

# Chapter 29

## Combining Indexing Units for Arabic Information Retrieval

**Souheila Ben Guirat**

*LISI: Laboratory of Computer Science for Industrial Systems, Carthage University, Tunisia & Jarir:  
Joint Group for Artificial Reasoning and Information Retrieval, Tunisia*

**Ibrahim Bounhas**

*LISI: Laboratory of Computer Science for Industrial Systems, Carthage University, Tunisia & Jarir:  
Joint Group for Artificial Reasoning and Information Retrieval, La Manouba University, Tunisia*

**Yahya Slimani**

*LISI: Laboratory of Computer Science for Industrial Systems, Carthage University, Tunisia & Jarir:  
Joint Group for Artificial Reasoning and Information Retrieval, La Manouba University, Tunisia*

### ABSTRACT

*Using either stems or roots as index terms offered considerable performance to Arabic Information Retrieval (IR) systems compared to the use of surface words for indexing. Many comparative works tried to find out the best from these two indexing approaches but until then, no of the two methods widely overtook the other. Each of the two index types performed better under different test circumstances in terms of recall and precision. In this paper, the authors propose a hybrid approach combining the two indexing units in a way they take the advantages from both of them and try to overcome their shortcomings. Then, based on some combining techniques, the authors assign a weight for each indexing unit and try to find out the best weighting values.*

### 1. INTRODUCTION

Choosing the indexing unit is yet a challenging problem in Arabic IR (Elayeb & Bounhas, 2015). When using surface words, precision reaches high levels, but we will report low recall rates because of the high derivation and agglutination of Arabic. Furthermore, this choice requires more resources in terms of storage space and processing time (Aljlal & Frieder, 2002). Thus, other types of indexing units gave better performance to Arabic IR systems like stems and roots. To better understand this problem, we present an example in section 1. We summarize our contribution in section 2.

DOI: 10.4018/978-1-5225-5191-1.ch029

## 1.1. Motivation

Let consider the Arabic root “kassama-مَسَق” and some related words (cf. Table 1). In one side, root-based methods will relate many derived words to the same form and this will cause ambiguity and reduce precision. If we consider the example, “inkassama-مَسَقن-ا” (was divided) and “akssama-مَسَق-ا” (swear) will be represented by the same index i.e. “kassama-مَسَق” (divide).

Stem indexes reduce ambiguity (Ayed, 2014; Bounhas et al., 2015), but it will from the other side, reduce recall, since some studies (Al-Kabi et al., 2011) showed that in most cases, morphological variants of words have somehow similar semantic interpretations and are not completely dissimilar and different word forms may bear similar meaning. When we apply a light stemming to the same example, “alinkissamat-تاماسقن-ال-ا” (the divisions) and “inkissam-مَسَقن-ا” (division) will have the same stem but their semantic relation with “kisma-قَمَسَق-” (division) will be ignored.

Indeed, light stemming methods offer less recall and more precision. However, lemmatization helps to achieve better recall rates, but reduces precision. Thus, combining these techniques seems promising as it realizes some compromise between precision and recall, in a way we ensure that “alinkissamat-تاماسقن-ال-ا” (the divisions) is equal to “inkissam-مَسَقن-ا” (division), not so far from “inkassama-مَسَقن-ا” (was divided), a little bit different from “iktassama-مَسَقن-ا” (share somebody in something), but not totally different from “kaassama-مَسَق-ا” (share something with somebody) and that “akssama-مَسَق-ا” (swear) and “istakssama-مَسَقن-ا” (conjure) are somehow related (cf. Figure 1).

## 1.2. Related Work

Many root extraction stemmers (A-Kabi et al., 2011; Al-Shawakfeh et al., 2010) were proposed in literature like the famous stemmer of Khoja & Garside (1999), which extracts roots by removing affixes then checking the remaining letters against a root list to avoid invalid roots. Khoja algorithm was the base of many other later works. Compared to surface based retrieval (Aljlal & Frieder, 2002), root based techniques performed better in terms of recall and resources economy.

Larkey (Larkey & Connell, 2001) and many other researchers (Aljlal & Frieder, 2002; Chen & Gey, 2002) focused on light stemming, which aims to extract stems through affixes truncation. These methods allowed better precision and reduced storage size and processing time, compared to indexing by surface words (Aljlal & Frieder, 2002).

A comparison between Khoja stemmer and Larkey light stemmer proved that these techniques still have some weaknesses (Handi et al., 2012). In (A-Kabi et al., 2011), the authors compared four root stemmers, e.g. Al-Mustafa (Al-Kabi & Al-Mustafa, 2006), Taghva et al. (2004), Al-Sarhan et al. (2003) and Rabab’ah et al. (2005).

Table 1. Some Arabic words related to the root “kassama-مَسَق”

مَسَقن-ا Inkassama	مَسَقن-ا istakssama	مَسَقن-ا iktassama	مَسَق-ا Akssama	مَسَق-ا kaassama	مَسَقن-ا inkissam	قَمَسَق kisma	تاماسقن-ال-ا alinkissamat
Was divided	Conjure	Share somebody in something.	Swear	Share something with somebody	Division	Division	The divisions

13 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:  
[www.igi-global.com/chapter/combining-indexing-units-for-arabic-information-retrieval/198571](http://www.igi-global.com/chapter/combining-indexing-units-for-arabic-information-retrieval/198571)

## Related Content

---

### The Effect of Stemming on Arabic Text Classification: An Empirical Study

Abdullah Wahbeh, Mohammed Al-Kabi, Qasem Al-Radaideh, Emad Al-Shawakfa and Izzat Alsmadi (2013). *Information Retrieval Methods for Multidisciplinary Applications* (pp. 207-225).  
[www.irma-international.org/chapter/effect-stemming-arabic-text-classification/75909](http://www.irma-international.org/chapter/effect-stemming-arabic-text-classification/75909)

### Template based Question Answering System Over Semantic Web

(2022). *International Journal of Information Retrieval Research* (pp. 0-0).  
[www.irma-international.org/article//299933](http://www.irma-international.org/article//299933)

### Need of Intelligent Search in Dynamic Social Network

Shailendra Kumar Sonkar, Vishal Bhatnagar and Rama Krishna Challa (2018). *Information Retrieval and Management: Concepts, Methodologies, Tools, and Applications* (pp. 144-157).  
[www.irma-international.org/chapter/need-of-intelligent-search-in-dynamic-social-network/198551](http://www.irma-international.org/chapter/need-of-intelligent-search-in-dynamic-social-network/198551)

### Automatic Ontology Construction using Conceptualization and Semantic Roles

Amita Arora, Manjeet Singh and Naresh Chauhan (2017). *International Journal of Information Retrieval Research* (pp. 62-80).  
[www.irma-international.org/article/automatic-ontology-construction-using-conceptualization-and-semantic-roles/181726](http://www.irma-international.org/article/automatic-ontology-construction-using-conceptualization-and-semantic-roles/181726)

### Integration of Multi-Class Service Paradigm With Generic Trust Mechanism for Innovation, Customization and Adaptability in MANETs

Nitin Khanna, Sandeep Singh, Anshu Bhasin and Kamal Malik (2022). *International Journal of Information Retrieval Research* (pp. 1-15).  
[www.irma-international.org/article/integration-of-multi-class-service-paradigm-with-generic-trust-mechanism-for-innovation-customization-and-adaptability-in-manets/300290](http://www.irma-international.org/article/integration-of-multi-class-service-paradigm-with-generic-trust-mechanism-for-innovation-customization-and-adaptability-in-manets/300290)