

satisfactory results or successful algorithms for generating Arabic word roots or morphological information.

This paper provides an algorithm that generates the root of Arabic words by applying permutations on the Arabic letters. The paper also provides another algorithm to build Arabic lexicon by generating different morphological information for these roots and then provides meanings and examples for each word from a holy Quran dataset [18] (the holy Quran is considered as a standard Arabic reference since around 15 century).

1.1 Contribution

This paper provides a contribution in the field of NLP by: providing a novel root generator algorithm for giving the roots of Arabic words using an approach that permutes and combines the Arabic alphabet. This way of generating the word root was used by Al Khalil Bin Ahmed in constructing his valuable and well-known lexicon before more than 1200 years ago.

Building a comprehensive electronic Arabic lexicon (Arabic language dataset) that can be used in Arabic language processing and for other purposes such as dictionaries, word processing, spelling checker.. etc.

1.2 Previous work

Most of Arabic language processing achievements are in the area of Arabic morphology. The research projects that work on Arabic stemming algorithms or root generating can be classified into two categories:

- Rule-Based stemming or root generating algorithm: This category is studied more than the other kinds of algorithms, because it provides more accurate results than other algorithms. Most of works on rule-based stemming algorithms depend on affix removal such as khoja stemmer and Tim Buckwalter morphological analyzer. In general there are many examples of recently research work under this category such as [7], [4],[1], [5] and [9].
- Classification or clustering algorithms: The stem in this category of algorithms is classified or clustered according to the stem type and language document. For example in [13] the association rules are used in clustering the stem of text. In [14] Stem is generated based on feature reduction technique which is used to categorize Arabic Text. In [21] algorithm based on Naive Bayes is developed to categorize Arabic document. In [15] the authors use classification with p-stemmer for Arabic text. Authors in [8] studied the effect of stemming on Arabic text classification.

In the area of lexicon development most of the previous work concentrated on morphological analysis either to provide a special corpus or documents set such as [2], [3],

[10] and [12] or in grammar or morphological computational such as [11],[17] , [16].

2. Proposed Algorithms

This paper presents an approach for generating Arabic word roots by emulating the Arabic letters and then generate the corresponding word, morphology, and adding the meaning to resulted Arabic lexicon. Consequently, two main algorithms have been initiated, one for generating roots and the other for developing lexicons.

2.1. Generating the Roots

Most of the roots of Arabic words are made up of three characters in form of (فعل) template, no roots less than three characters [20] exist in Arabic language. The roots that consist of two characters are usually has repeated characters or bugs. There are some roots, which consist of 4 or 5 characters, these types of roots are for special words which can be studied individually such as (دحرج. بعثر. زلزل.) as for 4 characters root, which is in form of (فعلل) template, and (سفرجل . زبرجد. غصنقر) as for 5 characters root, which is in form of (فعللل) template. Therefore, most of Arabic linguistics and morphological analysis are based on three-character root. Accordingly, our proposed algorithms work in generating three characters root. The root generator algorithm is formulated according to the following approach:

- 1- State the framework (the mathematical permutation and combination theory).
- 2- Apply the algorithm.
- 3- Present the results.
- 4- Test the results.

2.1.1. The framework for our proposed algorithm

Our proposed algorithm for generating the Arabic roots based on the mathematical permutation and combination theory. According to the permutation theory (equation (1)) if we want to arrange unrepeated three items out N elements then the result will be equal to $N \times (N-1) \times (N-2)$, and it is equal to $N \times N \times N$ if it has repeated items.

$$P_r^n = \frac{n!}{(n-r)!} \quad (1)$$

In Arabic language there are 28 alphabetical letters or characters (see table 1). The three characters in an Arabic root can be a combination or a permutation of any three letters without repeating of one character three times, repeating of two characters is allowed. Based on this

concept and on the permutation theory we can generate two types of correct Arabic roots:

- 1- Maximum Arabic word roots: the roots that can have adjacent repeated two characters. This kind of roots can be generated using the following formula:
 - a. Find all possible roots, the maximum three arrangement groups of characters (i.e. with one repeated character). This is equal to 21952 (28×28×28)
 - b. Remove the repeated one character root (28 roots: the number of Arabic alphabets), then the maximum number of roots that can be generated from 28 Arabic alphabetical letters will be 21924 roots (21952-28). These roots represent the maximum possible of all Arabic Language roots.
- 2- The strong Arabic roots: the roots that have no adjacent repeated two characters, these kinds of roots represent the roots of the common used Arabic words. To generate this kind of root we need to find a way for avoiding the occurrence of the adjacent repeated two characters, if we do that the total number of roots must be equal to 28×27×27= 20412.

Table 1: Arabic alphabetic

#	Arabic letter	English name	Pronounce
1	أ	Alef	A
2	ب	Baa'	Bah
3	ت	Taa'	Tah
4	ث	Thaa'	Thah
5	ج	Jeem	Jah
6	ح	Haa'	Hah
7	خ	Khaa'	Khah
8	د	Dal	Dah
9	ذ	Thal	Thah
10	ر	Raa'	Rah
11	ز	Zain	Zah
12	س	Seen	Sah
13	ش	Sheen	Shah
14	ص	Saad	Sua
15	ض	Daad	Dhah
16	ط	Taa'	Tah
17	ظ	Zaa'	Zhah
18	ع	Aain	Aaa'
19	غ	Ghain	Ghah
20	ف	Faa'	Fah
21	ق	Qaaf	Qah
22	ك	Kaf	Kah
23	ل	Lam	Lah
24	م	Meem	Mah
25	ن	Noon	Nah
26	هـ	HHaa'	HHah
27	و	Waw	Wah
28	ي	Yah	Yah

2.1.2 Root Generation algorithm

To generate all possible three characters root, either the maximum root or the strong root, we will use three nested FOR loops, as it shown in fFigure 2). For the maximum root we use the operation AND inside IF statement clause. For the Strong root we use the operation OR instead of AND. The results (the roots) are stored into structured database table because the number of roots is large and we are going to use this database table also to store electronic information about lexicon. We used visual basic.net and Microsoft access for implementing our algorithm. Figure 3) shows snapshot of actual code that used in generating our roots, while Figure 4) show snapshots for the database tables that used to store the maximum roots and strong roots respectively.

```

Array [أ, ب, ت, ث, .....]
For i=0 to 27
For j=0 to 27
For k= 0 to 27
If (Array[i] <> Array[j]) AND/OR ( Array[i] <> Array[k])
Then
Root= Array[i] & Array[j] & Array[k]
New Database record = root
End if
i=i+1
j=j+1
k=k+1
    
```

Figure 2: Root generation Algorithm

2.1.3 Root Generation results

The implementation of the maximum root generation algorithm provides all possible Arabic roots, 21924 roots, including adjacent repeated two characters, (figure 4-a) which is equal to the number that is calculated by the emulation equation (see section 2.1.1). All roots are made of three characters; no two roots have three repeated characters.

The roots that are generated by this algorithm (maximum root) can be suitable to use as a basis and framework for all Arabic language, because it generates all possible roots in Arabic language. The only point against this generator is that, it may generate roots that are not an actual Arabic word root, especially when the root includes adjacent repeated two characters, but we think this is an advantage because it may cover the extinct Arabic language, or the special Arabic language. (There are varieties of Arabic language according to varieties of regions or tribes). Accordingly, we suggest that the maximum roots can be used as reference or a useful Arabic language dataset.

```

Module1: Form1: Design: root
[Button] Click
Private Sub Button1_Click(ByVal sender As System.Object, ByVal e As System.EventArgs) Handles Button1.Click
    Dim str As String
    Dim r As Array = {"ا", "ب", "ت", "ث", "ج", "ح", "خ", "د", "ذ", "ر", "ز", "س", "ص", "ض", "ط", "ظ", "ع", "ف", "ق", "ك", "گ", "ل", "م", "ن", "ه", "و", "ي"}
    Dim i, j, k As Integer
    ProgressBar1.Maximum = 21900
    ProgressBar1.Step = 100

    For i = 0 To 27
        For j = 0 To 27
            For k = 0 To 27
                If r(i) = r(j) Or r(j) = r(k) Then
                    Str = ""
                Else
                    Con.Open()

                    Str = "Insert into alaini values("
                    Str += "'" & i & "' & "'" & j & "' & "'" & k & "'" & """"
                    Str += ",)"
                    Str += "'" & r(i) & "'" & "'" & r(j) & "'" & "'" & r(k) & "'" & """"
                    Str += ",)"
                    ProgressBar1.PerformStep()
                    ProgressBar1.Show()

                    Cmd = New OleDbCommand(Str, Con)
                    Cmd.ExecuteNonQuery()
                    Con.Close()

                    Dim = New OleDbDataAdapter("SELECT * FROM alaini ORDER BY ind", Con)
                    Dim.Fill(0, "alaini")
                    Con.Close()
                End If
            Next
        Next
    Next
End Sub
    
```

Figure 3: Snapshot of part of the root generation algorithm VB.Net code

The implementation of the strong root generator algorithm provides 20412 Arabic roots (figure 4-b) which is equal to the number that calculated by the emulation equation (see section 2.1.1). All roots are made of three characters without any three repeated characters or adjacent repeated two characters. The roots that generated by this algorithm has the characteristics of common and most used Arabic words. Accordingly, this root can be suitable for many of NLP process.

The suggested root generator algorithm is used for generating two kinds of roots, as it presented above. In addition, this algorithm can generate another kind of roots, for examples; if we remove one of the For Loops we can get two characters root, in the other hand if we add additional For-Loop we can get four characters root, and so on. Also, if we made any changes in the if-condition or if-operators we can get different kinds of roots, (figures (2) and (3)).

2.2. Developing the Arabic Lexicon

The proposed Arabic lexicon is developed in two steps: generating the morphology and providing the meaning.

ind	root
1	بزء
2	بزء
3	بزء
4	بزء
5	بزء
6	بزء
7	بزء
8	بزء
9	بزء
10	بزء
11	بزء
12	بزء
13	بزء
14	بزء
15	بزء
16	بزء
17	بزء
18	بزء
19	بزء
20	بزء
21	بزء
22	بزء
23	بزء
24	بزء
25	بزء

(A)

ind	root
1	ابء
2	ابء
3	ابء
4	ابء
5	ابء
6	ابء
7	ابء
8	ابء
9	ابء
10	ابء
11	ابء
12	ابء
13	ابء
14	ابء
15	ابء
16	ابء
17	ابء
18	ابء
19	ابء
20	ابء
21	ابء
22	ابء
23	ابء
24	ابء
25	ابء

(B)

Figure 4: Snapshot of part of the roots tables (a) part of the maximum root table (25 records out of 204120) (b) part of strong root table (25 records out 20412)

2.2.1 Generating the morphology

Arabic language is based on root [20]; hence, all Arabic linguistic and morphology analysis is carrying on the structural templates based on word root. Arabic morphology templates are providing the format in which words are structured (Table 2).

Table 2: An example of Arabic morphology templates

فعل	فاعل	فعلول	فعليل
مفعل	مفاعل	مفعول	مفاعيل
تفعل	يفعل	فعلال	يفتوعل
افعل	افعال	افتعال	يفتعل
استفعل	استفعال	مستفعل	يستفعل

The morphology generation algorithm generates different words from one root by inserting or adding letters to that root at certain positions according to the morphology template (Table 2). As it mentioned in section (2.1) the three characters root is of “فعل” template can be used as a basis for generating new words by adding characters at different positions. For example, if we want to generate a word in “فاعل” template then we will insert the letter “ا” at the second position, from right of the root string. Figure 5) shows how the algorithm works.

```

Open the root database table
While not end of the root table
    String = root
    Word1 in "فاعل" template = insert "ا" into
the second position of string
    Word2 in "فعلول" template = insert "و" into
the third position of string
    Word1 in "فعليل" template = insert "ا" into
the second position of string
    .....
    .....
    Update root set value = ....., .....
    Move to next record
Update table
End loop
    
```

Figure 5: Morphology generator algorithm

The morphology generation algorithm is implemented using VB.Net and MS-Access database. Four templates are used to generate four words, as an example, given that the same algorithm can be used for generating large number of words (all possible words can be generated). Figure 6) shows snapshot of part of the code that used in generating these words and Figure 7) shows snapshot of part of the database table where the roots and the words that are generated as a morphology from that roots are stored. This table presents the first stage in our proposed lexicon.

2.2.2 Providing the meaning

The last stage in developing our proposed Arabic lexicon is to provide meaning and examples for our generated roots and the words that structured from those roots. To do this we can use any searching algorithm that searches for the root or the structured word in one of the well-known Arabic datasets and then add the searching results to our proposed lexicon (the database records).

Figure 8) shows a snapshot of the code that we developed to search for the meaning of the generated words from one of holy Quran datasets [18]. The searching results are then added to our proposed lexicon, these results include all available data in the dataset, such as the Arabic meaning, English meaning, an example sentences include the word (the Quran verse that contains the word or the root), the pronouns ... etc. Figure 9) shows part from the resulted lexicon.

```

Con.Open()
Dad = New OleDbDataAdapter("SELECT * FROM alain1 ORDER BY ind", Con)
Dad.Fill(Dst, "alain1")

Con.Close()
Dim Str, s As String
Dim i, j As Integer
i = CurrentRow
j = Dst.Tables("alain1").Rows.count - 1
Try

Con.Open()
While i <> j

s = Dst.Tables("alain1").Rows(CurrentRow)("root")
Str = LSet(s, 1) & "ا" & Mid(s, 2)
Dst.tables("alain1").rows(CurrentRow)("faail") = Str
Str = "م" & LSet(s, 1) & "ا" & Mid(s, 2)
Dst.tables("alain1").rows(CurrentRow)("mfaail") = Str

Str = LSet(s, 2) & "و" & Mid(s, 3)
Dst.tables("alain1").rows(CurrentRow)("faool") = Str
Str = "م" & LSet(s, 2) & "و" & Mid(s, 3)
Dst.tables("alain1").rows(CurrentRow)("mfaool") = Str

Str = LSet(s, 2) & "ي" & Mid(s, 3)
Dst.tables("alain1").rows(CurrentRow)("faail") = Str
Str = "م" & LSet(s, 2) & "ي" & Mid(s, 3)
Dst.tables("alain1").rows(CurrentRow)("mfaail") = Str

Str = "ي" & s
Dst.tables("alain1").rows(CurrentRow)("yfal") = Str
    
```

Figure 6 :Snapshot of part from morphology generator VB.Net code

ind	root	faail	faool	faail	mfaail	mfaool	mfaail	yfal	tfal
3875	خزل	خزائل	خزول	خزول	مخزائل	مخزول	مخزول	يخزل	تخزل
14295	فعلل	فعلائل	فعلول	فعلول	مفعلائل	مفعلول	مفعلول	يفعلل	تفعلل
17764	لجج	لججال	لجول	لجول	ملججال	ملجول	ملجول	يلجج	تلجج
17343	لجج	لججال	لجول	لجول	ملججال	ملجول	ملجول	يلجج	تلجج
9288	صعز	صعزائل	صعزول	صعزول	مصعزائل	مصعزول	مصعزول	يصعز	تعصز
4195	خقل	خقائل	خقول	خقول	مخقائل	مخقول	مخقول	يخقل	تخقل
17386	لجف	لجفائل	لجفول	لجفول	ملجفائل	ملجفول	ملجفول	يلجف	تلجف
14215	فعض	فعضائل	فعضول	فعضول	مفعضائل	مفعضول	مفعضول	يفعض	تفعض
12861	عمزو	عمازوا	عمازوا	عمازوا	ممازوا	ممازوا	ممازوا	يمازوا	تمازوا
1757	نعب	نعبائل	نعبول	نعبول	منعبائل	منعبول	منعبول	ينعب	تنعب
12250	عاد	عادل	عادل	عادل	معادل	معادل	معادل	يعاد	تعاد
8036	شكك	شككائل	شككول	شككول	مشككائل	مشككول	مشككول	يشكك	تشكك
8937	شكك	شككائل	شككول	شككول	مشككائل	مشككول	مشككول	يشكك	تشكك
8941	شكك	شككائل	شككول	شككول	مشككائل	مشككول	مشككول	يشكك	تشكك
3894	خزص	خزصائل	خزصول	خزصول	مخزصائل	مخزصول	مخزصول	يخزص	تخزص
19630	هز	هزائل	هزول	هزول	مهزائل	مهزول	مهزول	يهز	تهز
17391	لحن	لحنائل	لحنول	لحنول	ملحنائل	ملحنول	ملحنول	يلحن	تلحن
17390	لحم	لحمائل	لحمول	لحمول	ملحمائل	ملحمول	ملحمول	يلحم	تلحم
14550	فوز	فوزائل	فوزول	فوزول	مفوزائل	مفوزول	مفوزول	يفوز	تفوز
17387	لحن	لحنائل	لحنول	لحنول	ملحنائل	ملحنول	ملحنول	يلحن	تلحن
12606	عصب	عصبائل	عصبول	عصبول	مصصبائل	مصصبول	مصصبول	يصصب	تعصص
3890	خز	خزائل	خزول	خزول	مخزائل	مخزول	مخزول	يخز	تخز
14206	فصع	فصعائل	فصعول	فصعول	مفصعائل	مفصعول	مفصعول	يفصع	تفصع
3885	خزج	خزجال	خزول	خزول	مخزجال	مخزول	مخزول	يخزج	تخزج
14538	فوت	فوتائل	فوتول	فوتول	مفوتائل	مفوتول	مفوتول	يفوت	تفوت

Figure 7: Snapshot of part of database table shows the roots and the words that generated as morphology for those roots

with 21,924 records, and the other with 20,412 records. Each record presents a root with its corresponding words, meaning, English translation ...etc (see figure (9)). In this paper, our concern was mainly to provide an algorithm that can be used to generate Arabic language dictionary including roots and other comprehensive lexicon.

References

- [1] Ababneh, Mohamad, et al. "Building an Effective Rule-Based Light Stemmer for Arabic Language to Improve Search Effectiveness." *International Arab Journal of Information Technology (IAJIT)* 9.4 (2012).
- [2] AbdelRaouf, Ashraf, et al. "Arabic Corpus Enhancement using a New Lexicon/Stemming Algorithm." *ICPRAM*. 2013.
- [3] AbdelRaouf, Ashraf, et al. "Building a multi-modal Arabic corpus (MMAC)." *International Journal on Document Analysis and Recognition (IJ DAR)* 13.4 (2010): 285-302.
- [4] Abu-Errub, Aymen, et al. "Arabic Roots Extraction Using Morphological Analysis." *International Journal of Computer Science* 11 (2014): 2.
- [5] Algarni, Mohammed, et al. "Simple Arabic Stemmer." *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*. ACM, 2014.
- [6] Ali Helmi Moussa, 1980 "A Comparison Study Between the Words Roots in Quran and Words Roots in Taj Alarooos" *Alam Alfikir*, Kuwait. Quran research special edition in 12- 1982.
- [7] Al-Kabi, Mohammed N., et al. "A novel root based Arabic stemmer." *Journal of King Saud University-Computer and Information Sciences* 27.2 (2015): 94-103.
- [8] Al-Kabi, Mohammed, Emad Al-Shawakfa, and Izzat Alsmadi. "The Effect of Stemming on Arabic Text Classification: An Empirical Study." *Information Retrieval Methods for Multidisciplinary Applications* (2013): 207.
- [9] Atwan, Jaffar, Masnizah Mohd, and Ghassan Kanaan. "Enhanced arabic information retrieval: Light stemming and stop words." *Soft Computing Applications and Intelligent Systems*. Springer Berlin Heidelberg, 2013. 219-228.
- [10] Bekkali, Mohammed, and Abdelmonaime Lachkar. "ARABIC TWEETS CATEGORIZATION BASED ON ROUGH SET THEORY." *International Journal of Computer Science & Information Technology* 6.6 (2014).
- [11] Boudelaa, Sami, and William D. Marslen-Wilson. "Morphological units in the Arabic mental lexicon." *Cognition* 81.1 (2001): 65-92.
- [12] Darwish, Kareem, Walid Magdy, and Ahmed Mourad. "Language processing for arabic microblog retrieval." *Proceedings of the 21st ACM international conference on Information and knowledge management*. ACM, 2012.
- [13] Haralambous, Yannis, Yassir Elidrissi, and Philippe Lenca. "Arabic Language Text Classification Using Dependency Syntax-Based Feature Selection." *arXiv preprint arXiv:1410.4863* (2014).
- [14] Harrag, Fouzi, Eyas El-Qawasmah, and Abdul Malik S. Al-Salman. "Stemming as a feature reduction technique for Arabic Text Categorization." *Programming and Systems (ISPS)*, 2011 10th International Symposium on. IEEE, 2011.
- [15] Kanan, Tarek, and Edward A. Fox. "Automated Arabic Text Classification with P-Stemmer". *Machine Learning, and a Tailored News Article Taxonomy*. Department of Computer Science, Virginia Polytechnic Institute & State University, 2015.
- [16] Khaliq, Bilal, and John Carroll. "Induction of Root and Pattern Lexicon for Unsupervised Morphological Analysis of Arabic." (2013): 1012-1016.
- [17] Maamouri, Mohamed, et al. "The penn arabic treebank: Building a large-scale annotated arabic corpus." *NEMLAR conference on Arabic language resources and tools*. 2004.
- [18] Mohamed Osman Hegazi, Anwer Hilal, and Mohammad Alhawarat 2015. "Fine-grained Quran Dataset" *International Journal of Advanced Computer Science and Applications(IJACSA)*, 6(12): 262-267, 2015.
- [19] Mohammad Alhawarat, Mohamed Hegazi and Anwer Hilal, 2015. "Processing the Text of the Holy Quran: a Text Mining Study" *International Journal of Advanced Computer Science and Applications(IJACSA)*, 6(2): 262-267, 2015.
- [20] Mustafa Ahmed Adel-Alim, 2010, "Arab characteristics between the ancient and modern" *United Arab Emirates University*,2010
- [21] Nehar, Attia, Djelloul Ziadi, and Hadda Cherroun. "Rational Kernels for Arabic Stemming and Text Classification." *arXiv preprint arXiv:1502.07504* (2015).



Mohamed Osman Hegazi has got his Ph.D. in Distributed Database from Sudan University of Science & Technology, Khartoum Sudan in 2004. He has worked as an associate professor in Alzaiem Alazhari University, Khartoum Sudan for 4 years, he then joined the College of Computer Engineering and Sciences of Prince Sattam Bin Abdulaziz University since 2013 as an associate professor of computer science.