

Author(s)	Nguyen, Qui V.
Title	Enhancing a Web Crawler with Arabic Search.
Publisher	Monterey, California: Naval Postgraduate School
Issue Date	2010-09
URL	<a href="http://hdl.handle.net/10945/7288">http://hdl.handle.net/10945/7288</a>

This document was downloaded on August 25, 2013 at 12:28:37



<http://www.nps.edu/library>

Calhoun is a project of the Dudley Knox Library at NPS, furthering the precepts and goals of open government and government transparency. All information contained herein has been approved for release by the NPS Public Affairs Officer.

**Dudley Knox Library / Naval Postgraduate School  
411 Dyer Road / 1 University Circle  
Monterey, California USA 93943**



<http://www.nps.edu/>



# **NAVAL POSTGRADUATE SCHOOL**

**MONTEREY, CALIFORNIA**

## **THESIS**

**ENHANCING A WEB CRAWLER WITH  
ARABIC SEARCH CAPABILITY**

by

Qui V. Nguyen

September 2010

Thesis Advisor:  
Second Reader:

Weilian Su  
John C. McEachen

**Approved for public release; distribution is unlimited**

THIS PAGE INTENTIONALLY LEFT BLANK

<b>REPORT DOCUMENTATION PAGE</b>			<i>Form Approved OMB No. 0704-0188</i>	
Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instruction, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188) Washington DC 20503.				
<b>1. AGENCY USE ONLY (Leave blank)</b>		<b>2. REPORT DATE</b> September 2010	<b>3. REPORT TYPE AND DATES COVERED</b> Master's Thesis	
<b>4. TITLE AND SUBTITLE</b> Enhancing a Web Crawler with Arabic Search Capability			<b>5. FUNDING NUMBERS</b>	
<b>6. AUTHOR(S)</b>				
<b>7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)</b> Naval Postgraduate School Monterey, CA 93943-5000			<b>8. PERFORMING ORGANIZATION REPORT NUMBER</b>	
<b>9. SPONSORING /MONITORING AGENCY NAME(S) AND ADDRESS(ES)</b> N/A			<b>10. SPONSORING/MONITORING AGENCY REPORT NUMBER</b>	
<b>11. SUPPLEMENTARY NOTES</b> The views expressed in this thesis are those of the author and do not reflect the official policy or position of the Department of Defense or the U.S. Government. IRB Protocol number ____ N.A. ____.				
<b>12a. DISTRIBUTION / AVAILABILITY STATEMENT</b> Approved for public release; distribution is unlimited			<b>12b. DISTRIBUTION CODE</b>	
<b>13. ABSTRACT (maximum 200 words)</b>  Many advantages of the Internet—ease of access, limited regulation, vast potential audience, and fast flow of information—have turned it into the most popular way to communicate and exchange ideas. Criminal and terrorist groups also use these advantages to turn the Internet into their new play/battle fields to conduct their illegal/terror activities. There are millions of Web sites in different languages on the Internet, but the lack of foreign language search engines makes it impossible to analyze foreign language Web sites efficiently. This thesis will enhance an open source Web crawler with Arabic search capability, thus improving an existing social networking tool to perform page correlation and analysis of Arabic Web sites. A social networking tool with Arabic search capabilities could become a valuable tool for the intelligence community. Its page correlation and analysis results could be used to collect open source intelligence and build a network of Web sites that are related to terrorist or criminal activities.				
<b>14. SUBJECT TERMS</b> Nutch, Lucene, Web Crawler, Information Retrieval in Arabic, Stemming in Arabic			<b>15. NUMBER OF PAGES</b> 121	
			<b>16. PRICE CODE</b>	
<b>17. SECURITY CLASSIFICATION OF REPORT</b> Unclassified	<b>18. SECURITY CLASSIFICATION OF THIS PAGE</b> Unclassified	<b>19. SECURITY CLASSIFICATION OF ABSTRACT</b> Unclassified	<b>20. LIMITATION OF ABSTRACT</b> UU	

NSN 7540-01-280-5500

Standard Form 298 (Rev. 2-89)  
Prescribed by ANSI Std. Z39-18

THIS PAGE INTENTIONALLY LEFT BLANK

**Approved for public release; distribution is unlimited**

**ENHANCING A WEB CRAWLER WITH ARABIC SEARCH CAPABILITY**

Qui V. Nguyen  
Lieutenant, United States Navy  
B.S., University of Texas at Dallas, 1999

Submitted in partial fulfillment of the  
requirements for the degree of

**MASTER OF SCIENCE IN ELECTRICAL ENGINEERING**

from the

**NAVAL POSTGRADUATE SCHOOL  
September 2010**

Author: Qui V. Nguyen

Approved by: Weilian Su  
Thesis Advisor

John C. McEachen  
Second Reader

R. Clark Robertson  
Chairman, Department of Electrical and Computer Engineering

THIS PAGE INTENTIONALLY LEFT BLANK

## **ABSTRACT**

Many advantages of the Internet—ease of access, limited regulation, vast potential audience, and fast flow of information—have turned it into the most popular way to communicate and exchange ideas. Criminal and terrorist groups also use these advantages to turn the Internet into their new play/battle fields to conduct their illegal/terror activities. There are millions of Web sites in different languages on the Internet, but the lack of foreign language search engines makes it impossible to analyze foreign language Web sites efficiently. This thesis will enhance an open source Web crawler with Arabic search capability, thus improving an existing social networking tool to perform page correlation and analysis of Arabic Web sites. A social networking tool with Arabic search capabilities could become a valuable tool for the intelligence community. Its page correlation and analysis results could be used to collect open source intelligence and build a network of Web sites that are related to terrorist or criminal activities.



THIS PAGE INTENTIONALLY LEFT BLANK

# TABLE OF CONTENTS

I.	INTRODUCTION.....	1
A.	MOTIVATION .....	1
B.	RESEARCH OBJECTIVES.....	1
C.	THESIS ORGANIZATION.....	2
II.	ARABIC INFORMATION RETRIEVAL .....	3
A.	INFORMATION RETRIEVAL .....	3
B.	THE CHALLENGES OF ARABIC IR .....	4
C.	RESEARCH IN ARABIC IR.....	5
D.	LIGHT STEMMER ALGORITHM.....	6
	1. Introduction.....	6
	2. The Algorithm .....	6
	a. Normalization.....	6
	b. Light Stemmers .....	7
	3. Results .....	8
E.	CHAPTER SUMMARY.....	9
III.	LUCENE AND NUTCH.....	11
A.	INTRODUCTION.....	11
B.	LUCENE.....	11
	1. Overview .....	11
	2. Indexing Process.....	11
	3. Analyzer .....	13
C.	NUTCH .....	13
	1. Architecture Overview .....	13
	2. Plug-In Architecture.....	14
D.	CHAPTER SUMMARY.....	15
IV.	ARABICANALYZER PLUG-IN DEVELOPMENT .....	17
A.	INTRODUCTION.....	17
B.	REQUIREMENT .....	18
C.	DEVELOPMENT PROCESS.....	18
	1. Implementation of the Light Stemmer Algorithm.....	18
	2. Development of <i>ArabicAnalyzer</i> Plug-in.....	18
	3. Creating Arabic Ngram profile .....	19
D.	CHAPTER SUMMARY.....	20
V.	EXPERIMENTAL SETUP .....	21
A.	PROBLEM STATEMENT .....	21
B.	HARDWARE AND SOFTWARE CONFIGURATION.....	21
C.	METHODOLOGY .....	22
D.	RESULTS AND DISCUSSION .....	23
	1. Terms Count.....	23
	2. Crawl Time.....	24

3.	Search Results .....	25
E.	CHAPTER SUMMARY.....	32
VI.	CONCLUSION .....	33
A.	SUMMARY .....	33
B.	FUTURE WORK.....	33
	APPENDIX A .....	35
	APPENDIX B .....	47
	APPENDIX C .....	57
	APPENDIX D .....	69
	APPENDIX E .....	81
	APPENDIX F .....	91
	LIST OF REFERENCES.....	101
	INITIAL DISTRIBUTION LIST .....	103

## LIST OF FIGURES

Figure 1.	String removed by light stemming. From [14] .....	7
Figure 2.	Monolingual 11-point precision results. From [14] .....	8
Figure 3.	Lucene indexing architecture. From [17] .....	12
Figure 4.	Nutch's architecture. From [18] .....	14
Figure 5.	The process of indexing a Web site .....	17
Figure 6.	<i>ArabicAnalyzer</i> plug-in architecture. From [19] .....	19

THIS PAGE INTENTIONALLY LEFT BLANK

## LIST OF TABLES

Table 1.	The uniterpolated average precision. From [14].....	9
Table 2.	The number of terms counts .....	23
Table 3.	Average crawl time of www.america.gov/ar/ .....	24
Table 4.	Average crawl time of www.bbc.co.uk/arabic/ .....	24
Table 5.	Average crawl time of www.addustour.com .....	24
Table 6.	Average crawl time of www.aawsat.com .....	24
Table 7.	Search terms.....	25
Table 8.	Search results of term “Economy” using ArabicAnalyzer .....	26
Table 9.	Search results of term “Economy” using NutchDocumentAnalyzer .....	27
Table 10.	Search results of term “The United States” using ArabicAnalyzer .....	28
Table 11.	Search results of term “The United States” using NutchDocumentAnalyzer .....	29
Table 12.	Search results of term “Democratic” using ArabicAnalyzer .....	30
Table 13.	Search results of term “Democratic” using NutchDocumentAnalyzer.....	31

THIS PAGE INTENTIONALLY LEFT BLANK

## **EXECUTIVE SUMMARY**

After more than eight years of the War on Terrorism, Improvised Explosive Devices (IEDs) have become the weapon of choice for the terrorist in Iraq and Afghanistan. IEDs accounted for the majority of casualties of Allied forces and civilians. One of the reasons for the proliferation of IEDs is the ease of access to training material available on the Internet. The Internet is a cheap, convenient, yet powerful tool to access a vast reservoir of information and knowledge. Unfortunately, the Internet also empowers technology-savvy terror networks and extremist groups to create IED education networks and distribute the IED know-how to their operatives and supporters quickly and efficiently.

One solution to counter this problem is a social networking tool that applies networking theory and social network analysis to identify terrorist IED education networks quickly. This tool would utilize an open source web crawler that could index Arabic websites into a searchable database for analyzing and querying to collect more actionable intelligence.

The Nutch project was selected as the search engine of choice for this social networking tool. Its transparency ranking information allows the users the ability to tailor the ranking to meet the user's specific requirements. Its versatile plug-in architecture provides extensibility, flexibility and maintainability.

To enable Nutch indexing of Arabic websites, an Arabic language analyzer needs to be added into Nutch's library. Multiple experiments were used to test the performance of the Arabic language analyzer with moderate results.

Overall, Nutch with an added Arabic analyzer would be a valuable tool improving an existing social networking tool to perform page correlation and analysis of Arabic websites. Its results could be used to identify IED education networks and to collect open source intelligence.



THIS PAGE INTENTIONALLY LEFT BLANK

## **ACKNOWLEDGMENTS**

I want to thank my family, Andrea and Rosie, for supporting me. Without their support and encouragement, I would not finish this thesis.

I also would like to thank Dr. Su for his guidance in the past year. His advice was valuable to the completion of this thesis.

THIS PAGE INTENTIONALLY LEFT BLANK

# **I. INTRODUCTION**

## **A. MOTIVATION**

Since its invention, the Internet has revolutionized communication. It enables people to exchange ideas and share information rapidly and cheaply. Unfortunately, its lack of regulation and pervasive communication also has turned it into the new tool for the tech-savvy terrorists: “Today, almost without exception, all major (and many minor) terrorist and insurgent groups have web sites” [1]. Many terror organizations such as Al-Qaeda actively use the Internet to recruit new members, solicit donations from sympathizers, and spread propaganda.

They also turn the Internet into their virtual training grounds, offering tutorials on building IEDs and planning attacks. These training materials are easily accessible to anyone with an Internet connection. This is the main contribution to the explosion of IED attacks in Iraq and Afghanistan. To counter the proliferation of IED technology, these IED education networks need to be identified, monitored and referred to sovereign authorities for further action as necessary.

One possible solution for this problem is a social networking tool that applies network science to identify the IED education network via the World Wide Web. In [2], network science is defined as the study of networks which “contrasts, compares, and integrates techniques and algorithms developed in disciplines as diverse as mathematics, statistics, physics, social network analysis, information science and computer science.” The social network tool would incorporate an open source web crawler that could index Arabic websites into a searchable database for analyzing and querying.

## **B. RESEARCH OBJECTIVES**

The research objectives of this thesis were to enhance a Web crawler engine with Arabic search capability that could index Arabic language websites proficiently, thus improving an existing social networking tool to perform page correlation and analysis of

Arabic websites. The newly enhanced Web crawler could help speed up the analytical process of the social networking tool to effectively identify IED education networks via the World Wide Web.

### **C. THESIS ORGANIZATION**

This thesis consists of six chapters. An overview of the motivation, objectives and thesis organization is provided in Chapter I. A brief discussion about information retrieval, a description of Arabic information retrieval challenges, stemming in Arabic and the light stemmer algorithm is contained in Chapter II. Lucene—a scalable Information Retrieval (IR) library; Nutch—an open source search engine; and Nutch’s plug-in architecture are introduced in Chapter III. The implementation process of the light stemmer algorithm into Lucene’s analyzer database, and development of the *ArabicAnalyzer* plug-in are discussed in Chapter IV. The performance of *ArabicAnalyzer* and *NutchDocumentAnalyzer* are compared in Chapter V. The summary of the thesis and future research recommendation are discussed in Chapter VI.

## **II. ARABIC INFORMATION RETRIEVAL**

### **A. INFORMATION RETRIEVAL**

The fast growth of the Internet accompanied with the explosion of data available via the World Wide Web has made the finding of useful information a tedious and difficult task. These difficulties have attracted renewed interested in Information Retrieval and its techniques.

Information Retrieval (IR) is the science of locating relevant documents in a large collection of documents. The retrieval process is influenced by queries supplied by the user's input, the indexing process and the natural language that is being indexing [3].

In [4], some popular IR classic strategies are the Vector Space Model, Probabilistic Retrieval, Language Model, and Inference Networks.

The Vector Space Model is a widely used retrieval strategy. In this model, both the query and each document are represented as vector in terms of space. A measure of similarity between the two vectors is computed.

In the Probabilistic Retrieval model, a probability based on the likelihood that a term will appear in a relevant document is computed for each term in the collection. For terms that match between a query and a document, the similarity measure is computed as the combination of the probabilities of each of the matching terms [4].

In the Language Model, a language model is inferred for each document; then the probability of generating the query according to each of these models is computed. Documents are then ranked according to these probabilities [5].

Inference Networks, also known as Bayesian networks, are used to model documents, the documents' contents and the query. It then uses this information to derive —“infer”—other relationships. The strength of this inference is then used as the similarity coefficient [4].

## **B. THE CHALLENGES OF ARABIC IR**

According to [6], there are over 200 million native Arabic speakers in the world and over 20 million people speaking it as a second language. There is also an exponential growth of Internet in speaking countries. From [7], the numbers of Internet users in Middle East countries alone have grown from 3 million in 2000 to 58 million in 2009. So, there is increasingly a demand for an Arabic IR as well, but Arabic poses many challenges for IR

First, Arabic has a very complex morphology system. In [8], the authors observed:

Arabic has two genders, feminine and masculine; three numbers, singular, dual and plural; and three grammatical cases, nominative, genitive, and accusative. A noun has the nominative case when it is a subject, accusative when it is the object of a verb, and genitive when it is the object of a preposition.

This would compound the complexity of any Arabic IR to deal with this morphology system.

Second, there are a lot of ambiguities in Arabic. One of the major contributions to this phenomenon is that orthographic variations are widespread in Arabic [9]. The authors gave an example that sometimes in combining HAMZA with ALEF (إ) or MADDA with ALEF (آ), the HAMZA (ء) or MADDA (ـ) is dropped, rendering it ambiguous to whether the HAMZA (ء) or MADDA (ـ) is present. Another contribution to the higher level of ambiguity is that sometimes vowels (diacritics) are omitted in written Arabic, which may change the meaning of the words. This uncertainty would affect the precision and recall of any Arabic IR.

Finally, the plural form of irregular nouns, broken plurals, is common in Arabic. A broken plural's form does not resemble its initial singular form. It does not obey normal morphological rules. Because of that, it is very difficult to design an algorithm to transform this kind of plural to singular form [9].

### C. RESEARCH IN ARABIC IR

Research on Arabic IR has focused on using word roots and stems as index terms. A stem is the remainder of the word after removing prefixes and suffixes. On the other hand, the root is the origin of the word that remains after removing nonessential characters, prefixes and suffixes. When using word roots as index terms, a linguistic knowledge and an understanding of the languages' morphology are needed. On the other hand, prior knowledge of the language is not required when using stems as index terms. In [10], the authors recognized that “stemming is one of many tools besides normalization that is used in information retrieval to combat the vocabulary mismatch problem.” As discussed in section 2b, Arabic is very difficult to stem, therefore, there were only a few available Arabic stemmers.

One of the earliest stemmers was the root-based stemmer proposed by Khoja and Garside. This stemmer removed all the stopwords, punctuation, and numbers. Then it peeled away prefixes and suffixes. After that, it matched the result against a list of patterns to extract the root. Finally, it matched the extracted root against a list of known “valid” roots. There are a few weaknesses in the Khoja stemmer. First, it can provide wrong solutions when removing prefixes and suffixes. It also can generate wrong roots for words that contain *EBDAL* [10], [11], [12].

Buckwalter's morphological analyzer is another useful stemmer. First, this stemmer converts the Arabic word into English letters. Then, it segments it into all probabilities of prefixes, stems, and suffixes. After that, it checks every probability with its build-in lexicon libraries (prefixes dictionary, stems dictionary and suffixes dictionary). If all the word elements (prefix, stem, suffix) are found in their respective libraries, three truth tables indicating their legal combination (prefixes-suffixes, prefixes-stems, and stems-suffixes) are used to determine whether they are compatible. If the word elements pass all three truth tables, the probability is valid. This stemmer provides highly reliable results, but its performance is slow [13].

The light stemmer is another approach for Arabic IR. Most light stemmers in [8], [14] are based on the same idea: extract stems by deleting the most frequent prefixes and



suffixes. These stemmers are not interested in producing the Arabic root. This thesis applies the light stemmer algorithm in [14] to enable the Web crawler with an Arabic search capability. A more detailed discussion is in the next section.

## **D. LIGHT STEMMER ALGORITHM**

### **1. Introduction**

The light stemmer allows for good information retrieval results without providing the correct morphological analyses [10]. Anyone can employ the light stemmer algorithm without the required language skills.

### **2. The Algorithm**

The stemmer has two parts: Normalization and Stemmer. The Normalization process is used to normalize the orthography—the writing system—of the queries and corpus. The stemmer removes suffixes using the light stemmer algorithm to extract the stems [14].

#### ***a. Normalization***

In [14], before stemming, corpus and queries are normalized as follows:

- (1). Convert to Windows Arabic encoding (CP12560).
- (2). Remove punctuation.
- (3). Remove diacritics (primary weak vowels).
- (4). Remove non letters.
- (5). Replace ٱ (ALEF with MADDA above), ْ (ALEF with HAMZA above), and ِ (ALEF with HAZA below) with ا (ALEF)
- (6). Replace final ع (ALEF MAKSURA) with هـ (YEH)
- (7). Replace final ة (TEH MARBUTA) with هـ (HEH)

**b. Light Stemmers**

After the corpus and queries are normalized, the stemmer is applied as follows:

- (1). Remove و (WAW) if the remainder of the word is three or more characters long.
- (2). Remove any of the definite articles if this leaves two or more characters.
- (3). Go through the list of suffixes once in the right to left order indicated in Figure 1, removing any that are found at the end of the word, if this leaves two or more characters.

	Remove from front	Remove Suffixes
Light1	ال، وال، بال، كال، فال	none
Light2	ال، وال، بال، كال، فال، و	none
Light3	“	ه، ة
Light8	“	ها، ان، ات، ون، ين، يه، ية، ه، ة، ي

Figure 1. String removed by light stemming. From [14]

Light1, Light2, Light3 and Light8 apply the same algorithm in the stemming process. The difference between them is the number of prefixes and suffixes that are removed in step 3 of the light stemmer's algorithm. In Light1, the Light Stemmer algorithm only removes five prefixes and no suffixes. In Light2, the Light Stemmer algorithm removes six prefixes and no suffixes. In Light3, the Light Stemmer algorithm removes six prefixes and two suffixes. In Light8, the Light Stemmer algorithm removes six prefixes and 10 suffixes.

### 3. Results

The authors in [14] compared the retrieval effectiveness of the light stemmer algorithm (Light8) and of a morphological analyzer (Khoja stemmer). Raw in Figure 2 means no normalization and stemming. From Figure 2, we see that the light stemmer outperforms Khoja stemmer and raw retrieval. From Table 1, we see that light stemmer improved over 90% in average precision from raw retrieval.

The authors concluded that stemming is very effective on Arabic IR. For monolingual retrieval, the light stemmer has demonstrated improvement of around 100% in average precision due to stemming and related processes.

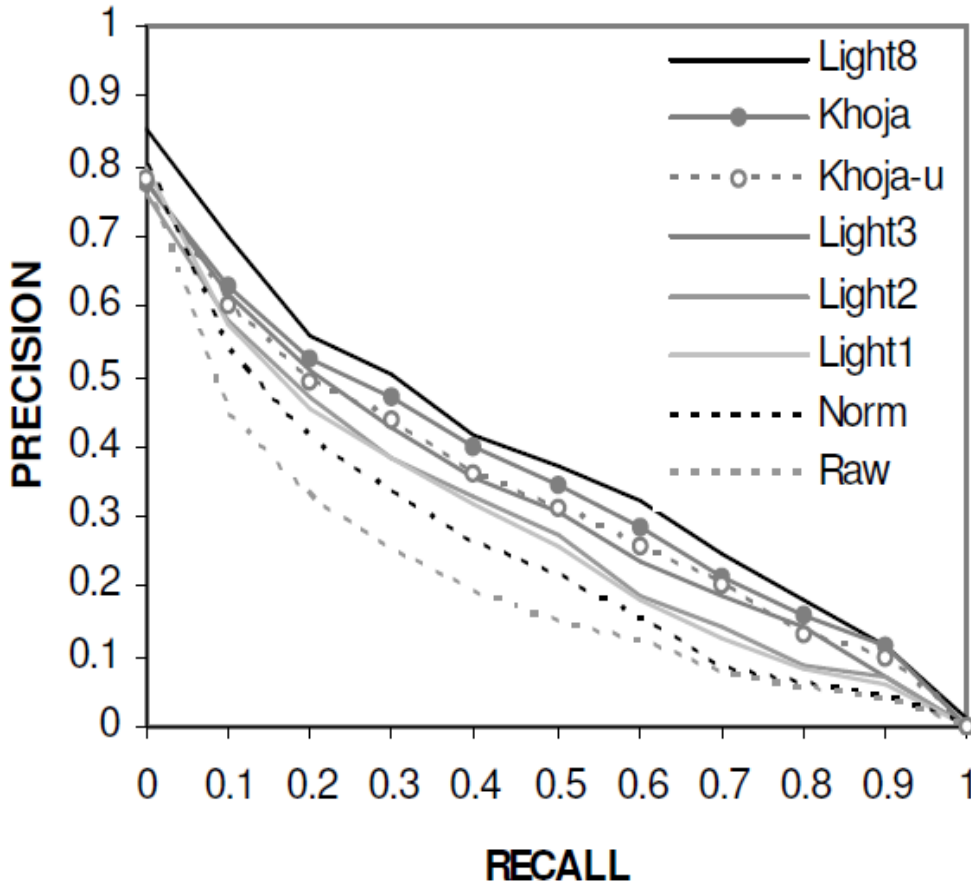


Figure 2. Monolingual 11-point precision results. From [14]

Table 1. The uniterpolated average precision. From [14]

<b>Stemmer</b>	<b>raw</b>	<b>khoja-u</b>	<b>khoja</b>	<b>light8</b>
<b>Av. Precision</b>	.194	.313	.341	.376
<b>Pct. Change</b>		61.7	76.2	94.3

## **E. CHAPTER SUMMARY**

In this chapter, the challenges of Arabic IR and past Arabic IR research were covered. Also discussed was the effectiveness of light stemmer in Arabic IR. In the next chapter, Lucene, Nutch and Nutch's plug-in architecture are introduced.

THIS PAGE INTENTIONALLY LEFT BLANK

### **III. LUCENE AND NUTCH**

#### **A. INTRODUCTION**

Lucene and Nutch, created by Doug Cutting, are two open-source software projects. According to [15], Lucene is a high performance, scalable Information Retrieval (IR) library that provides Java-based indexing and searching technology and advanced analysis/tokenization capabilities. On the other hand, Nutch is a search engine that was built on top of Lucene. Together, they can make a full-featured search engine that offers transparency into how Web sites are ranked, and an understanding of how a large search engine works [16].

#### **B. LUCENE**

##### **1. Overview**

Lucene is a software library that enables users to add indexing and searching capabilities to their application. Lucene can index and search any type of data as long as it can be converted into a text format. This means Lucene can be used to search Web pages, pdf files, and Microsoft® Word files because textual information can be extracted from them. With this feature, Lucene is the best toolkit for a search engine.

##### **2. Indexing Process**

Indexing is the process of converting text into an index, a data structure that improves the speed of data retrieval operations. The index is the fundamental component of Lucene.

From [16], to index data with Lucene, the data must be converted into a stream of plain text tokens, a format that Lucene can process. After that, Lucene prepares the data for indexing by breaking the stream of plain text into chunks or tokens and performing a number of operations on them. For instance, the tokens could be lowercase before

indexing, to make the search case-insensitive. This step is called analysis. After the input has been analyzed, it is ready to be added into the index. The Indexing process is illustrated in Figure 3.

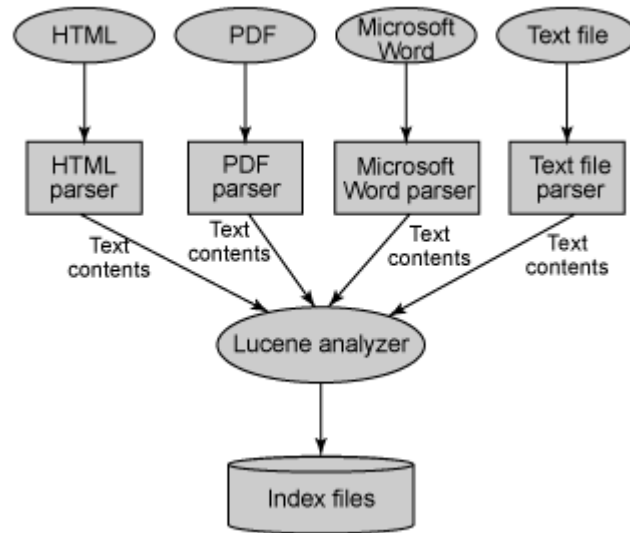


Figure 3. Lucene indexing architecture. From [17]

Lucene implements an innovative approach to maintaining the index—rather than maintaining a single index, Lucene builds multiple index segments and merges them periodically. Using segments allows a quick way to add new documents to the index by adding them to the newly created index segments and only periodically merging them with other existing segments. This process makes additions efficient because it minimizes physical index modifications.

Some IR libraries need to index the whole corpus again when new data is added to their index; Lucene does not need to do that because it supports incremental indexing. This means Lucene allows the contents of newly added documents be searchable immediately without indexing the whole corpus again [15].

### **3. Analyzer**

As discussed above, analysis is a very important step in the indexing process. It converts a field of text into the most fundamental indexed representation, terms. These terms are used to determine what documents match a query during searches.

An analyzer is an encapsulation of the analysis process. The analyzer's job is to process strings of text into a stream of tokens by performing any number of operations on them. Lucene includes several built-in analyzers that do a good job at analyzing English-based text. For analyzing non-English languages, specific language analyzers are needed. Lucene's core API provides building blocks to create custom language analyzers.

## **C. NUTCH**

### **1. Architecture Overview**

Nutch is a complete open-source Web search engine that can operate at one of three scales: local file system, intranet, or whole Web [15]. Nutch can be divided into two parts: the crawler and the searcher.

From [18], components of the crawler are WebDB, the fetch list, fetchers and updates. WebDB is a custom database that tracks every known page and relevant link. It maintains a small set of facts about each page, such as the last crawled date. Fetch lists are generated from WebDB. These lists contain the URLs that users want to download. The fetchers consume the fetch lists to produce the WebDB updates and the Web contents. The updates tell which page has changed since the last crawl. The contents are used to search. The WebDB-fetch cycle is designed to repeat forever, maintaining an up-to-date image of the Web.

Once the Web content is produced, Nutch can get ready to process queries using the searchers. First, the indexer processes the Web content of all terms and pages into an inverted index. The document set is divided into a set of index segments, each of which is fed into a single searcher process. Each searcher also draws upon the Web content



from earlier to provide a cached copy of any Web page. Finally, a pool of Web servers handles the interaction with users and contact searcher for results. A generic overview of Nutch's architecture is shown in Figure 4.

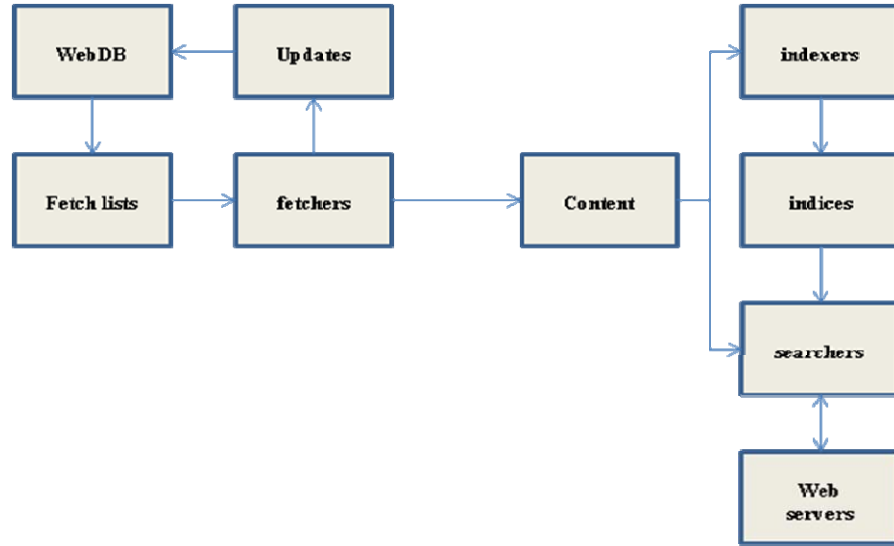


Figure 4. Nutch's architecture. From [18]

## 2. Plug-In Architecture

Nutch's plug-in system is based on the Eclipse 2.0 plug-in architecture. It provides a core service for controlling a set of tools working together to support programming tasks. After reviewing Eclipse's architecture from [19] and applying it to Nutch's plug-in system, we observe that the three most important components of Nutch's plug-in system are Extension, ExtensionPoints and Plug-in. The Extension class provides a way to add some new functions to a plug-in. It is defined by a plug-in that wants to extend its functionality to another plug-in. ExtensionPoints define an interface that must be implemented by the Extension. A plug-in, pluggable component, defines a number of extension-points that may allow it to be augmented by different kinds of extension.

This system is the mechanism of Nutch's extensibility. Users can contribute to the Nutch platform by wrapping their tools in plug-ins. The new plug-ins can add new processing elements to existing plug-ins, and Nutch provides a set of core plug-ins to assist the process.

## **D. CHAPTER SUMMARY**

In this chapter, the overview of Lucene's indexing process and analyzer were examined. The overview of Nutch's architecture and its plug-in system were also studied. In the next chapter, the implementation process of the light stemmer algorithm into Nutch is discussed.

THIS PAGE INTENTIONALLY LEFT BLANK

## IV. ARABICANALYZER PLUG-IN DEVELOPMENT

### A. INTRODUCTION

When Nutch finishes fetching a segment of Web sites, the *language-identifier* plug-in is called to identify the language of the Web sites and attach a language code to those Web sites. After that, the *Analyzerfactory* instantiates the *NutchAnalyzer* interface, which defines an extension point that associates with the specific language code. The *NutchAnalyzer* extension point is an abstract class that extends the Lucene Analyzer class, so that Lucene analyzers can be easily integrated as *NutchAnalyzer* plug-ins. The policy of the *Analyzerfactory* for finding the *NutchAnalyzer* extension to use is to return the first one that matches a specified language code. If none is found, then the default *NutchDocumentAnalyzer* is used. After *Analyzerfactory* identifies the right analyzer basing on the language code, the *NutchAnalyzer* calls the correct analyzer, in this case *ArabicAnalyzer*, from the Lucene analyzer library to index the Web site. The process of indexing a Web site is shown in Figure 5.

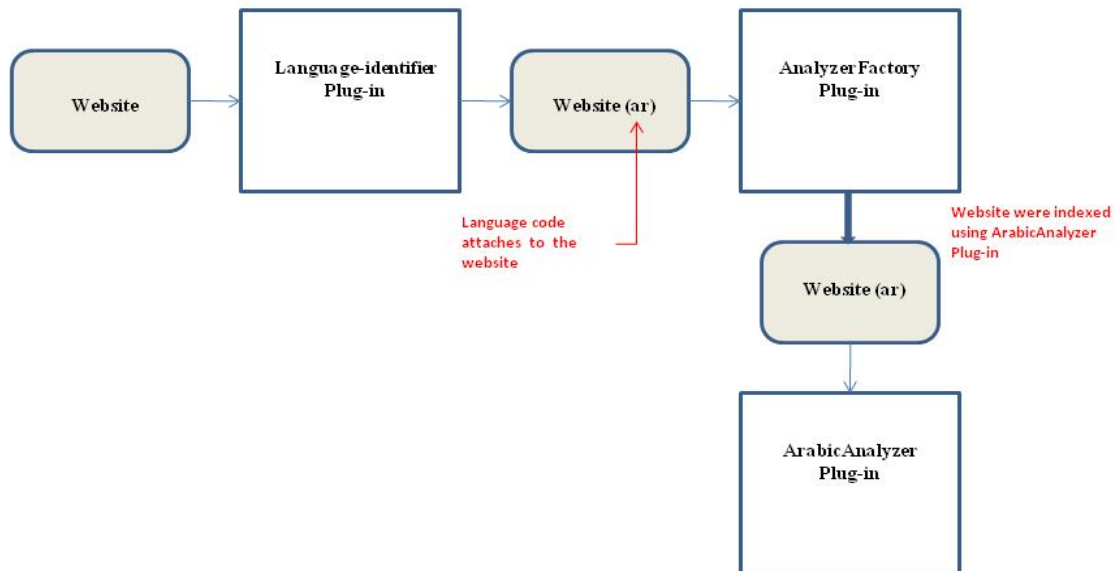


Figure 5. The process of indexing a Web site

## B. REQUIREMENT

To enable Nutch with Arabic-search capability, there are several tasks that need to be completed. First, the Lucene analysis library needs to be updated with the *ArabicAnalyzer* that implemented the light stemming algorithm. Secondly, an *ArabicAnalyzer* plug-in is needed for Nutch to be able to access the Lucene analysis library. Finally, an Arabic Ngram profile is needed to train Nutch how to recognize Arabic text.

## C. DEVELOPMENT PROCESS

### 1. Implementation of the Light Stemmer Algorithm

As stated above, the Lucene analysis library needs to be updated with the *ArabicAnalyzer*, which implements the light stemmer algorithm. The analysis package contains three primary files: *ArabicAnalyzer*, *ArabicNormalizationFilter*, and *ArabicStemFilter*.

The *ArabicAnalyzer* first creates a list of Arabic stop words that is based on the stoplist from <http://members.unine.ch/jacques.savoy/clef/index.html>. It uses the standard *Stopfilter* to filter out all the stop words from the token stream. The result is then fed into *ArabicNormalizationFilter*, which normalizes the orthography. The final result is then fed into the *ArabicStemFilter*, which stems the token stream using the light stemmer algorithm.

### 2. Development of *ArabicAnalyzer* Plug-in

The host plug-in is the *ArabicAnalyzer* class in Nutch. The *NutchAnalyzer*, a Nutch built-in extension point, defines the interface that must be implemented by the Nutch's *ArabicAnalyzer*. The extender plug-in is the *ArabicAnalyzer* from Lucene's analysis library that extends the functions of the Nutch's *ArabicAnalyzer*; in this case, the Lucene's *ArabicAnalyzer* enables the Nutch's *ArabicAnalyzer* to index Arabic text.

Basically, the Nutch's *ArabicAnalyzer* plug-in is a wrapper that sets the stages and makes it possible to run Lucene's *ArabicAnalyzer*. The *ArabicAnalyzer* plug-in architecture that was derived from [19] is shown in Figure 6.

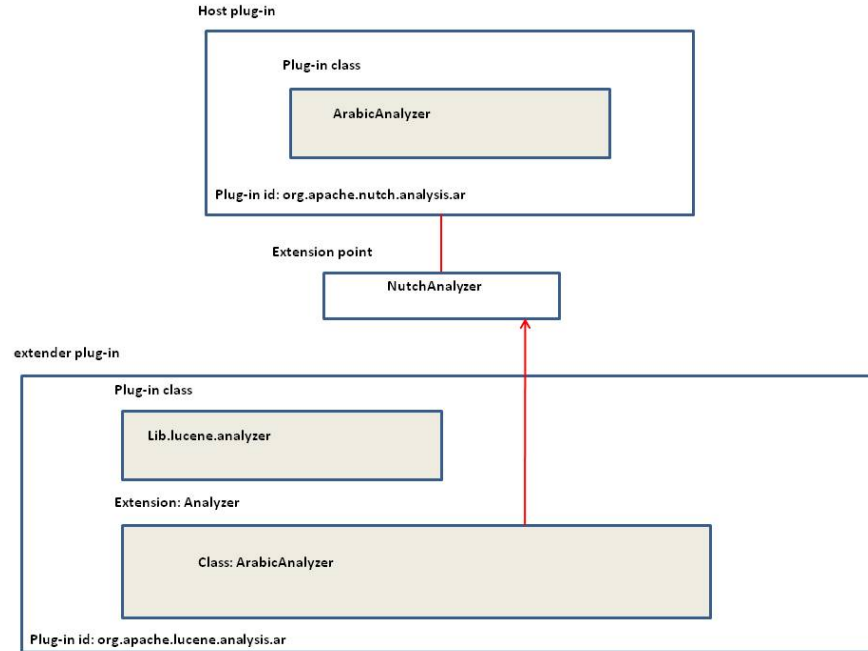


Figure 6. *ArabicAnalyzer* plug-in architecture. From [19]

### 3. Creating Arabic Ngram profile

Nutch uses the *language-identifier* plug-in in standard Nutch's library to create an Arabic profile based on the "1000 most frequent words" by Jacques Savoy from the Web site <http://members.unine.ch/jacques.savoy/clef/index.html>. This trains Nutch to "recognize" Arabic Web sites so that it can invoke the right analyzer to index the Web sites.

## **D. CHAPTER SUMMARY**

In this chapter, the *ArabicAnalyzer* plugin development process is discussed. The Lucene's analyzer library is enhanced with the *ArabicAnalyzer* that implements the light stemming algorithm. The Nutch's plug-in architecture is utilized to create the *ArabicAnalyzer* plug-in. The plug-in enables the Nutch search engine to index Arabic-language Web sites using the *ArabicAnalyzer* in the Lucene's analyzer library. In the next chapter, the performance of *ArabicAnalyzer* and *NutchDocumentAnalyzer* are compared.

## **V. EXPERIMENTAL SETUP**

### **A. PROBLEM STATEMENT**

These experiments will compare the result of Nutch when it used the default *NutchDocumentAnalyzer* with *ArabicAnalyzer* to analyze the same Web site.

The *NutchDocumentAnalyzer* separates the stream of tokens into individual terms without applying any filter. For example, the token stream “hello world” becomes “hello” “world” after *NutchDocumentAnalyzer* processes it. This study uses *NutchDocumentAnalyzer*’s index result as a baseline, because no term is discarded during indexing when using *NutchDocumentAnalyzer* [15].

On the other hand, the *ArabicAnalyzer* applies several filters when analyzing the stream of tokens. First, the token stream goes to *StopFilter*, which removes all the stop words in the custom-built stop words list. The result is then filtered again using *ArabicNormalizationFilter* to normalize the orthography. After that, the result again is filtered using *ArabicStemFilter*, which applies the light stemming algorithm. The final result is then stored into the index database.

### **B. HARDWARE AND SOFTWARE CONFIGURATION**

The platform used to conduct the experiments was a single Dell XPS M1530 laptop personal computer. This machine had an Intel Core 2 Duo CPU T9300 at 2.5 GHz with 4 GB of RAM and a 185 GB hard disk. The operating system used was Microsoft Windows Vista Home Premium with Service Pack 2.

Nutch 1.0 and Lucene 2.4.0 were used to implement the light stemmer algorithm and for all the experiments.



## C. METHODOLOGY

There were three experiments to collect data. The first experiment used Nutch to crawl eight Web sites with the depth of five and topN of 50. TopN determines the maximum number of pages that are retrieved at each level up to the depth. The Web sites are *alarabiya.net*, *aljazeera.net*, *alriyadh.com*, *addustour.com*, *aawsat.com*, *bbc.co.uk/Arabic/*, *arabic.cnn.com* and *america.gov/ar/*. Nutch only indexes the Web pages within these sites using *ArabicAnalyzer* and *NutchDocumentAnalyzer*.

The second experiment computes the average crawl time and its standard deviation. The crawler was set to crawl four out of the eight Web sites above 25 times each.

The third experiment compares the ranking of the top 10 pages after using the two algorithms to search for three different Arabic terms.

To disable *ArabicAnalyzer*, the following code was added into the property block of *nutch-site.xml* file in the *conf* folder so that *AnalyzerFactory* is forced to use *NutchDocumentAnalyzer* to index these sites by not specifying any analyzer:

```
<property>
<name>plugin.includes</name>
<value>protocol-http/urlfilter-regex/parse-(text/html/js)/index-
(basic/anchor)/query-(basic/site/url)/response-(json/xml)/summary-basic/scoring-
opic/language-identifier</value>
<description>Regular expression naming plugin directory names to
include. Any plugin not matching this expression is excluded.
</description>
</property>
```

To enable *ArabicAnalyzer*, the following code replaces the above code within the *nutch-site.xml* file. With *ArabicAnalyzer* on, the *AnalyzerFactory* uses it to index these sites:

```
<property>
<name>plugin.includes</name>
```

```

<value>protocol-http/urlfilter-regex/parse-(text/html/js)/index-
(basic/anchor)/query-(basic/site/url)/response-(json/xml)/summary-basic/scoring-
opic/language-identifier/analysis-ar</value>
<description>Regular expression naming plugin directory names to
include. Any plugin not matching this expression is excluded.
</description>
</property>

```

## D. RESULTS AND DISCUSSION

### 1. Terms Count

The first experiment shows that Nutch needs 20% to 37% fewer terms to index the same number of documents from the same Web site when it uses *ArabicAnalyzer*. The result also means that the *ArabicAnalyzer* plug-in is more efficient when searching its index database, because it searches fewer terms to locate the relevant documents. See Table 2 for the detailed breakdown of each Web site.

Table 2. The number of terms counts

Web sites	NutchDocumentAnalyzer (Terms count)	ArabicAnalyzer (Terms count)
arabic.cnn.com	24776	15827
alarabiya.net	21140	15806
alriyadh.com	20898	13163
aljazeera.net	18096	13658
bbc.co.uk/arabic/	16061	9957
america.gov/ar/	11435	7958
addustour.com	2888	2075
aawsat.com	1050	847

## 2. Crawl Time

The second experiment shows that Nutch takes longer to index the same Web sites when it uses *ArabicAnalyzer*. This result is expected, because there are more filters in *ArabicAnalyzer*: thus, it requires more processing power and time to index Web sites.

The results, as illustrated in Tables 3 to 6, also show that the crawl times fluctuated more when Nutch used *ArabicAnalyzer*.

Table 3. Average crawl time of [www.america.gov/ar/](http://www.america.gov/ar/)

	Average Crawl time (sec)	Standard Deviation (sec)
<b>NutchDocumentAnalyzer</b>	362.92	15.7
<b>ArabicAnalyzer</b>	375.2	25.32

Table 4. Average crawl time of [www.bbc.co.uk/arabic/](http://www.bbc.co.uk/arabic/)

	Average Crawl time (sec)	Standard Deviation (sec)
<b>NutchDocumentAnalyzer</b>	482.76	5.95
<b>ArabicAnalyzer</b>	546.64	37.05

Table 5. Average crawl time of [www.addustour.com](http://www.addustour.com)

	Average Crawl time (sec)	Standard Deviation (sec)
<b>NutchDocumentAnalyzer</b>	104.56	1.67
<b>ArabicAnalyzer</b>	105.12	2.38

Table 6. Average crawl time of [www.aawsat.com](http://www.aawsat.com)

	Average Crawl time (sec)	Standard Deviation (sec)
<b>NutchDocumentAnalyzer</b>	69.56	2.52
<b>ArabicAnalyzer</b>	70.2	2.84

### 3. Search Results

For the third experiment, the index database of the Web site [www.america.gov/ar/](http://www.america.gov/ar/) is used to collect search results data. The terms shown in Table 7 are used for the search.

Table 7. Search terms

Normal Form	Light Stemmer Form	Meaning
الاقتصاد	اقتصاد	Economy
أميركا	اميركا	The United States
الديمقراطية	ديمقراط	Democratic

The Light Stemmer forms are searched using the *ArabicAnalyzer*'s index database and the Normal forms are searched using the *NutchDocumentAnalyzer*'s index database.

When comparing the top 10 pages of the search term “economy,” the top seven pages are the same; for the search term “The United States,” all top 10 pages are the same; and for the search term “Democratic,” six pages are the same but with the ranking different. In all three cases, the search results from NutchDocumentAnalyzer have better ranking scores than the search results from ArabicAnalyzer.

By the title of the search results, one can conclude that their contents are related to the search terms. The two algorithms also hit a high mark on relevance of information that relates to the search terms. See Tables 8 through 13 for the breakdown.

Table 8. Search results of term “Economy” using ArabicAnalyzer

<b>Top 10 pages using ArabicAnalyzer</b>	<b>Score for Query</b>
<a href="http://www.america.gov/ar/econ.html">www.america.gov/ar/econ.html</a>	0.3486507
<a href="http://www.america.gov/ar/publications/books/outline-of-the-us-economy.html">www.america.gov/ar/publications/books/outline-of-the-us-economy.html</a>	0.12422927
<a href="http://www.america.gov/ar/econ/business.html">www.america.gov/ar/econ/business.html</a>	0.09217107
<a href="http://www.america.gov/ar/reviving_trade_ar.html">www.america.gov/ar/reviving_trade_ar.html</a>	0.033118278
<a href="http://www.america.gov/ar/publications/books.html#outline_economy">www.america.gov/ar/publications/books.html#outline_economy</a>	0.016127191
<a href="http://www.america.gov/ar/">http://www.america.gov/ar/</a>	0.003058498
<a href="http://www.america.gov/ar/multimedia/photogallery.html">http://www.america.gov/ar/multimedia/photogallery.html</a>	6.69E-04
<a href="http://www.america.gov/ar/publications/books.html">www.america.gov/ar/publications/books.html</a>	6.47E-04
<a href="http://www.america.gov/ar/publications/ejournalusa/1209.html">www.america.gov/ar/publications/ejournalusa/1209.html</a>	5.85E-04
<a href="http://www.america.gov/ar/index.html">www.america.gov/ar/index.html</a>	5.73E-04

Table 9. Search results of term “Economy” using NutchDocumentAnalyzer

<b>Top 10 pages using NutchDocumentAnalyzer</b>	<b>Score for Query</b>
<a href="http://www.america.gov/ar/econ.html">www.america.gov/ar/econ.html</a>	0.38501537
<a href="http://www.america.gov/ar/publications/books/outline-of-the-us-economy.html">www.america.gov/ar/publications/books/outline-of-the-us-economy.html</a>	0.13663794
<a href="http://www.america.gov/ar/econ/business.html">www.america.gov/ar/econ/business.html</a>	0.09989148
<a href="http://www.america.gov/ar/publications/books.html#outline_economy">www.america.gov/ar/publications/books.html#outline_economy</a>	0.01747075
<a href="http://www.america.gov/ar/">www.america.gov/ar/</a>	0.002472951
<a href="http://www.america.gov/ar/multimedia/photogallery.html">www.america.gov/ar/multimedia/photogallery.html</a>	5.41E-04
<a href="http://www.america.gov/ar/index.html">www.america.gov/ar/index.html</a>	4.64E-04
<a href="http://www.america.gov/ar/world/europe.html">www.america.gov/ar/world/europe.html</a>	4.64E-04
<a href="http://www.america.gov/ar/world/mideast.html">www.america.gov/ar/world/mideast.html</a>	4.64E-04
<a href="http://www.america.gov/ar/world/scasia.html">www.america.gov/ar/world/scasia.html</a>	4.02E-04

Table 10. Search results of term “The United States” using ArabicAnalyzer

<b>Top 10 pages using ArabicAnalyzer</b>	<b>Score for Query</b>
<a href="http://www.america.gov/ar/pages/footer/local/about-us.html">www.america.gov/ar/pages/footer/local/about-us.html</a>	0.1196895
<a href="http://www.america.gov/ar/publications/books-content/musliminamerica.html">www.america.gov/ar/publications/books-content/musliminamerica.html</a>	0.11078926
<a href="http://www.america.gov/ar/amlife.html">www.america.gov/ar/amlife.html</a>	0.105654
<a href="http://www.america.gov/ar/services/mobile.html">www.america.gov/ar/services/mobile.html</a>	0.042377986
<a href="http://www.america.gov/ar/multimedia/photogallery.html#/4110/mosques_ar/">www.america.gov/ar/multimedia/photogallery.html#/4110/mosques_ar/</a>	0.022628564
<a href="http://www.america.gov/ar/publications/books.html#beingmuslim">www.america.gov/ar/publications/books.html#beingmuslim</a>	0.015091554
<a href="http://www.america.gov/ar/multimedia/photogallery.html#/4110/religious_freedom_ar/">www.america.gov/ar/multimedia/photogallery.html#/4110/religious_freedom_ar/</a>	0.01136245
<a href="http://www.america.gov/ar/publications/books.html#governed">www.america.gov/ar/publications/books.html#governed</a>	0.01132718
<a href="http://www.america.gov/ar/multimedia/photogallery.html#/4110/islam_ar/">www.america.gov/ar/multimedia/photogallery.html#/4110/islam_ar/</a>	0.011314282
<a href="http://www.america.gov/ar/">www.america.gov/ar/</a>	0.003082759

Table 11. Search results of term “The United States” using NutchDocumentAnalyzer

<b>Top 10 pages using NutchDocumentAnalyzer</b>	<b>Score for Query</b>
<a href="http://www.america.gov/ar/pages/footer/local/about-us.html">www.america.gov/ar/pages/footer/local/about-us.html</a>	0.11997691
<a href="http://www.america.gov/ar/publications/books-content/musliminamerica.html">www.america.gov/ar/publications/books-content/musliminamerica.html</a>	0.11105819
<a href="http://www.america.gov/ar/amlife.html">www.america.gov/ar/amlife.html</a>	0.10594571
<a href="http://www.america.gov/ar/services/mobile.html">www.america.gov/ar/services/mobile.html</a>	0.042462345
<a href="http://www.america.gov/ar/multimedia/photogallery.html#/4110/mosques_ar/">www.america.gov/ar/multimedia/photogallery.html#/4110/mosques_ar/</a>	0.022691099
<a href="http://www.america.gov/ar/publications/books.html#beingmuslim">www.america.gov/ar/publications/books.html#beingmuslim</a>	0.01513323
<a href="http://www.america.gov/ar/multimedia/photogallery.html#/4110/religious_freedom_ar/">www.america.gov/ar/multimedia/photogallery.html#/4110/religious_freedom_ar/</a>	0.011393607
<a href="http://www.america.gov/ar/publications/books.html#governed">www.america.gov/ar/publications/books.html#governed</a>	0.011358418
<a href="http://www.america.gov/ar/multimedia/photogallery.html#/4110/islam_ar/">www.america.gov/ar/multimedia/photogallery.html#/4110/islam_ar/</a>	0.01134555
<a href="http://www.america.gov/ar/">www.america.gov/ar/</a>	0.002807639



Table 12. Search results of term “Democratic” using ArabicAnalyzer

<b>Top 10 pages using NutchDocumentAnalyzer</b>	<b>Score for Query</b>
<a href="http://www.america.gov/ar/global/democracy.html">www.america.gov/ar/global/democracy.html</a>	0.2665834
<a href="http://www.america.gov/ar/global.html">www.america.gov/ar/global.html</a>	0.16062789
<a href="http://www.america.gov/ar/publications/ejournalusa/608.html">www.america.gov/ar/publications/ejournalusa/608.html</a>	0.033635326
<a href="http://www.america.gov/ar/publications/ejournalusa/0110.html">www.america.gov/ar/publications/ejournalusa/0110.html</a>	0.030587077
<a href="http://www.america.gov/ar/democracy/global/index.html">www.america.gov/ar/democracy/global/index.html</a>	0.027611194
<a href="http://www.america.gov/ar/">www.america.gov/ar/</a>	0.002160408
<a href="http://www.america.gov/ar/multimedia/podcast.html">www.america.gov/ar/multimedia/podcast.html</a>	6.10E-04
<a href="http://www.america.gov/ar/publications/books.html">www.america.gov/ar/publications/books.html</a>	5.85E-04
<a href="http://www.america.gov/ar/amlife.html">www.america.gov/ar/amlife.html</a>	5.46E-04
<a href="http://www.america.gov/ar/publications/ejournalusa.html">www.america.gov/ar/publications/ejournalusa.html</a>	5.40E-04

Table 13. Search results of term “Democratic” using NutchDocumentAnalyzer

<b>Top 10 pages using NutchDocumentAnalyzer</b>	<b>Score for Query</b>
<a href="http://www.america.gov/ar/global/democracy.html">www.america.gov/ar/global/democracy.html</a>	0.29354417
<a href="http://www.america.gov/ar/global.html">www.america.gov/ar/global.html</a>	0.17680001
<a href="http://www.america.gov/ar/publications/ejournalusa/0110.html">www.america.gov/ar/publications/ejournalusa/0110.html</a>	0.033559922
<a href="http://www.america.gov/ar/democracy/global/index.html">www.america.gov/ar/democracy/global/index.html</a>	0.030395675
<a href="http://www.america.gov/ar/">www.america.gov/ar/</a>	0.002139257
<a href="http://www.america.gov/ar/multimedia/podcast.html">www.america.gov/ar/multimedia/podcast.html</a>	6.04E-04
<a href="http://www.america.gov/ar/amlife.html">http://www.america.gov/ar/amlife.html</a>	5.40E-04
<a href="http://www.america.gov/ar/amlife/people.html">www.america.gov/ar/amlife/people.html</a>	4.68E-04
<a href="http://www.america.gov/ar/econ.html">www.america.gov/ar/econ.html</a>	4.68E-04
<a href="http://www.america.gov/ar/multimedia.html">www.america.gov/ar/multimedia.html</a>	4.68E-04

A more detailed breakdown of the score of the top 10 pages using *ArabicAnalyzer* and *NutchDocumentAnalyzer* is shown in Appendices A through F.

## **E. CHAPTER SUMMARY**

In this chapter, the results of several experiments to compare the performance of *ArabicAnalyzer* and *NutchDocumentAnalyzer* were described. In the next chapter, the thesis summary and future work recommendations are discussed.

## VI. CONCLUSION

### A. SUMMARY

Arabic IR is a challenging problem because of the complexity of Arabic languages. Even though the light stemmer algorithm was not a perfect solution for Arabic IR problem, it showed improvement over other popular methods. The *ArabicAnalyzer* plug-in inherited the same strengths and weaknesses of the light stemmer algorithm. It also was not perfect, but it did show great promise in saving storage overhead.

The experiments completed in this thesis showed that there are advantages and disadvantages when implementing the *ArabicAnalyzer* plug-in. It is clear by looking at the data that, in general, the *ArabicAnalyzer* plug-in performed as well as the default setting. The query results were relevant to the search terms. It was observed that the plug-in ran slower than the default setting, but the speed issue could be overlooked since the data that this research was trying to gather did not have to be in real time. On the other hand, the *ArabicAnalyzer* plug-in would require at least 20% less memory for its index database, compared with the default setting: the savings in storage could become a major plus when indexing the Internet.

### B. FUTURE WORK

For future research, the plug-in needs to be integrated into the social networking tool and experiments need to be conducted to determine the recall, precision and relevance of the plug-in in the integration environment. The experiments should also help determine the strengths and weaknesses of the plug-in in such environments and recommend improvement.

THIS PAGE INTENTIONALLY LEFT BLANK

## APPENDIX A

This is the detail score for query of top 10 pages using *ArabicAnalyzer*.

Search Term: داصتقا (economy)

Page 1:

- boost = 0.22821301
- digest = 767d250a62c827c2bd330c0674546358
- lang = ar
- segment = 20100305180909
- title = داصتقالا - داصتقالا - America.gov
- tstamp = 20100305230954510
- url = http://www.america.gov/ar/econ.html

score for query: داصتقا

- 0.3486507 = (MATCH) sum of:
  - 0.18338637 = (MATCH) weight(anchor:داصتقا^2.0 in 15), product of:
    - 0.2879631 = queryWeight(anchor:داصتقا^2.0), product of:
      - 2.0 = boost
      - 4.075775 = idf(docFreq=5, numDocs=130)
      - 0.035326175 = queryNorm
    - 0.63683987 = (MATCH) fieldWeight(anchor:داصتقا in 15), product of:
      - 1.0 = tf(termFreq(anchor:داصتقا)=1)
      - 4.075775 = idf(docFreq=5, numDocs=130)
      - 0.15625 = fieldNorm(field=anchor, doc=15)
  - 6.6904654E-4 = (MATCH) weight(content:داصتقا in 15), product of:
    - 0.03756986 = queryWeight(content:داصتقا), product of:
      - 1.0635134 = idf(docFreq=121, numDocs=130)

- $0.035326175 = \text{queryNorm}$
- $0.017808065 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقا in 15})$ , product of:
  - $2.4494898 = \text{tf}(\text{termFreq}(\text{content: داصتقا})=6)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=15)$
- $0.16459529 = (\text{MATCH}) \text{weight}(\text{title: داصتقا}^{1.5} \text{ in 15})$ , product of:
  - $0.23745762 = \text{queryWeight}(\text{title: داصتقا}^{1.5})$ , product of:
    - $1.5 = \text{boost}$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{numDocs}=130)$
    - $0.035326175 = \text{queryNorm}$
  - $0.6931565 = (\text{MATCH}) \text{fieldWeight}(\text{title: داصتقا in 15})$ , product of:
    - $1.4142135 = \text{tf}(\text{termFreq}(\text{title: داصتقا})=2)$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{numDocs}=130)$
    - $0.109375 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=15)$

\*\*\*\*\*

Page 2:

- $\text{boost} = 0.16124225$
- $\text{digest} = \text{fdaa17fd08dfde3bb91a83a6d98afa04}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180909$
- $\text{title} = \text{يكريم ال داصتقا ال زجوم - Outline of the U.S. Economy - America.gov}$
- $\text{tstamp} = 20100305230918398$
- $\text{url} = \text{http://www.america.gov/ar/publications/books/outline-of-the-us-economy.html}$

**score for query: داصتقا**

- $0.12422927 = (\text{MATCH}) \text{sum of:}$

- $0.07335455 = (\text{MATCH}) \text{weight}(\text{anchor:داصتقا}^{2.0} \text{ in } 84)$ , product of:
  - $0.2879631 = \text{queryWeight}(\text{anchor:داصتقا}^{2.0})$ , product of:
    - $2.0 = \text{boost}$
    - $4.075775 = \text{idf}(\text{docFreq}=5, \text{numDocs}=130)$
    - $0.035326175 = \text{queryNorm}$
  - $0.25473595 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:داصتقا} \text{ in } 84)$ , product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:داصتقا})=1)$
    - $4.075775 = \text{idf}(\text{docFreq}=5, \text{numDocs}=130)$
    - $0.0625 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=84)$
- $9.948079\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:داصتقا} \text{ in } 84)$ , product of:
  - $0.03756986 = \text{queryWeight}(\text{content:داصتقا})$ , product of:
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.035326175 = \text{queryNorm}$
  - $0.026478883 = (\text{MATCH}) \text{fieldWeight}(\text{content:داصتقا} \text{ in } 84)$ , product of:
    - $5.0990195 = \text{tf}(\text{termFreq}(\text{content:داصتقا})=26)$
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=84)$
- $0.049879905 = (\text{MATCH}) \text{weight}(\text{title:داصتقا}^{1.5} \text{ in } 84)$ , product of:
  - $0.23745762 = \text{queryWeight}(\text{title:داصتقا}^{1.5})$ , product of:
    - $1.5 = \text{boost}$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{numDocs}=130)$
    - $0.035326175 = \text{queryNorm}$
  - $0.21005814 = (\text{MATCH}) \text{fieldWeight}(\text{title:داصتقا} \text{ in } 84)$ , product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{title:داصتقا})=1)$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{numDocs}=130)$



- $0.046875 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=84)$

\*\*\*\*\*

Page 3:

boost = 0.16781548

- digest = b4649130898e202ca38ef61b6b22b917
- lang = ar
- segment = 20100305180909
- title = قراچتلاو لامعأل - قراچتلاو لامعأل - America.gov
- tstamp = 20100305231006880
- url = http://www.america.gov/ar/econ/business.html

score for query: داصتقا

- $0.09217107 = (\text{MATCH})$  sum of:
  - $0.091693185 = (\text{MATCH})$  weight(anchor: داصتقا^2.0 in 16), product of:
    - $0.2879631 = \text{queryWeight}(\text{anchor}: \text{داصتقا}^2.0)$ , product of:
      - $2.0 = \text{boost}$
      - $4.075775 = \text{idf}(\text{docFreq}=5, \text{numDocs}=130)$
      - $0.035326175 = \text{queryNorm}$
    - $0.31841993 = (\text{MATCH})$  fieldWeight(anchor: داصتقا in 16), product of:
      - $1.0 = \text{tf}(\text{termFreq}(\text{anchor}: \text{داصتقا})=1)$
      - $4.075775 = \text{idf}(\text{docFreq}=5, \text{numDocs}=130)$
      - $0.078125 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=16)$
  - $4.7789037\text{E-}4 = (\text{MATCH})$  weight(content: داصتقا in 16), product of:
    - $0.03756986 = \text{queryWeight}(\text{content}: \text{داصتقا})$ , product of:
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.035326175 = \text{queryNorm}$

- $0.012720046 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقا in 16}),$   
product of:
  - $2.4494898 = \text{tf}(\text{termFreq}(\text{content: داصتقا})=6)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=16)$

\*\*\*\*\*

Page 4:

- $\text{boost} = 0.030659562$
- $\text{digest} = 4304d87a1d51187c1c1d0b2b4d1597a8$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181031$
- $\text{title} = \text{America.gov} - \text{قراچتلا داصتقا شاعنإ} - \text{قراچتلا داصتقا شاعنإ}$
- $\text{tstamp} = 20100305231127141$
- $\text{url} = \text{http://www.america.gov/ar/reviving\_trade\_ar.html}$

**score for query: داصتقا**

- $0.033118278 = (\text{MATCH}) \text{sum of:}$ 
  - $0.018338637 = (\text{MATCH}) \text{weight}(\text{anchor: داصتقا}^2.0 \text{ in } 104), \text{product of:}$ 
    - $0.2879631 = \text{queryWeight}(\text{anchor: داصتقا}^2.0), \text{product of:}$ 
      - $2.0 = \text{boost}$
      - $4.075775 = \text{idf}(\text{docFreq}=5, \text{numDocs}=130)$
      - $0.035326175 = \text{queryNorm}$
    - $0.06368399 = (\text{MATCH}) \text{fieldWeight}(\text{anchor: داصتقا in } 104),$   
product of:
      - $1.0 = \text{tf}(\text{termFreq}(\text{anchor: داصتقا})=1)$
      - $4.075775 = \text{idf}(\text{docFreq}=5, \text{numDocs}=130)$
      - $0.015625 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=104)$
  - $8.363082\text{E-}5 = (\text{MATCH}) \text{weight}(\text{content: داصتقا in } 104), \text{product of:}$

- $0.03756986 = \text{queryWeight}(\text{content: داصتقا}), \text{ product of:}$ 
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{ numDocs}=130)$
  - $0.035326175 = \text{queryNorm}$
- $0.002226008 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقا in 104}), \text{ product of:}$ 
  - $2.4494898 = \text{tf}(\text{termFreq}(\text{content: داصتقا})=6)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{ numDocs}=130)$
  - $8.544922\text{E-}4 = \text{fieldNorm}(\text{field}=\text{content}, \text{ doc}=104)$
- $0.014696008 = (\text{MATCH}) \text{weight}(\text{title: داصتقا}^{1.5} \text{ in 104}), \text{ product of:}$ 
  - $0.23745762 = \text{queryWeight}(\text{title: داصتقا}^{1.5}), \text{ product of:}$ 
    - $1.5 = \text{boost}$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{ numDocs}=130)$
    - $0.035326175 = \text{queryNorm}$
  - $0.06188897 = (\text{MATCH}) \text{fieldWeight}(\text{title: داصتقا in 104}), \text{ product of:}$ 
    - $1.4142135 = \text{tf}(\text{termFreq}(\text{title: داصتقا})=2)$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{ numDocs}=130)$
    - $0.009765625 = \text{fieldNorm}(\text{field}=\text{title}, \text{ doc}=104)$

\*\*\*\*\*

Page 5:

- $\text{boost} = 0.02675021$
- $\text{digest} = \text{aaf055c1e690c63cf69285f8ab04f499}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181330$
- $\text{title} = \text{بیتک - بیتک - America.gov}$
- $\text{tstamp} = 20100305231345369$
- $\text{url} = \text{http://www.america.gov/ar/publications/books.html\#outline\_economy}$

score for query: داصتقا

- 0.016127191 = (MATCH) sum of:
  - 0.016046308 = (MATCH) weight(anchor:داصتقا^2.0 in 80), product of:
    - 0.2879631 = queryWeight(anchor:داصتقا^2.0), product of:
      - 2.0 = boost
      - 4.075775 = idf(docFreq=5, numDocs=130)
      - 0.035326175 = queryNorm
    - 0.05572349 = (MATCH) fieldWeight(anchor:داصتقا in 80), product of:
      - 1.0 = tf(termFreq(anchor:داصتقا)=1)
      - 4.075775 = idf(docFreq=5, numDocs=130)
      - 0.013671875 = fieldNorm(field=anchor, doc=80)
  - 8.088332E-5 = (MATCH) weight(content:داصتقا in 80), product of:
    - 0.03756986 = queryWeight(content:داصتقا), product of:
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 0.035326175 = queryNorm
    - 0.0021528779 = (MATCH) fieldWeight(content:داصتقا in 80), product of:
      - 3.3166249 = tf(termFreq(content:داصتقا)=11)
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 6.1035156E-4 = fieldNorm(field=content, doc=80)

\*\*\*\*\*

Page 6:

- boost = 1.0000145
- digest = 0d5b023c802941ddb358071073a98833
- lang = ar
- segment = 20100305180856

- title = أمريكا - أول فصل - أول فصل - America.gov
- tstamp = 20100305230902835
- url = http://www.america.gov/ar/

**score for query: داصتقا**

- 0.0030584983 = (MATCH) sum of:
  - 0.0030584983 = (MATCH) weight(content: داصتقا in 0), product of:
    - 0.03756986 = queryWeight(content: داصتقا), product of:
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 0.035326175 = queryNorm
    - 0.08140829 = (MATCH) fieldWeight(content: داصتقا in 0), product of:
      - 2.4494898 = tf(termFreq(content: داصتقا)=6)
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 0.03125 = fieldNorm(field=content, doc=0)

\*\*\*\*\*

**Page 7:**

- boost = 0.23009512
- digest = 15d9ca5e7382f701cd03fb542ae3ab22
- lang = ar
- segment = 20100305180909
- title = أمريكا - أول فصل - أول فصل - America.gov
- tstamp = 20100305230915350
- url = http://www.america.gov/ar/multimedia/photogallery.html

**score for query: داصتقا**

- 6.6904654E-4 = (MATCH) sum of:
  - 6.6904654E-4 = (MATCH) weight(content: داصتقا in 38), product of:
    - 0.03756986 = queryWeight(content: داصتقا), product of:

- $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
- $0.035326175 = \text{queryNorm}$
- $0.017808065 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقا in 38}),$   
product of:
  - $2.4494898 = \text{tf}(\text{termFreq}(\text{content: داصتقا})=6)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=38)$

\*\*\*\*\*

Page 8:

- $\text{boost} = 0.22996004$
- $\text{digest} = \text{a0130240b4348578aa8a83e59187dfb3}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180909$
- $\text{title} = \text{ب تک - ب تک - America.gov}$
- $\text{tstamp} = 20100305231001279$
- $\text{url} = \text{http://www.america.gov/ar/publications/books.html}$

**score for query: داصتقا**

- $6.4706657\text{E-}4 = (\text{MATCH}) \text{sum of:}$ 
  - $6.4706657\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content: داصتقا in 73}), \text{product of:}$ 
    - $0.03756986 = \text{queryWeight}(\text{content: داصتقا}), \text{product of:}$ 
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.035326175 = \text{queryNorm}$
    - $0.017223023 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقا in 73}),$   
product of:
      - $3.3166249 = \text{tf}(\text{termFreq}(\text{content: داصتقا})=11)$
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=73)$

\*\*\*\*\*

Page 9:

- boost = 0.16516872
- digest = bc202e5e0f508e4291bb897eec7814dc
- lang = ar
- segment = 20100305180909
- title = 1209 - America.gov - ومناول مكحل
- tstamp = 20100305230941328
- url = http://www.america.gov/ar/publications/ejournalusa/1209.html

score for query: داصتقا

- 5.8529375E-4 = (MATCH) sum of:
  - 5.8529375E-4 = (MATCH) weight(content: داصتقا in 97), product of:
    - 0.03756986 = queryWeight(content: داصتقا), product of:
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 0.035326175 = queryNorm
    - 0.01557881 = (MATCH) fieldWeight(content: داصتقا in 97), product of:
      - 3.0 = tf(termFreq(content: داصتقا)=9)
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 0.0048828125 = fieldNorm(field=content, doc=97)

\*\*\*\*\*

Page 10:

- boost = 0.19712433
- digest = c25a22a11ab6bec420c26625155ced62
- lang = ar
- segment = 20100305180909
- title = America.gov - أول الفصل - أول الفصل

- tstamp = 20100305230929458
- url = http://www.america.gov/ar/index.html

**score for query: داصتقا**

- $5.7346845E-4$  = (MATCH) sum of:
  - $5.7346845E-4$  = (MATCH) weight(content:داصتقا in 30), product of:
    - $0.03756986$  = queryWeight(content:داصتقا), product of:
      - $1.0635134$  = idf(docFreq=121, numDocs=130)
      - $0.035326175$  = queryNorm
    - $0.015264055$  = (MATCH) fieldWeight(content:داصتقا in 30), product of:
      - $2.4494898$  = tf(termFreq(content:داصتقا)=6)
      - $1.0635134$  = idf(docFreq=121, numDocs=130)
      - $0.005859375$  = fieldNorm(field=content, doc=30)



THIS PAGE INTENTIONALLY LEFT BLANK

## APPENDIX B

This is the detail score for query of top 10 pages using *NutchDocumentAnalyzer*.

Search Term: داصتقالا (ecomony)

Page 1:

- boost = 0.22826105
- digest = c33a5dc3f7d8475491bfafcf91c8b283
- lang = ar
- segment = 20100307101102
- title = داصتقالا - داصتقالا - America.gov
- tstamp = 20100307151153574
- url = http://www.america.gov/ar/econ.html

**score for query: داصتقالا**

- 0.38501537 = (MATCH) sum of:
  - 0.1991137 = (MATCH) weight(anchor: داصتقالا^2.0 in 16), product of:
    - 0.29873407 = queryWeight(anchor: داصتقالا^2.0), product of:
      - 2.0 = boost
      - 4.2657595 = idf(docFreq=4, numDocs=131)
      - 0.035015345 = queryNorm
    - 0.6665249 = (MATCH) fieldWeight(anchor: داصتقالا in 16), product of:
      - 1.0 = tf(termFreq(anchor: داصتقالا)=1)
      - 4.2657595 = idf(docFreq=4, numDocs=131)
      - 0.15625 = fieldNorm(field=anchor, doc=16)
  - 5.40958E-4 = (MATCH) weight(content: داصتقالا in 16), product of:
    - 0.037221763 = queryWeight(content: داصتقالا), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)

- 0.035015345 = queryNorm
- 0.01453338 = (MATCH) fieldWeight(content:داصتقالا in 16), product of:
  - 2.0 = tf(termFreq(content:داصتقالا)=4)
  - 1.063013 = idf(docFreq=122, numDocs=131)
  - 0.0068359375 = fieldNorm(field=content, doc=16)
- 0.18536073 = (MATCH) weight(title:داصتقالا^1.5 in 16), product of:
  - 0.25088066 = queryWeight(title:داصتقالا^1.5), product of:
    - 1.5 = boost
    - 4.776585 = idf(docFreq=2, numDocs=131)
    - 0.035015345 = queryNorm
  - 0.7388403 = (MATCH) fieldWeight(title:داصتقالا in 16), product of:
    - 1.4142135 = tf(termFreq(title:داصتقالا)=2)
    - 4.776585 = idf(docFreq=2, numDocs=131)
    - 0.109375 = fieldNorm(field=title, doc=16)

\*\*\*\*\*

Page 2:

- boost = 0.16124225
- digest = 6120d6b7e6584b6a71b7d9990a68b952
- lang = ar
- segment = 20100307101102
- title = يڤريمال داصتقالا زجوم - Outline of the U.S. Economy - America.gov
- tstamp = 20100307151112494
- url = http://www.america.gov/ar/publications/books/outline-of-the-us-economy.html

**score for query: داصتقالا**

- 0.13663794 = (MATCH) sum of:
  - 0.07964548 = (MATCH) weight(anchor:داصتقالا^2.0 in 85), product of:

- $0.29873407 = \text{queryWeight}(\text{anchor: داصتقال}^2.0)$ , product of:
  - $2.0 = \text{boost}$
  - $4.2657595 = \text{idf}(\text{docFreq}=4, \text{numDocs}=131)$
  - $0.035015345 = \text{queryNorm}$
- $0.26660997 = (\text{MATCH}) \text{fieldWeight}(\text{anchor: داصتقال in 85})$ , product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{anchor: داصتقال})=1)$
  - $4.2657595 = \text{idf}(\text{docFreq}=4, \text{numDocs}=131)$
  - $0.0625 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=85)$
- $8.196751\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content: داصتقال in 85})$ , product of:
  - $0.037221763 = \text{queryWeight}(\text{content: داصتقال})$ , product of:
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.035015345 = \text{queryNorm}$
  - $0.022021394 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقال in 85})$ , product of:
    - $4.2426405 = \text{tf}(\text{termFreq}(\text{content: داصتقال})=18)$
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=85)$
- $0.056172792 = (\text{MATCH}) \text{weight}(\text{title: داصتقال}^1.5 \text{ in 85})$ , product of:
  - $0.25088066 = \text{queryWeight}(\text{title: داصتقال}^1.5)$ , product of:
    - $1.5 = \text{boost}$
    - $4.776585 = \text{idf}(\text{docFreq}=2, \text{numDocs}=131)$
    - $0.035015345 = \text{queryNorm}$
  - $0.22390243 = (\text{MATCH}) \text{fieldWeight}(\text{title: داصتقال in 85})$ , product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{title: داصتقال})=1)$
    - $4.776585 = \text{idf}(\text{docFreq}=2, \text{numDocs}=131)$
    - $0.046875 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=85)$

\*\*\*\*\*

Page 3:

- boost = 0.16784814
- digest = 2e923bcfb9409e9be88aad90198063bc
- lang = ar
- segment = 20100307101102
- title = America.gov - قراچتل او لامعأل - قراچتل او لامعأل
- tstamp = 20100307151205095
- url = http://www.america.gov/ar/econ/business.html

داصتقال: score for query:

- 0.09989148 = (MATCH) sum of:
  - 0.09955685 = (MATCH) weight(anchor:داصتقال^2.0 in 17), product of:
    - 0.29873407 = queryWeight(anchor:داصتقال^2.0), product of:
      - 2.0 = boost
      - 4.2657595 = idf(docFreq=4, numDocs=131)
      - 0.035015345 = queryNorm
    - 0.33326244 = (MATCH) fieldWeight(anchor:داصتقال in 17), product of:
      - 1.0 = tf(termFreq(anchor:داصتقال)=1)
      - 4.2657595 = idf(docFreq=4, numDocs=131)
      - 0.078125 = fieldNorm(field=anchor, doc=17)
  - 3.34631E-4 = (MATCH) weight(content:داصتقال in 17), product of:
    - 0.037221763 = queryWeight(content:داصتقال), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.035015345 = queryNorm
    - 0.008990197 = (MATCH) fieldWeight(content:داصتقال in 17), product of:
      - 1.7320508 = tf(termFreq(content:داصتقال)=3)

- $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
- $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=17)$

\*\*\*\*\*

Page 4:

- $\text{boost} = 0.02675021$
- $\text{digest} = 4eb9183dbdc405b0d40ef3c92da5ed66$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101458$
- $\text{title} = \text{بتيك - بتيك - America.gov}$
- $\text{tstamp} = 20100307151513307$
- $\text{url} = \text{http://www.america.gov/ar/publications/books.html\#outline\_economy}$

**score for query: داصتقالا**

- $0.01747075 = (\text{MATCH}) \text{ sum of:}$ 
  - $0.017422449 = (\text{MATCH}) \text{ weight}(\text{anchor: داصتقالا}^2.0 \text{ in } 81), \text{ product of:}$ 
    - $0.29873407 = \text{queryWeight}(\text{anchor: داصتقالا}^2.0), \text{ product of:}$ 
      - $2.0 = \text{boost}$
      - $4.2657595 = \text{idf}(\text{docFreq}=4, \text{numDocs}=131)$
      - $0.035015345 = \text{queryNorm}$
    - $0.058320932 = (\text{MATCH}) \text{ fieldWeight}(\text{anchor: داصتقالا} \text{ in } 81), \text{ product of:}$ 
      - $1.0 = \text{tf}(\text{termFreq}(\text{anchor: داصتقالا})=1)$
      - $4.2657595 = \text{idf}(\text{docFreq}=4, \text{numDocs}=131)$
      - $0.013671875 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=81)$
  - $4.8299826\text{E-}5 = (\text{MATCH}) \text{ weight}(\text{content: داصتقالا} \text{ in } 81), \text{ product of:}$ 
    - $0.037221763 = \text{queryWeight}(\text{content: داصتقالا}), \text{ product of:}$ 
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.035015345 = \text{queryNorm}$

- $0.0012976233 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقال in 81}),$   
product of:
  - $2.0 = \text{tf}(\text{termFreq}(\text{content: داصتقال})=4)$
  - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
  - $6.1035156\text{E-}4 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=81)$

\*\*\*\*\*

Page 5:

- $\text{boost} = 1.0000145$
- $\text{digest} = \text{eed4dd9817b50ffda0aef158be6e4c12}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101052$
- $\text{title} = \text{America.gov - داصتقال - داصتقال - داصتقال}$
- $\text{tstamp} = 20100307151057483$
- $\text{url} = \text{http://www.america.gov/ar/}$

**score for query: داصتقال**

- $0.002472951 = (\text{MATCH}) \text{sum of:}$ 
  - $0.002472951 = (\text{MATCH}) \text{weight}(\text{content: داصتقال in 0}), \text{product of:}$ 
    - $0.037221763 = \text{queryWeight}(\text{content: داصتقال}), \text{product of:}$ 
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.035015345 = \text{queryNorm}$
    - $0.06643831 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقال in 0}),$   
product of:
      - $2.0 = \text{tf}(\text{termFreq}(\text{content: داصتقال})=4)$
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.03125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=0)$

\*\*\*\*\*

Page 6:

- $\text{boost} = 0.23014276$
- $\text{digest} = 5f50883579dcc0acb85ff3052764f758$

- lang = ar
- segment = 20100307101102
- title = روصلاب ةيكريمأل ا ةياكحل ا - روص موبلأ - America.gov
- tstamp = 20100307151109851
- url = http://www.america.gov/ar/multimedia/photogallery.html

**داصتقالا: score for query**

- 5.40958E-4 = (MATCH) sum of:
  - 5.40958E-4 = (MATCH) weight(content:داصتقالا in 39), product of:
    - 0.037221763 = queryWeight(content:داصتقالا), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.035015345 = queryNorm
    - 0.01453338 = (MATCH) fieldWeight(content:داصتقالا in 39), product of:
      - 2.0 = tf(termFreq(content:داصتقالا)=4)
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.0068359375 = fieldNorm(field=content, doc=39)

\*\*\*\*\*

**Page 7:**

- boost = 0.19715214
- digest = e84ec632a6d47466f40d6beacbcfbdf7
- lang = ar
- segment = 20100307101102
- title = ةلوالا ةحفصلا - ةلوالا ةحفصلا - America.gov
- tstamp = 20100307151123237
- url = http://www.america.gov/ar/index.html

**داصتقالا: score for query**

- 4.636783E-4 = (MATCH) sum of:
  - 4.636783E-4 = (MATCH) weight(content:داصتقالا in 31), product of:



- $0.037221763 = \text{queryWeight}(\text{content: داصتقالا}), \text{product of:}$ 
  - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
  - $0.035015345 = \text{queryNorm}$
- $0.012457183 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقالا in 31}), \text{product of:}$ 
  - $2.0 = \text{tf}(\text{termFreq}(\text{content: داصتقالا})=4)$
  - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
  - $0.005859375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=31)$

\*\*\*\*\*

Page 8:

- $\text{boost} = 0.20153543$
- $\text{digest} = 67\text{da}72\text{f}899\text{f}80475\text{f}1\text{ae}62770921\text{f}1\text{bd}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101102$
- $\text{title} = \text{America.gov} - \text{ايساروأو ابوروأ} - \text{ايساروأو ابوروأ}$
- $\text{tstamp} = 20100307151127260$
- $\text{url} = \text{http://www.america.gov/ar/world/europe.html}$

**score for query: داصتقالا**

- $4.636783\text{E-}4 = (\text{MATCH}) \text{sum of:}$ 
  - $4.636783\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content: داصتقالا in 128}), \text{product of:}$ 
    - $0.037221763 = \text{queryWeight}(\text{content: داصتقالا}), \text{product of:}$ 
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.035015345 = \text{queryNorm}$
    - $0.012457183 = (\text{MATCH}) \text{fieldWeight}(\text{content: داصتقالا in 128}), \text{product of:}$ 
      - $2.0 = \text{tf}(\text{termFreq}(\text{content: داصتقالا})=4)$
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.005859375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=128)$

\*\*\*\*\*

Page 9:

- boost = 0.20091416
- digest = 9256cf74ef9b595d81f726d5f347898a
- lang = ar
- segment = 20100307101102
- title = ايقيرفأ لامشو طسوأل قرشل - ايقيرفأ لامشو طسوأل قرشل - America.gov
- tstamp = 20100307151148461
- url = http://www.america.gov/ar/world/mideast.html

score for query: داصتقالا

- 4.636783E-4 = (MATCH) sum of:
  - 4.636783E-4 = (MATCH) weight(content: داصتقالا in 129), product of:
    - 0.037221763 = queryWeight(content: داصتقالا), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.035015345 = queryNorm
    - 0.012457183 = (MATCH) fieldWeight(content: داصتقالا in 129), product of:
      - 2.0 = tf(termFreq(content: داصتقالا)=4)
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.005859375 = fieldNorm(field=content, doc=129)

\*\*\*\*\*

Page 10:

- boost = 0.20039715
- digest = adbc4a97340b57bcc256c62131041c4c
- lang = ar
- segment = 20100307101102
- title = ايسأ قرشو بونج - ايسأ قرشو بونج - America.gov
- tstamp = 20100307151119356

- url = <http://www.america.gov/ar/world/scasia.html>

**score for query:** داصتقالا

- $4.015572E-4 = (\text{MATCH})$  sum of:
  - $4.015572E-4 = (\text{MATCH})$  weight(content: داصتقالا in 130), product of:
    - $0.037221763 = \text{queryWeight}(\text{content: داصتقالا})$ , product of:
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.035015345 = \text{queryNorm}$
    - $0.010788237 = (\text{MATCH})$  fieldWeight(content: داصتقالا in 130), product of:
      - $1.7320508 = \text{tf}(\text{termFreq}(\text{content: داصتقالا})=3)$
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.005859375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=130)$

## APPENDIX C

This is the detail score for query of top 10 pages using *ArabicAnalyzer*.

Search Term: **الكريم** (The United States)

Page 1:

- boost = 0.15805063
- digest = 6b6baa67bd29d99a3cf293efeb2bc3e1
- lang = ar
- segment = 20100305181031
- title = America.gov - فووغ تود الكريم أعقوم لوح - فووغ تود الكريم أعقوم لوح
- tstamp = 20100305231050378
- url = http://www.america.gov/ar/pages/footer/local/about-us.html

**score for query: **الكريم****

- 0.1196895 = (MATCH) sum of:
  - 0.059957497 = (MATCH) weight(anchor:**الكريم**<sup>2.0</sup> in 68), product of:
    - 0.261373 = queryWeight(anchor:**الكريم**<sup>2.0</sup>), product of:
      - 2.0 = boost
      - 3.6703098 = idf(docFreq=8, numDocs=130)
      - 0.035606395 = queryNorm
    - 0.22939436 = (MATCH) fieldWeight(anchor:**الكريم** in 68), product of:
      - 1.0 = tf(termFreq(anchor:**الكريم**)=1)
      - 3.6703098 = idf(docFreq=8, numDocs=130)
      - 0.0625 = fieldNorm(field=anchor, doc=68)
  - 4.8168114E-4 = (MATCH) weight(content:**الكريم** in 68), product of:
    - 0.037867878 = queryWeight(content:**الكريم**), product of:
      - 1.0635134 = idf(docFreq=121, numDocs=130)

- 0.035606395 = queryNorm
- 0.012720046 = (MATCH) fieldWeight(content:الكريم in 68), product of:
  - 2.4494898 = tf(termFreq(content:الكريم)=6)
  - 1.0635134 = idf(docFreq=121, numDocs=130)
  - 0.0048828125 = fieldNorm(field=content, doc=68)
- 0.059250325 = (MATCH) weight(title:الكريم^1.5 in 68), product of:
  - 0.23934121 = queryWeight(title:الكريم^1.5), product of:
    - 1.5 = boost
    - 4.4812403 = idf(docFreq=3, numDocs=130)
    - 0.035606395 = queryNorm
  - 0.24755588 = (MATCH) fieldWeight(title:الكريم in 68), product of:
    - 1.4142135 = tf(termFreq(title:الكريم)=2)
    - 4.4812403 = idf(docFreq=3, numDocs=130)
    - 0.0390625 = fieldNorm(field=title, doc=68)

Page 2:

- boost = 0.16184442
- digest = 0f454ab63865ae2e08003bb23896bfad
- lang = ar
- segment = 20100305180909
- title = الكريم أ يف نوملسملا - Being Muslim in America - America.gov
- tstamp = 20100305231009575
- url = http://www.america.gov/ar/publications/books-content/musliminamerica.html

**score for query: الكريم**

- 0.11078926 = (MATCH) sum of:
  - 0.059957497 = (MATCH) weight(anchor:الكريم^2.0 in 72), product of:

- $0.261373 = \text{queryWeight}(\text{anchor:الكريم}^{2.0})$ , product of:
  - $2.0 = \text{boost}$
  - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
  - $0.035606395 = \text{queryNorm}$
- $0.22939436 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم in 72})$ , product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
  - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
  - $0.0625 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=72)$
- $5.5619737\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:الكريم in 72})$ , product of:
  - $0.037867878 = \text{queryWeight}(\text{content:الكريم})$ , product of:
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.035606395 = \text{queryNorm}$
  - $0.014687842 = (\text{MATCH}) \text{fieldWeight}(\text{content:الكريم in 72})$ , product of:
    - $2.828427 = \text{tf}(\text{termFreq}(\text{content:الكريم})=8)$
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=72)$
- $0.050275568 = (\text{MATCH}) \text{weight}(\text{title:الكريم}^{1.5} \text{ in 72})$ , product of:
  - $0.23934121 = \text{queryWeight}(\text{title:الكريم}^{1.5})$ , product of:
    - $1.5 = \text{boost}$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{numDocs}=130)$
    - $0.035606395 = \text{queryNorm}$
  - $0.21005814 = (\text{MATCH}) \text{fieldWeight}(\text{title:الكريم in 72})$ , product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{title:الكريم})=1)$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{numDocs}=130)$
    - $0.046875 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=72)$

\*\*\*\*\*

Page 3:

- boost = 0.23032264
- digest = ce4a12d589c1a56e886d5b6848609391
- lang = ar
- segment = 20100305180909
- title = أمريكا - أي حل - أي كريم أ ل أي حل - America.gov
- tstamp = 20100305230939904
- url = http://www.america.gov/ar/amlife.html

score for query: الكريم

- 0.105654 = (MATCH) sum of:
  - 0.10492562 = (MATCH) weight(anchor:الكريم^2.0 in 3), product of:
    - 0.261373 = queryWeight(anchor:الكريم^2.0), product of:
      - 2.0 = boost
      - 3.6703098 = idf(docFreq=8, numDocs=130)
      - 0.035606395 = queryNorm
    - 0.40144014 = (MATCH) fieldWeight(anchor:الكريم in 3), product of:
      - 1.0 = tf(termFreq(anchor:الكريم)=1)
      - 3.6703098 = idf(docFreq=8, numDocs=130)
      - 0.109375 = fieldNorm(field=anchor, doc=3)
  - 7.283851E-4 = (MATCH) weight(content:الكريم in 3), product of:
    - 0.037867878 = queryWeight(content:الكريم), product of:
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 0.035606395 = queryNorm
    - 0.019234907 = (MATCH) fieldWeight(content:الكريم in 3), product of:
      - 2.6457512 = tf(termFreq(content:الكريم)=7)

- $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
- $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=3)$

\*\*\*\*\*

Page 4:

- $\text{boost} = 0.15872316$
- $\text{digest} = \text{dcfeb490d3db633d16bfb0588d67076d}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180909$
- $\text{title} = \text{ربع نشيدي لايابوم فوغ تود الكريم ل اوجلا ةخسن - نشيدي لايابوم ةمدخ PDA - America.gov ةزهجأ}$
- $\text{tstamp} = 20100305230942596$
- $\text{url} = \text{http://www.america.gov/ar/services/mobile.html}$

**score for query: الكريم**

- $0.042377986 = (\text{MATCH}) \text{ sum of:}$ 
  - $4.8168114\text{E-}4 = (\text{MATCH}) \text{ weight}(\text{content:الكريم in 114}), \text{ product of:}$ 
    - $0.037867878 = \text{queryWeight}(\text{content:الكريم}), \text{ product of:}$ 
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.035606395 = \text{queryNorm}$
    - $0.012720046 = (\text{MATCH}) \text{ fieldWeight}(\text{content:الكريم in 114}), \text{ product of:}$ 
      - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=114)$
  - $0.041896306 = (\text{MATCH}) \text{ weight}(\text{title:الكريم}^{1.5} \text{ in 114}), \text{ product of:}$ 
    - $0.23934121 = \text{queryWeight}(\text{title:الكريم}^{1.5}), \text{ product of:}$ 
      - $1.5 = \text{boost}$
      - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{numDocs}=130)$
      - $0.035606395 = \text{queryNorm}$



- $0.17504844 = (\text{MATCH}) \text{fieldWeight}(\text{title:الكريم in 114}), \text{product of:}$ 
    - $1.0 = \text{tf}(\text{termFreq}(\text{title:الكريم})=1)$
    - $4.4812403 = \text{idf}(\text{docFreq}=3, \text{numDocs}=130)$
    - $0.0390625 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=114)$
- \*\*\*\*\*

Page 5:

- $\text{boost} = 0.04832446$
- $\text{digest} = 87c8a44e7bc9cb3221f6823da385f8dd$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181031$
- $\text{title} = \text{America.gov} - \text{روص لاب ةيكريم أةي الكحل} - \text{روص موبل}$
- $\text{tstamp} = 20100305231143209$
- $\text{url} =$   
 $\text{http://www.america.gov/ar/multimedia/photogallery.html\#/4110/mosques\_ar/}$

**score for query: الكريم**

- $0.022628564 = (\text{MATCH}) \text{sum of:}$ 
  - $0.02248406 = (\text{MATCH}) \text{weight}(\text{anchor:الكريم}^2.0 \text{ in } 50), \text{product of:}$ 
    - $0.261373 = \text{queryWeight}(\text{anchor:الكريم}^2.0), \text{product of:}$ 
      - $2.0 = \text{boost}$
      - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
      - $0.035606395 = \text{queryNorm}$
    - $0.08602288 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم in } 50), \text{product of:}$ 
      - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
      - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
      - $0.0234375 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=50)$
  - $1.4450435\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:الكريم in } 50), \text{product of:}$

- $0.037867878 = \text{queryWeight}(\text{content:الكريم}), \text{product of:}$ 
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $0.035606395 = \text{queryNorm}$
- $0.0038160137 = (\text{MATCH}) \text{fieldWeight}(\text{content:الكريم in 50}), \text{product of:}$ 
  - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $0.0014648438 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=50)$

\*\*\*\*\*

Page 6:

- $\text{boost} = 0.033444975$
- $\text{digest} = 7212084a79cd19adbfc07dc50d3c0ea4$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181031$
- $\text{title} = \text{America.gov} - \text{بتك} - \text{بتك}$
- $\text{tstamp} = 20100305231138931$
- $\text{url} = \text{http://www.america.gov/ar/publications/books.html\#beingmuslim}$

**score for query: الكريم**

- $0.015091554 = (\text{MATCH}) \text{sum of:}$ 
  - $0.014989374 = (\text{MATCH}) \text{weight}(\text{anchor:الكريم}^2.0 \text{ in } 75), \text{product of:}$ 
    - $0.261373 = \text{queryWeight}(\text{anchor:الكريم}^2.0), \text{product of:}$ 
      - $2.0 = \text{boost}$
      - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
      - $0.035606395 = \text{queryNorm}$
    - $0.05734859 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم in } 75), \text{product of:}$ 
      - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
      - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$

- $0.015625 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=75)$
- $1.0217999\text{E}-4 = (\text{MATCH}) \text{weight}(\text{content}:\text{الكريم} \text{ in } 75), \text{ product of:}$ 
  - $0.037867878 = \text{queryWeight}(\text{content}:\text{الكريم}), \text{ product of:}$ 
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.035606395 = \text{queryNorm}$
  - $0.002698329 = (\text{MATCH}) \text{fieldWeight}(\text{content}:\text{الكريم} \text{ in } 75), \text{ product of:}$ 
    - $3.4641016 = \text{tf}(\text{termFreq}(\text{content}:\text{الكريم})=12)$
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $7.324219\text{E}-4 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=75)$

\*\*\*\*\*

Page 7:

- $\text{boost} = 0.0420541$
- $\text{digest} = 3354b6239b6eb27b9d241073f88fc34e$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181031$
- $\text{title} = \text{America.gov} - \text{روص لاب ةيكريم أا ةيالكحلا} - \text{روص موبلأ}$
- $\text{tstamp} = 20100305231110244$
- $\text{url} =$   
 $\text{http://www.america.gov/ar/multimedia/photogallery.html\#/4110/religious\_freedom\_ar/}$

**score for query: الكريم**

- $0.01136245 = (\text{MATCH}) \text{sum of:}$ 
  - $0.01124203 = (\text{MATCH}) \text{weight}(\text{anchor}:\text{الكريم}^{2.0} \text{ in } 52), \text{ product of:}$ 
    - $0.261373 = \text{queryWeight}(\text{anchor}:\text{الكريم}^{2.0}), \text{ product of:}$ 
      - $2.0 = \text{boost}$
      - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
      - $0.035606395 = \text{queryNorm}$

- $0.04301144 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم in 52}),$   
product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
  - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
  - $0.01171875 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=52)$
- $1.20420285\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:الكريم in 52}),$  product of:
  - $0.037867878 = \text{queryWeight}(\text{content:الكريم}),$  product of:
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.035606395 = \text{queryNorm}$
  - $0.0031800114 = (\text{MATCH}) \text{fieldWeight}(\text{content:الكريم in 52}),$   
product of:
    - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.0012207031 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=52)$

\*\*\*\*\*

Page 8:

- $\text{boost} = 0.02675021$
- $\text{digest} = \text{d4493509fb1e3146c2003310c9b70cbd}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181330$
- $\text{title} = \text{America.gov - بتك - بتك}$
- $\text{tstamp} = 20100305231409034$
- $\text{url} = \text{http://www.america.gov/ar/publications/books.html\#governed}$

**score for query: الكريم**

- $0.01132718 = (\text{MATCH}) \text{sum of:}$ 
  - $0.01124203 = (\text{MATCH}) \text{weight}(\text{anchor:الكريم}^2.0 \text{ in 77}),$  product of:
    - $0.261373 = \text{queryWeight}(\text{anchor:الكريم}^2.0),$  product of:
      - $2.0 = \text{boost}$

- $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
- $0.035606395 = \text{queryNorm}$
- $0.04301144 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم in 77}),$   
product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
  - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
  - $0.01171875 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=77)$
- $8.514999\text{E-}5 = (\text{MATCH}) \text{weight}(\text{content:الكريم in 77}),$  product of:
  - $0.037867878 = \text{queryWeight}(\text{content:الكريم}),$  product of:
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.035606395 = \text{queryNorm}$
  - $0.0022486076 = (\text{MATCH}) \text{fieldWeight}(\text{content:الكريم in 77}),$   
product of:
    - $3.4641016 = \text{tf}(\text{termFreq}(\text{content:الكريم})=12)$
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $6.1035156\text{E-}4 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=77)$

\*\*\*\*\*

Page 9:

- $\text{boost} = 0.025411258$
- $\text{digest} = 6b7361561b7255632af783ca69a88410$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181330$
- $\text{title} = \text{America.gov} - \text{روص لاب ةيكريم أا ةيالكحلا} - \text{روص موبلأ}$
- $\text{tstamp} = 20100305231411435$
- $\text{url} = \text{http://www.america.gov/ar/multimedia/photogallery.html\#/4110/islam\_ar/}$

**score for query: الكريم**

- $0.011314282 = (\text{MATCH}) \text{sum of:}$ 
  - $0.01124203 = (\text{MATCH}) \text{weight}(\text{anchor:الكريم}^2.0 \text{ in 49}),$  product of:

- $0.261373 = \text{queryWeight}(\text{anchor:الكريم}^2.0)$ , product of:
  - $2.0 = \text{boost}$
  - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
  - $0.035606395 = \text{queryNorm}$
- $0.04301144 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم in 49})$ , product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
  - $3.6703098 = \text{idf}(\text{docFreq}=8, \text{numDocs}=130)$
  - $0.01171875 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=49)$
- $7.225217\text{E-}5 = (\text{MATCH}) \text{weight}(\text{content:الكريم in 49})$ , product of:
  - $0.037867878 = \text{queryWeight}(\text{content:الكريم})$ , product of:
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.035606395 = \text{queryNorm}$
  - $0.0019080068 = (\text{MATCH}) \text{fieldWeight}(\text{content:الكريم in 49})$ , product of:
    - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $7.324219\text{E-}4 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=49)$

\*\*\*\*\*

Page 10:

- $\text{boost} = 1.0000145$
- $\text{digest} = 0d5b023c802941ddb358071073a98833$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180856$
- $\text{title} = \text{America.gov} - \text{الوأل ءحفصلال} - \text{الوأل ءحفصلال}$
- $\text{tstamp} = 20100305230902835$
- $\text{url} = \text{http://www.america.gov/ar/}$

**score for query: الكريم**

- $0.0030827592 = (\text{MATCH})$  sum of:
  - $0.0030827592 = (\text{MATCH})$  weight(content:الكريم in 0), product of:
    - $0.037867878 = \text{queryWeight}(\text{content:الكريم})$ , product of:
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.035606395 = \text{queryNorm}$
    - $0.08140829 = (\text{MATCH})$  fieldWeight(content:الكريم in 0), product of:
      - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.03125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=0)$

## APPENDIX D

This is the detail score for query of top 10 pages using *NutchDocumentAnalyzer*.

Search Term: الأمريك (America)

Page 1:

- boost = 0.1580853
- digest = 65d01f780ed747de9fd07241fb39df44
- lang = ar
- segment = 20100307101231
- title = America.gov - فوغ تود الأمريك عقوقم لوح - فوغ تود الأمريك عقوقم لوح
- tstamp = 20100307151249455
- url = http://www.america.gov/ar/pages/footer/local/about-us.html

score for query: الأمريك

- 0.11997691 = (MATCH) sum of:
  - 0.060125146 = (MATCH) weight(anchor:الأمريك^2.0 in 69), product of:
    - 0.26155776 = queryWeight(anchor:الأمريك^2.0), product of:
      - 2.0 = boost
      - 3.6779728 = idf(docFreq=8, numDocs=131)
      - 0.035557326 = queryNorm
    - 0.2298733 = (MATCH) fieldWeight(anchor:الأمريك in 69), product of:
      - 1.0 = tf(termFreq(anchor:الأمريك)=1)
      - 3.6779728 = idf(docFreq=8, numDocs=131)
      - 0.0625 = fieldNorm(field=anchor, doc=69)
  - 4.8056475E-4 = (MATCH) weight(content:الأمريك in 69), product of:
    - 0.037797898 = queryWeight(content:الأمريك), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)



- 0.035557326 = queryNorm
- 0.01271406 = (MATCH) fieldWeight(content:أكريم in 69), product of:
  - 2.4494898 = tf(termFreq(content:أكريم)=6)
  - 1.063013 = idf(docFreq=122, numDocs=131)
  - 0.0048828125 = fieldNorm(field=content, doc=69)
- 0.0593712 = (MATCH) weight(title:أكريم^1.5 in 69), product of:
  - 0.23942009 = queryWeight(title:أكريم^1.5), product of:
    - 1.5 = boost
    - 4.488903 = idf(docFreq=3, numDocs=131)
    - 0.035557326 = queryNorm
  - 0.2479792 = (MATCH) fieldWeight(title:أكريم in 69), product of:
    - 1.4142135 = tf(termFreq(title:أكريم)=2)
    - 4.488903 = idf(docFreq=3, numDocs=131)
    - 0.0390625 = fieldNorm(field=title, doc=69)

\*\*\*\*\*

Page 2:

- boost = 0.16184442
- digest = be96f39b462a546d99cbfa50ba70c710
- lang = ar
- segment = 20100307101102
- title = أكريم في نوملسملا - Being Muslim in America - America.gov
- tstamp = 20100307151207698
- url = http://www.america.gov/ar/publications/books-content/musliminamerica.html

**score for query: أكريم**

- 0.11105819 = (MATCH) sum of:
  - 0.060125146 = (MATCH) weight(anchor:أكريم^2.0 in 73), product of:

- $0.26155776 = \text{queryWeight}(\text{anchor:الكريم}^{2.0})$ , product of:
  - $2.0 = \text{boost}$
  - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
  - $0.035557326 = \text{queryNorm}$
- $0.2298733 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم in 73})$ , product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
  - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
  - $0.0625 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=73)$
- $5.5490836\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:الكريم in 73})$ , product of:
  - $0.037797898 = \text{queryWeight}(\text{content:الكريم})$ , product of:
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.035557326 = \text{queryNorm}$
  - $0.014680931 = (\text{MATCH}) \text{fieldWeight}(\text{content:الكريم in 73})$ , product of:
    - $2.828427 = \text{tf}(\text{termFreq}(\text{content:الكريم})=8)$
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=73)$
- $0.050378136 = (\text{MATCH}) \text{weight}(\text{title:الكريم}^{1.5} \text{ in 73})$ , product of:
  - $0.23942009 = \text{queryWeight}(\text{title:الكريم}^{1.5})$ , product of:
    - $1.5 = \text{boost}$
    - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$
    - $0.035557326 = \text{queryNorm}$
  - $0.21041733 = (\text{MATCH}) \text{fieldWeight}(\text{title:الكريم in 73})$ , product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{title:الكريم})=1)$
    - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$
    - $0.046875 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=73)$

\*\*\*\*\*

Page 3:

- boost = 0.23039404
- digest = 8ed8fcd743ff1ce5d4c42db83fc549af
- lang = ar
- segment = 20100307101102
- title = أمريكا - أي حل - أي كريمة أم لا - America.gov
- tstamp = 20100307151136760
- url = http://www.america.gov/ar/amlife.html

score for query: الكريمة

- 0.10594571 = (MATCH) sum of:
  - 0.10521901 = (MATCH) weight(anchor:الكريمة^2.0 in 3), product of:
    - 0.26155776 = queryWeight(anchor:الكريمة^2.0), product of:
      - 2.0 = boost
      - 3.6779728 = idf(docFreq=8, numDocs=131)
      - 0.035557326 = queryNorm
    - 0.40227827 = (MATCH) fieldWeight(anchor:الكريمة in 3), product of:
      - 1.0 = tf(termFreq(anchor:الكريمة)=1)
      - 3.6779728 = idf(docFreq=8, numDocs=131)
      - 0.109375 = fieldNorm(field=anchor, doc=3)
  - 7.2669686E-4 = (MATCH) weight(content:الكريمة in 3), product of:
    - 0.037797898 = queryWeight(content:الكريمة), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.035557326 = queryNorm
    - 0.019225854 = (MATCH) fieldWeight(content:الكريمة in 3), product of:

- $2.6457512 = \text{tf}(\text{termFreq}(\text{content:الكريم})=7)$
- $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
- $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=3)$

\*\*\*\*\*

Page 4:

- $\text{boost} = 0.1587577$
- $\text{digest} = \text{a5795145f4a839cf52528d1b49e03bd1}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101102$
- $\text{title} = \text{ربع نشييدي إلي ابوم فوغ تود الكريم لأوجل أةخسن - نشييدي إلي ابوم فم د خ PDA - America.gov ةزهجأ}$
- $\text{tstamp} = 20100307151139253$
- $\text{url} = \text{http://www.america.gov/ar/services/mobile.html}$

**score for query: الكريم**

- $0.042462345 = (\text{MATCH}) \text{ sum of:}$ 
  - $4.8056475\text{E-}4 = (\text{MATCH}) \text{ weight}(\text{content:الكريم in 115}), \text{ product of:}$ 
    - $0.037797898 = \text{queryWeight}(\text{content:الكريم}), \text{ product of:}$ 
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.035557326 = \text{queryNorm}$
    - $0.01271406 = (\text{MATCH}) \text{ fieldWeight}(\text{content:الكريم in 115}), \text{ product of:}$ 
      - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=115)$
  - $0.04198178 = (\text{MATCH}) \text{ weight}(\text{title:الكريم}^{1.5} \text{ in 115}), \text{ product of:}$ 
    - $0.23942009 = \text{queryWeight}(\text{title:الكريم}^{1.5}), \text{ product of:}$ 
      - $1.5 = \text{boost}$
      - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$

- $0.035557326 = \text{queryNorm}$
- $0.17534778 = (\text{MATCH}) \text{fieldWeight}(\text{title:الكريم} \text{ in } 115), \text{ product of:}$ 
  - $1.0 = \text{tf}(\text{termFreq}(\text{title:الكريم})=1)$
  - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$
  - $0.0390625 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=115)$

\*\*\*\*\*

Page 5:

- $\text{boost} = 0.04832446$
- $\text{digest} = 22e534f031e9c7ac8682fcd4f86523e4$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101231$
- $\text{title} = \text{America.gov} - \text{روص لاب ةيكريم أا ةيالكحل} - \text{روص موبل أ}$
- $\text{tstamp} = 20100307151334977$
- $\text{url} =$   
 $\text{http://www.america.gov/ar/multimedia/photogallery.html\#/4110/mosques\_ar/}$

**score for query: الكريم**

- $0.022691099 = (\text{MATCH}) \text{sum of:}$ 
  - $0.02254693 = (\text{MATCH}) \text{weight}(\text{anchor:الكريم}^2.0 \text{ in } 51), \text{ product of:}$ 
    - $0.26155776 = \text{queryWeight}(\text{anchor:الكريم}^2.0), \text{ product of:}$ 
      - $2.0 = \text{boost}$
      - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
      - $0.035557326 = \text{queryNorm}$
    - $0.08620249 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم} \text{ in } 51), \text{ product of:}$ 
      - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
      - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
      - $0.0234375 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=51)$

- $1.4416942E-4 = (\text{MATCH}) \text{ weight}(\text{content:الكريم in 51})$ , product of:
  - $0.037797898 = \text{queryWeight}(\text{content:الكريم})$ , product of:
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.035557326 = \text{queryNorm}$
  - $0.0038142179 = (\text{MATCH}) \text{ fieldWeight}(\text{content:الكريم in 51})$ , product of:
    - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.0014648438 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=51)$

\*\*\*\*\*

Page 6:

- $\text{boost} = 0.033444975$
- $\text{digest} = 80c97402726fad635131db1bb29555be$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101231$
- $\text{title} = \text{America.gov} - \text{بتك} - \text{بتك}$
- $\text{tstamp} = 20100307151330985$
- $\text{url} = \text{http://www.america.gov/ar/publications/books.html\#beingmuslim}$

**score for query: الكريم**

- $0.01513323 = (\text{MATCH}) \text{ sum of:}$ 
  - $0.0150312865 = (\text{MATCH}) \text{ weight}(\text{anchor:الكريم}^2.0 \text{ in } 76)$ , product of:
    - $0.26155776 = \text{queryWeight}(\text{anchor:الكريم}^2.0)$ , product of:
      - $2.0 = \text{boost}$
      - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
      - $0.035557326 = \text{queryNorm}$
    - $0.057468325 = (\text{MATCH}) \text{ fieldWeight}(\text{anchor:الكريم in } 76)$ , product of:
      - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$

- $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
- $0.015625 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=76)$
- $1.01943166\text{E}-4 = (\text{MATCH}) \text{weight}(\text{content}:\text{الكريم} \text{ in } 76), \text{ product of:}$ 
  - $0.037797898 = \text{queryWeight}(\text{content}:\text{الكريم}), \text{ product of:}$ 
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.035557326 = \text{queryNorm}$
  - $0.0026970592 = (\text{MATCH}) \text{fieldWeight}(\text{content}:\text{الكريم} \text{ in } 76), \text{ product of:}$ 
    - $3.4641016 = \text{tf}(\text{termFreq}(\text{content}:\text{الكريم})=12)$
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $7.324219\text{E}-4 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=76)$

\*\*\*\*\*

Page 7:

- $\text{boost} = 0.0420541$
- $\text{digest} = 1e1bc6ad9ffbfcd82ea012b44610bed$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101231$
- $\text{title} = \text{America.gov} - \text{روص ل اب ةي ك ريم أ ل ا ةي الك حل ا} - \text{روص موب ل أ}$
- $\text{tstamp} = 20100307151307072$
- $\text{url} =$   
 $\text{http://www.america.gov/ar/multimedia/photogallery.html\#/4110/religious\_freedom\_ar/}$

**score for query: الكريم**

- $0.011393607 = (\text{MATCH}) \text{sum of:}$ 
  - $0.011273465 = (\text{MATCH}) \text{weight}(\text{anchor}:\text{الكريم}^2.0 \text{ in } 53), \text{ product of:}$ 
    - $0.26155776 = \text{queryWeight}(\text{anchor}:\text{الكريم}^2.0), \text{ product of:}$ 
      - $2.0 = \text{boost}$
      - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$

- $0.035557326 = \text{queryNorm}$
- $0.043101244 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم in 53})$ , product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
  - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
  - $0.01171875 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=53)$
- $1.2014119\text{E}-4 = (\text{MATCH}) \text{weight}(\text{content:الكريم in 53})$ , product of:
  - $0.037797898 = \text{queryWeight}(\text{content:الكريم})$ , product of:
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.035557326 = \text{queryNorm}$
  - $0.003178515 = (\text{MATCH}) \text{fieldWeight}(\text{content:الكريم in 53})$ , product of:
    - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.0012207031 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=53)$

\*\*\*\*\*

Page 8:

- $\text{boost} = 0.02675021$
- $\text{digest} = \text{df2eeaf879a60aaaddf7c8403cba7fa}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101458$
- $\text{title} = \text{America.gov} - \text{بتك} - \text{بتك}$
- $\text{tstamp} = 20100307151541037$
- $\text{url} = \text{http://www.america.gov/ar/publications/books.html\#governed}$

**score for query: الكريم**

- $0.011358418 = (\text{MATCH}) \text{sum of:}$ 
  - $0.011273465 = (\text{MATCH}) \text{weight}(\text{anchor:الكريم}^{2.0} \text{ in 78})$ , product of:
    - $0.26155776 = \text{queryWeight}(\text{anchor:الكريم}^{2.0})$ , product of:



- 2.0 = boost
- 3.6779728 = idf(docFreq=8, numDocs=131)
- 0.035557326 = queryNorm
- 0.043101244 = (MATCH) fieldWeight(anchor:أكريم in 78), product of:
  - 1.0 = tf(termFreq(anchor:أكريم)=1)
  - 3.6779728 = idf(docFreq=8, numDocs=131)
  - 0.01171875 = fieldNorm(field=anchor, doc=78)
- 8.495264E-5 = (MATCH) weight(content:أكريم in 78), product of:
  - 0.037797898 = queryWeight(content:أكريم), product of:
    - 1.063013 = idf(docFreq=122, numDocs=131)
    - 0.035557326 = queryNorm
  - 0.0022475494 = (MATCH) fieldWeight(content:أكريم in 78), product of:
    - 3.4641016 = tf(termFreq(content:أكريم)=12)
    - 1.063013 = idf(docFreq=122, numDocs=131)
    - 6.1035156E-4 = fieldNorm(field=content, doc=78)

\*\*\*\*\*

Page 9:

- boost = 0.025411258
- digest = 295971814b3454a9d44144054b5c194a
- lang = ar
- segment = 20100307101458
- title = روصلاب ةيكريمأل ةيالكحلأ - روص موبلأ - America.gov
- tstamp = 20100307151543423
- url = http://www.america.gov/ar/multimedia/photogallery.html#/4110/islam\_ar/

score for query: أكريم

- 0.0113455495 = (MATCH) sum of:

- $0.011273465 = (\text{MATCH}) \text{weight}(\text{anchor:الكريم}^2.0 \text{ in } 50)$ , product of:
  - $0.26155776 = \text{queryWeight}(\text{anchor:الكريم}^2.0)$ , product of:
    - $2.0 = \text{boost}$
    - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
    - $0.035557326 = \text{queryNorm}$
  - $0.043101244 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:الكريم} \text{ in } 50)$ , product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:الكريم})=1)$
    - $3.6779728 = \text{idf}(\text{docFreq}=8, \text{numDocs}=131)$
    - $0.01171875 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=50)$
- $7.208471\text{E}-5 = (\text{MATCH}) \text{weight}(\text{content:الكريم} \text{ in } 50)$ , product of:
  - $0.037797898 = \text{queryWeight}(\text{content:الكريم})$ , product of:
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.035557326 = \text{queryNorm}$
  - $0.0019071089 = (\text{MATCH}) \text{fieldWeight}(\text{content:الكريم} \text{ in } 50)$ , product of:
    - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:الكريم})=6)$
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $7.324219\text{E}-4 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=50)$

\*\*\*\*\*

Page 10:

- $\text{boost} = 1.0000145$
- $\text{digest} = \text{eed4dd9817b50ffda0aef158be6e4c12}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101052$
- $\text{title} = \text{America.gov - لىلوالا ءحفصلال - لىلوالا ءحفصلال}$
- $\text{tstamp} = 20100307151057483$
- $\text{url} = \text{http://www.america.gov/ar/}$

score for query: الكريّم

- $0.0028076388 = (\text{MATCH})$  sum of:
  - $0.0028076388 = (\text{MATCH})$  weight(content:الكريّم in 0), product of:
    - $0.037797898 = \text{queryWeight}(\text{content:الكريّم})$ , product of:
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.035557326 = \text{queryNorm}$
    - $0.07428029 = (\text{MATCH})$  fieldWeight(content:الكريّم in 0), product of:
      - $2.236068 = \text{tf}(\text{termFreq}(\text{content:الكريّم})=5)$
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.03125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=0)$

## APPENDIX E

This is the detail score for query of top 10 pages using *ArabicAnalyzer*.

Search Term: طارق مديد (Democratic)

Page 1:

- boost = 0.16689056
- digest = 6e1b0463970c5b60bb75636a698cf1b3
- lang = ar
- segment = 20100305180909
- title = طارق مديد - طارق مديد - America.gov
- tstamp = 20100305230951886
- url = http://www.america.gov/ar/global/democracy.html

**score for query:** طارق مديد

- 0.2665834 = (MATCH) sum of:
  - 0.15995954 = (MATCH) weight(anchor:طارق مديد^2.0 in 23), product of:
    - 0.30052778 = queryWeight(anchor:طارق مديد^2.0), product of:
      - 2.0 = boost
      - 4.2580967 = idf(docFreq=4, numDocs=130)
      - 0.035288982 = queryNorm
    - 0.5322621 = (MATCH) fieldWeight(anchor:طارق مديد in 23), product of:
      - 1.0 = tf(termFreq(anchor:طارق مديد)=1)
      - 4.2580967 = idf(docFreq=4, numDocs=130)
      - 0.125 = fieldNorm(field=anchor, doc=23)
  - 5.846775E-4 = (MATCH) weight(content:طارق مديد in 23), product of:
    - 0.037530307 = queryWeight(content:طارق مديد), product of:
      - 1.0635134 = idf(docFreq=121, numDocs=130)

- $0.035288982 = \text{queryNorm}$
- $0.01557881 = (\text{MATCH}) \text{fieldWeight}(\text{content:طارق مديد in 23}),$  product of:
  - $3.0 = \text{tf}(\text{termFreq}(\text{content:طارق مديد})=9)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=23)$
- $0.1060392 = (\text{MATCH}) \text{weight}(\text{title:طارق مديد}^{1.5} \text{ in 23}),$  product of:
  - $0.22539584 = \text{queryWeight}(\text{title:طارق مديد}^{1.5}),$  product of:
    - $1.5 = \text{boost}$
    - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
    - $0.035288982 = \text{queryNorm}$
  - $0.47045767 = (\text{MATCH}) \text{fieldWeight}(\text{title:طارق مديد in 23}),$  product of:
    - $1.4142135 = \text{tf}(\text{termFreq}(\text{title:طارق مديد})=2)$
    - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
    - $0.078125 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=23)$

\*\*\*\*\*

Page 2:

- $\text{boost} = 0.23113073$
- $\text{digest} = 5285dc46473be73851750b409de012a5$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180909$
- $\text{title} = \text{America.gov} - \text{يمل اعل يدحتل} - \text{يمل اعل يدحتل}$
- $\text{tstamp} = 20100305230921085$
- $\text{url} = \text{http://www.america.gov/ar/global.html}$

**score for query:** طارق مديد

- $0.16062789 = (\text{MATCH}) \text{sum of:}$ 
  - $0.15995954 = (\text{MATCH}) \text{weight}(\text{anchor:طارق مديد}^{2.0} \text{ in 22}),$  product of:

- $0.30052778 = \text{queryWeight}(\text{anchor:طارق م ي د}^{2.0})$ , product of:
  - $2.0 = \text{boost}$
  - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
  - $0.035288982 = \text{queryNorm}$
- $0.5322621 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:طارق م ي د in 22})$ , product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:طارق م ي د})=1)$
  - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
  - $0.125 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=22)$
- $6.683421\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:طارق م ي د in 22})$ , product of:
  - $0.037530307 = \text{queryWeight}(\text{content:طارق م ي د})$ , product of:
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.035288982 = \text{queryNorm}$
  - $0.017808065 = (\text{MATCH}) \text{fieldWeight}(\text{content:طارق م ي د in 22})$ , product of:
    - $2.4494898 = \text{tf}(\text{termFreq}(\text{content:طارق م ي د})=6)$
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=22)$

\*\*\*\*\*

Page 3:

- $\text{boost} = 0.031816483$
- $\text{digest} = \text{bba906c38386b2e71f42a4f7d365e8cb}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181031$
- $\text{title} = \text{America.gov - ةي طارق م ي د ل ا و ق ا و س أ ل ا - ةي طارق م ي د ل ا و ق ا و س أ ل ا}$
- $\text{tstamp} = 20100305231058196$
- $\text{url} = \text{http://www.america.gov/ar/publications/ejournalusa/608.html}$

**score for query:** طارق م ي د

- $0.033635326 = (\text{MATCH}) \text{sum of:}$

- $0.017495574 = (\text{MATCH}) \text{weight}(\text{anchor:طارق مريد}^{2.0} \text{ in } 98)$ , product of:
  - $0.30052778 = \text{queryWeight}(\text{anchor:طارق مريد}^{2.0})$ , product of:
    - $2.0 = \text{boost}$
    - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
    - $0.035288982 = \text{queryNorm}$
  - $0.058216166 = (\text{MATCH}) \text{fieldWeight}(\text{anchor:طارق مريد} \text{ in } 98)$ , product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:طارق مريد})=1)$
    - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
    - $0.013671875 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=98)$
- $2.33871\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:طارق مريد} \text{ in } 98)$ , product of:
  - $0.037530307 = \text{queryWeight}(\text{content:طارق مريد})$ , product of:
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $0.035288982 = \text{queryNorm}$
  - $0.006231524 = (\text{MATCH}) \text{fieldWeight}(\text{content:طارق مريد} \text{ in } 98)$ , product of:
    - $6.0 = \text{tf}(\text{termFreq}(\text{content:طارق مريد})=36)$
    - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
    - $9.765625\text{E-}4 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=98)$
- $0.015905881 = (\text{MATCH}) \text{weight}(\text{title:طارق مريد}^{1.5} \text{ in } 98)$ , product of:
  - $0.22539584 = \text{queryWeight}(\text{title:طارق مريد}^{1.5})$ , product of:
    - $1.5 = \text{boost}$
    - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
    - $0.035288982 = \text{queryNorm}$
  - $0.07056865 = (\text{MATCH}) \text{fieldWeight}(\text{title:طارق مريد} \text{ in } 98)$ , product of:
    - $1.4142135 = \text{tf}(\text{termFreq}(\text{title:طارق مريد})=2)$
    - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$

- $0.01171875 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=98)$

\*\*\*\*\*

Page 4:

- $\text{boost} = 0.11378951$
- $\text{digest} = \text{b8c157220365a4bf104bc045832885be}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180909$
- $\text{title} = \text{0110 - ةيطارقميدلا مظنلا يف مكحلا لقتني فيك :تاباختنا نم رثكأ - America.gov}$
- $\text{tstamp} = 20100305231013594$
- $\text{url} = \text{http://www.america.gov/ar/publications/ejournalusa/0110.html}$

**score for query:** طارقميد

- $0.030587077 = (\text{MATCH}) \text{ sum of:}$ 
  - $5.9466175\text{E-}4 = (\text{MATCH}) \text{ weight}(\text{content:طارقميد in 88}), \text{ product of:}$ 
    - $0.037530307 = \text{queryWeight}(\text{content:طارقميد}), \text{ product of:}$ 
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.035288982 = \text{queryNorm}$
    - $0.01584484 = (\text{MATCH}) \text{ fieldWeight}(\text{content:طارقميد in 88}), \text{ product of:}$ 
      - $4.358899 = \text{tf}(\text{termFreq}(\text{content:طارقميد})=19)$
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.0034179688 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=88)$
  - $0.029992415 = (\text{MATCH}) \text{ weight}(\text{title:طارقميد}^{1.5} \text{ in 88}), \text{ product of:}$ 
    - $0.22539584 = \text{queryWeight}(\text{title:طارقميد}^{1.5}), \text{ product of:}$ 
      - $1.5 = \text{boost}$
      - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
      - $0.035288982 = \text{queryNorm}$
    - $0.13306552 = (\text{MATCH}) \text{ fieldWeight}(\text{title:طارقميد in 88}), \text{ product of:}$



- $1.0 = \text{tf}(\text{termFreq}(\text{title:طارق مـيد})=1)$
- $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
- $0.03125 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=88)$

\*\*\*\*\*

Page 5:

- $\text{boost} = 0.028445216$
- $\text{digest} = \text{c04e43d37fb6f380a397373427882a1e}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305181151$
- $\text{title} = \text{America.gov - ريب عتلا ةيـرح مـل نون طاوم | مل اعلـا يـف ةيـطارق مـيدلا - ةمدقم}$
- $\text{tstamp} = 20100305231252420$
- $\text{url} = \text{http://www.america.gov/ar/democracy/global/index.html}$

**score for query:** طارق مـيد

- $0.027611194 = (\text{MATCH}) \text{ sum of:}$ 
  - $0.019994942 = (\text{MATCH}) \text{ weight}(\text{anchor:طارق مـيد}^{2.0} \text{ in } 14), \text{ product of:}$ 
    - $0.30052778 = \text{queryWeight}(\text{anchor:طارق مـيد}^{2.0}), \text{ product of:}$ 
      - $2.0 = \text{boost}$
      - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
      - $0.035288982 = \text{queryNorm}$
    - $0.06653276 = (\text{MATCH}) \text{ fieldWeight}(\text{anchor:طارق مـيد} \text{ in } 14), \text{ product of:}$ 
      - $1.0 = \text{tf}(\text{termFreq}(\text{anchor:طارق مـيد})=1)$
      - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
      - $0.015625 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=14)$
  - $1.181473\text{E-}4 = (\text{MATCH}) \text{ weight}(\text{content:طارق مـيد} \text{ in } 14), \text{ product of:}$ 
    - $0.037530307 = \text{queryWeight}(\text{content:طارق مـيد}), \text{ product of:}$ 
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.035288982 = \text{queryNorm}$

- $0.0031480505 = (\text{MATCH}) \text{fieldWeight}(\text{content:طارق مديد in 14}),$  product of:
  - $3.4641016 = \text{tf}(\text{termFreq}(\text{content:طارق مديد})=12)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $8.544922\text{E-}4 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=14)$
- $0.0074981037 = (\text{MATCH}) \text{weight}(\text{title:طارق مديد}^{1.5} \text{ in 14}),$  product of:
  - $0.22539584 = \text{queryWeight}(\text{title:طارق مديد}^{1.5}),$  product of:
    - $1.5 = \text{boost}$
    - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
    - $0.035288982 = \text{queryNorm}$
  - $0.03326638 = (\text{MATCH}) \text{fieldWeight}(\text{title:طارق مديد in 14}),$  product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{title:طارق مديد})=1)$
    - $4.2580967 = \text{idf}(\text{docFreq}=4, \text{numDocs}=130)$
    - $0.0078125 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=14)$

\*\*\*\*\*

Page 6:

- $\text{boost} = 1.0000145$
- $\text{digest} = 0\text{d}5\text{b}023\text{c}802941\text{ddb}358071073\text{a}98833$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180856$
- $\text{title} = \text{America.gov} - \text{ىل وأل ةحفصل} - \text{ىل وأل ةحفصل}$
- $\text{tstamp} = 20100305230902835$
- $\text{url} = \text{http://www.america.gov/ar/}$

**score for query:** طارق مديد

- $0.0021604078 = (\text{MATCH}) \text{sum of:}$ 
  - $0.0021604078 = (\text{MATCH}) \text{weight}(\text{content:طارق مديد in 0}),$  product of:
    - $0.037530307 = \text{queryWeight}(\text{content:طارق مديد}),$  product of:
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$

- $0.035288982 = \text{queryNorm}$
- $0.05756435 = (\text{MATCH}) \text{fieldWeight}(\text{content:طارق م ي د in 0}),$   
product of:
  - $1.7320508 = \text{tf}(\text{termFreq}(\text{content:طارق م ي د})=3)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $0.03125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=0)$

\*\*\*\*\*

Page 7:

- $\text{boost} = 0.22860475$
- $\text{digest} = 5a62cd3a20d5393ff5806fd92af1edef$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180909$
- $\text{title} = \text{تس الكدوب - تس الكدوب - America.gov}$
- $\text{tstamp} = 20100305230934690$
- $\text{url} = \text{http://www.america.gov/ar/multimedia/podcast.html}$

**score for query:** طارق م ي د

- $6.1011006\text{E-}4 = (\text{MATCH}) \text{sum of:}$ 
  - $6.1011006\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:طارق م ي د in 60}), \text{product of:}$ 
    - $0.037530307 = \text{queryWeight}(\text{content:طارق م ي د}), \text{product of:}$ 
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.035288982 = \text{queryNorm}$
    - $0.016256463 = (\text{MATCH}) \text{fieldWeight}(\text{content:طارق م ي د in 60}),$   
product of:
      - $2.236068 = \text{tf}(\text{termFreq}(\text{content:طارق م ي د})=5)$
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=60)$

\*\*\*\*\*

Page 8:

- $\text{boost} = 0.22996004$
- $\text{digest} = a0130240b4348578aa8a83e59187dfb3$

- lang = ar
- segment = 20100305180909
- title = بٲٲك - بٲٲك - America.gov
- tstamp = 20100305231001279
- url = http://www.america.gov/ar/publications/books.html

**score for query:** طارقم يد

- 5.846775E-4 = (MATCH) sum of:
  - 5.846775E-4 = (MATCH) weight(content:طارقم يد in 73), product of:
    - 0.037530307 = queryWeight(content:طارقم يد), product of:
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 0.035288982 = queryNorm
    - 0.01557881 = (MATCH) fieldWeight(content:طارقم يد in 73), product of:
      - 3.0 = tf(termFreq(content:طارقم يد)=9)
      - 1.0635134 = idf(docFreq=121, numDocs=130)
      - 0.0048828125 = fieldNorm(field=content, doc=73)

\*\*\*\*\*

Page 9:

- boost = 0.23032264
- digest = ce4a12d589c1a56e886d5b6848609391
- lang = ar
- segment = 20100305180909
- title = ةيكر يمأل ا ةاي حل - ةيكر يمأل ا ةاي حل - America.gov
- tstamp = 20100305230939904
- url = http://www.america.gov/ar/amlife.html

**score for query:** طارقم يد

- 5.4569903E-4 = (MATCH) sum of:
  - 5.4569903E-4 = (MATCH) weight(content:طارقم يد in 3), product of:
    - 0.037530307 = queryWeight(content:طارقم يد), product of:

- $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
- $0.035288982 = \text{queryNorm}$
- $0.0145402225 = (\text{MATCH}) \text{fieldWeight}(\text{content:طارق مديد in 3}),$   
product of:
  - $2.0 = \text{tf}(\text{termFreq}(\text{content:طارق مديد})=4)$
  - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
  - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=3)$

\*\*\*\*\*

Page 10:

- $\text{boost} = 0.2296042$
- $\text{digest} = \text{bc9c562d0a61b335f5a8730f14412dcb}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100305180909$
- $\text{title} = \text{America.gov} - \text{هي آسأ وي لانروج ي} - \text{هي آسأ وي لانروج ي}$
- $\text{tstamp} = 20100305230924025$
- $\text{url} = \text{http://www.america.gov/ar/publications/ejournalusa.html}$

**score for query:** طارق مديد

- $5.4010196\text{E-}4 = (\text{MATCH}) \text{sum of:}$ 
  - $5.4010196\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:طارق مديد in 86}), \text{product of:}$ 
    - $0.037530307 = \text{queryWeight}(\text{content:طارق مديد}), \text{product of:}$ 
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.035288982 = \text{queryNorm}$
    - $0.014391088 = (\text{MATCH}) \text{fieldWeight}(\text{content:طارق مديد in 86}),$   
product of:
      - $3.4641016 = \text{tf}(\text{termFreq}(\text{content:طارق مديد})=12)$
      - $1.0635134 = \text{idf}(\text{docFreq}=121, \text{numDocs}=130)$
      - $0.00390625 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=86)$

## APPENDIX F

This is the detail score for query of top 10 pages using *NutchDocumentAnalyzer*.

Search Term: ةي طارق م يدل (Democratic)

Page 1:

- boost = 0.16692342
- digest = c49d13e1fa4eb518258862a27800f398
- lang = ar
- segment = 20100307101102
- title = ةي طارق م يدل - ةي طارق م يدل - America.gov
- tstamp = 20100307151151020
- url = http://www.america.gov/ar/global/democracy.html

**score for query:** ةي طارق م يدل

- 0.29354417 = (MATCH) sum of:
  - 0.17619587 = (MATCH) weight(anchor: ةي طارق م يدل^2.0 in 24), product of:
    - 0.31401145 = queryWeight(anchor: ةي طارق م يدل^2.0), product of:
      - 2.0 = boost
      - 4.488903 = idf(docFreq=3, numDocs=131)
      - 0.03497641 = queryNorm
    - 0.5611129 = (MATCH) fieldWeight(anchor: ةي طارق م يدل in 24), product of:
      - 1.0 = tf(termFreq(anchor: ةي طارق م يدل)=1)
      - 4.488903 = idf(docFreq=3, numDocs=131)
      - 0.125 = fieldNorm(field=anchor, doc=24)
  - 5.458426E-4 = (MATCH) weight(content: ةي طارق م يدل in 24), product of:
    - 0.03718038 = queryWeight(content: ةي طارق م يدل), product of:

- $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
- $0.03497641 = \text{queryNorm}$
- $0.014680931 = (\text{MATCH}) \text{fieldWeight}(\text{content}: \text{ةيطارق م يدل} \text{ in } 24),$   
product of:
  - $2.828427 = \text{tf}(\text{termFreq}(\text{content}: \text{ةيطارق م يدل})=8)$
  - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
  - $0.0048828125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=24)$
- $0.116802454 = (\text{MATCH}) \text{weight}(\text{title}: \text{ةيطارق م يدل}^{1.5} \text{ in } 24),$  product of:
  - $0.23550858 = \text{queryWeight}(\text{title}: \text{ةيطارق م يدل}^{1.5}),$  product of:
    - $1.5 = \text{boost}$
    - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$
    - $0.03497641 = \text{queryNorm}$
  - $0.4959584 = (\text{MATCH}) \text{fieldWeight}(\text{title}: \text{ةيطارق م يدل} \text{ in } 24),$   
product of:
    - $1.4142135 = \text{tf}(\text{termFreq}(\text{title}: \text{ةيطارق م يدل})=2)$
    - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$
    - $0.078125 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=24)$

\*\*\*\*\*

Page 2:

- $\text{boost} = 0.23117816$
- $\text{digest} = 6f317cffadc06ecf85513b0eb565f1b8$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101102$
- $\text{title} = \text{America.gov} - \text{يملا عل يدحتلا} - \text{يملا عل يدحتلا}$
- $\text{tstamp} = 20100307151115329$
- $\text{url} = \text{http://www.america.gov/ar/global.html}$

**score for query:** ةيطارق م يدل

- $0.17680001 = (\text{MATCH}) \text{sum of:}$

- $0.17619587 = (\text{MATCH}) \text{ weight}(\text{anchor: إةيطارق مدي دل}^2.0 \text{ in } 23)$ , product of:
  - $0.31401145 = \text{queryWeight}(\text{anchor: إةيطارق مدي دل}^2.0)$ , product of:
    - $2.0 = \text{boost}$
    - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$
    - $0.03497641 = \text{queryNorm}$
  - $0.5611129 = (\text{MATCH}) \text{ fieldWeight}(\text{anchor: إةيطارق مدي دل} \text{ in } 23)$ , product of:
    - $1.0 = \text{tf}(\text{termFreq}(\text{anchor: إةيطارق مدي دل})=1)$
    - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$
    - $0.125 = \text{fieldNorm}(\text{field}=\text{anchor}, \text{doc}=23)$
- $6.041371\text{E-}4 = (\text{MATCH}) \text{ weight}(\text{content: إةيطارق مدي دل} \text{ in } 23)$ , product of:
  - $0.03718038 = \text{queryWeight}(\text{content: إةيطارق مدي دل})$ , product of:
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.03497641 = \text{queryNorm}$
  - $0.016248815 = (\text{MATCH}) \text{ fieldWeight}(\text{content: إةيطارق مدي دل} \text{ in } 23)$ , product of:
    - $2.236068 = \text{tf}(\text{termFreq}(\text{content: إةيطارق مدي دل})=5)$
    - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
    - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=23)$

\*\*\*\*\*

Page 3:

- $\text{boost} = 0.11378951$
- $\text{digest} = \text{ab333ad468abf764c43637fe53b7e4f7}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101102$
- $\text{title} = \text{إةيطارق مدي دل مظن لا يف مك حل لقتني فيك :تاباختنا نم رثكأ - 0110 - America.gov}$
- $\text{tstamp} = 20100307151214969$



- url = <http://www.america.gov/ar/publications/ejournalusa/0110.html>

score for query: **ةي ط ا ر ق م ي د ل ا**

- 0.033559922 = (MATCH) sum of:
  - 5.2319805E-4 = (MATCH) weight(content: **ةي ط ا ر ق م ي د ل ا** in 89), product of:
    - 0.03718038 = queryWeight(content: **ةي ط ا ر ق م ي د ل ا**), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.03497641 = queryNorm
    - 0.0140718855 = (MATCH) fieldWeight(content: **ةي ط ا ر ق م ي د ل ا** in 89), product of:
      - 3.8729835 = tf(termFreq(content: **ةي ط ا ر ق م ي د ل ا**)=15)
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.0034179688 = fieldNorm(field=content, doc=89)
  - 0.033036724 = (MATCH) weight(title: **ةي ط ا ر ق م ي د ل ا**^1.5 in 89), product of:
    - 0.23550858 = queryWeight(title: **ةي ط ا ر ق م ي د ل ا**^1.5), product of:
      - 1.5 = boost
      - 4.488903 = idf(docFreq=3, numDocs=131)
      - 0.03497641 = queryNorm
    - 0.14027822 = (MATCH) fieldWeight(title: **ةي ط ا ر ق م ي د ل ا** in 89), product of:
      - 1.0 = tf(termFreq(title: **ةي ط ا ر ق م ي د ل ا**)=1)
      - 4.488903 = idf(docFreq=3, numDocs=131)
      - 0.03125 = fieldNorm(field=title, doc=89)

\*\*\*\*\*

Page 4:

- boost = 0.028445216
- digest = 9212154ec8740ad77458648f74aa149c
- lang = ar
- segment = 20100307101343
- title = **امريكا.كوم - ريب عتلا ةير ح مهل نون طاوم | مل ا عل ا يف ةي ط ا ر ق م ي د ل ا - ةم د ق م**

- tstamp = 20100307151423606
- url = http://www.america.gov/ar/democracy/global/index.html

**score for query:** ةيطارق مدي دل ا

- 0.030395675 = (MATCH) sum of:
  - 0.022024484 = (MATCH) weight(anchor: ةيطارق مدي دل ا^2.0 in 15), product of:
    - 0.31401145 = queryWeight(anchor: ةيطارق مدي دل ا^2.0), product of:
      - 2.0 = boost
      - 4.488903 = idf(docFreq=3, numDocs=131)
      - 0.03497641 = queryNorm
    - 0.07013911 = (MATCH) fieldWeight(anchor: ةيطارق مدي دل ا in 15), product of:
      - 1.0 = tf(termFreq(anchor: ةيطارق مدي دل ا)=1)
      - 4.488903 = idf(docFreq=3, numDocs=131)
      - 0.015625 = fieldNorm(field=anchor, doc=15)
  - 1.12010006E-4 = (MATCH) weight(content: ةيطارق مدي دل ا in 15), product of:
    - 0.03718038 = queryWeight(content: ةيطارق مدي دل ا), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.03497641 = queryNorm
    - 0.0030126106 = (MATCH) fieldWeight(content: ةيطارق مدي دل ا in 15), product of:
      - 3.3166249 = tf(termFreq(content: ةيطارق مدي دل ا)=11)
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 8.544922E-4 = fieldNorm(field=content, doc=15)
  - 0.008259181 = (MATCH) weight(title: ةيطارق مدي دل ا^1.5 in 15), product of:
    - 0.23550858 = queryWeight(title: ةيطارق مدي دل ا^1.5), product of:
      - 1.5 = boost
      - 4.488903 = idf(docFreq=3, numDocs=131)

- $0.03497641 = \text{queryNorm}$
- $0.035069555 = (\text{MATCH}) \text{fieldWeight}(\text{title:إي طارقم يدل in 15}),$   
product of:
  - $1.0 = \text{tf}(\text{termFreq}(\text{title:إي طارقم يدل})=1)$
  - $4.488903 = \text{idf}(\text{docFreq}=3, \text{numDocs}=131)$
  - $0.0078125 = \text{fieldNorm}(\text{field}=\text{title}, \text{doc}=15)$

\*\*\*\*\*

Page 5:

- $\text{boost} = 1.0000145$
- $\text{digest} = \text{eed4dd9817b50ffda0aef158be6e4c12}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101052$
- $\text{title} = \text{America.gov - إىل وأل ةحفصل - إىل وأل ةحفصل}$
- $\text{tstamp} = 20100307151057483$
- $\text{url} = \text{http://www.america.gov/ar/}$

**score for query:**  $\text{إي طارقم يدل}$

- $0.0021392573 = (\text{MATCH}) \text{sum of:}$ 
  - $0.0021392573 = (\text{MATCH}) \text{weight}(\text{content:إي طارقم يدل in 0}), \text{product of:}$ 
    - $0.03718038 = \text{queryWeight}(\text{content:إي طارقم يدل}), \text{product of:}$ 
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.03497641 = \text{queryNorm}$
    - $0.05753726 = (\text{MATCH}) \text{fieldWeight}(\text{content:إي طارقم يدل in 0}),$   
product of:
      - $1.7320508 = \text{tf}(\text{termFreq}(\text{content:إي طارقم يدل})=3)$
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.03125 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=0)$

\*\*\*\*\*

Page 6:

- $\text{boost} = 0.22865272$
- $\text{digest} = \text{dc127d214554a59575782c318462f4e8}$

- lang = ar
- segment = 20100307101102
- title = تس الكدوب - تس الكدوب - America.gov
- tstamp = 20100307151128611
- url = http://www.america.gov/ar/multimedia/podcast.html

**score for query:** ةيط ارقم يدل ا

- 6.041371E-4 = (MATCH) sum of:
  - 6.041371E-4 = (MATCH) weight(content:ةيط ارقم يدل ا in 61), product of:
    - 0.03718038 = queryWeight(content:ةيط ارقم يدل ا), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.03497641 = queryNorm
    - 0.016248815 = (MATCH) fieldWeight(content:ةيط ارقم يدل ا in 61), product of:
      - 2.236068 = tf(termFreq(content:ةيط ارقم يدل ا)=5)
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.0068359375 = fieldNorm(field=content, doc=61)

\*\*\*\*\*

Page 7:

- boost = 0.23039404
- digest = 8ed8fcd743ff1ce5d4c42db83fc549af
- lang = ar
- segment = 20100307101102
- title = ةيط ارقم يدل ا - ةيط ارقم يدل ا - America.gov
- tstamp = 20100307151136760
- url = http://www.america.gov/ar/amlife.html

**score for query:** ةيط ارقم يدل ا

- 5.4035656E-4 = (MATCH) sum of:
  - 5.4035656E-4 = (MATCH) weight(content:ةيط ارقم يدل ا in 3), product of:
    - 0.03718038 = queryWeight(content:ةيط ارقم يدل ا), product of:

- $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
- $0.03497641 = \text{queryNorm}$
- $0.01453338 = (\text{MATCH}) \text{fieldWeight}(\text{content:} \text{ةيطارق م يدل} \text{ in } 3),$   
product of:
  - $2.0 = \text{tf}(\text{termFreq}(\text{content:} \text{ةيطارق م يدل})=4)$
  - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
  - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=3)$

\*\*\*\*\*

Page 8:

- $\text{boost} = 0.22700267$
- $\text{digest} = \text{c639ba79e6601f1242cee32b3ba640f4}$
- $\text{lang} = \text{ar}$
- $\text{segment} = 20100307101102$
- $\text{title} = \text{نكامل او سانل - نكامل او سانل - America.gov}$
- $\text{tstamp} = 20100307151120668$
- $\text{url} = \text{http://www.america.gov/ar/amlife/people.html}$

**score for query:** ةيطارق م يدل ا

- $4.6796253\text{E-}4 = (\text{MATCH}) \text{sum of:}$ 
  - $4.6796253\text{E-}4 = (\text{MATCH}) \text{weight}(\text{content:} \text{ةيطارق م يدل} \text{ in } 7), \text{product of:}$ 
    - $0.03718038 = \text{queryWeight}(\text{content:} \text{ةيطارق م يدل}), \text{product of:}$ 
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.03497641 = \text{queryNorm}$
    - $0.012586276 = (\text{MATCH}) \text{fieldWeight}(\text{content:} \text{ةيطارق م يدل} \text{ in } 7),$   
product of:
      - $1.7320508 = \text{tf}(\text{termFreq}(\text{content:} \text{ةيطارق م يدل})=3)$
      - $1.063013 = \text{idf}(\text{docFreq}=122, \text{numDocs}=131)$
      - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{doc}=7)$

\*\*\*\*\*

Page 9:

- $\text{boost} = 0.22826105$

- digest = c33a5dc3f7d8475491bfafcf91c8b283
- lang = ar
- segment = 20100307101102
- title = داصتقالا - داصتقالا - America.gov
- tstamp = 20100307151153574
- url = http://www.america.gov/ar/econ.html

**score for query:** ةيطارق م يدل ا

- 4.6796253E-4 = (MATCH) sum of:
  - 4.6796253E-4 = (MATCH) weight(content: ةيطارق م يدل ا in 16), product of:
    - 0.03718038 = queryWeight(content: ةيطارق م يدل ا), product of:
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.03497641 = queryNorm
    - 0.012586276 = (MATCH) fieldWeight(content: ةيطارق م يدل ا in 16), product of:
      - 1.7320508 = tf(termFreq(content: ةيطارق م يدل ا)=3)
      - 1.063013 = idf(docFreq=122, numDocs=131)
      - 0.0068359375 = fieldNorm(field=content, doc=16)

\*\*\*\*\*

Page 10:

- boost = 0.23110132
- digest = 3c7f5c1dc4d604ef275f72043bf8cfc1
- lang = ar
- segment = 20100307101102
- title = secondary Multimedia - ةي م ا ل ع ل ل اس و - America.gov
- tstamp = 20100307151158462
- url = http://www.america.gov/ar/multimedia.html

**score for query:** ةيطارق م يدل ا

- 4.6796253E-4 = (MATCH) sum of:
  - 4.6796253E-4 = (MATCH) weight(content: ةيطارق م يدل ا in 38), product of:

- $0.03718038 = \text{queryWeight}(\text{content:}\text{ةي طارق م ي د ل ا}), \text{ product of:}$ 
  - $1.063013 = \text{idf}(\text{docFreq}=122, \text{ numDocs}=131)$
  - $0.03497641 = \text{queryNorm}$
- $0.012586276 = (\text{MATCH}) \text{fieldWeight}(\text{content:}\text{ةي طارق م ي د ل ا} \text{ in } 38), \text{ product of:}$ 
  - $1.7320508 = \text{tf}(\text{termFreq}(\text{content:}\text{ةي طارق م ي د ل ا})=3)$
  - $1.063013 = \text{idf}(\text{docFreq}=122, \text{ numDocs}=131)$
  - $0.0068359375 = \text{fieldNorm}(\text{field}=\text{content}, \text{ doc}=38)$

## LIST OF REFERENCES

- [1] B. Hoffman, "The use of the Internet by Islamic Extremists," in Testimony presented to the House Permanent Select Committee on Intelligence, p. 4, May 4, 2006.
- [2] K. Börner, S. Sanyal, and A. Vespignani, "Network science," in *Annual Review of Information Science & Technology*, Vol. 41, B. Cronin, ed., pp. 537–607, Information Today, Inc./American Society for Information Science and Technology, Medford, NJ, 2007.
- [3] H. M. Harmanani, W.T. Keirouz, and S. Raheel, "A rule based extensive stemmer for information retrieval with application to Arabic," *Int. International Arab J. of Inform. Tech.*, vol. 3, no. 3, pp. 265–272, July 2006.
- [4] D. A. Grossman and O. Frieder, *Information Retrieval Algorithm and Heuristics*, 2<sup>nd</sup> Ed. Springer, Norwell, MA, 2004.
- [5] W. B. Croft, Ed., *Advances in Information Retrieval: recent research from the Center for Intelligent Information Retrieval*. Kluwer Academic Publishers, Norwell, MA, 2003.
- [6] R. El Gamal, "Arabic speakers, a dying breed in the Arab world?" *Kuwait Times*, December 27, 2007.[Online]. Available: [http://www.kuwaittimes.net/read\\_news.php?newsid=MzA5NTkzOTI1](http://www.kuwaittimes.net/read_news.php?newsid=MzA5NTkzOTI1) (accessed: May 17, 2009).
- [7] "Internet Usage Statistics The Internet Big Picture World Internet Users and Population Stats." [Online]. Available: <http://www.internetworldstats.com/stats.htm>. (accessed Feb 10, 2009).
- [8] A. Chen and F. Gey, "Building an Arabic Stemmer for Information Retrieval," in *Proc. 11th Text Retrieval Conf.*, 2002, pp. 631–640.
- [9] J. Xu, A. Fraser and R. Weischedel, "Empirical Studies in Strategies for Arabic Retrieval," in *Proc. 25th Annual International Conf. Research and Development in Information Retrieval(SIGIR 2002)*, 2002, pp. 269–274.
- [10] L. S. Larkey and M.E Connell, "Arabic Information Retrieval at UMass in TREC-10," in *Proc. 10<sup>th</sup> Text Retrieval Conf.*, 2001, pp. 562–570.
- [11] R. Sonbol and N. Ghneim, "Arabic Morphological Analysis: a New Approach," in *Information and Communication Technologies: From Theory to Applications*, pp. 1–6, 2008.



- [12] K. Taghva, R. Elkhoury and J. Coombs, "Arabic Stemming Without a Root Dictionary," Information Science Research Institute (ISRI), 2005. [Online]. Available: <http://www.isri.unlv.edu/publications/isripub/Taghva2005b.ps>. (accessed: May 20, 2009).
- [13] T. BuckWalter, "Qamus: Arabic lexicography," 2003. [Online]. Available: <http://www.qamus.org/> (accessed: March 10, 2009).
- [14] L. S. Larkey, L. Ballesteros and M.E. Connell, "Improving stemming for Arabic Information Retrieval: Light Stemming and Co-occurrence Analysis," in *SIGIR 2002*, pp. 275–282, 2002.
- [15] O. Gospodnetic and E. Hatcher, *Lucene In Action*. Manning Publications Co., Greenwich, CT, 2005.
- [16] T. White, "Introduction to Nutch, Part 1: Crawling," January 10, 2006. [Online]. Available: <http://today.java.net/pub/a/today/2006/01/10/introduction-to-nutch-1.html> (accessed: May 10, 2009).
- [17] D. P. Zhou, "Delve inside the Lucene indexing mechanism," June 27, 2006. [Online]. Available: <http://www.ibm.com/developerworks/library/wa-lucene/> (accessed: May 10, 2009).
- [18] M. Cafarella and D. Cutting, "Building Nutch: Open Source Search," May 5, 2004. [Online]. Available: <http://queue.acm.org/detail.cfm?id=988408> (accessed: May 17, 2009).
- [19] A. Bolour, "Notes on the Eclipse Plug-in Architecture," July 3, 2003. [Online]. Available: [http://www.eclipse.org/articles/Article-Plug-in-architecture/plugin\\_architecture.html](http://www.eclipse.org/articles/Article-Plug-in-architecture/plugin_architecture.html) (accessed: May 17, 2009).

## INITIAL DISTRIBUTION LIST

1. Defense Technical Information Center  
Ft. Belvoir, Virginia
2. Dudley Knox Library  
Naval Postgraduate School  
Monterey, California
3. Chairman, Code EC  
Department of Electrical and Computer Engineering  
Naval Postgraduate School  
Monterey, California
4. Dr. Weilian Su  
Department of Electrical and Computer Engineering  
Naval Postgraduate School  
Monterey, California
5. Prof. John C. McEachen  
Department of Electrical and Computer Engineering  
Naval Postgraduate School  
Monterey, California
6. Qui Nguyen  
Lieutenant, United States Navy  
Naval Postgraduate School  
Monterey, California