# ENHANCING THE ACCURACY OF SONBOL'S ARABIC ROOT EXTRACTION ALGORITHM

Nisrean Thalji[1], Nik Adilah Hanin[1], Zyad Thalji[2] and Sohair Al-Hakeem[3]

## ABSTRACT

*Root extraction is an important primary process in most Arabic applications, such as information retrieval systems, text mining, text classifiers, question answering systems, data compression, indexes, spelling checkers, text summarization and machine translation. Any weaknesses of root extraction will affect negatively the performance of these applications. Sonbol's Arabic root extraction algorithm achieves high accuracy of performance and gives new classification for Arabic's letters which minimizes the affix ambiguity. The comparison and testing of the existing Arabic root extraction algorithms on unify datasets shows that they still need some enhancements. Arabic root extraction is mainly based on using patterns, where as much as the algorithm has patterns as much as the accuracy is better. In this study, we improve Sonbol's Arabic root extraction algorithm, by enhancing its rules and increasing its patterns. We use 4320 patterns to extract the roots, which is the largest patterns' list extracted by Thalji's corpus. We test the new algorithm on Thalji's corpus that contains 720,000 word-root pairs. This corpus is mainly built to test and compare Arabic root extraction algorithms. The new algorithm is compared with Sonbol's Arabic root extraction algorithm. The algorithm of Sonbol et al. achieves an accuracy of 68%, whereas the new algorithm achieves an accuracy of 92%.*

## 1. INTRODUCTION

Arabic language is one of the most used Semitic languages. Semitic languages are spoken in a number of regions that were common in distant times in many regions of Africa and Asia over many decades. Some of these languages are not used now, such as Akkadian, Assyrian and Babylonian and some languages are still used nowadays, such as Arabic, Hebrew and Syriac. Semitic languages are a branch of the Afro-Asiatic language family originating in the Middle East (Bennett, 1998) [1].

In Arabic, vowels are used to ensure the exact meanings of words. If the word is non-vocalized, in many cases it is an ambiguous word and then we need to read the sentence and sometimes the whole text to understand the exact meaning. These vowels are written above or under the letter. Table 1 shows vowels in Arabic and corresponding letter/s in English. Some words in non-vocalized texts may have more than one meaning (ambiguous words). So, they have different roots. For example, the non-vocalized Arabic word "والدين/WALDN" has three possible words "وَالَدَيْن/WALEDAYN", "وَالدّيْن/WA ADDAYN" and "وَالدّين/WA ADDEEN". And then, the possible roots are "وَلَد/WALAD" (son), "دَيْن/DAYN" (debt) and "دِين/DEEN" (religion). Another example is the non-vocalized Arabic word "كتب/KTB" which has three possible interpretations: "كَتَبَ/ KATABA" (he wrote), "كُتِبَ/ KUTIBA" (has been written), and "كُتُبٌ/ KUTUBUN" (books) [2]. We converted the Arabic letters and words into Latin characters (uppercase), so that the reader who is not familiar with Arabic texts can read it with ease. This way, the reader will be able to read words as the way they sound phonetically.

Root extraction is an important primary process in most Arabic applications, such as information retrieval systems, text mining, text classifiers, question answering systems, data compression, indexes, spelling checkers, text summarization and machine translation.

---

1. N. Thalji and N. Hanin are with Department of Computer Engineering, School of Computer and Communication Engineering, University Malaysia Perlis, Malaysia. Emails: nnthalji1980@gmail.com, adilahhanin@unimap.edu.my
2. Z. Thalji is with Department of Management Information System, Imam Abdulrahman Bin Faisal University, Kingdom of Saudi Arabia. Email: zythalji@iau.edu.sa
3. S. Al-Hakeem is with Department of Computer Science, Ajloun National University, Jordan. Email: drsohair@gmail.com

"Enhancing the Accuracy of Sonbol's Arabic Root Extraction Algorithm", N. Thalji, N. Hanin, Z. Thalji and S. Al-Hakeem.

Table 1. Vowels in Arabic and corresponding letter/s in English.

| No. | Vowels in Arabic | Corresponding letter/s in English |
|:---:|:---:|:---:|
| 1 | ◌ | E |
| 2 | ◌ | O |
| 3 | ◌ | A |
| 4 | ◌ | No letter |
| 5 | ◌ | En |
| 6 | ◌ | Un |
| 7 | ◌ | An |
| 8 | ◌ | Duplicate the letters |

Therefore, many Arabic root extraction algorithms are presented in many different studies that tried to algorithm of Sonbol, Ghneim and Desouki [3]. The algorithm of Sonbol et al. comes after the algorithm of Khoja and Garside [4], which is well-known in extracting Arabic roots. Khoja and Garside algorithm's accuracy amounted to 95% when they tested their algorithm with their own dataset. Sonbol et al. tried to improve Arabic root extraction algorithms in order to increase the percentage of accuracy. The algorithm of Sonbol et al. accuracy amounted to 98% when they tested their algorithm using their own dataset.

Al-Shawakfa, Al-Badarneh, Shatnawi, Al-Rabab'ah and Bani-Ismail [5] made a comparison study of existing Arabic root extraction algorithms. This comparison included the algorithm of Sonbol et al., Khoja and Garside algorithm and other algorithms. This comparison was conducted in a unified dataset, in order to evaluate these algorithms fairly. Khoja and Garside algorithm's accuracy was 34%, whereas Sonbol algorithm's accuracy was 24%. Variance in the accuracy values is due to the differences of datasets that were used in the testing process. The study by Al-Shawakfa et al. revealed that existing Arabic root extraction algorithms still need more improvement. It also presented some weaknesses of Khoja and Garside's algorithm and the algorithm of Sonbol et al. In this study, we continue completing Sonbol's work by improving their algorithm.

There are three different approaches to extract the roots of the word; rule-based approach, lookup table approach and statistics-based approach. Recently, most of the root extraction algorithms are rule-based approach [5]. This approach mainly has two parts; the lists of affixes (roots and patterns) and the rules. Each Arabic root extraction algorithm tried to enhance the lists and/or the rules.

However, most algorithms suffer from the following problems:

- There are neither standard dataset nor complete lists of Arabic affixes, patterns and roots. Each work collects or generates their own dataset or lists. Most of the lists which contains the affixes, patterns, roots and lists that they used in each work are not publicly available. They just wrote samples of these lists. As a result, every time a new root extraction method is proposed, the researchers need to collect their own data or generate their own list independently. This will cause overlapped works, where each work is trying to solve the same problem instead of improving each other's work, which resulted a waste of time and resources. In addition, the lists used might have significant difference in terms of number of words, which will make it difficult for researchers to fairly compare the performance of existing works. Therefore, in recent works, the researchers tried to extend these lists by adding new contents [6].

- Arabic has a complex structure, which makes it difficult to extract the roots [7]. All Arabic root extraction algorithms suffer from affixes' ambiguity, so that it is difficult to distinguish between affix letters and root letters.

This work focus on solving these problems. The structure of this paper is organized as follows; in Section 2, different related previous studies and their drawbacks are discussed. Section 03 describes the proposed methodology which includes the details of each process. Section 0 explains the experimental implementation of our algorithm and its evaluation process. Section 5 concludes the main points of the paper and gives some future directions.

161

Jordanian Journal of Computers and Information Technology (JJCIT), Vol. 04, No. 03, December 2018.

## 2. PREVIOUS STUDIES

In this section, we give a brief overview of previous rule-based Arabic root extraction algorithms. Khoja and Garside and Garside algorithm [4] is a very popular rule-based Arabic root extraction algorithm.

Khoja and Garside and Garside algorithm reported 96% accuracy of their stemmer using newspaper text. Al-Shalabi [8] presented an Arabic root extraction algorithm, which is a rule-based algorithm that is used to extract trilateral roots of Arabic words. This algorithm has been tested on a corpus of 72 abstracts 10582 words from the Saudi Arabian National Computer Conference, where its accuracy was about 92%.

Al-Kabi and Al-Mustafa algorithm [9] is based on affix removal. They tested their algorithm on small data sets containing 1,827 words. The system failed to analyze 55 words, since their patterns are unknown. This failure was mostly due to foreign (Arabized) words. The system enables to analyze the rest 1,772 words and achieved 91% of accuracy.

Sonbol, Ghneim and Desouki [3] algorithm is a rule-based root-extraction algorithm, the principal idea of which is based on encoding Arabic letters with a new code that preserves morphologically useful information and simplifies its capturing toward retrieving the root. They conducted their experiments using two different corpuses. The first corpus consisted of lists of word-root pairs 167,162 pairs. The second corpus was a collection of 585 Arabic articles from different categories (policy, economy, culture, science and technology and sport). This corpus consisted of 377,793 words. In general, the accuracy was about 96%-98%.

Another work is Ghwanmeh, Al-Shalabi, Kanaan, Khanfar and Rabab'ah algorithm [10], which is a rule-based algorithm used to find trilateral Arabic roots. According to Ghwanmeh et al.., their algorithm has only failed to analyze words that are normally foreign, irregular or do not have trilateral roots. A corpus of 242 abstracts from the Proceedings of Saudi Arabian National Computer conference in machine-readable form was used in the testing procedure. The set of abstracts was chosen randomly from the corpus for analysis. The results obtained showed that the algorithm extracts the correct roots with an accuracy rate up to 95%. Many algorithms have been conducted under this type, like the Kchaou and Kanoun algorithm [11], El-Defrawy, El-Sonbaty and Belal algorithm [12] and Ayedh and Guanzheng algorithm [13].

Also, many morphological analyzers have been conducted to properly provide maximum morphological information on Arabic words, such as the proclitic, the prefix, the lemma, the suffix, the stem, the root, the enclitic, the tag and the pattern. One of them is MADAMIRA [14], a morphological analyzer that provides many information on Arabic words. MADAMIRA combines two morphological analysis systems; MADA [15], [16] and [17] and AMIRA [18].

In addition, Al-Khalil Morphological System 2 [19] is a recent morphological analyzer that provides many information on Arabic words. It deals with vocalized and non-vocalized Arabic words. It overcomes many errors of the previous system Al-Khalil Morphological System 1.

In general, all rule-based Arabic roots extraction algorithms share seven processes, which are: normalization, removing prefixes and suffixes, matching the word against patterns, extracting the roots from the patterns, comparing the roots with the roots' list, returning the extracted roots and finally making enhancement to extract the correct roots, as shown in Figure 1.

The difference between one algorithm and the others is in the details of each process. Also, every algorithm has different lists of prefixes, suffixes, roots and patterns. The main problem is in process two and process four. In process two, in many cases, the algorithms remove the matched prefixes and suffix letters, but these letters are part of the root. So, the result is a wrong root. Our proposed solution to this problem is done by not removing prefixes or suffixes and collecting more rules to reduce affix ambiguity. In process four, the algorithms match the word against patterns to extract the root. The main problem is the limited number of patterns that are collected till now. The extraction accuracy will improve if the algorithm can test as many as existing word patterns. Our proposed solution to this problem is using longer pattern lists. The proposed algorithm uses Thalji's pattern [6], which is the most recent list. This list is automatically generated from most of the Arabic dictionaries and contains (4320) patterns, which is the longest list discovered until now. These patterns contribute to enhancing the accuracy of Arabic root extraction algorithms.

Process 1 : Normalization

Process 2 : Remove prefixes and suffixesz

Process 3 : Match the word against patterns

Process 4 : Extract the roots from the patterns

Process 5 : Compare the roots with the roots' list

Process 6 : Return the extracted roots

Process 7 : Make enhancement to the extracted roots

Figure 1. Main processes in rule-based Arabic root extraction algorithms.

## 3. METHOD

The root is the base form of the word that gives the main meaning of the word. In this section, the methodology for the proposed Arabic root extraction algorithm is explained. The proposed algorithm is an enhancement of the algorithm of Sonbol et al. by increasing the rules and the lists to find all possible roots of the word.

### 3.1 Normalization

The normalization steps for the algorithm of Sonbol et al. are as follows:
- Removing the kasheeda symbol ("_").
- Removing the diacritics.
- Replacing the Hamza's forms (ء, آ, إ, ؤ, ئ) with the letter (أ).
   In this section, we extend Sonbol's normalization process by the following steps:
- Removing the punctuations.
- Removing the non-letters.
- Duplicating any letter that has the Shaddah: "ّ" symbol.

### 3.2 Encoding

In this step, Sonbol coded the Arabic letters based on six symbols {O, P, S, PS, U, A}, representing six groups of letters each of which shares certain characteristics:

**O**: Original letters. These letters are surely part of the root. They are: ("ث", "THAA"), ("ج", "JEEM"), ("ح", "HAA"), ("خ", "KHAA"), ("د", "DAL"), ("ذ", "THAL"), ("ر", "RAA"), ("ز", "ZAY"), ("ش", "SHEEN"), ("ص", "SAD"), ("ض", "THAD"), ("ظ", "DAA"), ("ط", "TAA"), ("ع", "AYEN"), ("ق", "GAF"), ("غ", "GAYEN"). This means that if the word contains one or more of these letters, these letters should be part of root's letters.

**P:** Prefix letters "ب, ف, س, ل" (BAA, FAA, SEEN, LAM). These letters should be treated as part of the root word if they appear in a different part of the word other than the prefix part. Otherwise, these letters are considered ambiguous letters (can be part of the root word or added letters to the root word). If these letters are ambiguous letters, the algorithm initially considers them as prefix letters. There is a possibility for these letters to become (O) letters (root's letters) in the next steps.

**S:** Suffix letter ("ه", Haa). This letter should be treated as part of the root word if it appears in a different part of the word other than the suffix part. Otherwise, this letter is considered as an ambiguous letter. If this letter is an ambiguous letter, the algorithm initially considers it as suffix letter. There is a possibility for this letter to become (O) letter (root's letter) in the next steps.

**PS:** Prefix-Suffix letters "ك, م, ن" (KAF, MEEM, NOON). These letters can appear only on both sides of the word; i.e., in the suffix part or in the prefix part. These letters should be treated as part of the root

163

Jordanian Journal of Computers and Information Technology (JJCIT), Vol. 04, No. 03, December 2018.

word if they appear in a different part of the word other than the prefix-suffix part. Otherwise, these letters are considered ambiguous (can be part of the root word or added letters to the root word). If these letters are ambiguous letters, the algorithm initially considers them as prefix-suffix letters. There is a possibility for these letters to become (O) letters (root's letters), (P) prefix letters or (S) suffix letters in the next steps.

**U:** Uncertain letters ("ت, ي, و, ا" (TAA, YAA, WAW, ALEF)). These letters can appear anywhere in the word. It is not possible to verify whether these letters are part of the root word letters. Several cases are associated with these letters as they may change, omit or convert from one letter to another during the derivation process following well known Arabic rules "EBDAL and EALAL". For example, the Arabic word"قيْل/KEEL" (It was said) and its root "قول", the letter "و" is converted into "ي" during the derivation process.

**A:** Added letter ("ة "only (TAA MARBUTA)). This letter is always considered an additional letter. This letter is always deleted.

### 3.3 Some Improvements to the Last Coding

In this section, Sonbol et al. added some improvements to the last coding, by applying the following conditions:

- The letter "ب" (BAA) is a prefix letter if it is situated in the first three letters; otherwise, it is an original letter (part of the root word letter). This rule means that if the letter "ب" (BAA) is situated among the first three letters, it is an ambiguous letter. On the other hand, it will be part of the root word if it appears in place other than the first three letters of the word. For example, with the words " بأمرك, وبالصــعيد"(BEAMRK, WBLSA"EED), the letter "ب" (BAA) is situated among the first two letters. So, it's an ambiguous letter. Initially, the algorithm considers it as a prefix letter. Then in the next steps, it may change to (O) letter (root's letter). Another example is the word "افبالباطل/AFBLBATL"; this word has the letter "ب" (BAA) appearing twice in the word. The first "ب" (BAA) letter is an ambiguous letter, as it appears as the third letter in the word. Initially, the algorithm considers it as a prefix letter. Then, in the next steps, it may change to (O) letter (root's letter). The second "ب"(BAA) will be considered part of the root word, as it appears as the sixth letter of the word.

- The letter "ف" (FAA) is a prefix letter if it appears among the first two letters; otherwise, it is an original letter. This means that if the letter "ف" (FAA) appears among the first two letters of the word, it is an ambiguous letter. If the letter "ف" (FAA) appears in a place other than the first two letters, it is part of the root word. Initially, the algorithm considers it as prefix letter. For example, in the words " فهد, فاســتأجرتها, افغير" (FHD/FASTAJRTHA, AFGYR), the letter "ف" (FAA) appears among the first two letters of the word. So, it's an ambiguous letter. Initially, the algorithm considers it as a prefix letter. It may change to (O) letter (root's letter) in the next steps of the algorithm following certain rules. Another example; the words "متفائلين/ MTFAELEEN" and "تسافهت / TSAFHT", "وأتفكت/ WETFKT"; the letter "ف" (FAA) is considered to be part of the root word as it appears in places other than the first two letters of the word.

- The letter "س" (SEEN) is a prefix letter if it is followed by one of the letters " أ, ن, ي, ت " (HAMZA, NOON, YAA AND TAA); otherwise, it is part of the root word. For example, the words "ســتبقى/ STBKA, ســنتواجد /SNTWAJD, ســينجلي /SYNJLE"; the letter "س" (SEEN) is an ambiguous letter, as it is preceding one of the letters " أ, ن, ي, ت " (HAMZA, NOON, YAA AND TAA). Initially, the algorithm considers it as a prefix letter. It may change to (O) letter (root's letter) in the next steps of the algorithm. Another example is with the words "كالأســد, مأســور" (KALASD, MASOR); the letter "س" (SEEN) is considered part of the root word, as it is not preceding one of the letters " أ, ن, ي, ت " (HAMZA, NOON, YAA AND TAA).

- The letter "ل" (LAM) is considered a prefix letter if it appears among the first five letters of the word; otherwise, it is part of the root word.

- The letter "ه" (HAA) is considered a suffix letter if it appears among the last three letters of the word; otherwise, it is part of the root word.

- The letter "ك" (Kaf) is considered a prefix letter if it appears among the first three letters of the word; otherwise, it is a suffix letter.

If the algorithm of Sonbol et al. finds three O-Letters (or more) in the encoded word, these letters are considered root letters and the algorithm is terminated. However, in this work, the enhancement of the algorithm of Sonbol et al. is to continue searching for other possible roots and for longer roots (more than three letter root word).

## 3.4 Applying Transformation Rules

In this section, the algorithm of Sonbol et al. applies transformation rules between groups to obtain a maximum number of original letters. Transformation rules are mentioned below:

- R1) Change each (P) after (O) to (O).

For example, with the word "ضيوف/ DYUF", "ض" (DAA) is an (O) letter, "ف" (FAA) is a (P) letter; in this case, (P) comes after (O). So, "ف" (FAA) is changed to (O) letter, which means that it should be one of the root's letters. Until now "ض, ف" (FAA, DAA) are part of the root's letters.

- R2) Change each (S) before (O) to (O).

For example, with the word "الهداية /ALHDAYH", "د/ DAA" is an (O) letter, "ـهـ" is a (S) letter; in this case, (S) comes before (O). So, "ـهـ/HAA" is changed to (O) letter, which means that it should be one of the root's letters. Until now "د, ـهـ/ DAA, HAA" are part of the root's letters.

- R3) Change each (PS) before (P) to (P).

For example, with the word "كالسيوف/KALSOYUF", "كـ/KAA" is a (PS) letter, "ـسـ/SEEN" is a (P) letter; in this case, (PS) comes before (P). So, "كـ/KAA" is changed to (P) letter, which means that it should be one of the root's letters or prefix letters, but not a suffix letter.

- R4) Change each (PS) before (O) to (P).

For example, with the word "منتقمون/ MNTKMON", "ـنـ/ NOON" is a (PS) letter, "ـقـ/ KAA" is an (O) letter; in this case, (PS) comes before (O). So, "ـنـ/NOON" is changed to (P) letter, which means that it should be one of the root's letters or prefix letters, but surely not a suffix letter.

- R5) Change each (PS) after (S) to (S).

For example, with the word "منتهك/ MNTHK", "ك/KAF" is a (PS) letter, "ـهـ/HAA" is a (S) letter; in this case, (PS) comes after (S). So, "ـك/KAF" is changed to (S) letter, which means that it should be one of the root's letters or suffix letters, but surely not a prefix letter.

- R6) Change each (PS) after (O) to (S).

For example, with the word "علمك/ELMK", "ك/ KAAF" is a (PS) letter, "ـعـ/A'A" is an (O) letter; in this case, (PS) comes after (O). So, "ك/ KAAF" is changed to (S) letter, which means that it should be one of the root's letters or suffix letters, but surely not a prefix letter.

- R7) Change each (P) after (S) to (O).

For example, with the word "التهبت/ELTHBT", "ب/ BAA" is a (P) letter, "ـهـ/ HAA" is an (S) letter; in this case, (P) comes after (S). So, "ب/ BAA" is changed to (O) letter, which means that it should be one of the root's letters. Until now, "ـهـ, ب" (BAA, HAA) are root's letters.

- R8) Change each (S) before (P) to (O).

For example, with the word "البهتان/ ALBHTAN", "ـهـ/HAA" is an (S) letter, "ب/BAA" is a (P) letter; in this case, (S) comes before (P). So, "ـهـ/HAA" is changed to (O) letter, which means that it should be one of root's letters. Until now, "ـهـ, ب" (BAA, HAA) are the root's letters.

As in the previous step, if the algorithm of Sonbol et al. has three O-letters in the encoded word, these letters are considered root letters and the process will terminate here. However, the enhancement of the algorithm of Sonbol et al. is to continue searching for other possible roots and for longer roots, with more than three root length.

## 3.5 Extracting All Possible Patterns of the Word

The algorithm of Sonbol et al. uses the idea of traditional algorithms, but with the aid of the encoded

word. Traditional algorithms store lists of Arabic prefixes, suffixes and patterns. These algorithms delete prefixes and suffixes, then use the pattern to extract the root from the reminder. The enhancement of this process is done by using larger lists' content and not removing the prefixes or suffixes, but applying the patterns. It's worth to mention that the presented algorithm assumes that patterns are composed of three elements: prefix, stem and suffix. One of the problems that were experienced by most of the previous algorithms is that they delete clitics before comparing with patterns. And in many cases, these clitics are parts of roots and not clitics. For example, with the words "كـالحون, التقى"(Altka, kalehon), removing "كال, ال" (AL, KAL) will give these roots "حون, تقى"(TKA, HON), ignoring other possible roots "كلح, لقى"(LKA, KLH). In this section, we use the pattern's list of Thalji's corpus that was automatically extracted [6]. Up until now, this corpus contains the largest list of 4,320 patterns, which is the most appropriate list to be used in this work. Previous algorithms have used short lists that were manually collected. In addition, they did not publicly publish all the lists' contents. Thalji's patterns are listed in appendix A, so that future researchers can benefit from them.

In this step, we compare the word with the Thalji's list of patterns and return all matched patterns. For example, for a word like "فهد/ FHD", the algorithm found two original letters, "د ,ــهــ"(DAA, HAA). Next, the word is compared to the list of patterns and all matched patterns were returned. The word "فهد/ FHD" matches the pattern"فعل / FA'L". The word "فهد/ FHD" is the first possible root. Another example is the word "البحر/ ALBAHER", where the algorithm just finds two original letters, which are"ر, حــ " (RAA, HAA). The word is compared to the list of patterns and all matched patterns were returned. The word "البحر/ ALBHR" matches the pattern "الفعل/ ALFA'L". The word "بحر/ BHR" is the first possible root. Also, this word matches the patterns "افعل, فعلل"(FA"LL/ EFA"L), then "ابحر, البحر" (ALBHR/ ABHR) are also possible roots.

## 3.6 Extracting All Possible Roots for the Word

All possible roots are found by matching the words against the list patterns. All the possible roots that match the patterns are extracted after finding all possible patterns.

## 3.7 Solving the Problems with Ealal Rules and Ebdal Rules

When we have a weak letter (ALEF, YAA AND WAW), we replace this letter with the two other letters and check if the result is a valid root. If so, we add this root to the possible roots. For example, in the word "قـال /KAL", the algorithm replaces "ا/ ALEF" with "ي/ YAA" and "و/ WAW". So, "قول, قيـل" (KEEL, KAWL) are possible roots.

## 3.8 Minimizing Possible Roots by Comparing them with Roots' List

In this section, we use the roots' list of Thalji's corpus that was automatically extracted from most well-known Arabic dictionaries. It is the largest roots' list found till now with 12,000 roots. This list is longer than the list that is used in Ababneh, Al-Shalabi, Kanaan and Al-Nobani stemmer [20]. It has about 11,347 roots. The distribution of roots for these two different lists is shown in Table 2.

Table 2. The distribution of the roots for two different lists.

| Roots | List of Thalji's corpus | List of Ababneh, Al-Shalabi, Kanaan and Al-Nobani stemmer |
|---|---|---|
| Two-letter roots | 500 | 115 |
| Three-letter roots | 7912 | 7198 |
| Four-letter roots | 3180 | 3739 |
| Five-letter roots | 360 | 295 |
| Six-letter roots | 48 | 0 |

The presented algorithm uses Thalji's list to minimize the possible roots, whereas the algorithm of Sonbol et al. used a short list of roots. For example, in the word "البحر/ ALBHR", the possible roots are " بحر, ابحر, البحر" (BHR, ALBHR, ABHR), while the roots "ابحر, البحر" (ALBHR, ABHR) are excluded, because they are not found in the roots' list.

# 4. EXPERIMENT AND EVALUATION

In this section, the presented algorithm is compared with other algorithms with the same approach, which is the rule-based approach. These algorithms are Khoja and Garside's Arabic root extraction algorithm and Sonbol's Arabic root extraction algorithm. In addition, the presented algorithm is also compared with one of the most recent morphological analyzer systems, which is Al-Khalil Morphological System 2. Al-Khalil Morphological System 2 gives maximum morphological information of Arabic words, such as the proclitic, the prefix, the lemma, the suffix, the stem, the root, the enclitic, the tag and the pattern. A complete comparison was conducted between the algorithms on Thalji's corpus in terms of accuracy. Thalji's corpus is an automatic corpus that is built from ten old Arabic dictionaries; this corpus is mainly built to test and fairly compare Arabic root extraction algorithms. This corpus contains 720,000 word-root pairs, which helps to avoid the interference of a human expert normally needed to verify the correct roots of each word used in the testing or comparison process. Moreover, this corpus has more than 4,320 types of words derived from 12,000 roots. Therefore, the list used in this experiment is more comprehensive compared to previous works.

The result of testing shows that the accuracy of Khoja and Garside's algorithm was 63%, the accuracy of the algorithm of Sonbol et al. was 68%, the accuracy of Al-Khalil Morphological System 2 was 75%, whereas the accuracy of the presented algorithm was 92%. Figure 2 shows the performance accuracy of all compared algorithms.



Figure 2. Accuracy of the algorithm of Khoja and Garside, the algorithm of Sonbol et al., the algorithm of Al-Khalil and the presented algorithm.

The main problems of Khoja and Garside's algorithm are that it does not consider many roots, prefixes, suffixes and patterns. It suffered from affix ambiguity problems. In addition, it returned just one solution for non-vocalized words, ignoring other possible solutions. Besides, it replaced a vowel letter with the letter "و" that sometimes returns a root that is not related to the derivation word. Finally, it produced wrong roots being unsuccessful to extract roots for derivation words that contain the "ابدال/ EBDAL" rule.

The main problems of the algorithm of Sonbol et al. arise if the root does not contain any constant letter, if the root does not start with a constant letter or if the root contains only one constant letter. Also, it does not consider many roots, prefixes, suffixes and patterns. In addition, it returned just one solution for non-vocalized words, ignoring other possible solutions.

The main problems of Al-Khalil Morphological System 2 are that it failed to analyze some words, which were about 25% of the input words. Table 3 shows a sample of these words. For example, the words

167

Jordanian Journal of Computers and Information Technology (JJCIT), Vol. 04, No. 03, December 2018.

"ارصـــــدت, مقـاصـــر, مـدارح, جـارودة, شـــــجير" (ARSDT, MKASR, MDARJH, JARODH, SHJER) are straightforward to find the roots, because they contain three original letters, but the algorithm fails to analyze them. And in some cases, it returns non-acceptable roots. For example, with the word "التقى/ALTKA", the generated roots are "وقى, لقي"(WKA,LKE), where the root "وقى/WKE" is not an acceptable root, because "ل/LAM" letter cannot be an infix letter. Also, the algorithm matches the word to the wrong pattern, which is "فعاة/ FA"AH". In addition, the algorithm fails to find all possible roots of non-vocalized words, like the word "التقى/ ALTKA", where it doesn't return the possible root "تقى/TKA".

The main problems of the proposed algorithm are that it fails to extract the root of derivation words with one letter length. In Arabic language, there are some few derivation words with one letter length, like "ق, ر,ع"(KE, RE, A'E). These derivation words are derived from a weak root with a length of three letters and these weak letters are deleted during the derivation process.

Table 3. A sample of unanalyzed words by Al-Khalil Morphological System 2.

| سخريا | البضيع | دواعب | ارصدت | مقاصر | مضمار | مدارج | قضيم | عواطف | عفير |
|---|---|---|---|---|---|---|---|---|---|
| الجريم | الجبابرة | الجارودة | الشجير | التواجد | التماجيد | التقدمية | التساخين | البهيم | البنود |
| الرجيع | الرجوليه | الدبسه | الخليف | الخبيص | الحوارد | الحكومات | الحريصه | الجورب | الجفلى |
| الطوارف | الضراغمه | الصناديد | الشريم | السميد | السبله | الزهرة | الزغاليل | الرعاديد | الرحيق |
| زهيره | رصيد | دهريا | جحوشا | بالمناقيش | المغبرة | المراشد | اللطيم | العنزي | العريكة |

Another case in which the proposed algorithm still fails is to find the root of derivation words as in the word "درهم/ DERHAM". The algorithm produces these roots "درر, درأ,دري,دور, ودر"(WDR, DWR, DRE, DRA, DRR). In this derivation word's matter, the algorithm finds two constant letters in the derivation word and tries to find the third constant letter in order to produce trilateral roots. However, the proposed algorithm is stopped to continue looking for the fourth one.

The proposed algorithm and Al-Khalil Morphological System 2 produce more than one possible root of the derivation words. In contrast, Khoja and Garside's algorithm and the algorithm of Sonbol et al. produce just one root. In this section, the proposed algorithm and Al-Khalil Morphological System 2 are evaluated in terms of the average of possible roots per word and the number of processed words per second. The result is summarized in Table 4.

Table 4. Comparison between the proposed algorithm and Al-Khalil Morphological System 2.

| The algorithm | The average of possible roots per word | The number of processed words per second |
|---|---|---|
| The proposed algorithm | 3 | 101 |
| Al-Khalil Morphological System 2 | 5 | 105 |

## 5. FUTURE WORK

The presented algorithm particularly contributes to enhancing the algorithm of Sonbol et al. by increasing its rules and extending its lists' contents by using Thalji's lists. The presented Arabic root extraction algorithm is compared with Khoja and Garside's Arabic root extraction algorithm, Sonbol's Arabic root extraction algorithm and Al-Khalil Morphological System 2. The testing and comparing processes are conducted on Thalji's corpus, where the result of testing shows that the accuracy of Khoja and Garside's algorithm was 63%, whereas the accuracy of the algorithm of Sonbol et al. was 68%, the accuracy of Al-Khalil Morphological System 2 was 75%. The presented algorithm achieved an accuracy of 92%.

In future, we plan to enhance the accuracy of the presented algorithm, overcome some weakness points and enhance the result to return just the exact root word. In order to implement this, the system must have the ability to understand the whole sentence or sometimes the whole paragraph.

"Enhancing the Accuracy of Sonbol's Arabic Root Extraction Algorithm", N. Thalji, N. Hanin, Z. Thalji and S. Al-Hakeem.

# REFERENCES

[1]     W. Abo Thuaaib, History of Sematic Languages, Lebanon: Darul Kalam for Pub. and Printing, 2016.

[2]     A. Al-Taani and S. A. Al-Rub, "A Rule-based Approach for Tagging Non-vocalized Arabic Words," The International Arab Journal of Information Technology, vol. 6, no. 3, pp. 320-328, 2009.

[3]     R. Sonbol, N. Ghneim and M. S. Desouki, "Arabic Morphological Analysis : A New Approach," Information and Communication Technologies: From Theory to Applications, Proc. of the IEEE 3rd International Conference, pp. 1-6, 2008.

[4]     S. Khoja and R. Garside, "Stemming Arabic Text," Computing Department, Lancaster Univ., UK, 1999.

[5]     E. Al-Shawakfa, A. Al-Badarneh, S. Shatnawi, K. Al-Rabab'ah and B. Bani-Ismail, "A Comparison Study of Some Arabic Root Findings," Journal of the American Society for Information Science and Technology, vol. 61, no. 5, pp. 1015-1024, 2010.

[6]     N. Thalji, N. A. Hanin, Y. Yacob and S. Al-Hakeem, "Corpus for Test, Compare and Enhance Arabic Root Extraction Algorithms," International Journal of Advanced Computer Science and Applications, vol. 8, no. 5, pp. 229-236, 2017.

[7]     M. Sawalha and E. Atwell, "Comparative Evaluation of Arabic Language Morphological Analyzers and Stemmers," Proc. of COLING 22nd Inter. Conference on Comptational Linguistics, pp. 107-110, 2008.

[8]     R. Alshalabi, "Pattern-based Stemmer for Finding Arabic Roots," Information Technology Journal, pp. 38-43, 2005.

[9]     M. N. Al-Kabi and R. Al-Mustafa, "Arabic Root-based Stemmer," Proceedings of the International Arab Conference on Information Technology, 2006.

[10]    S. Ghwanmeh, S. Rabab'Ah, R. Al-Shalabi and G. Kanaan, "Enhanced Algorithm for Extracting the Root of Arabic Words," Proc. of the 6th International Conference on Computer Graphics, Imaging and Visualization, pp. 388-391, 2009.

[11]    Z. Kchaou and S. Kanoun, "Arabic Stemming with Two Dictionaries," IEEE International Conference on Innovations in Information Technology, pp. 688-691, 2008.

[12]    M. El-Defrawy, Y. El-Sonbaty and N. Belal, "A Rule-based Subject-correlated Arabic Stemmer," Arabian Journal for Science and Engineering, vol. 41, no. 8, pp. 2883-2891, 2016.

[13]    A. Ayedh and T. Guanzheng, "Building and Benchmarking Novel Arabic Stemmer for Document Classification," Journal of Computational and Theoretical Nanoscience, vol. 13, no. 3, pp. 1527-1535, 2016.

[14]    A. Pasha, M. Al-Badrashinyy, M. Diaby, A. El Kholy, R. Eskander, N. Habash, M. Pooleery, O. Rambow and R. Roth, "MADAMIRA: A Fast, Comprehensive Tool for Morphological Analysis and Disambiguation of Arabic," Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC'14), Japan, 2014.

[15]    N. Habash and O. Rambow, "Arabic Tokenization, Part-of-Speech Tagging and Morphological Disambiguation in One Fell Swoop," Proceedings of the 43rd Annual Meeting of Association for Computational Linguistics, pp. 573-580, Association for Computational Linguistics, Michigan, 2005.

[16]    N. Habash, O. Rambow and R. Roth, "MADA+ TOKAN: A Toolkit for Arabic Tokenization, Diacritization, Morphological Disambiguation, POS Tagging, Stemming and Lemmatization," Proc. of the 2nd International Conference on Arabic Language Resources and Tools (MEDAR), Egypt, 2009.

[17]    N. Habash, R. Roth, O. Rambow, R. Eskander and N. Tomeh, "Morphological Analysis and Disambiguation for Dialectal Arabic," Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, New Orleans, 2013.

[18]    M. Diab, K. Hacioglu and D. Jurafsky, "Automated Methods for Processing Arabic Text: from Tokenization to Base Phrase Chunking," Arabic Computational Morphology: Knowledge-based and Empirical Methods, Kluwer/Springer, 2007.

[19]    M. Boudchiche, A. Mazroui, M. Bebah, A. Lakhouaja and A. Boudlal, "Al-Khalil Morphological System 2: A Robust Arabic Morpho-syntactic Analyzer," Journal of King Saud University-Computer and Information Sciences, vol. 29, no. 2, pp. 141-146, 2017.

[20]    M. Ababneh, R. Al-Shalabi, G. Kanaan and A. Al-Nobani, "Building an Effective Rule-based Light Stemmer for Arabic Language to Improve Search Effectiveness," Int. Arab Jour. of IT vol. 9, no. 4, 2012.

169

Jordanian Journal of Computers and Information Technology (JJCIT), Vol. 04, No. 03, December 2018.

[21] K. Taghva, R. Elkhoury and J. Coombs, "Arabic Stemming without a Root Dictionary," Proc. of the IEEE International Conference on Information Technology: Coding and Computing, pp. 152-157, 2005.

[22] M. Sawalha and E. Atwel, "Corpus Linguistics Resources and Tools for Arabic Lexicography," Proceedings of the Workshop on Arabic Corpus Linguistics (UCREL), 2011.

[23] K. Mezher and O. Nazlia, "A Backpropagation Neural Network to Improve Arabic Stemming," Journal of Theoretical and Applied Information Technology , vol. 82, no. 3, pp. 385-394, 2015.

[24] G. Kanaan, R. Al-Shalabi and M. Sawalha, "Full Automatic Arabic Text Tagging System," Proceedings of the International Conference on Information Technology and Natural Sciences , pp. 258-267, 2003.

[25] E. Al-Shammari and J. Lin, "A Novel Arabic Lemmatization Algorithm," Proceedings of the 2nd Workshop on Analytics for Noisy Unstructured Text Data (ACM), pp. 113-118, 2008.

[26] H. M. Al-Serhan, R. Al Shalabi and G. Kannan, "New Approach for Extracting Arabic Roots," Proceedings of the Arab Conference on Information Technology, pp. 42-59, 2003.

[27] M. N. Al-Kabi, S. A. Kazakzeh, B. M. Abu Ata, S. A. Al-Rababah and I. M. Alsmadi, "A Novel Root-based Arabic Stemmer," Journal of King Saud University-Computer and Information Sciences, pp. 94-103, 2015.

[28] A.-K. N. Al-Kabi, "Towards Improving Khoja Rule-based Arabic Stemmer," Proc. of the IEEE Jordan Conference on Applied Electrical Engineering and Computing Technologies (AEECT), pp. 1-6, 2013.

[29] S. Al-Fedaghi and F. S. Al-Anzi, "A New Algorithm to Generate Arabic Root-pattern Forms," Proceedings of the 11th National Computer Conference and Exhibition, 1989.

[30] F. Abu Hawas and K. E. Emmert, "Rule-based Approach for Arabic Root Extraction: New Rules to Directly Extract Roots of Arabic Words," Journal of Computing and Information Technology, vol. 22, no. 1, pp. 57-68, 2014.

[31] K. Abainia, S. Ouamour and H. Sayoud, "A Novel Robust Arabic Light Stemmer," Journal of Experimental and Theoretical Artificial Intelligence, vol. 29, no. 3, pp. 557-573, 2017.

[32] B. Abuata and A. Al-Omari, "A Rule-based Stemmer for Arabic Gulf Dialect," Journal of King Saud University-Computer and Information Sciences, pp. 104-112, 2015.

[33] S. A. Yousif, V. Samawi, I. Elkabani and R. Zantout, "The Effect of Combining Different Semantic Relations on Arabic Text Classification," World of Computer Science and Information Technology Journal, pp. 112-118, 2015.

[34] G. Kanaan, R. Al-shalabi and M. Sawalha, "Improving Arabic Information Retrieval Systems Using Part of Speech Tagging," Information Technology Journal, pp. 32-37, 2005.

## APPENDIX

فعل ,الفعلة ,يفعل ,فعلا ,مفعول ,المفعول ,والفعلة ,الفعل ,والفوعل ,الفوعل ,الفوعلا ,والفوعلا ,فيعلا ,فعولا ,فعول ,فوعل ,وفوعلا ,الفعول ,فعلة ,وفعول
,وفاعلون ,فاعلة ,فواعل ,وفاعلات ,الفعيلة ,وفعلت ,فعلت ,الفاعل ,المفعل ,المفعل ,والفعال ,الفعال ,فعالة ,الفعلاة فعلون ,فعلاة ,فعيل ,فعال ,والتفعل ,تفعل ,مفعلين
,ومفعل ,مفعل ,الفعل ,والفعل ,والمفعل ,فعلته ,فعلت ,المفاعل ,فعالها ,فعال ,مفعولة ,والمفاعل ,مفاعل ,لفعلة ,وفعيلك ,يفاعلك ,يفاعلك ,وتفاعله ,بفعلها ,فعيلها
,والتفاعل ,التفعل ,والافتعال ,يفتعل ,فعيلية ,وفيعلة ,وفيعلة ,والفيعل ,فيعل ,وفعالا ,وفعل ,وتفعل ,الفعالية ,والفعالة ,ويفاعلها ,فاعل ,وللفاعل ,الفيعلة ,وفيعلت
,فاستفعلتها ,فعوال ,فعوالة ,ففعل ,بفعوالة ,والفواعل ,,يفعلون ,وتفعلت ,تتفعل ,الفيعل ,وفيعل ,الفعلان ,والفعلان ,فعلى ,فعالى ,والفاعل ,وفعلان
,وفاعلة ,وفعولت ,وفعالة ,فاعل ,الفواعل ,وفاعل ,فواعلا ,الفعلا ,الفعلا ,الفيعالة ,بفيعالة ,الفياعيل ,فيعال ,وفيعلتها ,فياعلة ,وفياعل ,بالفيعل ,فعيلا ,وتفاعلت
,بالفاعلين ,متفعلا ,والفعول ,متفعلا ,والفعيلة ,فعيلة ,الفيعله ,الفيعله ,تفعلوا ,تفعلت ,وفعيلت ,والانفعال ,بالفعل ,بالمفعول ,مفعلا ,انفعلت ,وانفعلت ,تفعلا ,وافتعل
,والفعيل ,فعلنا ,وافتعلت ,افتعالا ,وفعلة ,تفعله ,يفعله ,كالفعيل ,الفعيل ,والمتفعل ,والمفعول ,انفعل ,فياعل ,فعالا ,فعالا ,بفعالته ,الفعالة ,وفعلا ,فعلها ,فعليه ,فعيلته
,وفعاليل ,فعولهم ,ينفعل ,وانفعل ,فتفعل ,فانفعلت ,لانفعاله ,وفعال ,مفاعلا ,مفاعلا ,ففعله ,يفعلنه ,وفعيل ,والفاعلة ,والفوعلة ,وفعلاوات ,تفاعل ,وتفاعل ,المفعلة
,بالمفاعل ,ومتفعل ,وفعولته ,بالفعول ,والفعلولة ,الفعاليل ,والفعاليل ,ومفعال ,فانفعل ,منفعل ,ومنفعل ,فعله ,ومنفعل ,الفعلات ,فعلان ,فاستفعلته ,تفعيلا ,يفعلهن
,مفعلة ,ومفعلة ,والتفعيل ,التفعيل ,مفعولات ,والمنفعل ,استفعلته ,ليفعل ,مفعال ,لفاعل ,لفعل ,فعلي ,فعلي ,فيفعل ,بالفعيل ,وانفعال ,مفعله ,والمفعلات
,وفعلات ,افتعل ,والمفعلة ,وفعله ,واليفعيل ,فاعلته ,فعلات ,بفعل ,الفوعلة ,افعل ,يفوعل ,مفعلات ,يفتعله ,والفعولة ,يفتعلها ,وفعيلتك ,فعيله ,وفعيلا
,فعلن ,وفواعله ,وفواعل ,فعلل ,وفعللة ,وافعل ,والفعل ,والفعل ,فعلك ,وافتعلك ,فينفعل ,يفعلوا ,بالفعال ,والفاعلان ,كفعيل ,وفاعله ,الافتعال ,كفعل ,كالفعل
,فعالات ,والمفاعلة ,يفاعلها ,وفعيلته ,وتفعلهم ,ويفعل ,فعلاتهم ,لمفعل ,لمفعول ,فعوله ,والفعلوة ,وفعلوتان ,فعلوتين ,فعلوا ,الفعالي ,واتفعله ,ومفتعل ,مفتعل
,والفعلات ,وافعلتها ,الفعلى ,وتفاعلوا ,وفاعلته ,ففعلته ,وافعلت ,يفتعلوا ,وفاعلت ,يفاعله ,فعلوا ,والمفعال ,وفعلته ,ومفعلان ,ومفعلانة ,فعلتها ,تفتعل
,مفعولا ,ومفعول ,وفعالين ,الفاعول ,فاعول ,فعال ,ففاعل ,وفعلها ,ومفعال ,ومفعل ,تفعلني ,مفعاله ,وتفعيل ,والفعيلي ,والفعلى ,بالفعلة ,بالتفعيل ,ومفاعيل
,يتفعل ,وفعالته ,والمفتعل ,بالمفتعل ,والافعيلال ,وافتعلتها ,وافتعلت ,تتفعل ,والفعلي ,الفعلي ,بالفعلي ,فالفعول ,وفعلي ,وفعلي ,مفاعيل ,مفتتعل ,افتعلت ,المفاعلة ,افتعلوا
,تفاعلوا ,افتعل ,وتفعلتها ,والفيعول ,الفيعول ,وفعلا ,وفعولا ,المستفعل ,واستفعلت ,والفعيله ,ومفعولة ,ويفعلن ,متفعلة ,فتفتعل ,فعلتين ,فعاليل
,بمفعال ,وفعليه ,فعيلان ,ويفتعلان ,ويتفاعلان ,ففعلتم ,فعلانا ,وفعلى ,استفعلت ,واستفعلت ,يتفاعلان ,بفاعلة ,ففاعله ,ففاعله ,مفاعلة ,مفاعله ,يفعلان ,واليفعول ,يفاعيل
,ويفعول ,بفعله ,بيفاعيل ,اليفاعيل ,بالتفاعيل ,تفعيل ,وفعالية ,كالفعلة ,وفعلتهم ,التفعيلا ,مفاعلها ,المتفعلون ,تفعلهم ,مفعلان ,بفعالة
,وفعلهم ,فعالي ,مفتعلون ,المفتعل ,وتفعلكم ,فعالاه ,بفعالي ,وفوعل ,فافتعلوا ,الفعلانيون ,ويفعلون ,فعولة ,فعلناه ,فاعلين ,الفعولا ,فاعلتها ,وفعلني

"Enhancing the Accuracy of Sonbol's Arabic Root Extraction Algorithm", N. Thalji, N. Hanin, Z. Thalji and S. Al-Hakeem.

,بالفعلات ,والفعلول ,بفعلول ,وفعلك ,بمفعلين ,فعيلات ,فاعلات ,وفعلها ,وفعلنا ,فعلاه ,وفعلها ,مفتعلا ,مفتعلا ,المفاعيل ,يتفاعل ,المفاعل ,بتفعيل ,وتفعيل ,ومنفعله
,افتعال ,افتعله ,فعلاوات ,تفعلان ,والمفعولة ,والاستفعال ,الاستفعال ,فاعلون ,استفعل ,واستفعلته ,والافعالة ,ومفاعلته ,وافتعال ,والفعلات ,يفتعلون ,وافتعلي ,فالفعل
,مستفعل ,بفعول ,فعليات ,بالفعلا ,فعلتك ,كتفعل ,لفعالة ,وفعاله ,وفعاله ,لفعال ,كفعال ,ومستفعل ,فعاللة ,فعلني ,فعاول ,الفعاول ,الفعلانة ,والفعولة ,والفعولي
,وفعولة ,لتفعل ,وفعلكم ,فاعلهم ,وفعلون ,فعلين ,الفعلون ,للفعال ,تفعلم ,مفاعلة ,مفاعلوك ,لمفاعلة ,والفعالي ,وافعولوا ,الفعلات ,يفعلك ,فعلولة ,وفعلول
,فعلناها ,لمفعلة ,والفعله ,فعليل ,الفعلية ,فعلة ,والفعلن ,المفتعلة ,فافتعل ,يستفعل ,وافعوللت ,وافعولوا ,يتفعلن ,فعلهم ,فعلتان ,فانفعلا ,وفعللنا
,فيعول ,نفعل ,وفعالتا ,كفعلك ,وفعالتها ,فاعلت ,وفعلتهم ,بمفاعيل ,الفعلني ,والفعلنة ,فعلانه ,فعلية ,الفاعلة ,الفاعلين ,فاعلي ,فواعلها ,الفاعلتين
,المفعلا ,وفعلانها ,فعليها ,فعلانات ,لفعلانات ,وفعلتها ,افتعاله ,افتعاله ,بالمفعلة ,وفعلولة ,والفعول ,الفعاويل ,فعاويل ,وتفعيلي ,انفعال ,فتفعلوا ,الفعلين ,فيفتعل
,افتعالي ,وفعولها ,افعلت ,المفعلات ,يفعلنا ,افتعلته ,يفعللها ,الفعول ,وفعلولة ,والفعلول ,فعلولا ,افعولا ,يفعول ,وافعول ,افعولا ,وافعول ,وافعولا ,افعوال ,وافعول ,والمفعالين ,المفعالين ,الفعالين ,فالمفاعل
,لفعول ,الفعلوة ,وفعيله ,فعالتي ,فعيلتي ,مفعول ,مفعلوني ,الفعلية ,وافعول ,يفعول ,افعول ,افعول ,فافعلوها ,فافعلوها ,الفعلولة ,كالمفعيل ,وفوعلة ,المتفعل ,للفاعل
,فافتعله ,والافتعالة ,فعولات ,وفعاول ,فعيلون ,بفاعل ,كالفعلان ,والفعلان ,والفعالي ,الفعلى ,للفعلى ,والفعلانة ,مفعوله ,للفعلى ,فوعلانية ,فيعلاني ,فيعلان ,تفعيلهما ,للفعول
,وافتعلوا ,وفعيلي ,والفعلاوة ,والفعاول ,الفاعول ,لفعلان ,يفعولة ,كفعلة ,والفعلاة ,كالفعلاة ,الفعيلة ,وافعيلال ,بمفعول ,وافعيلال ,تفعيلهما ,ومفعلات ,بالفواعل ,وفعوله
,يفعلها ,افعلا ,لنفعلا ,اليفعول ,بالمفعل ,وافعال ,افعيلالا ,ومفتعلا ,وفاعلي ,المفعال ,وافتعله ,وفاعلتك ,وفعلوه ,كالمفعل ,للفعل ,والمفعالة ,فعولته
,مفعالة ,والفاعلية ,فوعلة ,بفعلهم ,فعلتني ,تفعليني ,فعيلهم ,فتعلل ,الفعيول ,فعاييل ,فعيولون ,مفعيل ,المفعيل ,افعوال ,والفعلان ,الفعلين ,مفعلون
,تفعلونني ,ليفعلوا ,لبفاعلوا ,الافعال ,يفعلهما ,التفاعل ,افعلوا ,ويفتعلها ,فيفعلها ,وفعله ,وفعله ,افعله ,وافعله ,والمفعلتان ,افعالا ,وافتعالا ,فعليلا ,والفعليلة
,يتفعل ,تتفعل ,كفعليل ,الفعليل ,لافعل ,والفعلوية ,فعولاة ,وفعولية ,الفعولية ,افعول ,افعولا ,وافعول ,وفعالى ,بفعال ,فاعلها ,والفعلى ,ففعلناها
,وفعلين ,بالفعلوية ,تفعيلة ,بالفعالة ,تفعلك ,بفعلة ,تفعلته ,فعلول ,فالفعيل ,افتعلهما ,ففعلنا ,لمتفعل ,وتفعلته ,وتفعلته ,وفعلات ,وفعلاني ,وفعلانية ,فعلاوه ,فتفعله
,يفعلني ,والفعلاتان ,لفعلي ,والتفعل ,تفعلل ,بالفعيلي ,مفاعلي ,وفعلانا ,المفعولة ,وفعيلات ,لمتفعل ,تفعلي ,ففعلها ,لبتفاعل ,بمفاعل ,فالفعال ,وتفعله
,والفعلوت ,لتستفعل ,يتفعلون ,وفاعلها ,الفعلوان ,فعلوان ,والفعلوانة ,الفعلوانة ,الفعلوانة ,يتفاعلني ,فعيلي ,فعيلك ,المفعلون ,ومفعله ,والفعل ,ويفعلها ,معفالا
,فيعلت ,كالفعال ,والفعلون ,فعلوة ,كفعلان ,والفعلوة ,وتفعلت ,ففعلت ,وتفعلت ,بالفعلى ,بالفعلين ,فياعيل ,فيعالة ,بفعلات ,مفتعلات ,ففعلا ,وفعالية ,وفعلية ,وفعليت ,والفعاليت
,الفعليون ,فعلني ,فعلناة ,وفاعولة ,وففعلت ,وتفعلنا ,فعلاني ,فعلاني ,فعليته ,المستفعلة ,فاستفعلوا ,وتفعلوا ,فعليا ,والفعلانية ,والفاعلي ,فعلله ,والتفعالة ,تفعالة
,بمفعول ,الفعاليات ,فعاليته ,فعلو ,لفعلك ,يفعلنوها ,ليفعلوها ,وفعلوها ,لتفعلها ,واليفعل ,واليفعل ,وافتعال ,وافتعال ,والفاعولة ,والمفعا ,فيعله ,واليفعلة ,يفعلات ,واليفعلات ,بالياعل ,واليفعلات ,وافيعيل ,وافيعيل
,افيعللي ,افتعلناهم ,وفعلوان ,للفعلان ,وفعلا ,يفعولا ,يفاعلى ,يفاعلات ,لفعلا ,الفعالي ,وفعلتان ,فعلانتن ,فعلاتهم ,فتعلاة ,فعللا ,فعلا ,الفعال ,الفعلة ,فعلة
,فعلية ,والفعلة ,والفعال ,تفعللا ,الفعاللة ,مفعلل ,والمفعل ,مفعلل ,فعللت ,وفعللت ,الفعيل ,وفعلة ,المفعلل ,الفعلال ,وفعللته ,فعاللها ,فعللوا ,افعللا ,افعلالا
,ومفعل ,وافعل ,فعللا ,فعلوا ,وفعلا ,الفعلالا ,والفعاللة ,المتفعل ,فعلني ,فعلليون ,الفعليون ,الفعيلول ,والفيعلول ,فيعلول ,والفيعلول ,الفعليلة ,وفعللون ,وفعللت
,والفعليلية ,مفعلا ,فعيللانة ,فعيللان ,وفعلني ,وفعالة ,وفعالا ,فعيللي ,والفعلني ,والفعللي ,الفعللان ,وفعلان ,وفعللا ,بالفعلل ,افعللالا ,يفعل ,فافعل ,ويفعل
,والفعللان ,فعلولها ,تفعللت ,فعالله ,الفعللان ,الفعالي ,فتعلللت ,كفعللات ,فتعللت ,افعللت ,كالفعلل ,فعللا ,وفعللا ,فعللها ,وفعللها ,الفعللا ,التفعلل ,الفعللا ,متفعللا ,فعللتها
,الفعللوت ,الفعللوه ,والفعلله ,وفعال ,فعللون ,الفعلول ,الفعلول ,الفعللة ,وفعللة ,بفعللا ,وفعلال ,والفعللى ,والفعللى ,الفعللية ,الفعلالة ,فعلالة ,فعللا ,وتفعلل ,الفعللية ,فعلليات
,لمفعلل ,مفعلله ,الفعللل ,فعللل ,الفعللى ,فعللي ,وفعللى ,فعللا ,بفعللا ,وفعلالة ,فعلليل ,فعلليل ,الفعلليل ,والفعلليل ,فعلليه ,فعلليله ,فعلوله ,فعللها ,للفعللل ,والفعللل ,للفعللل
,والفعلللة ,فعلللون ,وفعلللات ,افعللل ,الفعللل ,الفعلللة ,فعلللات ,فعللللى ,فعللللة ,والفعللللانة ,الفعللللى ,فعللللت ,والافعال ,فعيلتها ,وبالفعال ,والفاعله ,فعلانون
,بفعوال ,والفعلتان ,الفعلتين ,فوعلت ,فواعيل ,وفعلوها ,وفعلوها ,المفعولين ,الفعيلا ,فتفعنل ,فتفعنل ,تفتعله ,وفاعلني ,والفعليك ,والفعاليك ,فعلكة ,فاعله
,وفاعلناه ,فعيلي ,الفعلان ,بفاعول ,مفعلها ,للفعلة ,وفعلها ,وفعلنا ,وفعلتا ,فعلتيهم ,فوعلتا ,ويفتعلوا ,مفاعله ,بالمفعال ,وفاعلا ,وفيعلا ,وفيعلا ,بفعال ,وافعلته ,وافعلته ,افتعلنا
,وتفاعلنا ,وتفتعل ,فعلتي ,كالفعلال ,للفاعلة ,الفاعلات ,تفعلي ,والفعلية ,بمستفعل ,يفعلهم ,بمستفعل ,بمفعله ,والمفعلي ,مفعلي ,المفعلي ,افتعالك ,فعلاها
,تفعلونا ,التفعال ,فعلو هم ,للمفتعل ,بفعالها ,تفعلن ,لفيعل ,فعلوه ,وفاعلوه ,بفعيل ,فالفعلة ,تفعليا ,فوعلته ,وافتعليل ,المفعولا ,وتفعلتهم ,فاعلك ,فعالتك
,لمفعال ,استفعله ,وفعلناهم ,والمفاعيل ,وافعالت ,الفاعله ,فعلتنا ,وفعلناهم ,وفعلناهم ,فيفعلكم ,ونفعل ,المفعلين ,منفعلا ,كتفعيلك ,منفعلة ,والمفعلان ,الافعيل
,وفيعلين ,وفيعلون ,والمستفعل ,التفاعيل ,الفياعل ,والفعيلان ,الفعيلان ,بفاعله ,بفعلان ,للمفعول ,فالمفعل ,والمفعل ,الفعالات ,وكفعل ,ففعلان ,كالفعالة ,وفيعلان
,والتفعال ,والفعايل ,فعالت ,والمفعيل ,وافتعل ,كالتفاعل ,والفيعلون ,تفاعلا ,كالمفاعلة ,لافعلين ,كفاعل ,مفوعلا ,وافعلى ,يتفاعلون ,وفعالاك ,الفعليت
,والفعاله ,تفاعلت ,متفاعل ,والمتفاعل ,تستفعلوا ,افعالت ,الفعلتان ,وفعلناها ,وفعلناها ,فعلان ,بافعلين ,بفعلان ,لبفتعل ,وفعلنا ,وافعيل ,بفعلين ,يفاعلون ,والتفاعيل ,لاستفعالنه ,مستفعلات ,الفعلن ,يفعال
,فالمفعولة ,والفعلون ,الاستفعال ,وتفعلي ,تفعلة ,بالفعالين ,بفعلة ,لبفتعل ,والفتعال ,الفعالية ,وتفعلة ,الفعيلية ,فعلالا ,وفعلله ,والفعللوة
,وفعللوات ,الفعللول ,لبتفعلل ,الفيعلال ,فعلته ,بالفعاليل ,وفعلالة ,والفعالة ,الفعلالة ,فتفعلل ,افعللي ,والمفعللي ,والفعيلى ,والفعيلي ,افعيلت ,افعيالا ,مفعيلة
,يتفعول ,تفعوا لا ,متفيعل ,يتفيعل ,الفعولات ,وفيعلا ,بفاعلات ,بفاعلات ,افوعل ,متفعله ,فاعلناه ,مستفعلة ,والتفعلة ,مفتعلة ,مفتعلة ,بالفاعلة ,وفاعلتهم ,تفعلون ,فاعلين
,فعالهم ,مستفعلا ,مفعولي ,بالمفاعلين ,المفاعلين ,يفعلونني ,الفعيلي ,وفعيلي ,الفعيلي ,بالفعلتين ,وفيعيلي ,والفيعلانة ,والفيعلانة ,بفيعلا ,ومفعوله ,فواعله ,فافتعلت ,وافتعاله
,والفعلولية ,وفعلايا ,فعالي ,بفياعل ,بفيعل ,فيعيل ,بمفعلات ,وفعالي ,المفاعلا ,المفاعل ,وفعالي ,وفعلا ,فعلوتا ,وتفيعن ,تفيعلا ,فاعيل ,وفعلى ,فعالي ,ففعلهم ,ففعلهم ,افتعلها ,افتعلها ,فتنفعل ,والفعوا لا ,فافعلوا
,يتفعلوا ,والفياعله ,الفياعله ,الفعلو ,وفعولي ,وفعولي ,والفوعلات ,فيفعلون ,فعلاوي ,وفعلاوية ,وفعلونا ,افعوال ,وفعلونا ,بتفعول ,فعلانيهم ,وفعليهم ,وفعلانيه ,وفعلانيه ,وفعلانيه ,والفعلوي
,ليفعله ,ومفتعله ,لفاعلة ,وبنوفعال ,والفعلانة ,وفعلوت ,فعيللة ,وافعلناه ,ووافعلتياه ,وافعلتياه ,فاعتلاله ,فافتعلتها ,بالتفاعل ,ونستفعل ,تفاعله ,تفاعله ,الفعولي
,ففيعل ,مفعولان ,ولفعلة ,الافعلال ,وبفعلل ,لتتفعلل ,المفعللا ,وفاعل ,فاعلا ,الفعالل ,الفعللان ,يفعللان ,فعللانة ,الفعلله ,منفعلي ,بالمفاعيل ,لفعيل ,والفعلني
,فاعل ,الانفعال ,بمفعلة ,والفاعلات ,وفيعيلن ,وفيعلول ,وافتعلوه ,استفعلوه ,وتفعلمي ,فعيلاه ,وتفعلها ,فعيلاه ,وفياعيل ,وفياعيل ,لفعلانه ,والافعال ,والافعيل
,ولافعيل ,وتفتعلون ,نتفعل ,وفاعلنا ,الفعلون ,كالفعول ,فاعلتي ,والمفعلون ,يفعلوه ,والفعليت ,فعلت ,والفعليت ,وفعلى ,وفعلانه ,والفاعلون ,لفعله ,وفعلناه ,والفعلني
,والفيعلان ,فيعلانة ,ومفعو لا ,وفعالتي ,وفاعلي ,وفعالتك ,متفاعلا ,الفاعل ,وافعلتك ,فافافعل ,افتعاله ,لفعيلى ,وافعولت ,وفعلون ,فعلان ,وفعلان ,ومفتعلات ,لفعلن ,وفعالتها
,وبفعل ,ومفعلون ,مفاعلك ,الفعيلات ,وفعيل ,بفعيل ,مفعولون ,فاعلنا ,فاعلتنا ,بالمفعالين ,وافعله ,وافعله ,ومسفوعل ,ومفوعل ,وفعالته ,للمفعل ,فعلاتها ,بالفعالات
,واستفعلناها ,فعالتها ,فويعلان ,لفعلوا ,كمفعل ,استفعلوا ,الفواعيل ,يفعول ,وفعلها ,ومفعلها ,المفتعلة ,وفعلهم ,افتعال ,تفعلتها ,مفاعلته ,فعالك ,فعلتكم ,تستفعل
,ومفاعليها ,مستفعلين ,الفاعلتان ,فعلوت ,وفعلتمو ,الفعلوت ,فعلان ,وفعلتها ,وفعلها ,وافتعلها ,فيستفعله ,يستفعلوا ,وافعلتها ,وفاعلتنا ,تفعيلك ,والفوعلل ,الفوعلل ,المفعللة
,الفعللي ,الفعللات ,وفعلليل ,ومفعلة ,المفعلتل ,فعلوت ,فعلتلا ,افعللوا ,افعلتلا ,فعليلا ,فعليلا ,ويتفعلل ,فعلاوات ,وفعلاوات ,والافعلال ,متفعلل ,وفعللة ,وفعلية ,كفعلول
,لفعالة ,فعلولا ,فاعلني ,الفعليل ,فسيفعلون ,وافعوال ,فيعللتهم ,فياعلا ,افعلال ,افعلوا ,وافعلوا ,والفعاول ,والفعاول ,والفعلوا ,بالتفعل ,والفعلاوان ,للتفعيل ,وبالمفعلة ,يتمفعلون
,مفيعلان ,ومفعالين ,والمفعله ,فتعللها ,والفعلين ,والفعلين ,افعيل ,والفيعلي ,والفيعلي ,بالتفعل ,فاعلول ,فعلول ,فاعلول ,ففعلا ,ففعلا ,فعللها ,فعللها ,فعلها ,وتفعلت ,وتفعلت ,بفعلة ,يفعللوا ,بالفعلال ,فليفعل
,فاعوله ,كالفعلات ,لفعيلة ,تفعلهم ,بفاعوله ,والفوعلي ,وافعولي ,افتعلتك ,فتفعلها ,وفعلوها ,افعلوها ,وفعالات ,وفعالات ,وفعال ,كتفعيل ,للفعيل ,وافعالي ,وانفعالي ,ليفعلني ,وفعلته ,وفعله
,والفوعليل ,لفعلنا ,وفعلانه ,وفعلانه ,مفعلته ,تفاعيل ,والفويعلة ,والفويعلة ,تفعلها ,وتفتعلها ,وبفعلها ,كالمفعلة ,نفتعله ,فعولتك ,ويتفعلونه ,واعتلال ,افاعتلال ,لفعالات
,ومفعيل ,فاعلو ها ,تفوعله ,وفعلنيه ,والفعلاو ,والافعلان ,والفعلان ,وفعللة ,مفاعيله ,مفاعيله ,فاعلان ,وافعلان ,الفعلاني ,والفعلاني ,فعلول ,سيفعلني ,افعلان ,فاعولة ,فعلاته ,واستفعله ,المتفاعل
,مفتعلين ,متفاعلين ,فعيلتنا ,وفاعول ,والفيعوال ,وفعلو ,ويتفعل ,فاعو لا ,بالفاعل ,ففعلوا ,الفاعل ,وفعلتم ,ويفعلونه ,فعلهما ,فعلهما ,يستفعلون ,ومفاعلاتي ,ومفاعلة
,استفعال ,مستفعلك ,والفاعلى ,بفعلتين ,بالفعلين ,وبالفعل ,للمفاعل ,فعلانيا ,مفعلات ,فعلانيا ,فعيلا ,كالفعلني ,فعلاتك ,فعلانيا ,فعلاتك ,وبفعلون ,متفاعلون ,فعلهما ,يستفعلون ,تفعللها ,الفعاله ,والفعلالية
,كفعلالية ,والفعلاني ,الفعلاني ,الفعلانية ,فعللين ,فعلاليل ,فعلانيا ,مفعلات ,فعلانيا ,فعيلا ,كالفعلني ,فعلاتك ,فعلانيا ,فعلاتك ,ومفعلاتك ,وافعلناها ,وافتعلنا ,تفاعلنا
,تفعيله ,نفعله ,والتفعلى ,المفعلى ,المفعل ,المنفعل ,الفعاييل ,الفعيلك ,والفعلكة ,فعلانها ,فعلانها ,تمفعل ,فعوله ,مفعيلا ,وفعلياوان ,والمفعولات ,فعليان ,فعليون ,فعليان ,فعليون

171

Jordanian Journal of Computers and Information Technology (JJCIT), Vol. 04, No. 03, December 2018.

فعولكم ,والفعيول ,فعيول ,لفعول ,فعيول ,وانفعلوا ,وافعلها ,فانفعلوا ,لفيعلن ,ويفعلني ,فمفعل ,فاعلناهم ,ففعلناهم ,الفاعلية ,فالتفعيل ,وفاعوله ,للتفعل ,بفعلتيه

,بفاعلو ,فعولى ,الفاعلا ,بالفعولة ,وتفعلوه ,بمفتعل ,فعاليلي ,فعليله ,فعلليلة ,وتفعلني ,الفعولا ,والفوعال ,فوعالة ,والفوعلان ,افعولة ,الفياعلة

,يفعلانه ,والفعلو ,للفاعول ,فاتفعل ,لفعولهم ,ومفعولها ,واستفعلنهم ,الفيعال ,وفيعالة ,والفعاليان ,افعيالها ,الفعله ,فويعل ,وفويعل ,فافتعلها ,كفعول

,وافيعللول ,يفعلول ,فعلاله ,بالافعالة ,والافعالتين ,الافعالة ,لافتعالك ,كفعلتي ,يفتعل ,بفتعل ,فعالتيه ,كفعالة ,فعللله ,اليفعل ,كالفيعال ,المفيعل ,الفعالة ,الفيعلان

,فاعلوهم ,فعلنه ,يتفعله ,بفعلي ,وتفاعيل ,لمفعولة ,باقتعالك ,لافتعالها ,وفعتل ,لفعلة ,فعالتكم ,فالعلة ,يبفعل ,فويعلون ,وفويعلات ,الفويعلة ,فافتعلنا ,ومتفعلة

,واستفعالي ,فوعالا ,وفعلللة ,متفاعلان ,فعلاتة ,الفعلانية ,وفعلونة ,الفيعلي ,لتفعيله ,وفيعليات ,متفعلاتها ,يفتعلان ,الافعال ,الافعيلال ,بالفعيلة

,يتفعلونها ,فيعلي ,والفوعيل ,والفعولل ,فتعل ,فعللولة ,فعلنا ,تفيعل ,يفيعل ,فعلانيه ,كمفعال ,والفاعل ,والفاعل ,بفوعل ,لتفاعل ,الافتعال ,فعايل ,والفيعال ,الفواعيل

,فعولك ,يفاعلونها ,الفعوايل ,فعوايلات ,فتفعول ,ففعلتكم ,مفاعلون ,للمفاعل ,التفعلة ,الفاعلي ,بفوعل ,لتفاعل ,الافتعال ,وتفاعل ,فعايل ,والفيعال ,وافتعللته

,للمفعلات ,فعالاة ,فاع ,والفاع ,والفاع ,الفعلول ,الفعلول ,فعول ,الفعلولية ,فافعللت ,بفعللين ,ففعلله ,وافعلليته ,وافعلليني ,فاعليل ,ويفعل ,وافتعلها ,افعلن

,فعلتا ,كالفعللي ,يفعلونه ,والفاعولي ,نفعلها ,وتفاعلته ,تفاعلها ,ويفتعلك ,فواعلهن ,وتفاعللوا ,فعولوا ,ولفعيل ,وفاعلوا ,والتفعلي ,بالمفعيل ,وتمفعلت ,وتفعليك ,كالفعالة ,والفاعل

,فعلوات ,فيعالا ,بالفعالى ,وافعالات ,فعتل ,وفعلاه ,وفعت ,فالتفعل ,بفاعلين ,الفاعل ,افعلني ,افعلنا ,بافتعال ,الفاعل ,ففعلة ,فعلتم ,مفعلتك ,فعالان ,افعلي ,وافعلي

,وافاعتلال ,واستفعلوا ,فيفتعلون ,بفعلك ,فعتل ,وفعلاه ,وفعت ,فالتفعل ,بفاعلين ,الفاعل ,افعلني ,افعلنا ,بافتعال ,الفاعل ,ففعلة ,فعلتم ,مفعلتك ,فعالان ,فاعيل ,فعالان

,افعلتها ,والفعلوان ,فعلويها ,افعلالها ,والتفعول ,تفعولة ,لبفعلوهم ,مفعلاني ,المفعلاني ,كالمفعلاني ,ويتفعل ,والفواعلة ,وافعلاني ,وافعلاة ,الفاعة

,الفاعات ,فاعة ,وفاعة ,فعوا ,فيعوا ,الفيع ,فيع ,والفياع ,بالفياع ,تفيعت ,والفياعى ,فاعت ,تفيع ,فيعنا ,وفيعة ,وتيفيعت ,التفييعة ,والفاعة ,فاعه ,بفاعتها

,والفواع ,فوعه ,تفويعا ,فالفعيلة ,فيعلى ,ومتفعل ,كفعالي ,كالمستفعل ,الفاعولة ,واستفعلني ,وفاعلتن ,والفاعلاتين ,فاعولتها ,فاعللة ,الفيعلاني ,ويفعللها

,بالمفاعلة ,فعيلين ,ليفعلوك ,استفعلوا ,فالاستفعال ,فاعلا ,وافعاليلهن ,الفعوال ,وفعاولة ,ومفاعلهم ,الفعلولي ,فالفعلي ,لتفتعله ,لمفتعل ,وفعيلها

,لفعوله ,بفعلانه ,وبفعالته ,فعيلو ,تفعيلكم ,مفيعيل ,وللفعيل ,لفعيله ,افعلوه ,افعلوه ,لتفعلكم ,بالمفعلات ,فعلوت ,فعليك ,بفعليه ,للفعالة ,والمفعولي ,فلتفتعل

,فاستفعل ,وفعلانهم ,فعلانهم ,بالفعلان ,فالفعولة ,للفعليين ,فعيل ,الفيعلانات ,فعلواني ,فعلان ,المتفعلين ,فواعلهم ,فاعلتكم ,لفعلتكم ,الفعلاوات ,الفاعلان

,بمفعله ,فاعلنيها ,مفيعلة ,بالفيعال ,مفعلتهم ,فافعلهما ,فعيلاات ,الفوعالي ,الفوعال ,كالفعاليل ,وفواعيل ,والفعاليات ,مفعلتان ,مفاعيلهم

,بفاعال ,المفيعيل ,المتفاعلون ,مفاعلهم ,افعلوني ,والفعلوتي ,والفعلوى ,لفليفعلها ,وفاعلتان ,وتفعيلات ,مفعالي ,الفعلولين ,بالفعلول ,بالفعللة

,فعلانك ,للمفعلة ,وفعلبك ,المفاعلون ,الفعلاتوني ,الفيعلون ,فيعلون ,الفعلاني ,فعلتكم ,ويفعلوا ,تتفاعل ,كالمتفعل ,الفاعولة ,فعيلتين ,المفعلية ,لتفعلنها

,بفعلاهم ,المفاعلات ,بفعولي ,نفعلك ,الفوعيل ,فعلالو ,فعلانك ,الفعلاتين ,الفعلني ,فيعلانا ,كالفيعل ,لتفعلني ,فتفعلي ,مفعلتي ,وافعال ,بالفعال ,بالفعلية ,للفعلية ,لتفعلي

,مفعلهم ,وينفعل ,ويفعلك ,فتفاعلت ,لنفعلنك ,استفعلكم ,الفاعلون ,استفعلنك ,الفيعل ,لتفعلني ,فتفعلي ,مفعلتي ,وافاعل ,بفعلال ,بالفعلية ,للفعلية ,وافتعلوهم

,وافتعلهم ,وفعولى ,فعيلتك ,بالافتعال ,افتعالها ,وانفعال ,ومفاعلي ,مفيع ,واليفاعيل ,اليفعولي ,فعلاوي ,والمفعولية ,ولفعل ,تفعلوها ,وفعولية ,وتفاعلها

,الفاعلون ,الفعليلية ,ويتفاعل ,كمفعول ,والتفعلي ,وفعلهما ,للفاعلين ,المنفعلة ,الفعلال ,وفيعول ,فتعلل ,فاعتل ,فاعتل ,بتفعليها ,فعولي ,وفاعلين ,استفعلنا

,يتفعلها ,فليتفعل ,الفاعولات ,وافعليه ,المتفعلة ,ومفعلين ,وفعلوكم ,الفعوليون ,تفعيلية ,الفعيال ,لفعلهن ,الفعلايلات ,تنفعلان ,وفعلنه ,فعلنت

,الفاعلة ,فتفعلوهم ,فيتعلل ,وافتعليلة ,افتعليلة ,بافاعتلال ,لفعولته ,فنفعله ,بفعلته ,للاستفعال ,الفاعيلاات ,بتفاعيل ,الفاعيلان ,الفعيلاوان ,تفعلوا ,فافعلوه

,ومفعلاني ,فعلوك ,نفاعلهم ,وينفعل ,بمفعيل ,لتمفعل ,كفعلاني ,وفعيلاه ,مفاعليه ,مفاعليه ,الفعليان ,بالفعليان ,فعليانية ,فافعلها ,فليفعلوا ,ويفعلهم ,وبفعيلهم

,بمفعاله ,فعلانهما ,فتفاعلا ,فالتفاعل ,فافعلن ,بالفاعول ,والمتفعلة ,ومفعلوه ,فعلاه ,مفعلكم ,مفعاله ,يفعيل ,ولتفعل ,ستفعل ,وفعلن ,لفعلين ,بفعلتها ,التفعيلية ,وفعولتهن

,والفاعولات ,والفعيلل ,ففعلن ,وفعلتك ,فواعلة ,والفاوعل ,والفعلانة ,الفعاليين ,لافتعال ,بمفتعلها ,بفواعل ,للمفعلين ,وفعاليات ,والفعيلية ,والفعلويه ,وفعلويه ,ويفاعل

,لفعلها ,بفعوله ,وبفعلك ,فاعولي ,تستفعلها ,فاعلتهم ,فالفواعل ,افتعلهم ,فيعوله ,الفعلان ,وفعلانية ,بالفعلان ,كالفاعل ,المتفاعلان ,الفاعل ,ففعلتهم ,والمفيعل ,والمفعل ,الفعلنة ,وافتعللته ,والنفعل

,واستفعلها ,وبالمفعال ,وبالمفتعل ,تفاعلهم ,وافتعلني ,وفعلانية ,بالفعلان ,كالفاعل ,المتفاعلان ,الفاعل ,ففعلتهم ,والمفيعل ,والمفعل ,الفعلنة

,فعلنته ,وافتعلله ,فوعلات ,مفتعلان ,فالمفاعلة ,الفاعلون ,وفعالين ,فعلوي ,فعلان ,فعلاتهن ,والفيعلاني ,وافتعلا ,مفعلان ,والفعليات ,ففعالها ,امفعل

,فتفعلته ,يتفعللون ,فعليهم ,والفعلانين ,ومفتعلة ,فعيولة ,وافتعالات ,فعاليان ,فعاليان ,الفعاليين ,الفعللين ,وفعلوك ,تفعلني ,وفعلوهن ,ويفعلونها ,تفعلونهم ,مفتعلنا ,مفتعله

,تفاعلته ,مفعلها ,بمفاعله ,فاعلتللها ,فتعللت ,فعالنة ,فاعتلال ,تفعلنا ,بافتعل ,وبفعله ,وبمفعالها ,لتفعلتهم ,فعليلة ,فعيللتان ,افعلنيا ,وتفاعلا ,ومفعلي ,وفعلنتها

,اليفعيل ,بفعولها ,وفعليات ,للمفاعلة ,لفاعلون ,وفعلوي ,الفعتلان ,وفعلاتهم ,مفعليها ,وافعلو هن ,بالفاعلات ,وفعلنا ,ولفعلنا ,فوعلي ,وافعيل ,فعلللتان ,استفعلتم ,كفعله ,فعوليا ,فاعلتانات

,مفيعيلة ,مفيعلات ,تمفعلت ,والفوعول ,الفعليبانة ,لفعالته ,فعالته ,مفعليها ,وفعليها ,وافعلو هن ,بالفاعلات ,وفعلنا ,ولفعلنا ,فوعلي ,وافعيل ,فعلللتان ,استفعلتم ,كفعله ,فعوليا ,فاعلتانات

,والفعلتان ,وفاعلهم ,وفاعلان ,وفاعلان ,فيعلنه ,وبفعلته ,انفعلوا ,والافعولة ,لتفعيلي ,وتفعيلي ,فالفعلياني ,انفعاله ,يفتعلهما ,ويفوعل ,مفعلاة ,بفعلاي ,وفعيلتهم ,فعلوها ,وفاعولى

,لفوعلا ,وافتعلاة ,التفعول ,وافتعلوها ,ليتفعل ,بفعلى ,وللفعل ,والفعلوك ,والفعلاية ,والفعلاية ,فوعلية ,وتفعلتني ,فعلوكة ,وفعلوكة ,تفعاليل ,والافعل

,والافيعل ,والفعولى ,وسناوفعلت ,ومفتعلون ,بالافعال ,وفعليته ,لمفتعلون ,فعولنة ,والفعلاوات ,وتفعل ,وفعيلون ,لفعليها ,واستفعلهم ,استفعلهم ,ونفعلا

,الفعلانات ,يتفعللن ,فعاعلا ,ويفتعلون ,فعاعلا ,مفتعلتان ,والفعيلية ,وفعلوة ,فالمفعول ,والفعلانية ,والفوعلانية ,والفيعلانية ,واليفتعول ,الفعليه

,افتعلت ,بفعلة ,افعولا ,والفاياعيل ,الفعيلية ,بفعلته ,لفعلات ,وتفعيله ,والفعيولة ,والفعيولة ,والفعاويل ,والفعاويل ,والفعيولون ,والفعيولون ,والافعوال ,وافعولت ,وفعولنه ,وافعولته

,وتفعولنه ,الافعوال ,وافاعلل ,فالفاعلة ,كالفعالي ,فعالكم ,لفعلول ,المفتعلات ,الفعلاوة ,وفعلان ,فعليكما ,الفعولين ,والفتعلل ,والفتعلل ,مفعلت ,مفعلت ,فيفاعلك

,لتفعلن ,متفاعلات ,والمتفعلات ,فعلوانة ,والمتفعلات ,فانفتعلت ,بفعلي ,مفعالين ,وفعليلة ,الفعتل ,وافتعلتهم ,وتفعلتم ,وتفعالة ,وتفعيلة ,والمفعلاة ,الفعلنى ,ومفعالا ,الفعاليتان

,وفعالتهم ,الفعيلللين ,فعلوها ,فوعلوا ,والفويعلان ,وفعلو ,اليفعل ,وفتعل ,فعاليون ,فعلوت ,وفعوالا ,متفعلتان ,وتفعالة ,وتفعيلة ,والمفعلاة ,الفعلنى ,ومفعالا ,الفعاليتان

,نفاعلكم ,وافيعلي ,وافيعال ,وافعلوانه ,وفعلوانه ,فعلوانا ,فعالاك ,الفتعلتان ,الفتعلتان ,وفتعللة ,وفعيال ,والفعاليتان ,كالمفعول ,والمفوعل ,والمفوعل ,بفعلتهم ,الافعال ,ومفعول

,فعاليكما ,كافتعال ,والفاعيل ,فعييل ,لفعالا ,ومفعيلة ,الفعليات ,الفعولوهو ,والمفعول ,والمفعول ,وفعلان ,وفوعلان ,وليفعل ,فعاليت ,وبفعلة ,ومفاعله ,ومستفعلة ,والتفيعل

,فعلانكو هو ,لتفعلهم ,لمفتعلا ,والفوعلى ,وفيااعيله ,بالفعيلة ,فتعاللة ,فتعالل ,وتفعيلا ,وتفعيلا ,فيعتلل ,وافعتل ,والفتل ,فعلولنا ,فاعلولة ,والفعلللتان ,فاعولي ,والمفعلن ,فتعلله ,مفعوللا

,فاعلاوات ,والفاعيلة ,فعاليك ,بالفعيلية ,فعتعاللة ,فتعاللة ,فتعالل ,وتفعيلا ,فيعتلل ,وافعتل ,والفتل ,فعلولنا ,فاعلولة ,والفعلللتان ,فاعولي ,والفعلللتان ,وفعيبلوا ,فليفعله ,المتفيعلون ,الفيعلانيات

,فعتله ,فعلاي ,وفعيلاي ,تفعلهما ,وفعولوا ,وفعلاي ,تفعلهما ,وتفعلنه ,والمتفاعلات ,تنتفيع ,لفعلتها ,الفوعلي ,والفيعيل ,لتفعلت ,وتفعلات ,وتفعلت ,المفيعلة ,وفعليون ,لفعولي

,فعوليون ,مفعولين ,يفعيله ,وفعيله ,وفعيله ,لفيعال ,الفعلانين ,الفعيلة ,فعلون ,للفعلية ,وفعليان ,ومفعالة ,افاعيل ,لمفعله ,والفيعلى ,الفعلوى ,فعلاك ,وفعلوتي ,افعيلاله

,مفوعلة ,والمفوعلة ,وكالمفعولة ,كالمفعولة ,والفعلاه ,وافعلاه ,وافعليتاه ,ووافعليتاه ,فعلينة ,وفعلنية ,يافعال ,وفعلنته ,وفعلنته ,يفيعلون ,وفتاعل ,المفعولون ,الفعلانيات

,فافعلت ,والفعليان ,وافتعالت ,افيعلاتا ,وفيعلى ,والفعلواني ,الفيعلى ,فعلونة ,الفعلانة ,الفعلونة ,وافعلنا ,الفعليا ,الفعليا ,فعلتنك ,فيعلتك ,فعلتيها ,الفعلوتان ,والفعاولة ,وفعلنة

,وفعلتان ,لفتعلل ,والفعللى ,والافعلان ,وافعلانية ,وفتعلة ,فتعلة ,وفعلاه ,وفعلاه ,وفعلاه ,والفعلوللان ,والفعلولي ,الفعلولي ,واقتعلني ,لبفعلي ,مفعلية ,مفعليا ,لبفعلي ,مفعلية ,مفعليا ,بفعيله ,تنمفعل ,نتمفعل

,المفيعلان ,ومفيعلاناتها ,فعواله ,وفعللوا ,بفعلله ,والمتفعلل ,المفعللل ,المفعللل ,فعلللى ,مفعلللة ,لتفعيلهم ,وفعلويلة ,وفعللنك ,افتعلنك ,وفعلانك ,والفعيليات ,وتفوعل

,لفاعلين ,وفعيول ,والفعالينة ,الفوعلية ,فيفتعله ,بفعيلتهم ,بفعيلته ,وتفعالها ,لمفعيل ,وتفعالها ,لمفعيل ,وتفعالها ,لمفعيل ,وفعليته ,والمفعلى ,والمفعلى ,يفوعله ,ومفعلى ,لتفعله ,فاعلتان ,اليفعلية

,للفعالي ,فعلوته ,فيعلل ,وفاعلتي ,لينفعل ,مفعلى ,مفعلى ,الفعلتين ,وينفعال ,فعلالة ,فيعللها ,افعلها ,افعلها ,ويتفيع ,وفعلتكم ,بفواعلك ,بفواعله ,باستفعاله ,استفعلك

,لمفاعلتهم ,وتتفعل ,فيفعلونها ,مفعولتان ,مفعولتان ,وفعلوها ,الفوعليل ,الفوعليل ,وفعلاة ,وفعللين ,وفعللين ,الفعلويل ,الفعلويل ,الفعلللون ,وفعللللون ,وفعلللون ,الافعلاون ,وتفعلنته

,لمتفعل ,وافعللى ,والفعللى ,والفعللللت ,والفعلولات ,مفعللل ,الفعلاللة ,الفعلاللة ,لفعلاني ,وفياعلة ,فعيلا ,وتمفعل ,واستفعالك ,لبفاعل ,بفعالين ,والفعيلات ,افاعلوا

,ومفو عل ,بمفاعلة ,الفعتلي ,الفعاتلة ,والمفعوتل ,مفعوتل ,بفعيلته ,وبفعيلته ,يفعلونك ,الافعلال ,وفاعلناهم ,للفعيلي ,وفعلتهم ,فعلاتهم ,فعلاوة ,وفعلاو ,وفعلاوون ,وفعلاوون

,فاعلو ,مفعلانة ,لتفعالة ,ويتفعلها ,وافتعلته ,وفاعول ,والفاعول ,وكفعيلة ,كفعله ,الفوعلان ,كفعاله ,والفوعلان ,لفعال ,ومفعال ,لفعلال ,الفوعلال ,الفعلللل ,والفعلللل ,وتفعلنا ,وتفعللنا ,الفعللنا ,وفعلتني

,الفعاليه ,وفعلنت ,فعولا ,الفيعوال ,بمنفعل ,وفعيلان ,وفعالتان ,لبفعلة ,والفعالاة ,والفعالات ,فالمفعله ,والمفعله ,المتفعلات ,وتفعلى ,يفعلنكم ,لافعلنكم ,والتفعيل ,وافاعول

والينفعل, افتعلتها, افعلانية, وفعلتيه, وفعلته, والمفاعلك, ومفاعلك, الفعلنون, وفعلنين, وتفعلله, وافعلليت, بفعاليله, بفعاليه, وفعلوله, والفعلانة, لفعلل, الفعلال, الفعلالان, يفعللل, تفوعل, وتفيعلت, اليفعلة, وفاعلو هم, نستفعل, فعياله, والمفعلةتكون, وافتعلتنا, فتفاعلون, ولافتعاله, وفعللول, الفيعالون, الفو على, لفعيال, فعيلكم, الفعيللية, بمفيعل, الافعلة, فعولهن, والفيعالة, يفعلت, والتفعلت, افعلة, الافعلين, وفعيلين, لتفاعلها, تفتعلان, يفتعلونها, الفتعلة, الفيعلية, كفاعلة, واليفاعل, فعالليون, واليفاعل, كفاعلة, وفعوللة, لفوعلة, الفعلان, افتعلله, الفعيلتان, كفعلي, وافتعلت, لافتعاله, والافاعلة, افتعلنا, المفتعلل, والفعلانة, بافتعلل, افعيلي, والفعويل, ومفيعلا, مفيعلا, مفيعلون, مفعلك, فعليلته, افعللة, المفعللي, والافعلاه, الفعلاه, بفوعله, وبفعاله, وتفيعلة, بفاعله, الفاعلى, والفعيللان, كفاعل, تفعلين, فعيالة, وافعاللت, فعيلاك, الفاعولي, الفاعولي, ومفعتل, لفعتل, الفعلاينة, واستفعلناه, متفعللة, وفعللى, الفو عيلة, يتفعلى, بافعالي, مستفعله, لفعلانها, الفتعلل, والفلان, بالتفعلى, تفعلتا, تفعلية, لفيعلي, لفو علاني, لفعيلهم, فعلاية, والمفعللية, بالفعليل, لفاعلان, ويفعيل, التفعالة, لافعلة, واليفعولة, وافعلليت, افتعلتهم, بفعيلتها, بفعيلتها, والفاعلات, الفعلي, والمفعولون, والمفعولون, مفاعلين, تفعاله, وانفعلنا, فعلل, والتفعلوت, تفعلوت, بفعالتين, ويتفعلون, والفوعلية, بافتعليلة, وفو عال, افعوللوا, تفعلى, مافعلة, النفعول, مافعلة, كفعلولة, لفعيلي, وفعالاهما, الفعليين, والتفعيلة, كالفاعلية, وفعلتة, فعيلتة, لمفعلا, يتفعلني, للفعلول, فعلولة, وفعلولة, سيفعلون, وفعلاه, وفعلولوا, فعيلتي, وفعللاني, ففعللاني, الفعاتيل, الفعاتيل, الافاعل, والافاعل, بفعلانه, فافتعلوه, فتفعلهم, وفاعلني, الافعل, فيتفعل, يتفعلتون, وبالمفعللة, وفعلالا, الفعليللة, الفعيلات, الفو علان, الافعيلل, الافعيلل, وافاعلا, وفاعلا, فعالياتها, الافاعين, والفاعليات, الفاعليات, الافاعل, والفعلال, وافو عل, وافو علالا, فاعلوك, وفعلتنى, فاعاول, الفياعول, فاعيال, فاعيال, فيعالها, مفاعلتي, مفاعلتي, فيعولة, بفعاليها, وبمفعل, والافاعلة, الافعلان, لمفاعلتها, فيفعلله, فعلياوان, وفعلياوي, وفعلياوى, فعلياوين, فعللان, فعلات, وافعلل, بمفعلل, وفعللاه, لافعلله, وفعللان, فعوللان, وفعيللان, بتفعيلك, وبلفعلل, والفعللين, فعلللاة, فالفياعل, مفعلن, ومفعالي, فعلانات, فعللنا, فعللتين, وفعول, ومتفاعل, بالفواعيل, الافعول, والافاعيل, بالفعللول, وفعيليل, تفتعلنى, مفعولها, بالفعولل, فعلتهما, وافعلالا, كمستفعل, افعلى, بالفوعلة, مفعلكم, وفوعلا, بمفعلنا, فافعلي, الفتعللان, تفعلوني, الفعلت, ففعلتة, نفعال, وفعللية, تفاعلتم, ونفعله, فعاليلا, فعيالا, وفعيالة, تفعالها, كالفعلول, وفعلوهم, وفعلوهم, فو عليا, الفعالك, الفعلوك, وفعالك, والمفعولات, فاعلوه, والافعول, افعيلا, الفعاوييلات, الفعللالة, ومفعلته, وفاعلتني, بفعلتي, واتفعل, واتفعلالا, والاتفعلال, فاعليه, نفتعل, بالتفعال, والافعيل, وفعاليا, الفعللتين, مفعللتان, بتفعلوتها, الفعاة, فو عى, وفعوت, وفعيته, فعيا, فعا, يتفاعلن, الفعاليا, ويتفعتل, اتفعلل, متفعلات, فعلنتها, مفعلنا, فعليانة, ونستفعلا, افتعالتها, وتفاعيله, واستفعالك, استفعال, انفعالهم, فالفاعول, فتعليله, تفعلنيهم, تفاعلن, فعلنية, فعلنيه, وفعوتل, افعالها, الفواعلا, لفعلت, فعليلته, كالفعيله, وتفعلا, ومفاعلون, وافعولاه, وافعولاه, يفعلته, ففعلوهم, وفعلتهما, وفعلتهما, ففعلناه, يستفعلك, يفعلونهم, يفعلونهم, يستفعلان, فعلا, بفعلا, بمفعاليه, مفمفاعل, وتفاعلوه, وفواعلة, فمتفعل, فعيالها, فاعلتهم, فعالاتهم, يتفاعله, الفعاولة, فعاولة, وفاعله, فاستفعلني, كفعلات, مفعالات, يستفعلان, يستفعلاه, وفعاليه, المفعلتين, فعيلاي, فنفعل, الفاعلينا, واستفعلتك, بفعليها, ونفتعل, بفعلتك, فعالهن, يستفعلون, لتتفعال, يستفعلن, وافعلا, وافعل, وففعلي, وفعلاه, فيفاعل, لنفعله, وفيعولي, وفياعيلي, فياعيلك, يتفعلان, يفتعلن, بفعيلها, وفو علوا, وفو علوا, فعليبه, مفعلته, فعلتموه, ومفتعلهم, مستفعلو, بمفعالكم, ومفاعلكم, ومفاعلته, فاعلاه, فوا اعلهما, ويفاعلهم, واستفعلنا, والنفعلت, تفاعلاه, بفواعيلها, بفيعالها, ويفاعلني, تفاعلتا, فعلاتك, ومفاعلها, فتفاعل, بفيعلانته, فياعيله, وفوا اعلهم, فاعلتيها, افتعلوه, فمفعلة, بمفعولة, وتفاعلاه, ويفاعلن, ففاعلت, نستفعلون, ويفعلونهم, وفيعالتيه, فيعالات, بفاعلته, بفوا اعلها, بالفعيله, فيعولي, مفعالهم, والفعلك, وفعلتين, وفعلتين, فافعله, بالفوعال, فالمفعلة, فتفاعلاه, وبمفاعلهم, افعاليل, ليستفعل, وفوعلية, واستفعلك, يفعلله, وفعولهم, يفاعلني, بمستفعلين, الفو عله, ويفاعليه, فعلاوين, وافعلن, مفعلنة, ومفعولتها, وفعليهما, اقتفعلني, فوا اعلي, فاقتعلناه, فعلالى, ومفعولو هو, فاعلية, تفاعلون, ومفاعلات, وبفعول, نتفاعل, بالفو عل, وتفاعلوني, لبالمفعال, فعلتاهما, وفعلتاه, بفعولا, للفعاله, مفعلانه, بمفعالهم, والمفعلتين, والمفعلتين, وافاعله, المفعلانية, وفعلانة, فعالنة, لفعلهم, وفعولات, وفعولات, فعالياها, تفعلونه, وفعلاتاته, فعلولك, وفعليها, فتفاعلوه, الفاعلية, تفعيلاته, بمفعولين, وفوعلته, وفعيلتها, تنتفعلان, يتفعلهم, ويتفعلهم, ومتفيعل, فوا اعلكم, وافتعللني, كفاعلته, وفعلاته, وفاعلو ها, فعلاكما, ومفعالها, والفاعلتين, والفاعلتين, فعليه, فعلتني, ففعلتني, مفعليه, فتفعلا, فاعلوني, ومفتعلي, ومتفاعلة, فو علها, ومفعلتهم, وفعيلتي, لفعلتهم, وفعاليها, وفعاليها, فعياله, فيعالي, فياعيلي, فاستفعله, يفعلكم, فاعلوني, ومفتعلي, وفويعلي, وفويعلي, وفعلتي, وفاعلاه, وفاعلك, كالمفعلي, بفياعليه, وفعلاتها, بفعلانه, فاعلاته, لاستفعل, بفعاليل, ومفتعلك, وافتعلتك, وافتعلتك, يفتعلني, يفتعلهم, يفتعلهم, فيعلين, ومفاعيله, لفاعله, فعولان, فواعلك, فتفعيل, وتفاعلناه, مفاعلنا, وفعلاتها, بفعلانتيه, بمفعلته, استفعلتها, لمتفاعل, فاستفعلا, ففاعلتهم, ففاعلتهم, وفاعلتهم, ومتفعله, كالمفاعل, وفعلتك, وفعولتك, ومفعلتي, وفعلانيهم, وفعلتهن, فعلاتكم, ومتفعلاته, مفعلتكم, تستفعلان, تفعلتنا, وافتعلهم, فاقتعلهم, تفعلتنا, ويفتعلنه, ويتفعلنه, وبفعاله, وفعالتين, فعالتين, فعالتا, نفعلهما, يتفاعلاه, الفعاليت, المفاعلية, يتفعلونه, وتفعلهما, متفعلين, والفعولات, فعلتمو ها, فعلتمو ها, ولفعلك, بفعلكما, ويفعلة, اليفعلات, فعلهية, ومفعلهم, مفعوللة, وفعللتهم, مفاعلهن, مفتعلك, ومتفعلك, كالفياعل, وتفاعلني, فعلاهما, فعلاهما, ففعلناهما, فاعلها, فعلناهن, وفعاليلهم, وفعاليلهم, فتعللوا, فاعتلله, فعلاتي, بفعللها, استفعلهم, فتعاللوا, والفعلتين, فوعال, كالفعيلة, ومفعلتها, ويفعلن, ومفعاله, بفعالى, بفيعول, وفعيليلة, اليفعلية, والتفعلية, التفعلية, بتفاعل, ومفعليه, وفعلانون, وفعلاتنه, ويفتعلهم, وفاعلهم, تفاعيلها, لمنفعل, وبمفعولة, منفعله, متفعلون, وبالمفاعل, بالفعولاة, بالمفاعل, لمفاعل, لاستفعلت, والفعاتيل, تفو علا, بفعاليه, ستفعله, ومتفعلا, ومتفاعلا, فو علا, لمفاعلته, وبفعلات, التفعلات, وافعلة, بالفو علي, ففاعلته, وفيتعال, مفعلهن, ومفاعلهن, وفاعله, ويفاعله, وافو علة, وافوعلي, ويفعلي, بمفعليه, لفعالها, وتفعوله, فاعتلالها, نتنفعله, ونفتعله, فعتال, الفتعالل, لنفعلن, فعلناك, فتعلات, وفتعلها, بالفعله, وفعلتها, فعليتان, فعالينه, ويفتعلونها, فاعولها, اعولها, افتعلهن, واستفعلوه, يستفعلونه, يفاعيلها, اليفعلى, استفعلتهم, فاعليها, المفاعله, يفعلو, للافتعال, الفوعل, الفوعل, ومفتعال, بمفتعال, تفعلنه, يفاعلنا, واستفعلوهم, تفعلانه, التفاعلية, واستفعلتها, ولافتعلت, فاقتعلني, بفاعلتي, وفيعلة, افتعله, وفعلم, وفعلتم, يالفعلة, والمفعلان, وبفعلان, والمستفعلين, فاستفعلها, وبفوا اعل, اقتعلاى, يستفعلها, وساتفعله, ومفاعلتك, وتفاعلوها, وافعلوا, لبفعله, لمستفعل, ومفعولك, فاقتفعلاني, فعلناكم, الفاعلار, الفاعلان, الفعللانية, ويالفعلة, ونفعلل, فعللتني, الفعلالل, كافتعل, ففعال, لانفعالها, وفعلانى, افعللت, والنفعلة, النفعلة, فعلاوها, استفعلوها, فعلاهم, وفعلاهم, كتفاعيل, الفيعولا, فيفيعلا, فاعولته, وفعلاويون, وافعيلالا, افعيلال, الفعتلة, النفاعل, لفعلا, لبفتعلن, لفعلا, فعللك, فعللك, كفعلل, والمفعللة, فاستفعلى, وفعلاني, فلتفعل, فياعلها, فياعلانيا, كالمفعلل, واستفعاله, وفعلهن, والتفعيلا, فعللكم, فعلالكم, ليستفعلها, ويستفعلون, فعلانكم, نفعلهم, فافعال, اليتفعل, لفاعول, والفعتل, بيفعل, اسفعل, فعلينان, فعاتلة, فاعولك, بفاعولي, فعويل, بالفيعول, مفعلو ها, يفاعلنهن, بتفاعلهم, وفاعلهم, كالفعلال, بالفيعلي, بالفيعلي, وفعيلات, وفعيلات, ففاعيل, الفعلوتين, بالفعلوة, والفعلوان, والفعلان, تفعلنهم, فعلتمو هم, افاعلتم, كالفعاله, الفعالي, الفعالى, يفاعلنهن, وفعلين, كالفعللة, افتعلهن, بفعولة, وفعلنانة, النفعل, مفعلتين, مفاعيلها, بالفعالين, لفعولها, وللفاعول, فالمتفعل, والفاعولية, والفاعلولى, تفاعيلي, تفاعيله, فعلواي, للفعلية, وفعيلية, وفعلولوه, بفعوله, الفعلولى, فعللال, الفعلالية, وفعلاهما, وفعلهما, ومفعلتك, الفعلالي, ففاعلني, وفعلله, وفعللوه, وفعلهم, وافعلهم, فتفعلها, نتفاعله, المفتعلون, متفعلها, لافتعالهم, انفعالها, بفواعيل, الفو اعلات, المنفعل, والفعالا, والفعاال, فتفعللت, بفعللاتها, بالفعاله, تفاعلك, يتفاعلوا, الفعيالة, وتفاعلونه, والفعيلانة, يتفاعلونه, وتفاعلنه, بفعلانها, انفعلا, الفعلليات, فعللتا, تفعلينا, فتفعلني, وللفعال, لافعيلاله, ففعلك, فففعلها, وفيعلا, لبفعلوه, فعولتها, وفعلتها, وانفعلتا, وافعلنا, الفعلاه, تفعلوهم, الفعال, واليفعليين, اليفعليين, فعلويه, يفاعلا, لفعلل, فعوللى, والفعوللى, الفعوللى, فو علتي, وفاعولهما, وفاعولهما, وفاعلوهم, وفيعمل, ويفعلكم, ويفعلكم, بفعاليلهم, الفعالت, وفاعلو, بفعلتكم, فعوليين, فالمفتعل, بالفعالة, لافعلال, لافعال, الفعلال, كفعاته, الفعاللهم, وفعلو هن, فعلو هن, فعلكم, بفعلكم, ينفعلوا, فتفعللوا, فعليانات, فالفعلات, ياعلوا, وياعلوا, نفاعل, فعلاتي, وفعللتيه, بالمفعولة, فعيلتيه, ولنفعلنه, يفعلونا, بفعلالة, ومفعلن, ومفعلنة, فيفعلن, والفعلولى, الفعللولى, الفعللول, الفعللول, الفعاليك, كفعلى, فعلانتهم, فوعلتها, وفيعلالا, المفوعل, يفعلنها, متفاعلنا, فو عل, والمفعلا, فيفعلهم, ولفعله, و والفعلا, والفعيلا, وفياعلك, وفياعيلك, ففعلى, فعلتة, وفعلى, فو علان, فوعلن, فوعلن, وافتعلاه, فعلتاها, المفاعيلا, لبيفعلك, والفعلى, فالفعلى, فعيلال, لبيفعلكم, ولتفعلكم, لنفعلكم, ويفعلونكم, ومفعلة, بفو عالة, وفيعلته, وفيعلته, كفعلتنه, فويعلين, نفعول, لبيفعلها, ويفاعلهن, كالفعلالات, ينفعلان, للفعيلات, الفعيلون, بالفعيلات, الفعللا, والفيعلوة, بفاعلي, والفيعلولة, يفتعلونه, كفعلها, وافعللل, لبفعلها, فعلنا, تستفعلي, كالافتعال, بافعلال, تفعلناك, تفعلاني, ففعليل, انفعلن, والفعلولة, الفعلانى, انفعلن, والفعلولانية, فعلاوانتك, فعلوانتك, الفعلوانة, فعلاني, فعلالاني, نفو عل, بالمفعلين, نفو عل, الفعاليان, وفعاليلها, وفعاليلها, فعلتكم, يفتعلكم, وبفعالكم, وتفعلو ها, والفيعللى, الفيعللى, تفتعلنا, ويفتعلو ها, ويفتعلوه, فعلليلات, فتيفعل, والفيعلولة, وفعيلولة, يفتعلونه, كفعلها, والفعلللل, وافعللل, تستفعلي, وبيفعلانها, ففعولها, كالافتعان, وتفعلوها, الفيعللى, تفتعلنا, فاعلوا, لبفعلن, المستفعلون, المستفعلون, ولفعول, فاعلتين, بفعول, فالفعلة, فالفعللة, والفعلاله, كالمفاعيل, الفعلال, وفعلاه, والفعلال, وتفعلوان, بفعالك, كالمفتعل, يفتعلوه, وفعيللات, فتيفعل, والفعلللل, للفعللل, الفعلللل, وفعيلن

فعيلنا ,يفعلياه ,ويفعلايه ,وتفاعلك ,سنستفعلهم ,المفعالا ,فعتلة ,فعلايه ,لفعللوا ,فوعول ,كفاعله ,وتفاعلتها ,وافاعلتها ,فوعللات ,ففعلتان ,كفعلله ,ففعلها,

بالفعلال ,يفعيلا ,فمفعالا ,الفوعلين ,لبفاعله ,والفعللون ,الفعلونهن ,يفعلونهن ,كالفواعل ,وبالفعليا ,والفعليا ,بتفعول ,بالتفعول ,ومفيعل ,وكمفعل ,كافعله ,للفعولية,

للتفعلة ,المفعالة ,فعيلتيها ,ويفعلات ,لمفعلان ,كالتفعيلة ,فعلاتنا ,فعالنا ,فعولتهم ,فيفعلوك ,ففعول ,فمفعول ,والفعلوتى ,فعلانيون ,كالفعلى ,الفعلاوان,

الفعليلا ,وتفعيلهم ,وفاعليه ,والفيوعل ,للمفعلى ,ممفعل ,النفعل ,بنفعل ,لتفعيل ,ففعلاته ,وفويعلة ,كفعيلة ,التفعلى ,ففعلان ,وفعاليتان ,وفعاليان ,باليفاعيل,

للفعيلتين ,افعلون ,افتعلنها ,والمفاعلا ,ففعاله ,وللفعالة ,فتفعله ,بمفعالة ,المتفاعلة ,وفاعولها ,وفاعولنا ,بالفاعولة ,فاعليين ,فاعلاتين ,بفيعلان ,ومفوعلة ,المفوعلة,

ففعالا ,والممفعل ,فممفعل ,الممفعل ,فعلكما ,فعلياته ,فالفعلن ,المفعلوها ,بتفعلوها ,افعللانا ,والافعللان ,مفعللين ,فافعللوا ,والفعلنية ,وفعالاها ,للفعلي,

تفعلتي ,والمفعلاني ,بفوعلا ,فعلياوون ,فعلياوي ,فعلياواي ,وافعللتا ,وافعللتا ,فعاليا ,وافعليها ,كافعل ,بفعلله ,وافعللتا ,فعاليتا ,وتفعلتا ,الفيعليين ,ففعللها ,لتفعلنلها ,لتفعللت,

بفعالاتهم ,بفعالال ,تفعلناه ,يفعللك ,لبفعللنك ,فو على ,ويفعلنا ,وتفاعل ,تفعلوه ,فعلينا ,وفعلتاها ,وفعالياها ,وبفعولته ,فيعولا ,الفعليتين ,الفعليتين ,الفعليتا,

بالفعيليات ,وفعلولته ,والافعلل ,ويتفعله ,التفعيلات ,كالتفعيلات ,والمفعيول ,فاعلتم ,الفعلالا ,فاعيني ,افتعللته ,والفعللي ,وفعلولها ,ومفعللا ,بتفعيله,

افعللاه ,لفعلك ,ففواعله ,وافتعل ,وافتعلنها ,فاعلیا ,كالفعلوت ,الفعلولا ,وفعلتيها ,وافتعلن ,كيفعلون ,مفتعيل ,المفيعلون ,والمفاعلين ,افتعللوني ,والفعلللي,

الفعلكة ,بالفاعليين ,بالفعللي ,فعللان ,وفاعال ,وفاعالات ,وفاعال ,فعللالي ,فاعلن ,وفعالاه ,الفعللالون ,الفعلانا ,ففعلونا ,المفعلتان ,تفعلتك ,بفوعال ,والفعللا ,ولفاعل,

فعللات ,فافعلوني ,وتفعلوت ,فالمفعولات ,مفعلاه ,فو علاه ,فعاويله ,فعاويلات ,فعللالا ,افعللت ,فعويلا ,للتفاعيل ,الفعللول ,الفعللالول ,يفتعلون ,بفعللالون ,فيعلنا ,بالفعللية,

فعلليه ,ومتفعلين ,فافعليها ,ومستفعلات ,مفاعلان ,فتتفعل ,ومفعولهما ,يستفعلني ,فاعلاك ,ليفعلون ,وافعليها ,فالافتعال ,واليفعلان ,ففاعلهم ,مفاعلونا,

الفويعل ,الفعيللا ,وفعليله ,الفيعيل ,فعولنا ,فعليلي ,يستفعلونني ,ويستفعلونها ,فعلولي ,بفعيلتين ,بالمتفعل ,فاقفعلن ,مفاعلتكم ,سيفتعل ,فاعلتاها ,بفعالكا,

ففعلتا ,بالمفعله ,الفيعولان ,فاعللل ,تفعلتم ,فياعول ,بالفياعيل ,والمفيعلة ,افعيلانا ,والمفتعلات ,المتفاعلات ,بمتفاعله ,فعتيل ,فعتلا ,وفعلاتة ,الفعليول,

للفعلتين ,وافاعلوا ,بمفعلكم ,سيفعله ,فعاليا ,تفعلتني ,افعللا ,تفعليها ,بفعيلات ,المفعلاينا ,كالفعلي ,ففعلنا ,المفتعلان ,بفعلينا ,فعلكا ,فيفعلونه ,افتعلك,

وفعلوني ,وبفعيلته ,مفعولاته ,فيعلها ,كالتفعال ,والمفعلى ,بالتفعلل ,الفعليك ,الفيعلين ,لانفعالهما ,ليفتعلنكم ,مفعللله ,افعلللا ,بالفيعلاني ,فالفيعل,

فاعلاهم ,والفعلللي ,فعللان ,بالفياعلة ,بالفعللي ,ولافتعال ,نفعلا ,ببتفعل ,الفعلليون ,الفعلليون ,فعيتلوها ,بفعللي ,والفياعلا ,الفياعلا ,فليفعلانه ,كفعللة ,بالمفعللة ,فافعلنا,

وفعليلية ,وفعلايه ,تفعلت ,بالفيعلان ,لفعلنا ,فعلنا ,سيفعلها ,يستفعله ,بفعيلتكم ,يستفعله ,وافعللين ,بالفعيلين ,نفيعيل ,فتيعيل ,تفيعيل ,فوعالات ,وافعللوا ,وافعللوا ,فعليا ,الفعلة ,والفعاة,

فوعي ,وفعوي ,فعوي ,فعلانينا ,فعلانيا ,متفعلينا ,وفعلكم ,فاقفعلا ,لفعولتهما ,كالتفعيل ,وتفعللها ,كيفاعيل ,بيفعول ,وافعلوها ,كالفاعلة ,يتفاعلونها,

يفاعلونهم ,ففاعلتم ,فوعلتم ,فعلاواه ,الفعللوتات ,وفعللوه ,فعلانان ,فعلانان ,فيفعلوا ,الفوعلاني ,ومفعولى ,فيفتعلان ,وبفعالي ,فعوليتهم ,تستفعلني ,لفعالي,

للمفعل ,وفعالتاها ,ومفعولك ,ومفعللا ,وفعولان ,بفعالتها ,المستفعلات ,فعوللانة ,وفعيللانة ,ففعلال ,والفعوللة ,الفعوالة ,كالفاعول ,سنفعل ,لتفعلوه,

وفعلتموهم ,والفيعالية ,وفو علانة ,وفيعلانة ,وفيعلاني ,فو عالني ,فعالوك ,ففعلون ,فالفعلون ,مفتعلهم ,باليفعول ,فعلیتة ,كفعلیت ,وفعلنياتها ,الفعلناة,

المفاعلي ,مفاعليان ,لبفعل ,بفعلى ,والفعيللتان ,بفعلويه ,الفعلويهان ,والفعلويهون ,بفعلونا ,تفعلونا ,ليفعلوني ,فعللهما ,فعللهما ,والفعولن ,والفعولن ,تفعلله ,يتفعلكم ,فعولين,

والمفعلانيات ,وافعيلاله ,وافعلوه ,والمتفعول ,فالفعيلان ,الفعلاوات ,استفعلوهم ,كالفعلين ,كتفعيله ,فو علك ,والفعليل ,والفعلون ,فاقفعلون ,فعيلتهما ,فعلانه ,وفيعلول ,افعلات,

ففعلنها ,بفعالبهما ,والفعيللة ,استفعاله ,والفعالا ,الفيعولة ,وكفعلي ,لفياعل ,الفعلتن ,فاعلاتها ,واعلاتها ,فيعليتي ,فالفعلان ,ومفيعلانا ,ومفيعلانها ,مفيعلانات,

فعليكم ,المفعلت ,افعوللت ,مفعوللينا ,فعلللك ,يفعلليني ,فالمفعال ,والتفعلول ,بالاستفعال ,ليفعلن ,ويتمفعلون ,يتفعلن ,فالفعالية ,بالفعالية ,بيفتعل ,فمفعلات,

فعاليها ,والفعلية ,بالمنفعل ,وفعولتي ,وفعالّه ,الفعالل ,وفعالية ,وفعالهم ,وفعللهم ,ففعللتها ,وفعالاته ,وفعالاته ,وفعللا ,وفعولا ,فالمفاعلا ,وفعولا ,فيعلتهم ,وفيعلاتهم ,تتفعله,

فيعيال ,تتفعلون ,فاعتللهم ,لمنفعلها ,فاعلتنا ,ونفعلة ,ونفعلة ,ونفعل ,الفيعولا ,فاعلوا ,واعلوا ,فعلانة ,فعالية ,فيعلله ,فيعلله ,والفعلولا ,والفعلوا ,الفعللالان ,الفعللالین ,فيعللها ,فتعتلله ,كفاعول,

الفاعولية ,ليفعتله ,بالفاعلية ,فعلنات ,والفعتال ,ولنفعلا ,فعيل ,تفعلا ,فعلتلت ,الفاعولا ,ماتفعل ,المتفيعل ,فتتعلل ,كتفيعل ,امفعل ,بالفياعل ,مفاعیلن,

وفاعلانن ,فاعليان ,وفعلاك ,فيفعللون ,كالفعيلا ,كالفعلا ,لبتفعلوا ,والمستفعلة ,والفاعلاني ,فعلاتي ,فعلانيان ,ففعللاني ,فعلانيتان ,بفعلیهم ,فالفعليات ,مفاعلكم,

الفعولاة ,الفعولاة ,وفعولك ,للمتفاعل ,فعلانين ,بانفعال ,كفواعل ,سيفعلك ,والمفعلي ,فيفعلو ها ,فيفعلو ها ,فالفعلك ,فتتفاعل ,فتتفاعل ,مفعلاتهم ,مفعلاته ,وفعللها ,فعللاته ,فاففعلهم ,كافعللت,

افعللها ,فيفاعلونهم ,يفتعلنه ,فمفاعله ,بالمفعال ,والفعلانية ,والفعولي ,والانفعال ,تفتعلون ,فالمفعلات ,يستفعلا ,بفعالللة ,وفعاللى ,وفعاللي ,بفعاللنا ,بفعولك ,ففعليه,

كفعليه ,وفيعوال ,الفعلالية ,فعلليلته ,والبفعيلة ,وفعلو ,وفعلية ,استفعلها ,ففاعلوهم ,فيفعلوه ,تفعلنها ,بفعلاك ,ليفعلاك ,بفعلت ,بافتعلت ,لتفعيلها ,والفعتلان ,افعلانا,

فافتعلوهم ,وبالمفعل ,وفعالون ,لفعيلكم ,ففعللهم ,لنفعل ,لفعلونا ,لنفعل ,لفعلونا ,وفيعللانة ,وفيعلللة ,فعولهما ,بفعيال ,فيعولية ,لاتفعله ,والمفعالا ,والمفعالا ,فعيلتيه ,ففعالیل,

والمفاعله ,بالفعيول ,فعاليتهم ,فتفاعله ,كافتعلال ,وتفعلاتي ,افتعيل ,فتعللنا ,وفيعليل ,ويفعلول ,لفواعل ,الفعلايلا ,وافعلال ,الفعلايلا ,تفتعلني ,فتعليل,

وتفعلهن ,كتفعالة ,افيعلل ,افتوعيل ,فعلتناه ,فاستفعلنا ,لفاعتلال ,بفعلللتيه ,افتعللة ,وافعيلانها ,وفتعله ,يفعلو ها ,وفعلوها ,بفعاليهم ,كفعلت ,تفاعلينا ,الفتعللل,

الفعللال ,فاعولنا ,الفعوليين ,والفعيلاوان ,وفعلوهن ,و استفعلتم ,ويستفعلن ,فاستفعلتم ,فتفعلنا ,وبالفعلة ,بتفعبلهم ,فيفعلهما ,فتعللال ,فاتفعلت ,اتفعلت,

والفاعلول ,افتعلوني ,وللفعلات ,فافعلها ,وفعيلين ,بفعيلين ,فعياينى ,فالفواعلا ,وتفتعله ,بمفاعلته ,الايفعلان ,سيفعل ,فاعلناكم ,لمفعليه ,يستفتعل ,افتعالهم,

المتفاعلين ,فاعلكم ,افعلونا ,افتعلونا ,وفعيلهم ,وفعيله ,وفعلتنا ,فعولتنا ,فعيلنك ,فليفاعلها ,فيفعلون ,وسفعل ,يستفعل ,وسفعل ,سفعلا ,كتفاعل ,فليتفاعل ,فتفاعلوا ,فتفاعلو ها,

تفاعلو ها ,الفعليون ,فاعلتني ,تفاعلة ,وتفعلان ,لاستفعالهم ,والافعيال ,الافعيال ,فعايلة ,فعايلات ,وفعايلات ,ويفاعلات ,ففعلايل ,يفاعلاوات ,لبفعال ,الفاعليل,

ويفعولة ,مستفعلان ,يفعلنك ,فتفاعلته ,فتفاعلته ,لفعلكم ,فياعلهم ,لتفتعلنها ,افتعلتم ,بفعالتك ,وافعلها ,فيستفعلون ,تفتعلون ,لتفعتلون ,وافتعالي ,كالفيعلان ,وافتعالي ,الفعليل,

فاعليلة ,تفعيلاتهم ,بالتفعيلات ,تفعيلات ,ونفاعل ,فو علون ,فعيلهن ,والافعلة ,وفعلها ,وفعلاها ,مفعلونه ,فتفعلتني ,ففعلو هما ,ومفعولته ,فعلالانة ,مفعلى ,فعالكا,

بالفعالك ,وفعلكي ,الفعلواني ,بفعالنهم ,فاعلاتي ,التمفعل ,لتفتعل ,وفواعلات ,لفعلل ,فاعلناها ,فلتفعله ,لمفعلهم ,وفعلهم ,وفيعلها ,تفعللون ,فافتعلو ها ,فتفعلوني,

الفيعالاني ,فعلانن ,فعللن ,فعللها ,وفعللنان ,الفواعلين ,الفواعلين ,ومفيعلا ,بالمفيعلات ,ومفيعال ,فعايلها ,فعايلكم ,وفعلوا ,فيعلوا ,والفعايل ,الفعايل ,تتفعلوا ,وليفعلنهم ,فعللان,

فاعلو هن ,فعلاكم ,تفاعلو هن ,يافعلى ,فستفعل ,لمفعولون ,وفعولتهن ,لفعولتهن ,وفعلتهن ,لمفعلين ,بفاعليه ,لتفعلوا ,وليفعلوا ,ولتفعلوا ,لفعلكما ,افعلكم ,افعلوك,

افعلو هم ,افعلتني ,افعلتم ,فافعلني ,لافعلتم ,لافعلوك ,لافعلناكم ,لافعلناكم ,نفعلكم ,وافعلتهم ,وافعلك ,وافعلن ,وافعلون ,يفعلوكم ,وافعلون ,مفعليهم ,لفاعلوا ,يافعل ,لبفاعلوكم,

يفاعلونك ,كالمفعلين ,ففعلناهن ,لفعلناه ,لنفعلها ,وفعلناكم ,ولنفعلك ,ولنفعله ,ونفعله ,ونفعلهم ,ليفعلنكم ,وسيفعلها ,وسيفعلها ,وفعولهما ,وفعولهما ,وتفاعلون ,تفعلونهما ,وسيفعل,

فافعلهم ,وافعلهم ,ففاعلناها ,وتفعلونه ,يفاعلكم ,يافعلنا ,يافعلتنا ,يافعلتي ,وسيفعلها ,فيستفعلهم ,لنفعلهم ,وافعلو هم ,وفعلو هم ,لمفعولون ,وسيفعلون ,وسيفعلون ,والفاعلات,

فتفعلوه ,لنفعلنكم ,لتفعلنا ,وتستفعلوا ,وتستفعلون ,ولنفعلنهم ,ويستفعلهم ,يفعلاكم ,يفعلاكم ,فعلكن ,تفاعلو هم ,فعلتموني ,لافتعلتم ,ليستفعلنهم ,ويستفعلكم,

بفعالكم ,ببفعلهن ,فعلناهما ,والمنفعلة ,سنفعلهم ,سيفعلهم ,ليفعلهم ,ونفعلهم ,فتفعلون ,فستفعلون ,ستفعلكم ,ونفعلكم ,كفعلكم ,كفعلكم ,فللفعل ,لتفعلة ,لفعلى ,فعلی,

ونفعلك ,فبفعل ,فكفعل ,فلفعل ,بفعولكم ,بفعولهم ,نفعلن ,ويفعلا ,تفعلونها ,ففعللناك ,لفعلناك ,افعلهما ,يتفاعلا ,افعلهما ,بالفعلن ,بفعلتنا ,للفعلن ,تفعلانه,

بفعلنا ,بفعالاتي ,وليفعلن ,مفعلوا ,وبالفعول ,وبفعوله ,وبفعولي ,ولنفعلن ,ولنفعلن ,وفعلانن ,والمفعلين ,فياعيلهم ,للفياعيل ,لفيعيل ,لبفعيل ,بفعيل ,ليفعلن ,بفعلهما ,فعاليلهم ,فعاليلون,

مفتعيلا ,ولنفعلنكم ,مفتعلون ,ولتفعلن ,يفاعلي ,يفعلونكم ,والمفعلين ,فياعيلهم ,فياعيلهم ,بفعالاتهم ,وفعالتهما ,التفعلة ,لفياعيل ,ففعالة ,بفعالكم ,بفاعلهم,

والفعيلين ,والمتفعلين ,بمفعلي ,ليفعلنا ,استفعلوني ,فيفاعله ,فبفاعله ,والمستفعلين ,انفعلتم ,فعلتموهن ,انفعلتم ,فعلتموهم ,يفتعلوكم ,وفاعلو هن ,وفاعلو هن ,لفعاله,

يافعالي ,تستفعلوه ,ويستفعلونك ,يستفعلونك ,ولفعال ,فتفعلونها ,فلفعلتهم ,لتفاعلوا ,لفعلتهم ,افتعلوكم ,افتعلموه ,افتعلموهم ,يفتعلوكم ,وفاعلوهن,

يامفعل ,تفعلونهن ,فيتفعلون ,و استفعلكم ,فاستفعلوه ,و استفعلي ,ويتستفعلوا ,ويستفعلونه ,ستفعلوا ,لفعالين ,لنفعلهم ,ففعيلين ,ففعيلا ,الفعيلا ,لفعلتا ,لفعلتم ,لمفعول,

تفاعلونهم ,فلفاعلوكم ,فلبفاعل ,فاعلوكم ,يفاعلوكم ,يفاعلونكم ,وفاعله ,بفاعلوا ,وليفعتلوا ,وليفعتلوا ,بمفعلهم ,فتفعلوا ,فينفعلوا ,فعولكما ,فعولكم ,لمنفعلون,

لفعولكم ,متفعلكم ,منفعلون ,وتفعلك ,وفعولهن ,ينفعلون ,بفعيلي ,للمتفعلين ,ففعلكم ,ففعلتمو ,ففعلتمه ,فتفعلون ,وبفعلهم ,وللفاعلين ,وللفاعلين ,بفعلاته ,لفعلاته ,وفيعللوا,

.وليتفعلوا ,ويبتفعلوا ,فافتعلوا ,لفعلناهم ,ويفعلونكم ,ولتفاعلتم ,وتفاعلتم ,لتفاعلوكم ,يفاعلنك ,استفعلوكم ,لافتعل ,وبفعلت ,وبفعله ,فلفعله ,فسيفعلونها ,يستفعله ,فافعلوهن

**ملخص البحث:**

إنّ استخراج الجـــذور عمليـــة أساســية مهمـــة فـــي غالبيـــة تطبيقـــات اللغـــة العربيـــة، مثـــل: أنظمـــة استرجاع المعلومـــات، وتحليـــل النصـــوص، وتصـــنيف النصـــوص، وأنظمـــة الإجابـــة عـــن الأســـئلة، وضـــغط البيانـــات، والفهرســـه، والتـــدقيق الإملائـــي، وتلخـــيص النصـــوص، والترجمـــة الآليـــة. وإن أي ضـــعفٍ فـــي استخراج جـــذور الكلمـــات مـــن شـــأنه أن يؤثر سلباً على الأداء في تلك التطبيقات.

تُحقـــق خوارزميـــة سُـــنبُل لاستخراج الجـــذور دقـــة عاليـــة فـــي الأداء وتقـــدّم تصـــنيفاً جديـــداً للأحـــرف العربيـــة مـــن شـــأنه أن يقلـــل كثيـــراً مـــن الغمـــوض المتعلـــق بالبادئـــات واللاحقـــات. وتُبـــين المقارنـــة والفحـــص للخوارزميـــات المتاحـــة لاستخراج الجـــذور فـــي اللغـــة العربيـــة أن تلك الخوارزميات لا تزال في حاجة إلى بعض التحسينات.

يعتمـــد استخراج الجـــذور فـــي اللغـــة العربيـــة علـــى استخدام الاوزان، إذ تـــزداد دقـــة العمليـــة بزيـــادة عـــدد الاوزان. فـــي هـــذه الدراســـة، يجـــري تحســـين خوارزميـــة سُـــنبُل لاستخراج الجـــذور مـــن خـــلال تحســـين قواعـــدها وزيـــادة أوزانهـــا. وتـــم استخدام 4320 وزن لاستخراج الجـــذور، وتعـــد هـــذه القائمـــة أطـــول قائمـــة أوزان تـــم استخراجها مـــن مجموعـــة قواعـــد بيانـــات "ثلجـــي". كمـــا تـــم فحـــص الخوارزميـــة الجديـــدة باستخدام تلـــك المجموعـــة مـــن قواعـــد البيانـــات التـــي تحتـــوي علـــى 720000 زوج مـــن الجـــذور والكلمـــات. والجـــدير بالـــذكر أن مجموعـــة قواعـــد البيانـــات المـــذكورة إنّمـــا بُنيـــت لفحـــص خوارزميـــات استخراج الجذور ومقارنتها.

لقـــد تمـــت مقارنـــة الخوارزميـــة الجديـــدة بخوارزميـــة سُـــنبُل لاستخراج الجـــذور، وكانـــت النتيجـــة أنّ خوارزميـــة سُـــنبُل حققـــت دقـــة بلغـــت 68%، فـــي حـــين بلغـــت دقـــة الخوارزميـــة المحسَّنة الجديدة 92%.