

Arabic Text Preprocessing for the Natural Language Processing Applications

مُعالجة النصوص العربية لغايات تطبيقات المُعالجة الآلية للغات الطبيعية

Arafat Awajan

عَرَفات عَوْجان

Computer Science Department, Princess Sumaya University for Technology,

P O Box 1438, Al- Jubeiha 11941, Fax (00962) 65347295, Amman, Jordan

P. O. Box: 249, Abha, Saudi Arabia.

E-mail: awajan@psut.edu.jo

ABSTRACT: A new approach for preprocessing vowelized and unvowelized Arabic texts in order to prepare them for Natural Language Processing (NLP) purposes is described. The developed approach is rule-based and made up of four phases: text tokenization, word light stemming, words' morphological analysis, and text annotation. The first phase preprocesses the input text in order to isolate the words and represent them in a formal way. The second phase applies a light stemmer in order to extract the stem of each word by eliminating the prefixes and suffixes. The third phase is a rule-based morphological analyzer that determines the root and the morphological pattern for each extracted stem. The last phase produces an annotated text where each word is tagged with its morphological attributes. The preprocessor presented in this paper is capable of dealing with vowelized and unvowelized words, and provides the input words along with relevant linguistics information needed by different applications. It is designed to be used with different NLP applications such as machine translation, text summarization, text correction, information retrieval, and automatic vowelization of Arabic text. **Key words:** *Arabic Text Preprocessing, Stemming, Morphological Analysis, Text Annotation, Part of speech tagging.*

المستخلص: تهدف الورقة إلى وصف أسلوب جديد لمعالجة النصوص العربية المُشكَّلة وغير المُشكَّلة من أجل تهيئتها للاستعمال في تطبيقات المُعالجة الآلية للغات الطبيعية. بُني الأسلوب الجديد على قواعد محددة مسبقاً تتكون من أربعة مراحل، هي فصل المكون الأساس في النص (الكلمة)، فصل جذور الكلمات، التحليل الصرفي للكلمات، وإضافة توصيف للكلمات على النص. تعالج المرحلة الأولى النص بغرض عزل الكلمات وإعادة تمثيلها بطريقة معيارية، ويخضع النص في المرحلة الثانية إلى معالج يقوم باستخراج جذور كلمات النص وذلك بإزالة الإضافات السابقة واللاحقة لها، وتشمل المرحلة الثالثة محللاً صرفياً مبنياً على قواعد محددة، يستخرج الجذر والنمط الصرفي لكل كلمة، أما المرحلة الأخيرة فتضيف توصيفات على النص تشمل الخصائص الصرفية لكل كلمة. تَمَكَّن الطريقة المُقترحة في هذه الورقة من التعامل مع النصوص العربية المُشكَّلة وغير المُشكَّلة وتنتج لكل كلمة من النص مجموعة من المعلومات اللغوية الضرورية للعديد من التطبيقات. تم تصميم المعالج بحيث يمكن استخدامه مع الكثير من تطبيقات المُعالجة الآلية للغات الطبيعية مثل ترجمة وتلخيص وتصحيح النصوص، إلى جانب استخراج المعلومات والتشكيل الآلي للنصوص العربية. **كلمات مدخلية:** *مُعالجة، النصوص العربية، جذور الكلمات، تحليل صرفي، حاشية النص، توصيف الكلمات، تطبيقات، مُعالجة آلية، لغات طبيعية.*

INTRODUCTION

The Natural Language Processing (NLP) is one of the most important and evolving fields of

investigation in computer science and artificial intelligence. NLP deals with the creation of programs that are capable of processing and understanding human languages. A natural

language is a very complicated phenomenon, and its study involves many levels of analysis related to the phonology, morphology, syntactical rules, and semantics of the language (Allen, 1995; Manning and Schutze, 2000).

NLP covers a wide range of useful applications, i.e. machine translation, text summarization, text correction, document analysis, human-machine interaction and information retrieval. Although the objectives of each application determine the processing techniques and the transformations to be applied on the original texts and the order in which these transformations should be applied, the first and most critical step is the preprocessing of the input text.

The text preprocessing is a core task in the natural language processing procedure. It aims at creating an intermediate form from the input text based on the extraction of words, the morphological analysis, and the text annotation. Many researches related to the preprocessing of natural language texts have been published, mainly for the European languages. These works cover the tokenization of text e.g., (Grefenstette and Tapanainen, 1994), the morphological analysis of words e.g., (Antworth, 1994), and the part of speech tagging of words e.g., (Jurafsky and Martin, 2000).

Few researches have been published on the subject of the Arabic text preprocessing for the NLP applications. Published papers cover mainly the morphological analysis of Arabic words (Al-Sughaiyer and Al-Kharashi, 2004). In addition, there are fewer papers, which treated the tagging of words (Khoja, *et al.* 2001) and the preparation of the text for text understanding. These papers generally ignore the presence of diacritics in Arabic texts or limit the analysis to words generated from 3-letter roots (Beesley, 1996; Larkey, *et al.* 2002). Some of these papers are based on the generalization of concepts used for European languages to the case of the Arabic language (De Roeck and Al-Fares, 2000); (Larkey, *et al.* 2002).

This paper presents a new approach of preprocessing the Arabic texts based on the features of the Arabic language and is able to analyze Arabic words as they appear in real texts. The developed approach is capable of

dealing with vowelized, unvowelized or partially vowelized Arabic words. It transforms the input texts into a new format that is more appropriate and adequate for the different NLP applications. In addition to the original text, the new format contains additional information at the word level. The main purpose of the new format is to make the NLP applications faster and more accurate.

The approach consists of preprocessing the input text in four phases. The first phase, called the tokenization of the text, preprocesses the input text in order to detect and isolate the words. The second phase is a light stemmer that eliminates prefixes and suffixes in order to extract the kernel words or the stems. The third phase is the morphological analysis of words; it consists of a rule-based morphological analyzer that decomposes each word into its basic morphological components; the root and the morphological pattern. The fourth and last phase is the text's annotation that tags the words by adding morphological attributes in order to facilitate their analysis. Figure 1 describes schematically the main phases of the text preprocessing.

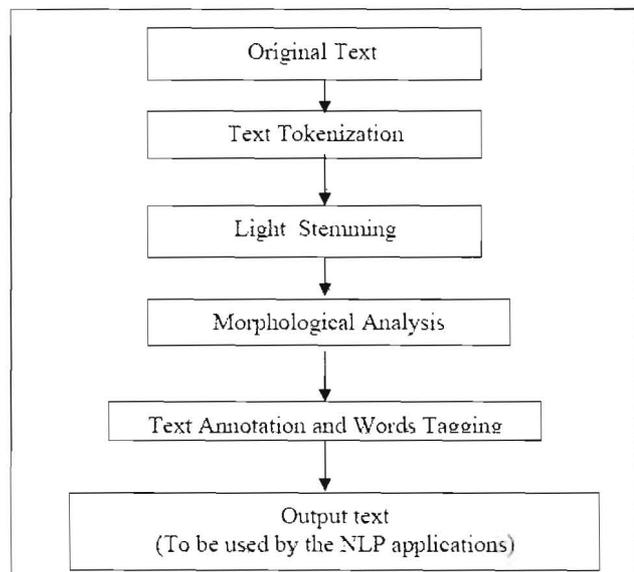


Fig. 1. Main Phases of Text Preprocessing.

BACKGROUND

Words in the Arabic language may be classified into two categories: derivative words, and non-derivative words. The derivative Arabic words are generated from basic entities called roots or radicals according to a predefined list of

standard patterns called morphological patterns or balances. These morphological patterns represent the major spelling rules of Arabic words. The non-derivative words include two sub-categories: fixed words and foreign words. The fixed Arabic words are a set of words that do not obey the derivation rules. These words are generally functional words like pronouns, prepositions, conjunctions, question words and the like. The foreign words are nouns borrowed from foreign languages.

The majority of Arabic words belong to the category of derivative words. All kinds of words (verbs, nouns, adjectives and adverbs) can be generated from roots according to the standard patterns. The pattern associated to a word determines its various attributes such as gender (masculine/feminine), number (singular/plural), and tense (past, present, and imperative). Based on the above, a derivative Arabic word can be represented by its root along with its morphological pattern. For example, the word (لاعبون) in English «players» is generated from the root (to play: لعب) according to the pattern (فاعلون). This pattern indicates that the word is a noun, its gender is masculine, and its number is plural. The final meaning of this Arabic word will be players (root (play): Attributes (noun; plural; masculine)).

Nevertheless, the morphological analysis of Arabic words presents many challenges that must be considered. The first challenge is the fact that some letters of the root may be dropped or modified during the derivation of words from roots. The second challenge is that many affixes can be attached to the beginning of the word (prefixes) and the end of the word (suffixes). These affixes may be formed from one or more letters. The third challenge is that the Arabic language words are written without short vowels. Different diacritical marks are used to replace the short vowels. The number of these marks is eight diacritics: three diacritical marks to indicate the short vowels (َ , ِ , ُ), three double diacritic marks which combine the single ones (ِ , ُ , ُ), one diacritical mark to indicate the absence of vowelization (ْ), and a single diacritical mark to indicate the duplicate occurrence of a consonant (ّ). These marks play a vital role in determining

the possible meaning of the word. Actually, two different patterns may have the same sequence of consonants, but they differ from each other in terms of the diacritical marks.

According to the extent of using the diacritical marks, Arabic texts may be classified into three different categories: unvowelized, partially vowelized, and fully vowelized texts. The first category represents the texts without diacritics such as many of the typed, printed texts and newspapers. The second category represents the texts partially vowelized such as textbooks and scientific texts. The last category represents the fully vowelized Arabic texts, in which every consonant is followed by a diacritical mark such as the Holy Quran text, children books and literature texts.

METHODOLOGY

Text Tokenization

Tokenization is the process of isolating word-like units from a text (Grefenstette and Tapanainen, 1994). In addition to words, text documents often contain white spaces, punctuation marks, and a number of mark-ups that indicate font changes, text subdivision, and special characters. The aim of the tokenization phase is to detect and isolate the individual words by eliminating these additional components. For the purposes of this work, it is assumed that an Arabic word is a sequence of Arabic letters and diacritical marks without separators (space or punctuation marks). The detection of individual words is based on a simple text-matching algorithm. A new word is a sequence of Arabic letters and diacritical marks starting by a letter and ended by a letter or diacritical marks. The detection of a separator or special character marks the end of the word.

Word Representation

In order to deal with the different forms that an Arabic word can take (i.e., unvowelized, partially vowelized, and fully vowelized words), an additional diacritical mark called “EXTRA-SEKOUN” is introduced in this paper. This special diacritical mark “EXTRA-SEKOUN”, represented by a period in the different examples given in this paper, is used

to replace the missed diacritical marks in a word.

A rule-based procedure "Check_Diacritics" takes the string of characters forming a word and checks the presence of diacritic marks after each consonant. It produces the new presentation of the word "New Word" by replacing the missed diacritical marks by the special character EXTRA-SEKOUN. It calls a function "IsDiacritic" defined for finding out the diacritical marks. The function "IsDiacritic" returns TRUE if a character is one of the diacritical marks of the Arabic language. The Prolog-style description of the rule "Check_Diacritics" is presented as follows:-

```
//Stopping case when the new presentation
is finalized.
```

```
Check_Diacritics ([ ], New Word, False).
```

```
//The last consonant is not followed by a
diacritic; an EXTRA-SEKOUN //mark is
then added
```

```
Check_Diacritics ([ ], New Word, True):-
Check_Diacritics ([ ],['.'|New Word],
False).
```

```
//Detection of a consonant at the first call
of the rule or after a diacritic.
```

```
Check_Diacritics([Head|Tail], New Word,
False):-
```

```
    Check_Diacritics (Tail, [Head|T1],
    True).
```

```
//Detection of a diacritic.
```

```
Check_Diacritics ([Head|Tail], New Word,
True):-
```

```
    Is Diacritic (Head),
    Check_Diacritics (Tail, [Head | New
    Word], False).
```

```
//Adding EXTRA-SEKOUN if 2
consecutive consonants are detected
```

```
Check_Diacritics ([Head|Tail], New Word,
True):-
```

```
    Check_Diacritics ([Head|Tail],
    ['.',New Word], False).
```

The word is then represented by a list of characters New Word with the following format:

$$[C_1 V_1, C_2 V_2, \dots, C_n V_n]$$

where C_i is a consonant and V_i is one of the diacritical marks including the EXTRA-SEKOUN mark replacing the missed diacritical marks in the input word.

Word Light Stemming

In the Arabic language, like many other languages, some lexical elements can be attached to a word in order to add new information to the word or to form a sentence or a part of sentence. Examples of these additive parts, called affixes, are the conjunctions (ex.: ف), the prepositions (ex.: ل), the pronouns (ex.: هم) and the article (ex.: ال). The number of affixes is limited, and they may be added at the beginning of the word (prefixes), or at the end of the word (suffixes). A word may have up to two prefixes and up to three suffixes (Al-Sughaiyer and Al-Kharashi, 2004).

The process of extracting the kernel word or stem from the original text by eliminating the suffixes and prefixes is called stemming. As opposed to the English language, the removal of prefixes and suffixes from an Arabic word does not usually reverse the meaning of the word (Abu Salem, *et al.* 1999; Xu, *et al.* 2002; Moukdad, 2006). The stemmer used is developed based on the light stemming approaches described in (Darwish, 2002; Larkey, *et al.* 2002) with additional restrictions in the list of strippable prefixes and suffixes. The light stemming refers to a process of stripping off a small set of prefixes and/or suffixes, without trying to neither deal with infixes nor recognize patterns and roots. On one hand, this approach reduces the risk of root consonant loss that can be produced if a heavy stemming is used. On the other hand, the drawback of the use of light stemmer will be corrected by the morphological analyzer that will be applied on the stems in the next phase.

The stemmer decomposes the input string into the additive parts (prefixes, suffixes) and the stem. The decomposition is realized by applying a set of identification rules that test the first characters and last characters of the word against the possible and strippable additive parts. The word is decomposed into three parts. The first part is in the list of possible prefixes and the third part is in the list of possible suffixes. The first and third parts may be null. The light stemmer developed in this paper is based on the following assumptions: 1) a word is composed of three parts: prefix, stem and suffix; 2) any of the additive parts (prefixes and suffixes) may be empty; 2) the stem of the word has at least three letters, any word with less than

four letters will be left without decomposition; 4) a prefix can have 0 to 3 letters and exist in the list of prefixes, and a suffix can have 0 to 3 letters and exist in the list of suffixes; and 5) only the consonants of the word (letters) are taken into account for the purpose of this decomposition.

The process of word-decomposition and suffix-removal is repeated until one of the following conditions is verified: the number of the letters in the word is less than or equals to 3, and there are no prefixes, nor suffixes detected. The final output of the light stemmer takes the following structure: [Prefix1][Prefix2] stem [Suffix1][Suffix2][Suffix3].

An explicit list of strippable affixes is provided in a table. They are classified according to their type (prefix/suffix) and their length (1, 2 and 3 letters). This restricted list contains mainly conjunctions, prepositions, pronouns and the article. Table 1 shows examples from this list. The affixes used to determine the person, number and gender are not included in this table and therefore not strippable in this phase. The priority of detection and removal is given for the three-letter affixes over the two-letter affixes, for the two-letter affixes over the one-letter affixes, and for prefixes over suffixes. However, some words contain letters that may be detected as prefixes or suffixes. To solve this problem, we will consider the results and feedback of the morphological analyzer to correct this type of error. If the morphological analyzer fails in detecting the morphological structure of the stem, the last action taken by the light stemmer will be disregarded.

Morphological Analysis

The proposed morphological analyzer is a rule-based technique, designed to identify the morphological structure of vowelized and unvowelized Arabic words. The morphological analyzer processes the extracted stems in order to determine their roots and patterns (morphological knowledge). As the affixes used to determine the person, number and gender are not removed in the previous phase, the list of morphological patterns used to generate verbs and nouns is extended to include new computational patterns generated from the classical patterns by adding these affixes. Table 2 shows examples of the extended morphological patterns represented by their list of consonants.

Stem Decomposition

The stem $[C_1 V_1, C_2 V_2, \dots, C_n V_n]$ is split into two lists: the first one LC contains the sequence of consonants $[C_1, C_2, \dots, C_n]$ and the second one LV contains the sequence of diacritical characters $[V_1, V_2, \dots, V_n]$. Table 3 illustrates this representation for the three situations of Arabic texts where the "EXTRA-SEKOUN" character marked by a dot is used to replace the missed diacritical marks in the original word.

The recursive procedure Decompose performs the decomposition of the word in the two lists: LC and LV. The procedure Decompose may be presented by the following Prolog-style description: // Base (Stopping) case when the decomposition is terminated

Decompose ([], LC, LV).

// Detection of a diacritic

Decompose ([Head|Tail], _ , LV):-

IsDiacritic (Head),

Decompose (Tail , _ , [Head|LV]).

// Detection of a consonant

Decompose ([Head|Tail], LC , _):- Decompose (T , [Head|LC] , _).

The list LC of consonants represents the letters of the word's root and the letters added to the root to form the stem according to a standard pattern.

Root Representation

In order to extract the root of a stem, the list LC can be represented by the following general description:

$$[X_1 [X_2[X_3]]] R_1 [Y_1] R_2 [Y_2] R_3 [[Y_3] R_4 [Y_4] R_5] [Z_1 [Z_2[Z_3]]]$$

where X_1, X_2, X_3 represent a prefix of maximum 3 letters, Z_1, Z_2, Z_3 represent a suffix of maximum three letters and Y_1, Y_2, Y_3, Y_4 represent the possible infixes.

The maximum number of letters considered for the prefixes and suffixes takes into account that some of these components may not be detected by the light stemmer used before. The slots $R_1, R_2, R_3, R_4,$ and R_5 represent the letters of the root used to generate the word. This representation allows the manipulation of words generated by all kinds of roots (3-letter roots, 4-letter roots, and 5-letter roots).

The three examples in Table 3 share the same list of consonants LC. This list LC contains

a prefix with one consonant $X_1 = \text{"ي"}$, a suffix with two consonants $Z_1 = \text{"و"}$, $Z_2 = \text{"ن"}$, and a 3-letter root R_1, R_2, R_3 , where $R_1 = \text{"ذ"}$, $R_2 = \text{"هـ"}$, and $R_3 = \text{"ب"}$. Table 4 shows additional examples illustrating the decomposition of the list of consonants LC into prefixes, suffixes, infixes and root letters.

Table 1. Examples from the List of Strippable Affixes.

Types of Affixes	Group	Examples of Affixes	Examples(Words)
Prefixes	P_G1 (one letter)	و، ف، ب، ك، ل	فسمع
	P_G2 (two letters)	أل	الكتاب
	P_G3 (three letters)	وال، فال، كال	كالعصفور
Suffixes	S_G1 (one letter)	ي، ك، هـ	كتابي
	S_G2 (two letters)	هم، هن، ها، كن	يعلمهم
	S_G3 (three letters)	هما	رافقهما

Table 2. Examples from the List of Extended Morphological Pattern.

Standard Pattern	Person	Number	Masculine	Feminine	Example
Verb standard pattern: فعل	First	Singular	فعلت	فعلت	ذهبت
		Dual	فعلتا	فعلتا	ذهبتا
		Plural	فعلنا	فعلنا	ذهبتنا
	Second	Singular	فعلت	فعلت	ذهبت
		Dual	فعلتما	فعلتما	ذهبتما
		Plural	فعلتم	فعلتن	ذهبتن
	Third	Singular	فعل	فعلت	ذهبت
		Dual	فعلتا	فعلتا	ذهبتا
		Plural	فعلوا	فعلن	ذهبوا
Noun standard pattern: فاعل		Dual	فاعلتان	فاعلتان	لاعتان
		Plural	فاعلون	فاعلات	لاعون

Table 3. Decomposition of Stems.

Word	Case	List of Consonants LC	List of Diacritics LV
يَنْقُبُونَ	Fully vowelized	[ي ذ ه ب و ن]	[َ ُ َ َ َ َ]
يَنْقُبُونْ	Partially vowelized	[ي ذ ه ب و ن]	[َ َ َ َ َ]
ينقبون	Unvowelized	[ي ذ ه ب و ن]	[.]

Table 4. Decomposition of the List of Consonants.

Input Word	List of Consonants	Root R1R2R3	Prefix X1X2X3	Infix Y1 Y2	Suffixes Z1Z2Z3
سيدرسون	[س ي د ر س و ن]	[درس]	[س ي]	[]	[و ن]
دارسون	[د ا ر س و ن]	[درس]	[]	[ا]	[و ن]
مدارس	[م د ا ر س]	[درس]	[م]	[ا]	[]

Morphological Pattern Representation

Each one of the morphological patterns is represented by a list L of characters with the same structure as we proposed for the words. The slots of the root letters are marked by (*); and may be replaced by any consonant. For example, the morphological pattern “بَفَعَلُونَ” is represented by the list [بَ * * * وَ نَ]; and decomposed into two lists: the list of consonants LC (بَ * * * وَ نَ) and the list of diacritical marks LV (وَ نَ). This partition of consonants and diacritics reduces significantly the number of patterns to be tested. The characters “*” represent slots where consonants can be inserted to form a real word.

The morphological patterns can be re-grouped in classes according to their list of consonants. The patterns of the same class share the same list of consonants and they are different among one another in terms of the lists of diacritical marks. Table 5 shows an example of three different patterns of the same class; these patterns have the same list LC and have different lists of diacritical marks LV. The set of patterns will be represented by the set of consonant lists LC, where we associate with each entry all the possible and correct combinations of diacritical marks LV. The couplet LC and LV will determine the morphology of the word.

Root and Pattern Identification

The step of identification of the root and pattern is realized by two recursive procedures: FindPattern and FindRoot. The recursive procedure FindPattern (LC(word), LC(pattern)) receives the list of consonants of the stem and returns the corresponding standard pattern. The recursive procedure Find Root (LC(word), LC(Pattern), LC(Root)) receives the word and its pattern and extracts the root by comparing the two entries and applying the rules relating the pattern to the root. The following Prolog style code represents these two rule-based recursive procedures:

```
//Base (Stopping) case when the decomposition
is terminated
Find Pattern ([ ], [ ]).
//Skip the letters of the stem located in the root
letters slots of the pattern
Find Pattern ([Head|Tail1], [ '*'|Tail2]):- Find
Pattern (Tail1, Tail2).
//Finding the letters of the standard pattern
Find Pattern ([Head|Tail1], [Head|Tail2]):- Find
Pattern (Tail1, Tail2).
//Base-Stopping) case when the Root is completely
detected
Find Root ([ ], [ ], Root).
// Skip of the letters added to the root according
to the standard pattern
Find Root ([Head|Tail1], [Head|Tail2], Root):-
Find Root (Tail1, Tail2, Root).
// Find the root letters, corresponding to slots ‘*’
in the morphological pattern
Find Root ([Head1|Tail1], [ '*'|Tail2], Root):-
Find Root (Tail1, Tail2, [Head1|Root]).
```

Text Annotation

The text annotation is based on a categorical approach that uses the constituent parts of the words generated by the previous phases. The words are automatically categorized and classified into predefined categories. Three main categories are defined as: the derivative words, non-derivative words, and undefined words. The derivative words are the words that the system is able to extract their roots, patterns, prefixes, and suffixes. The pattern will determine for this kind of words, additional features such as type (name/verb), gender, number, etc. The non-derivative words are Arabic words that are not generated from standard pattern such as pronouns, prepositions, conjunctions, question words, foreign names and the like. Those words will be stored in predefined tables. The last category “undefined” word class is used to categorize the words that the system

Table 5. Grouping Patterns According to their List of Consonants.

Pattern	List of Consonants LC	List of Diacritical Marks LV
بَفَعَلُونَ	[بَ * * * وَ نَ]	[وَ نَ]
بُفَعَلُونَ	[بَ * * * وَ نَ]	[وَ نَ]
بُفَعَلُونَ	[بَ * * * وَ نَ]	[وَ نَ]

categories. Table 6 shows the main categories and subcategories used for the purposes of this study.

A categorical grammar notation is used to describe the results of the last phase. This notation determines the category and subcategory of each word. For the derivative words, their constituent parts (the root, morphological pattern, and the prefixes and suffixes produced by the light stemmer) are also attached to this notation. Additional attributes are attached to the derivative words based on the features and attributes of their standard pattern discovered by the morphological analyzer. These attributes are.

1. For the nouns: (Gender(M/F), Number(Singular/Dual/Plural), (Definite/Un-definite))
2. For the verbs: (Gender(M/F), Number(Singular/Dual/Plural), Person (First/Second/Third), Time (Present/Past))

EXPERIMENTS AND EVALUATION

Implementation

To evaluate the performance and accuracy of the developed approach, a prototype of the proposed preprocessor has been implemented and tested on real Arabic texts. The implementation includes the different tables needed in the different phases of the preprocessing. A table entitled (ROOT) is defined to represent the roots in the language to generate the derived words; 672 three-letter roots and 41 four-letter roots are stored in this table for the purposes of this experiment. A Table entitled (PATTERN) aims at representing the standard and extended patterns, whereby each entry contains the list of consonants of the pattern along with all the possible combinations of diacritical marks. A table entitled (AFFIXES) contains the strippable affixes in the light stemming phase. A table entitled (NON-DERIVATIVE) is defined and

Table 6. Word Categories and Subcategories.

	Level 1	Level 2	Level 3
Word	Derivative (1)	Noun (1.1)	Constituent parts: root pattern Attributes: Gender Number Definite or Un-definite.
		Verb (1.2.)	Constituent parts: root pattern Attributes: Gender Number Person Time
	Non Derivative (2)	Fixed Words(2.1)	Pronoun (2.1.1)
			Preposition (2.1.2)
			Conjunction (2.1.3)
			Question (2.1.4)
	Foreign Words (2.2)		
	Proper Names (2.3)		
Undefined (3)			

Table 7. Examples of the Results of the Annotation Phase.

Word	Category	Constituent Parts	Attributes
اللاعبون The players	(1.1) Derivative words, Noun	Light Stemming Phase Prefix: ال (Article the) Suffix: Morphological Analyzer Phase Root: لعب Pattern: فاعلون	Gender (M) Number (P) Definite
تأكلها She eats it	(1.2) Derivative word Verb	Light Stemming Phase Prefix: -- Suffix: ها (Pronoun it) Morphological Analyzer Phase Root: أكل Pattern: تفعل	Gender (F) Number (S) Person (Third) Time (Present)
أين Where	(2.1.4) Non Derivative word Fixed word Question	--	--
أوروبا Europe	(2.2) Non Derivative Foreign word	--	--

contains the fixed words and tools of the language. In addition, the foreign words that the Arabic language has borrowed from other languages as well as the proper nouns have been listed.

In the Arabic language, as well as in many other languages, a word may belong to more than one category. The implemented prototype is designed to be able to produce all the possible and detected categories of the input words.

Data Sets

The performance of the proposed preprocessing technique have been empirically tested and evaluated on real data sets. Three different data sets were used to represent the three cases of texts: vowelized, partially vowelized and unvowelized texts. The texts were collected from different resources. Manual treatments were necessary in some cases to eliminate or insert the diacritical marks to create these three different categories of sets. Table 8 shows the details of the data sets. It is noticed that despite

the fact that the number of fixed words is limited, their frequency in real text is very high.

Performance Measurement

In order to evaluate the quality and accuracy of the results, four performance measures have been used: the recall, the precision, the error rate, and the category precision. They are defined as:

$$\text{Precision} = C / N$$

$$\text{Recall} = C / TN$$

$$\text{Error -Rate} = E / TN$$

$$\text{Category-Precision} = TC / TD$$

where: TN is the total number of words in the data set, E is the number of words incorrectly categorized, C is the number of words correctly categorized, N is the number of words categorized by the system ($N = C + E$), TC is the total number of categories correctly detected, and TD is the total number of categories detected. The last measure indicates the ratio of correct categories detected at the word level.

Table 8. Data Sets.

Data Set Category	Number of Words	Number of Derivative Words	Number of Non-Derivative Words
Vowelized text	1180	542 (46%)	638
Partially vowelized text	1467	820 (56%)	647
Unvowelized texts	1452	650 (45%)	802

Table 9. Results of the Preprocessor on the three Data Sets.

Data Set category	Precision	Recall	Error Rate	Category Precision
Vowelized text	0.97	0.88	0.027	#1.0
Partially vowelized text	0.93	0.85	0.062	0.94
Unvowelized texts	0.91	0.83	0.073	0.87

Results

The results of the preprocessor are controlled manually to check their correctness. A word is considered correctly categorized if all the categories associated to this word are correct. A word is considered incorrectly categorized if any of the categories associated to the word is incorrect.

As the results depend directly on the above-mentioned tables, the words that the preprocessor might fail to identify and categorize due to the fact that their roots or pattern are not included in the predefined tables are not taken into consideration in the performance analysis of the prototype. The results of the evaluation test for the three data sets are shown in Table 9.

In general, the preprocessor has achieved the best results in the fully vowelized texts, as there was a very low rate of inappropriate results. The incorrectly categorized words are generally the words which have some of their letters mixed up with the affixes, and then stripped during the light stemming phase. The Error rate is higher in the case of the partially vowelized and unvowelized texts, mainly because of the missed diacritical marks especially the diacritical mark (') used to indicate the duplicate occurrence of a consonant.

CONCLUSION

We have developed, successfully, a four-phase approach to the task of preprocessing Arabic texts, text tokenization, light stemming, morphological analysis, and text annotation. These four phases are applied in sequence and collaborate in order to transform an Arabic text into a new format designated to be used for the NLP applications.

Our approach works on Arabic texts as they appear in real world documents. This approach can be used as a part of a larger application of NLP. The output of the preprocessing steps determine whether the word is generated from a root or not as well as a set of attributes that can be used by the NLP application to determine the part of speech of the word in addition to its possible meanings.

The proposed technique gives accurate results for the fully vowelized Arabic texts. In the absence of diacritical marks, it produces a list of possible morphological patterns, in general between 1 to 5 patterns for each word that share the same consonants. However, they are different from the others in terms of the list of diacritics.

In spite of the good, accurate and general results obtained by our approach, it has been noticed that it requires many tables and rules. It is assumed that an expansion of the proposed preprocessing techniques in the direction of using the syntactical information will help in getting more results that are accurate.

REFERENCES

- Abu Salem, H, Al-Omari, M, and Even, M** (1999) Stemming Methodologies over Individual Query Words for an Arabic Information Retrieval System. *Journal of the American Society for Information Science* **50** (6): 524-529.
- Al-Sughaiyer, IA, and Al-Kharashi, IA** (2004) Arabic Morphological Analysis Techniques: A Comprehensive Survey. *Journal of the American Society for Information Science and Technology* **55** (3): 189-213.
- Allen, J** (1995) *Natural Language Understanding*. The Benjamin-Cummings Publishing Company, Redwood City, USA, pp1-654.
- Antworth, E** (1994) *Morphological Parsing with a Unification-based Word Grammar*. A paper presented at North Texas Natural Language Processing Workshop, May 23, University of Texas, Arlington, Texas. (Available: <http://www.sil.org/pckimmo/ntnlp94.html>).
- Beesley, KR** (1996) Arabic Finite-State Morphological Analysis and Generation. In: *Proceedings of the 16th International Conference on Computational Linguistics*, August 9-15, ACL, East Stroudsburg, USA, pp 89-94.
- Darwish, K** (2002) Building a Shallow Arabic Morphological Analyzer in One Day. In: *Proceeding of the 40th Annual Meeting of the Association for Computational Linguistics*, ACL, University of Pennsylvania, Philadelphia, USA, pp47-54.
- De Roeck, A, and Al-Fares, W** (2000) A Morphological Sensitive Clustering Algorithm for Identifying Arabic Roots. In: *Proceeding of the 38th Annual Meeting of the Association for Computational Linguistics*, ACL, Hong Kong, China, pp199- 206.
- Grefenstette, G and Tapanainen, P** (1994) What is a Word, What is a Sentence: Problems of Tokenization. In: *Proceeding of the 3rd Conference on Computational Lexicography and Text Research*. Complex, July 7-10, Budapest, Hungary, pp79- 87.
- Jurafsky, D, and Martin JH** (2000) *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, Prentice Hall PTR, Upper Saddle River NJ, USA, pp1- 934.
- Khoja, S, Garside, R, and Knowles, G** (2001) A Tag set for the Morphosyntactic Tagging of Arabic In: *Proceedings of the Corpus Linguistics*. 30 March- 2nd. April, Lancaster University, Peter Lang, Lancaster, UK, pp1-642. (Available: www.archimedes.fas.harvard.edu/mdh/arabic/CL2001.pdf).
- Larkey, LS, Ballesteros, L, and Connell, ME** (2002) Improving Stemming for Arabic Information Retrieval: Light Stemming and Co-occurrence Analysis. In: *Proceeding of the 25th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. 11- 15 August, Tampere, Finland, pp275-282.
- Manning, CD, and Schutze, H** (2000) *Foundation of Statistical Natural Language Processing 2nd. ed.*, MIT Press, USA, pp1- 620.
- Moukdad, H** (2006) Stemming and Root-based Approaches to the Retrieval of Arabic Documents on the Web. *Webology*, **3** (1): 1-27. (Available: www.webology.ir/2006/v3n1/a22.html).
- Xu, J, Farsner, A, and Weischedel, R** (2002) Empirical Studies in Strategies for Arabic Retrieval. In: *Proceeding of the 25th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'02)*. August 11-15, University of Tampere, Tampere, Finland, pp269 – 274.
- Ref. No. (2456)
Rec. 07/ 10/ 2007
In- revised form: 18/ 02/ 2008