

The best match algorithm, BM25, with the combined or general stoplist was the best performing function for retrieval in the Arabic language. The performance of a general stoplist or a combined list was relatively close. The use of any of them is recommended but the general stoplist is certainly preferred if we are dealing with a different corpus.

The Kullback-Leibler Divergence Model had problems performing with stopwords. Further investigation with the model could reveal the extent of this problem especially when dealing with different smoothing algorithms and variant query lengths.

Developing a new weighting algorithm that could use the characteristics of the Best Match algorithm, BM25, and combining it with language specific characteristics may lead to further improvements. For instance, using the syntactical structure of the Arabic sentence in calculating the term weight with the BM25 could improve upon the efficiency of this algorithm.

Lemur comes equipped with several applications; this study has used only a small percentage of these applications. Experimenting with the other applications provided by the toolkit to determine their performance in Arabic is another area that could be explored, for instance the use of feedback in retrieval, summarization...etc.

References

- [1] M. Abdul-Rauf, *Arabic for English Speaking Students*, Al- Saadawi Publications, 1996
- [2] A. Al-Raghi, *Grammatical Application*, Dar El-Nahda Al-Arabia Lelteba'ah Wa Al-nashr, 1981.
- [3] A. Alshehri, Optimization and Effectiveness of N-Grams Approach for Indexing and Retrieval in Arabic Information Retrieval Systems, Ph. D Thesis, School of Information Sciences University of Pittsburgh, 2002
- [4] B. Brunner (Ed.), Time Almanac 2005 with Information Please, <http://www.infoplease.com/ipa/A0775272.html>, 2005
- [5] A. Chen, F. Gey, Translation Term Weighting and Combining Translation Resources in Cross-Language Retrieval, *Tenth Text REtrieval Conference (TREC 2001)*, http://trec.nist.gov/pubs/trec10/papers/berkeley_trec10.pdf, 2002.
- [6] C. Fox, A Stop List for General Text. *SIGIR Forum*, Vol. 24, No. 1-2, 1990, pp.19-35.
- [7] D. Hull, Using Statistical Testing in the Evaluation of Retrieval Experiments, *Proceedings of the 16th annual international ACM SIGIR conference on Research and Development in Information Retrieval*, 1993, pp. 329-338.
- [8] S. Khoja, R. Garside, Stemming Arabic Text, <http://www.comp.lancs.ac.uk/computing/users/khoja/stemmer.ps>, 1999.
- [9] R. R. Korfhage, *Information Storage and Retrieval*, John Wiley, 1997.
- [10] L. S. Larkey, L. Ballesteros, M. E. Connell, Improving Stemming for Arabic Information Retrieval: Light Stemming and Co-occurrence Analysis, *Proceedings of the 25th Annual International AC SIGIR Conference on Research and Development in Information Retrieval*, 2002, pp. 275-282.
- [11] L. S. Larkey, M. E. Connell, Arabic Information Retrieval at UMass in TREC-10, *Tenth Text REtrieval Conference (TREC 2001)*, http://trec.nist.gov/pubs/trec10/papers/UMass_TREC10_Final.pdf, 2002
- [12] LDC, TREC 2002 Arabic Topic Development and Assessment, http://www ldc.upenn.edu/Projects/TREC/TR_EC_2002/index.html, 2002a.
- [13] LDC, Arabic Newswire Part 1, <http://www ldc.upenn.edu/Catalog/LDC2001 T55.html>, 2002b
- [14] Lemur Toolkit Team, Lemur Toolkit Overview & Documentation, <http://www.lemurproject.org>, 2003
- [15] G. Salton, C. Buckley, Term- Weighting Approaches in Automatic Text Retrieval. *Information Processing & Management*, Vol. 24, No. 4, 1998, pp.513-523.
- [16] G. Salton, M. McGill, *Introduction to Modern Information Retrieval*, McGraw-Hill Book Company, 1983.
- [17] J. Savoy, Y. Rasolofo, Report on the TREC-11 Experiment: Arabic, Named Page and Topic Distillation Searches. *Eleventh Text REtrieval Conference (TREC 2002)*, http://trec.nist.gov/pubs/trec11/papers/uneuch_atel.pdf, 2003.
- [18] P. Schauble, *Multimedia Information Retrieval: Content-based Information Retrieval from Large Text and Audio Databases*, Kluwer Academic Publishers, 1997.
- [19] S. Siegel, N. J. Castellan, *Nonparametric Statistics of the Behavioral Sciences*, 2nd edition, McGraw-Hill Book Company, 1988.
- [20] TREC, Query Set for Arabic Retrieval at the TREC-2001 Cross-Language Information Retrieval Track,

http://trec.nist.gov/data/topics_noneng/Arabic_topics.txt (Arabic);
http://trec.nist.gov/data/topics_noneng/english_topics.txt (English), 2002a

- [21] TREC, Relevance Judgments for Arabic Retrieval at the TREC-2001 Cross-Language Information Retrieval Track,
http://trec.nist.gov/data/qrels_noneng/xlingual_t10qrels.txt, 2002b
- [22] E. M. Voorhees, D. Harman, Overview of TREC 2001. *Tenth Text Retrieval Conference, TREC 2001*, 2002, pp. 13-16.
- [23] C. Zhai, Dirichlet Prior Smoothing Parameter. Personal Communication, 2003



Dr. Ibrahim Abu El-Khair

Received in 2000, the MLIS, from the School of Library and Information Science (currently School of Information Studies) - University of Wisconsin, Milwaukee, USA. In 2003, he received a Ph.D. in Information Science, from the

School of Information Sciences, University of Pittsburgh, USA. Currently he is an assistant professor at the Department of Library and Information Science, Faculty of Arts, Minia University, Egypt.

Appendix A. General Stoplist

| | | | | | | | |
|---------|----------|----------|---------|----------|---------|---------|---------|
| بأى | بألتين | إنكم | أى | أمامنا | اولئك | الغير | انها |
| بأي | باللذان | إنكما | أي | أمامه | اي | القادم | اتناء |
| بأياً | باللذين | إنما | أياً | أمامها | اياه | اللاتي | اجل |
| بأية | باللواتي | إننا | أيان | أمامهم | ايضا | اللاحق | احدا |
| بأيها | بالنسبة | إنني | أية | أمامهما | اين | الللتان | احدى |
| بأيهم | | إنه | أيضا | أمامهن | ايها | الللتين | احيانا |
| بأيهما | بامكان | إنها | أيضاً | أمامي | آخر | للذان | اخرى |
| بأيهن | بان | إنهم | أين | أن | أبدا | للذين | اخيرا |
| بأحدى | باناه | إنهما | أينما | أنا | أثناء | | اذ |
| بإذا | باولئك | إنهن | أيها | أنت | أجل | اللواتي | اذا |
| بإلا | بآخر | إني | أيهم | أنتم | أحد | المقبل | اذن |
| بإياك | بأحد | إياك | أيهما | أنتما | أحياناً | الممكن | ازاء |
| بإياكم | بأشياء | إياكم | أيهن | أنتن | أخرى | المنصر | استمرا |
| بإياكما | بأقل | إياكما | إحدى | أنك | أخيرا | م | ر |
| بإياكن | بألا | إياكن | إذ | أنكم | أخيراً | النحو | اصبح |
| بإياه | بأن | إياه | إذا | أنكما | أزاء | الى | اصبحت |
| بإياها | بأنا | إياها | إذا | أنكن | أشياء | اليه | اكثر |
| بإياهم | بأنك | إياهم | إزاء | أنما | | اليها | الا |
| بإياهما | بأنكم | إياهما | إطلاقاً | أننا | أصبح | اليهم | الان |
| بإياهن | بأنكما | إياهن | إلا | أنني | أقل | اما | الآن |
| بإياى | بأنكن | إياى | إلى | أنه | أكثر | امام | الأمام |
| ببضع | بأننا | بئس | إلي | أنها | ألا | امس | الأمر |
| ببضعة | بأنني | بالأمام | إليك | أنهم | ألست | ان | الأن |
| ببعض | بأنه | بالأمر | إليكم | أنهما | ألستم | انا | الإطلاق |
| ببعضها | بأنها | بالإضا | إليكما | أنهن | ألستما | انت | البعض |
| ببعضهم | بأنهم | فة | إليكن | | ألستن | انتم | التي |
| بتلك | بأنهما | بالتالي | إلينا | أنى | ألسن | انك | التي |
| بحديث | بأنني | بالتأكيد | إليه | أنى | أليس | انكم | الجاري |
| بدلاً | بأوآخر | بالتى | إليها | أو | أليست | انكن | الحالي |
| بدون | بأولئك | بالذي | إليهم | أوآخر | أليسوا | انما | الخ |
| بدوننا | بأولاء | بالذين | إليهما | أولئك | أم | اننا | الذان |
| بدونه | بأولائك | بالرغم | إليهن | أولا | أما | انني | الذى |
| بدونها | بأولائكم | بالضبط | إما | أولاء | أمام | انه | الذي |
| بدونهم | بأولائكم | بالغير | | أولائك | أمامك | انهم | الذين |
| بدونهما | ا | بالقول | إن | أولائكم | أمامكم | انهما | الرغم |
| بدونهن | بأولائك | باللاتي | إننا | أولائكما | أمامكما | انهن | السابق |
| بذا | ن | باللتان | إنك | أولائكن | أمامكن | او | السواء |

| | | | | | | | |
|---------|---------|---------|--------|--------|--------|--------|----------|
| بذاك | بك | بينكن | حولها | دون | شيئان | عليكم | غيره |
| بذلك | بكافة | بينما | حولهم | دونك | شيئين | عليكما | غيرها |
| بذو | بكل | بيننا | حولهما | دونكم | ضدك | عليكن | غيرهم |
| بذي | بكم | بينه | حولهن | دونكما | ضدكم | علينا | غيرهما |
| برغم | بكما | بينها | حولي | دوننا | ضدكما | عليه | غيرهن |
| بسبب | بكن | بينهم | حيث | دونه | ضدكن | عليها | غيري |
| بسوى | بكيف | بينهما | حيثما | دونها | ضدنا | عليهم | فاذ |
| بشأن | بل | بينهن | حين | دونهم | ضده | عليهما | فاذا |
| بشكل | بلا | بيني | حينئذ | دونهما | ضدها | عليهن | فاكثر |
| بشيء | بلى | تحتة | حيناً | دونهن | ضدهم | عما | فالآن |
| بشيئاً | بما | تقريباً | حينذاك | ذا | ضدهما | عن | فالآن |
| بشيئان | بماذا | تقريباً | حينما | ذات | ضدهن | عنا | فالتي |
| بشيئين | بمتى | تقول | حينه | ذاتك | ضدي | عند | فالذي |
| بصورة | بمزيد | تكن | حينها | ذاتكما | ضدين | عندئذ | فالذين |
| بضع | بمزيداً | تكون | خارجاً | ذاته | ضرورة | عندك | فالغير |
| بضعة | بمفرده | تكونوا | خاصاً | ذاتها | ضرور | عندكم | فالقول |
| بعد | بمن | تلك | خاصة | ذاتهم | ي | عندكما | فالاتي |
| بعندئذ | بن | تلكم | خصو | ذاتهما | ضرور | عندما | فاللتان |
| بعدة | بنا | تلكما | صاً | ذاتهن | ياً | عنده | فاللتين |
| بعدم | بنحو | تماماً | خصيد | ذاك | ضمن | عندها | فاللذان |
| بعدها | بنسبة | ثم | صا | ذلك | طالما | عندهم | فاللذين |
| بعض | به | ثمة | خلا | ذلكم | طويل | عندهما | فاللواتي |
| بعضاً | بهؤلاء | جدا | خلال | ذلكما | طويلاً | عندهن | فان |
| بعضها | بها | جداً | خلاله | ذو | طويلة | عناك | فانك |
| بعضهم | بهاتان | جيدا | خلف | ذي | طويله | عنكم | فاننا |
| بغض | بهاتين | حاشا | خلفك | ربما | ظل | عنكما | فانه |
| بغير | بهذا | حالما | خلفكم | رغم | عام | عنم | فانها |
| بغيرك | بهذان | حالياً | خلفكما | رغماً | عامة | عنه | فانهم |
| بغيركم | بهذه | حالياً | خلفكن | رقم | عبر | عنها | فاولئك |
| بغيركما | بهذي | حنماً | خلفنا | سواء | عدا | عنهم | فأحد |
| بغيركن | بهذين | حتى | خلفه | سواءاً | عدة | عنهما | فأقل |
| بغيرنا | بهل | حسب | خلفها | سوف | عدم | عنهن | فأكثر |
| بغيره | بهم | حوالي | خلفهم | سوى | عدمه | عني | فألا |
| بغيرها | بهما | حول | خلفهما | شانه | عديدة | غير | فأما |
| بغيرهم | بهن | حولك | خلفهن | شأنه | عسى | غيرك | فأن |
| بغيرهما | بين | حولكم | خلفي | شنتي | على | غيركم | فأنا |
| بغيرهن | بينك | حولكن | دائماً | شيء | علي | غيركما | فأنت |
| بغيري | بينكم | حولنا | دائماً | شيئاً | عليّ | غيركن | فأنتم |
| | بينكما | حوله | داخلاً | شيئاً | عليك | غيرنا | فأنتما |

| | | | | | | | |
|----------|----------|---------|---------|---------|----------|---------|---------|
| فأنتن | فانها | فبك | فحيث | فحنأ | ففوقه | فكثير | فلديهن |
| فأنه | فإنهم | فبكل | فحيثما | فحنذ | ففوقها | فكثيراً | فلذا |
| فأنهم | فإنهما | فبكم | فحين | فحنذذ | ففوقهم | فكذلك | فلذاك |
| فأنى | فانى | فبكما | فحينئذ | فحنذك | ففوقهما | فكل | فلذلك |
| فأولئك | فيايك | فبكن | فحيناً | فحنذكم | ففوقهن | فكلا | فلذي |
| فأولاء | فيايكم | فبما | فحينذاك | فحنذكما | ففي | فكلانا | فلست |
| فأولائك | فيايكما | فبماذا | فحينما | فحنذما | ففيك | فكلاهما | فلستم |
| فأولائكم | فيايكن | فبنا | فحينه | فحنذه | ففيكم | فكلتا | فلستما |
| فأولائكم | فيايه | فبنسبة | فحينها | فحندها | ففيكن | فكلكم | فلستن |
| ا | فياها | فبهؤلاء | فخلا | فحندهم | ففيما | فكلنا | فلسن |
| فأولائك | فياهم | فبها | فخلال | فحندهما | ففيها | فكله | فلسوف |
| ن | فياهما | فبهاتان | فدائماً | فحندهن | ففيه | فكلها | فلعدم |
| فأى | فياهن | فبهاتين | فذا | فحنك | ففيها | فكلهم | فلعل |
| فأيان | فياى | فبهذا | فذاك | فحنكم | ففيهم | فكلهن | فلقد |
| فأين | فبئس | فبهذان | فذلك | فحنكما | ففيهما | فكلينا | فلك |
| فأينما | فبالتى | فبهذه | فذو | فحنه | ففيهن | فكليهما | فلكل |
| فإذ | فبالذي | فبهذين | فذي | فحنها | فقبل | فكم | فلكلا |
| فإذا | فبالذين | فبهم | فسواء | فحنهم | فقد | فكما | فلكلتا |
| فإلا | فبالغير | فبهما | فسواء | فحنهما | فقدماً | فكي | فلكم |
| فإلى | فبالقول | فبهن | فسوف | فحنهن | فقط | فكيف | فلكما |
| فإلى | فباللاتى | فبين | فسوى | فحنى | فقلت | فكيلا | فلكن |
| فإليك | فباللتان | فبينك | فطالما | فغير | فقول | فلا | فلكنك |
| فإليكم | فباللتين | فبينكم | فعدا | فغيرك | فكالتى | فلأحد | فلكنه |
| فإليكما | فباللذان | فبينكما | فعدة | فغيركم | فكالذي | فلأنه | فلكنهم |
| فإليكن | فباللذين | فبينكن | فعدم | فغيركما | فكالذين | فلأولئك | فلكنهما |
| فإلينا | فباللوات | فبينما | فعلما | فغيركن | فكالقول | فلأحدى | فلكنهن |
| فإليه | ي | فبيننا | فعلى | فغيرنا | فكاللاتى | فلإنه | فلكى |
| فإليها | فبالنسبة | فبينه | فعلياً | فغيره | فكاللتان | فلبئس | فلكليلا |
| فإليهم | فبالولئك | فبينها | فعليك | فغيرها | فكاللتين | فلتلك | فلم |
| فإليهما | فبالألا | فبينهم | فعليكم | فغيرهم | فكاللذان | فلدى | فلما |
| فإليهن | فبالولئك | فبينهما | فعليكما | فغيرهما | فكاللذين | فلدى | فلماذا |
| فإما | فبتلك | فبينهن | فعلیکن | فغيرهن | فكاللوا | فلديك | فلماذا |
| فإن | فبحيث | فبيني | فعلينا | فغيري | تى | فلديكم | فلن |
| فإننا | فبذا | فتحت | فعليه | ففوق | فكان | فلديكما | فلنا |
| فإنك | فبذاك | فتاك | فعلوها | ففوقك | فكانك | فلدينا | فله |
| فإنكم | فبذلك | فثم | فعليهم | ففوقكم | فكانه | فلديه | فلهؤلاء |
| فإنكما | فبذى | فجأة | فعليهما | ففوقكما | فكانهم | فلديها | فلها |
| فإننا | فبعد | فجأة | فعليهن | ففوقكن | فكانهما | فلديهم | فلهاتان |
| فإنه | فبعده | فحاشا | فعن | ففوقنا | فكانهن | فلديهما | فلهاتين |

| | | | | | | | |
|---------|----------|---------|---------|----------|--------|--------|----------|
| لإياهما | لأمامي | كنت | كإياكن | كالتى | فهم | فممكن | فلهتان |
| لإياهن | لأن | كنتم | كإياه | كالذى | فهما | فمعنا | فلهتين |
| لإياى | لأنا | كنتما | كإياها | كالذين | فهنا | فمعه | فلهذا |
| لبئس | لأنك | كهؤلاء | كإياهم | كالقول | فهناك | فمعها | فلهذان |
| لبعض | لأنكم | كهاتين | كإياهما | كاللتي | فهو | فمعهم | فلهذه |
| لنتك | لأنكما | كهذا | كإياهن | كاللتان | فهو | فمعهما | فلهذين |
| لدى | لأنكن | كهذه | كإياى | كاللتين | فهى | فمعهن | فلهم |
| لدى | لأننا | كهذي | كبيرا | كاللذان | فوق | فمعي | فلهما |
| لديك | لأنني | كهذين | كتك | كاللذين | فوقك | فمما | فلهن |
| لديكم | لأنه | كونه | كثير | كاللواتي | فوقكم | فمن | فلو |
| لديكما | لأنها | كونها | كثيرا | كان | فوقكما | فمنا | فلولا |
| لدينا | لأنهم | كونوا | كثيراً | كانا | فوقكن | فمنذ | فلولاك |
| لديه | لأنهما | كي | كذا | كانت | فوقنا | فمنك | فلولاكم |
| لديها | لأنى | كيف | كذاك | كانتا | فوقه | فمنكم | فلولاكما |
| لديهم | لأواخر | كيلا | كذلك | كانوا | فوقها | فمنكما | فلولاكن |
| لديهما | لأولئك | لئلا | كذو | كأحد | فوقهم | فمنكن | فلولانا |
| لديهن | لأولاء | لا | كسوى | كان | فوقهما | فمننا | فلولاه |
| لذا | لأولئك | لابد | كغير | كانك | فوقهن | فمنه | فلولاها |
| لذاك | لأولائكم | لان | ككل | كانكم | فى | فمنها | فلولاهم |
| لذلك | لأولائكم | لانه | كل | كاننا | في | فمنهم | فلولاهما |
| لذو | ما | لانها | كلا | كانه | فيك | فمنهما | فلولاهن |
| لذي | لأولئك | لانهم | كلانا | كانها | فيكم | فمنهن | فلولاى |
| لست | ن | لاولئك | كلاهما | كانهم | فيما | فمني | فليس |
| لستم | لأى | لاي | كلتا | كانهما | فيما | فمهما | فليست |
| لستما | لأى | لآخر | كلكم | كانهن | فيه | فحن | فليسوا |
| لستن | لأياً | لأحد | كلما | كانى | فيها | فهؤلاء | فما |
| لسن | لأية | لأمام | كلنا | كأولاء | فيهم | فهاتان | فماذا |
| لسوف | لأيها | لأمامك | كله | كأولائك | فيهما | فهاتين | فما عدا |
| لسوى | لأيهم | لأمامكم | كلها | كأولائك | فيهن | فهأنت | فمتى |
| لعدم | لأيهما | لأمامكم | كلهم | م | فيومئذ | فهأنتم | فمثل |
| لعل | لأيهن | ا | كلهن | كأولائك | قبل | فهأنذا | فمثلاً |
| لغى | لإحدى | لأمامكن | كلينا | ما | قبله | فهتان | فمثلاً |
| لغير | لإياك | لأمامنا | كليهما | كأولائك | قبلها | فهتين | فمدام |
| لقد | لإياكم | لأمامه | كم | ن | قد | فهذا | فمدة |
| لك | لإياكما | لأمامها | كما | كأى | قديماً | فهذان | فمع |
| لكل | لإياكن | لأمامهم | كماذا | كإحدى | قريباً | فهذه | فمعاً |
| لكلا | لإياه | لأمامهم | كمن | كإياك | كافة | فهذي | فمعك |
| لكلنا | لإياها | ا | كن | كإياكم | كافياً | فهذين | فمعكم |
| لكم | لإياهم | لأمامهن | كنا | كإياكما | كالآن | فهل | فمعكما |

| | | | | |
|---------|---------|--------|---------|---------|
| ورائكن | مننا | مثلا | لهاتين | لكما |
| ورائهم | منه | مثلاً | لهتان | لكن |
| ورائهما | منها | مثلما | لهتين | لكنك |
| ورائهن | منهم | مثله | لهذا | لكننا |
| يا | منهما | مثلها | لهذان | لكنه |
| يبدو | منهن | مثلهم | لهذه | لكنها |
| يكن | مني | مدة | لهذي | لكنهم |
| يكون | مهما | مدى | لهذين | لكنهما |
| يكونوا | نحن | مرة | لهم | لكنهن |
| يلي | نظرا | مزيد | لهما | لكني |
| يمكن | نعم | مزيداً | لهن | لكي |
| يمكنه | هؤلاء | مطلقاً | لو | لكيلا |
| يومئذ | هاتان | مع | لولا | للأمم |
| | هاتين | معا | لولاك | للأمر |
| | هاذين | معاً | لولاكم | للتي |
| | هامة | معظم | لولاكم | للذي |
| | هأنت | معك | لولاكن | للذين |
| | هأنتم | معكم | لولانا | للغاية |
| | هأنذا | معكما | لولاه | لللاتي |
| | هذا | معكن | لولاها | لللتان |
| | هذان | معنا | لولاهم | لللتين |
| | هذه | معه | لولاهما | للذان |
| | هذي | معها | لولاهن | للذين |
| | هذين | معهم | لولاى | للواتي |
| | هكذا | معهما | لي | للمزيد |
| | هل | معهن | ليس | لم |
| | هم | معي | ليست | لما |
| | هما | مم | ليسوا | لماذا |
| | هن | مما | ليكون | لمدة |
| | هنا | ممكناً | مؤكداً | لمذا |
| | هناك | ممكناً | ما | لمزيد |
| | هنالك | ممن | مادام | لمزيداً |
| | هو | من | ماذا | لمن |
| | هي | منا | مازال | لن |
| | وراء | منذ | مازالت | لنا |
| | وراءه | منك | ماعد | له |
| | ورائك | منكم | ماهو | لهؤلاء |
| | ورائكم | منكما | متى | لها |
| | ورائكما | منكن | مثل | لهاتان |

Appendix B. Corpus-Based Stoplist

| | | | | | |
|-------|--------|-------|--------|----------|----------|
| هو | كانون | شخصا | ايام | الدولية | ابر |
| هي | كأس | شرق | ايضا | الذي | ابو |
| وح | كبير | صباح | أن | الذي | اجتماع |
| وزارة | كرة | صحيفة | باسم | الذين | اجل |
| وزير | كل | صفر | بان | السابق | احد |
| وكالة | كما | ضمن | برس | الساعة | اخرى |
| يتبع | لا | عام | بسبب | السبت | اذا |
| يجب | لدى | عاما | بشكل | السلطات | ارا |
| ينكر | لقاء | عبد | بطولة | السلطة | اربعة |
| يكون | لكرة | عبر | بعد | الشرطة | اط |
| يمكن | لكن | عدة | بعض | العاصمة | اطار |
| ينا | للامم | عدد | بن | العسكرية | اطلاق |
| يول | لم | عدم | به | العمل | اعادة |
| يون | لن | عشر | بها | الف | اعلن |
| يونيو | له | عشرة | بيان | القدم | اغو |
| | لها | على | بين | اللجنة | اف |
| | لوكالة | علي | تشرين | الماضي | افب |
| | مؤتمر | عليه | تم | المباراة | اكثر |
| | ما | عليها | ثلاثة | المتحدث | اكذ |
| | مار | عمان | ثم | المجلس | الا |
| | ماي | عن | جمت | المجموعة | الاتفاق |
| | مايو | عند | جميع | المرحلة | الاثنين |
| | مباراة | عندما | جنوب | المصدر | الاحد |
| | مجموعة | غدا | جهة | المقبل | الاخيرة |
| | مدينة | غير | حال | المقرر | الاراضي |
| | مساء | فان | حاليا | الملك | الاربعاء |
| | مصادر | فبر | حتى | النار | الاسبوع |
| | مصدر | فرانس | حزيران | النهائي | الان |
| | مع | فسب | حول | الوزراء | الانباء |
| | مليون | في | حيث | الوضع | الاول |
| | من | في | حين | الوقت | الاولى |
| | منذ | فيه | خلال | الى | التعاون |
| | منها | فيها | دورة | اليوم | التي |
| | موا | قال | دولار | اما | التي |
| | موسع | قبر | دولة | امام | الثانية |
| | نحو | قبرص | دون | امس | الثلاثاء |
| | نفسه | قبل | ديس | ان | الجمعة |
| | نقطة | قد | ديسك | انه | الجيش |
| | نهاية | قدم | ذكرت | انها | الحالي |
| | نوف | قرار | ذلك | او | الحدود |
| | هان | قمة | رئيس | اوك | الحزب |
| | هذا | قوات | زيارة | اول | الحكم |
| | هذه | كان | سبت | اي | الخميس |
| | هناك | كانت | سنوات | ايار | الدولة |