

Systematic Verb Stem Generation for Arabic *

Jim Yaghi

DocRec Ltd.,
34 Strathaven Place, Atawhai,
Nelson, New Zealand.
jim@docrec.com

Sane M Yagi

Department of English,
University of Sharjah,
Sharjah, U.A.E.
saneyagi@yahoo.com

Abstract

Performing root-based searching, concordancing, and grammar checking in Arabic requires an efficient method for matching stems with roots and vice versa. Such mapping is complicated by the hundreds of manifestations of the same root. An algorithm based on the generation method used by native speakers is proposed here to provide a mapping from roots to stems. Verb roots are classified by the types of their radicals and the stems they generate. Roots are moulded with morphosemantic and morphosyntactic patterns to generate stems modified for tense, voice, and mode, and affixed for different subject number, gender, and person. The surface forms of applicable morphophonemic transformations are then derived using finite state machines. This paper defines what is meant by 'stem', describes a stem generation engine that the authors developed, and outlines how a generated stem database is compiled for all Arabic verbs.

1 Introduction

Morphological parsers and analysers for Arabic are required to dissect an input word and analyse its components in order to perform even the simplest of language processing tasks. The letters of the majority of Arabic words undergo transformations rendering their roots unrecognisable. Without the root, it is difficult to identify a word's morphosemantic template, which is necessary for pinpointing its meaning, or its morphosyntactic pattern, which is essential for realising properties of the verb, such as its tense, voice, and mode, and its subject's number, gender, and person. It is fundamental that an analyser be able to reverse the transformations a word undergoes in order to match the separated root and template with the untransformed ones in its database. Unfortunately, defining rules to reverse transformations is not simple.

Research in Arabic morphology has primarily focused on morphological analysis rather than stem generation.

Sliding window algorithms (El-Affendi, 1999) use an approximate string matching approach of input words against lists of roots, morphological patterns, prefixes, and suffixes. Algebraic algorithms (El-Affendi, 1991), on the other hand, assign binary values to morphological patterns and input words, then perform some simple algebraic operations to decompose a word into a stem and affixes. Permutation algorithms (Al-Shalabi and Evens, 1998) use the input word's letters to generate all possible trilateral or quadrilateral sequences without violation of the original order of the letters which is then compared with items in a dictionary of roots until a match is found. Linguistic algorithms (Thalouth and Al-Dannan, 1990; Yagi and Harous, 2003) remove letters from an input word that belong to prefixes and suffixes and place the remainder of the word into a list. The members of this list are then tested for a match with a dictionary of morphological patterns.

The primary drawback of many of these techniques is that they attempt to analyse using the information found in the letters of the input word. When roots form words, root letters are often transformed by replacement, fusion, inversion, or deletion, and their positions are lost between stem and affix letters. Most attempts use various closest match algorithms, which introduce a high level of uncertainty. In this paper, we define Arabic verb stems such that root radicals, morphological patterns, and transformations are formally specified. When stems are defined this way, input words can be mapped to correct stem definitions, ensuring that transformations match root radicals rather than estimate them.

Morphological transformation in our definition is largely built around finite state morphology (Beesley, 2001) which assumes that these transformations can be represented in terms of regular relations between regular language forms. Beesley (2001) uses finite state transducers to encode the

* The authors wish to thank the anonymous reviewers of this article as their suggestions have improved it significantly.

intersection between roots, morphological patterns, and the transformation rules that account for morphophonemic phenomena such as assimilation, deletion, epenthesis, metathesis, etc.

In this paper, a description of the database required for stem generation is presented, followed by a definition of stem generation. Then the database together with the definition are used to implement a stem generation engine. This is followed by a suggestion for optimising stem generation. Finally, a database of generated stems is compiled in a format useful to various applications that the conclusion alludes to.

In the course of this paper, roots are represented in terms of their ordered sequence of three or four radicals in a set notation, i.e., {F,M,L,Q}. When the capitalised Roman characters F, M, L, and Q are used, they represent a radical variable or place holder. They stand for First Radical (F), Medial Radical (M), Last Radical in a trilateral root (L), and Last Radical in a quadrilateral root (Q).

For readability, all Arabic script used here is followed by an orthographic transliteration between parentheses, using the Buckwalter standard¹. Buckwalter's orthographic transliteration provides a one-to-one character mapping from Arabic to US-ASCII characters. With the exception of a few characters, this transliteration scheme attempts to match the sounds of the Roman letters to the Arabic ones. The following list is a subset of the less obvious transliterations used here: َ (@), ُ (Y), َ (a), ِ (i), ُ (u), َ (o), and َ (~).

2 Stem Generation Database

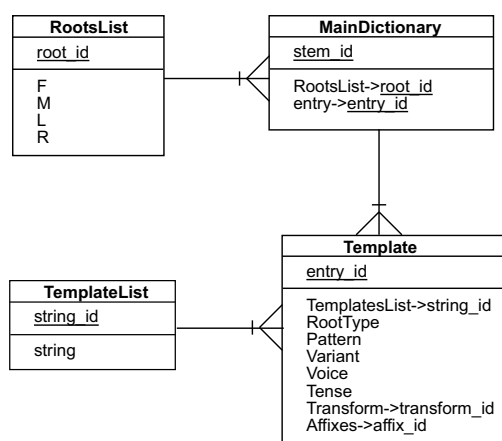


Figure 1: The stem generation database tables and their relations.

¹The complete table of orthographic transliteration may be found at <http://www.qamus.org/transliteration.htm>

Arabic stems can be generated if lists of all roots and all morphological patterns are provided. It is necessary that this data be coupled with a database that links the roots with their morphological patterns (or templates) so that only valid stems are generated for each root. The roots in this database may be moulded with morphosemantic and morphosyntactic patterns to generate intermediate form stems. The stems may then be transformed into final surface forms with a number of specific morphophonemic rules using a finite state transducer compiling language.

Figure 1 shows a summary of the stem generation tables and their relations. The *RootsList* table contains all verb roots from the popular Arabic dictionary, Al-Waseet, (Mustapha et al., 1972), with F, M, L, and Q representing the table fields for up to four radicals per root. A root identifier is used to link this table to the *Template* table. The *Template* table lists all morphosemantic and morphosyntactic patterns used to generate stems from roots of a certain type. This table also specifies the syntactic properties of stems (voice and tense) generated by using the template entry. The *MainDictionary* table links the *RootsList* and *Template* tables together and specifies which entries apply to which roots.

Stems generated with these tables are unaffixed stems. The *affix_id* field links each entry to a subject pronominal affix table that uses transformation rules generating affixed stems. Although object pronominal affixes are not dealt with in this paper, they are generally agglutinating in nature and therefore cause no morphophonemic alterations to a stem. They can be added for generation or removed for analysis without affecting the stem at all.

Affixation and transformation rules are both specified using PERL regular expressions (Friedl, 2002). Regular expressions (Regexp) is an algebraic language that is used for building finite state transducers (FSTs) that accept regular languages. In the next section, Regexp is used to perform morphophonemic transformations and to generate affixed forms of stems. If generated stems are to be useful for root extraction and morphological analysis, it is essential at every stage of generation to be able to track exactly which letters are members of the root radical set, which belong to the template, and what transformations occur on the untransformed stem producing the final surface form.

3 Definition of Stem Generation

In order to be useful in analysis applications, Arabic stems need to be in a surface form which will only undergo agglutinating changes for any further mor-

phological modification. Stems should be defined in terms of the root radicals, morphosemantic and morphosyntactic template letters, and morphophonemic alterations. By doing so, inverting stem transformations becomes trivial. We require the automatic stem generator to always be aware of the origin of each of the letters in stems it generates and to be able to distinguish between letters in the original radical set or in the template string. The stem generator may then be used to compile a complete list of all affixed stems from database roots while retaining all transformation information. The resulting list of stems may then be turned into a searchable index that holds the complete morphological analysis and classification for each entry.

Since originally Arabic words can have a maximum of four root radicals, a root radical set R is defined in terms of the ordered letters of the root as follows:

$$R = \{r_F, r_M, r_L, r_Q\} \quad (1)$$

In the database, pattern, root, variant, and voice-tense ids identify a particular morphological pattern s . Templates are used to generate a stem from a root. The text of s is defined in terms of the letters and diacritics of the template in sequence ($x_1 \dots x_l$) and the radical position markers or place holders (h_F , h_M , h_L , and h_Q), that indicate the positions that letters of the root should be slotted into:

$$s = x_1 x_2 \dots h_F \dots h_M \dots h_L \dots h_Q \dots x_n \quad (2)$$

Stem Generator (SG) uses regular expressions as the language for compiling FSTs for morphophonemic transformations. Transformation rules take into account the context of root radicals in terms of their positions in the template and the nature of the template letters that surround them. Transformations are performed using combinations of regular expression rules applied in sequence, in a manner similar to how humans are subconsciously trained to process the individual transformations. The resulting template between one morphophonemic transformation and the next is an intermediate template. However, in order to aid the next transformation, the transformed radicals are marked by inserting their place holders before them. For example, $h_F \tilde{r} h_M \tilde{s} h_L \tilde{m}$ (**FraMsaLma**) is an intermediate template formed by the root radical set $R = \{r, s, m\}$ ($\{r, s, m\}$) and the morphological pattern $s = h_F \tilde{r} h_M \tilde{s} h_L \tilde{m}$ (**FaMaLa**).

To create the initial intermediate template i_0 from the radical set R and morphological pattern s , a

function $\text{Regexp}(String, SrchPat, ReplStr)$ is defined to compile FSTs from regular expressions. The function accepts in its first argument a string that is tested for a match with the search pattern ($SrchPat$) in its second argument. If $SrchPat$ is found, the matching characters in $String$ are replaced with the replace string ($ReplStr$). This function is assumed to accept the standard PERL regular expression syntax.

A function, $\text{CompileIntermediate}(R, s)$, accepts the radical set R and morphological pattern s to compile the first intermediate template i_0 . A regular expression is built to make this transformation. It searches the morphological pattern text for radical place holders and inserts their respective radical values after them. Since Regexp performs substitutions instead of insertions, replacing each marker with itself followed by its radical value is effectively equivalent to inserting its radical value after it. Let p be a search pattern that matches all occurrences of place holders h_F , h_M , h_L , or h_Q in the morphological pattern, then an initial intermediate form i_0 may be compiled in the following manner:

$$\begin{aligned} i_0 &= \text{CompileIntermediate}(R, s) \\ &= \text{Regexp}(s, p, pR_p) \\ &= \{x_1 \dots h_F r_F \dots h_M r_M \dots h_L r_L \dots h_Q r_Q \dots x_n\} \end{aligned} \quad (3)$$

Let $T = \{t_1 \dots t_m\}$ be the transformation rules applied on each intermediate template to create subsequent intermediate templates. Transformation rules are defined as:

$$t_j = (SrchPat_j, ReplStr_j) \quad (4)$$

A second function $\text{Transform}(i, t)$ is required to perform transformations. A subsequent intermediate template i_{j+1} is the recursive result of transforming the current intermediate template i_j with the next rule t_{j+1} . Each transformation is defined as:

$$\begin{aligned} i_{j+1} &= \text{Transform}(i_j, t_{j+1}) \text{ for } 0 \leq j < m \\ &= \text{Regexp}(i_j, SrchPat_{j+1}, ReplStr_{j+1}) \end{aligned} \quad (5)$$

At any point in the transformation process, the current transformed state of radicals (R') and template string (s') may be decomposed from the current intermediate template as follows:

$$\text{CompileIntermediate}^{-1}(i_j) = (R', s') \quad (6)$$

To turn final intermediate template i_m into a proper stem, a regular expression is built that deletes the place holders from the intermediate template. To do this with a regular expression, the place holders matched are replaced with the null string during the matching process as follows:

$$\text{Regexp}(i_m, p, \text{null}) \quad (7)$$

Basic stems are only modified for tense and voice. Additional morphosyntactic templates or affixation rules further modify proper stems for person, gender, number, and mode. Affixation rules are regular expressions like transformation rules. However, these rules modify final intermediate templates by adding prefixes, infixes, or suffixes, or modifying or deleting stem letters. They require knowledge of the radical positions and occasionally their morphophonemic origins. Adding affixes to a stem operates on the intermediate template which retains the necessary information.

Let a be the affixation rule that is being applied to a certain intermediate template:

$$a = (\text{SrchPat}, \text{ReplStr}) \quad (8)$$

Now using the function Transform that was defined earlier, affixes are added to i_m to produce the intermediate affixed template i_{m+1} :

$$\begin{aligned} i_{m+1} &= \text{Transform}(i_m, a) \\ &= \text{Regexp}(i_m, \text{SrchPat}, \text{ReplStr}) \end{aligned} \quad (9)$$

To convert for output i_{m+1} to an affixed stem, one may remove place holders using the following:

$$\text{Regexp}(i_{m+1}, p, \text{null}) \quad (10)$$

With this definition, generated stems are described by intermediate templates. Intermediate templates retain knowledge of the current state of template and radical letters without losing the ability to recall their origins. This algorithm, therefore, would avoid guesswork in the identification of root radicals. Automatic rule-based stem generation and analysis are both facilitated by this feature of intermediate templates.

4 Stem Generation Engine

A stem generation engine may be built on the basis of the definition just advanced. The three components, *Stem Transformer*, *Affixer*, and *Slotter*, applied in sequence, make up SG. *Stem Transformer* applies the appropriate transformation rules to the morphological pattern, *Affixer* adds specific affixes to the transformed template; and *Slotter* applies the radicals to the transformed affixed template to produce the final affixed stem.

SG begins with a stem ID from the *MainDictionary* table as input to *Stem Transformer* (See Figure 1). The root and entry associated with the stem ID are used to identify the radicals of the root, the morphological pattern string, a list of transformation rules, and an affix table ID.

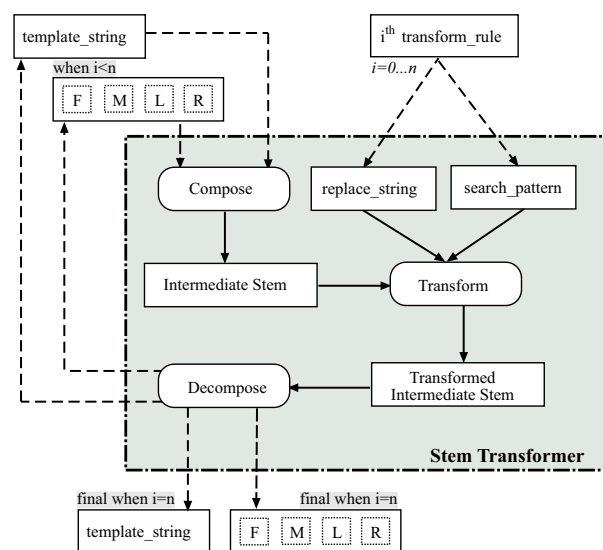


Figure 2: Stem Transformer

Stem Transformer applies transformation rules that are localised to the root radicals and letters of the template in the contexts of one another. To prepare the template and root for transformation, the engine begins by marking radicals in the template.

Stem Transformer is applied incrementally using the current radical set, the template string, and one transformation rule per pass, as in Figure 2. The output of each pass is fed back into Stem Transformer in the form of the j^{th} -rule-transformed template string and radicals, along with the $(j+1)^{\text{th}}$ transformation rule. When all rules associated with the template are exhausted, the resultant template string and radicals are output to the next phase.

To illustrate, assume the morphological pattern $s = | h_{\text{F}} \text{ } \text{ } h_{\text{M}} \text{ } h_{\text{L}} (\text{AiFotaMaLa})$, the radical set $R = \{ \text{ذ}, \text{ك}, \text{ر} \} (\{ @, \mathbf{k}, \mathbf{r} \})$, and the transformation rule set $T = \{ 1, 12 \}$.

Stem Transformer generates a proper stem using the following steps:

Equation 3 above creates the initial intermediate template when passed the radical set and morphological template, thus producing:

$$\begin{aligned}
 i_0 &= \text{CompileIntermediate}(R, s) \\
 &= \text{ } \overset{\text{F}}{\text{h}} \overset{\text{M}}{\text{t}} \overset{\text{L}}{\text{t}} \overset{\text{R}}{\text{h}} \\
 &\quad (\text{AiF@taMkaLra})
 \end{aligned}$$

The first transformation rule $t_1 = 1, t_1 \in T$ is a regular expression that searches for a t (t) following h_F and replaces t (t) with a copy of r_F . To transform i_0 into i_1 with rule t_1 , Equation 5 is used, thus producing:

$$\begin{aligned}
 i_1 &= \text{Transform}(i_0, t_1) \\
 &= \text{ } \overset{\text{F}}{\text{h}} \overset{\text{M}}{\text{t}} \overset{\text{L}}{\text{t}} \overset{\text{R}}{\text{h}} \\
 &\quad (\text{AiF@o@aMkaLra})
 \end{aligned}$$

Next, a gemination rule $t_2 = 12, t_2 \in T$ is applied to i_1 . The gemination regular expression searches for an unvovelled letter followed by a vowelled duplicate and replaces it with the geminated vowelled letter. Once more, Equation 5 is used to make the transformation:

$$\begin{aligned}
 i_2 &= \text{Transform}(i_1, t_2) \\
 &= \text{ } \overset{\text{F}}{\text{h}} \overset{\text{M}}{\text{t}} \overset{\text{L}}{\text{t}} \overset{\text{R}}{\text{h}} \\
 &\quad (\text{AiF@~aMkaLra})
 \end{aligned}$$

To obtain the proper stem from the intermediate template, the final intermediate template i_2 may be substituted into Equation 7:

$$\begin{aligned}
 \text{Stem} &= \text{Regexp}(i_2, p, null) \\
 &= \text{ } \overset{\text{F}}{\text{h}} \overset{\text{M}}{\text{t}} \overset{\text{L}}{\text{t}} \overset{\text{R}}{\text{h}} \\
 &\quad (\text{Ai@~akara})
 \end{aligned}$$

To summarise, the final output of Stem Transformer is a root moulded into a template and a template-transformed radical set. These outputs are used as input to the affixation phase which succeeds stem transformation. Affixer, applied iteratively to the product of Stem Transformer, outputs 14 different subject-pronominally affixed

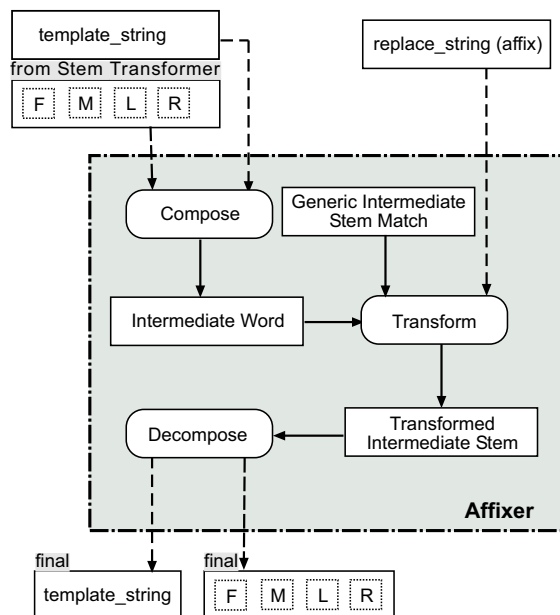


Figure 3: The Affixer Phase

morphosyntactic forms for every input except the imperative which only produces 5. There are 9 different tense-voice-mode combinations per subject pronominal affix, so most roots produce 117 affixed stems per dictionary entry. Affixer is run with different replace strings that are specific to the type of affix being produced. It modifies copies of the transformed stem from the previous phase, as in Figure 3. Using the example cited shortly before, Affixer is passed the last intermediate template i_m and the affix regular expression a . In this example, a is a regular expression that searches for $h_L r_L$ and replaces it with $h_L r_L \overset{\text{t}}{\text{t}}$ ($Lr_L \overset{\text{t}}{\text{t}}$); this corresponds to the past active third person feminine singular affix.

Now applying Equation 9 produces:

$$\begin{aligned}
 i_3 &= \text{Transform}(i_2, a) \\
 &= \text{ } \overset{\text{F}}{\text{h}} \overset{\text{M}}{\text{t}} \overset{\text{L}}{\text{t}} \overset{\text{R}}{\text{h}} \overset{\text{t}}{\text{t}} \\
 &\quad (\text{AiF@~aMkaLrato})
 \end{aligned}$$

In the last stage of stem generation, *Slotter* replaces the place holders in the transformed template with the transformed radical set, producing the final form of the affixed stem. For the example, the result of applying Equation 10 is:

$$\text{Regexp}(i_3, p, null) = \text{ } \overset{\text{F}}{\text{h}} \overset{\text{M}}{\text{t}} \overset{\text{L}}{\text{t}} \overset{\text{R}}{\text{h}} \overset{\text{t}}{\text{t}} \\
 (\text{Ai@~akarato})$$

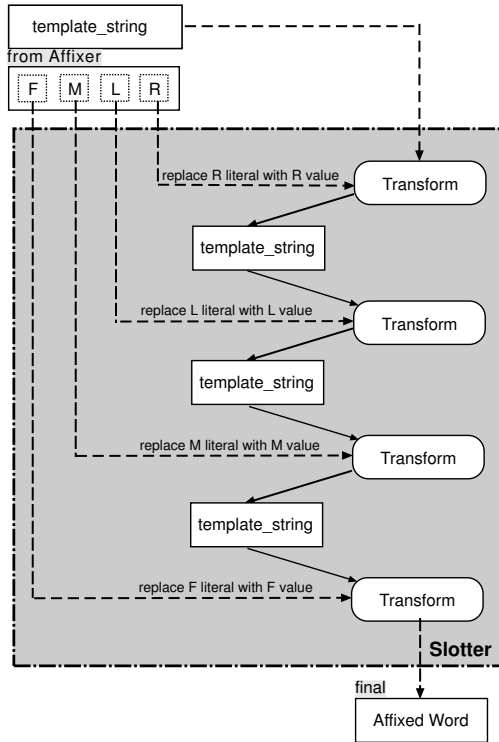


Figure 4: The Slotter Phase

5 Optimisation

Data produced for the use of SG was designed initially with no knowledge of the actual patterns and repetitions that occur with morphophonemic and affix transformation rules. In fact, SG is made to create stems this way: A root is added to a morphosemantic template, then morphosyntactic templates are applied to it, inducing in some patterns morphophonemic transformation. However, while this may be useful in many language teaching tools, it is extremely inefficient. The original data was used to discover patterns that would allow stems to be created in an optimal manner.

Following the classification in Yaghi (2004), there are 70 verb root types associated with 44 theoretically possible morphological patterns. There is an element of repetition present in the classification. In addition, the *Template* table lists sequences of rules that operate on morphological patterns in a manner similar to how native speakers alter patterns phonemically. These rules could be composed into a single FST that yields the surface form.

For example, in the previous section, the morphophonemic transformation rule set $T = \{1, 12\}$ could have been written into one rule. In its non-optimised form the rule duplicates r_F in place of ت (t) creating intermediate form $h_F \text{ت} \text{ت} h_M \text{ك} h_L \text{ر}$ (AiF@o@aMkaLra) and then deletes the first of the

duplicate letters and replaces it with a gemination diacritic that is placed on the second repeat letter. The resulting surface form is اَكْر (Ai@~akara). Instead, one rule could achieve the surface form by replacing the letter ت (t) in the template with a geminated ت (@) yielding the same result.

Compiling separate regular expressions for each transformation rule is costly in terms of processing time especially when used with back-references, as SG does. Back-references group a sub-pattern and refer to it either in the search pattern or substitute string. Such patterns are not constant and are required to be recompiled for every string they are used with. It is desirable, therefore, to minimise the number of times patterns are compiled. To optimise further, the transformation may be made on the morphological pattern itself, thus producing a sound surface form template. This procedure would eliminate the need to perform morphophonemic transformations on stems.

Each template entry in the *Template* table (see Figure 1) is given a new field containing the surface form template. This is a copy of the morphological pattern with morphophonemic transformations applied. A coding scheme is adopted that continues to retain letter origins and radical positions in the template so that this will not affect affixation. Any transformations that affect the morphological pattern alone are applied without further consideration. The coding scheme uses the Roman characters **F**, **M**, **L**, and **Q** to represent place holders in the templates. Each place holder is followed by a single digit indicating the type of transformation that occurs to the radical slotted in that position. The codes have the following meanings: 0=no alteration, 1=deletion, 2=substitution, 3=gemination. If the code used is 2, then the very next letter is used to replace the radical to which the code belongs.

Take for example, the *Template* table entry for the root type 17 (all roots with $F=\text{و}$ (w) and $L=\text{ي}$ (y)), its morphological pattern $h_F \text{ت} h_M \text{ك} h_L \text{ر}$ (AiFotaMaLa), and its variant (ID 0). The morphophonemic transformation rules applied to the template are $T=\{20,12,31,34,112\}$. These rules correspond to the following:

- 20=change r_F to a duplicate of the next letter ت (t)
- 12=geminate duplicate letters
- 31=delete diacritic after the ي (y) in position h_L
- 34=convert ي (y) to ا (A)
- 112=convert ا to ي (Y)

Surface Form		r_F2	r_M0	r_L2		(Ai F2t~a M0a L2Y)
Affix				r_L2		(L2yotumaA)
Combined Result		r_F2	r_M0	r_L2		(Ai F2t~a M0a L2yotumaA)

Table 1: Surface form template aligned with an affix entry rule.

The surface form template can be rewritten as $h_F2h_Mh_L2$ (AiF2t~aM0aL2Y). This can be used to form stems such as ائدى (Ait~adaY) by slotting the root {و, د, ي} ({w,d,y}).

The affix tables use a similar notation for coding their rules. Every affix rule indicates a change to be made to the surface form template and begins with a place holder followed by a code 0 or 2 unless the rule redefines the entire template in which case the entry begins with a 0. Radical place holders in affix rules define changes to the surface form template. These changes affect the template from the given radical position to the very next radical position or the end of the template, whichever is first.

Affix rules with code 0 following radical place holders signify that no change should be made to that section of the surface form template. However, a code 2 after a place holder modifies the surface form template in that position by replacing the letter that follows the code with the rest of that segment of the rule. Affix rules using code 2 after place holders override any other code for that position in the surface form template because affixation modifies morphophonemically transformed stems.

Creating affixed stems from templates and affixes formatted in this way becomes far more optimal. If a surface form template was specified as $h_F2h_Mh_L2$ (AiF2t~aM0aL2Y) and it was to be combined with the affix rule r_L2 (L2yotumaA) then SG simply needs to align the affix rule with the surface form template using the place holder symbol in the affix rule and replace appropriately as in Table 1.

With the resulting affixed surface form template SG may retain the radicals of the original root where they are unchanged, delete radicals marked with code 1 and 3, and substitute letters following code 2 in place of their position holders. If the example above is used with the root {و, د, ي} ({w, d, y}), the final stem is: ائدى (Ait~adayotumaA, meaning "the two of you have accepted compensation for damage").

To use the original regular expression transformations would take an average of 18000 seconds to produce a total of 2.2 million valid stems in the database. With the optimised coding scheme, the time taken is reduced to a mere 720 seconds; that is

4% of the original time taken.

6 Generated Stem Database Compiler

هي	[انفعل - 0]	{وق ي}
انقت	Ait~aqa0	ماضي معلوم Past Active
انقت	Aut~uqiya0	ماضي مجهول Past Passive
تنقي	tat~aqaY	مضارع مرفوع معلوم Present Indicative Active
تنقي	tut~aqaY	مضارع مرفوع مجهول Present Indicative Passive
لن تنقي	tat~aqaY	مضارع منصوب معلوم Present Subjunctive Active

Figure 5: Output from the Stem Generation CGI

Once the dictionary database has been completed and debugged, an implementation of SG generates for every root, template, and affix the entire list of stems derived from a single root and all the possible template and affix combinations that may apply to that root entry. The average number of dictionary entries that a root can generate is approximately 2.5. Considering that each entry generates 117 different affixed stems, this yields an average of approximately 300 affixed stems per root. However, some roots (e.g., {ب,ت,ك} ({k,t,b})) produce 13 different entries, which makes approximately 1,500 affixed stems for each of such roots.

The generated list is later loaded into a B-Tree structured database file that allows fast stem search and entry retrieval.

A web CGI was built that uses the Stem Generation Engine to produce all affixed stems of any given root. A section of the results of this appears in Figure 5.

7 Conclusions

In this paper, we have discussed our attempt at imitating the process used by Arabic speakers in generating stems from roots. We formulated a definition of the process, facilitating an encoding of Arabic stems. The encoding represents stems in terms of their components while still allowing a simple mapping to their final surface forms. A stem's components are a root, morphosemantic and morphosyntactic templates, and any morphophonemic alterations that the stem may have undergone. In doing so, the problem has been reduced to the much smaller task of obtaining stems for the words subject to analysis, and then matching these against the surface forms of the pre-analysed stems. The encoding retains most of the information essential to stem generation and analysis, allowing us to trace the various transformations that root radicals undergo when inflected. Root extractors and morphological analysers can match an input word with a defined verb stem, then use the information in the definition to determine with certainty the stem's root and morphological pattern's meaning. The authors intend to use a similar strategy to define stems for Arabic nouns.

Mapping from words to defined stems is now much easier. The stem generation algorithm here attempts to produce a comprehensive list of all inflected stems. Any verb may be found in this list if some simple conjoin removal rules are first applied. Conjoins are defined here as single letter conjunctions, future or question particles, emphasis affixes, or object pronominal suffixes that agglutinate to a verb stem. Because conjoins may attach to a verb stem in sequence and without causing any morphological alteration, extracting stems from Arabic words becomes similar to extracting stems from English words. In fact, many of the Arabic word analysis approaches reviewed in the introduction to this paper would yield more accurate results if applied to stem extraction instead of root extraction. It would become possible to use for this purpose conventional linguistic, pattern matching, or algebraic algorithms.

The dictionary database described here can be used to form the core of a morphological analyser that derives the root of an input word, identifies its stem, and classifies its morphosemantic and morphosyntactic templates. An analyser based on these principles may be used in many useful applications, some of which are detailed in Yaghi (2004). Example applications include root, lemma based, and exact word analysis, searching, incremental searching, and concordancing.

References

- S. S. Al-Fedaghi and F. S. Al-Anzi. 1989. A New Algorithm to Generate Arabic Root-Pattern Forms. In *Proceedings of the 11th National Computer Conference and Exhibition*, pages 391–400, Dhahran, Saudi Arabia, March.
- Riyad Al-Shalabi and Martha Evens. 1998. A Computational Morphology System for Arabic. In *Proceedings of the COLING/ACL98*, pages 66–72, Montréal, Québec, Canada, August.
- Kenneth R Beesley. 2001. Finite-State Morphological Analysis and Generation of Arabic at Xerox Research: Status and Plans in 2001. In *ARABIC Language Processing: Status and Prospects*, Toulouse, France, July. Arabic NLP Workshop at ACL/EACL 2001.
- Mohammed A. El-Affendi. 1991. An Algebraic Algorithm for Arabic Morphological Analysis. *The Arabian Journal for Science and Engineering*, 16(4B).
- Mohammed A. El-Affendi. 1999. Performing Arabic Morphological Search on the Internet: A Sliding Window Approximate Matching Algorithm and its Performance. Technical report, CCIS Report King Saud University.
- Jeffery E. F. Friedl. 2002. *Mastering Regular Expressions*. O'Reilly, 2nd edition, July.
- Lama Hamandi, Rached Zantout, and Ahmed Guesoum. 2002. Design and Implementation of an Arabic Morphological Analysis System. In *Proceedings of the International Conference on Research Trends in Science and Technology 2002*, pages 325–331, Beirut, Lebanon.
- Ibrahim Mustapha, Ahmed H. Al-Zayat, Hamid AbdelQadir, and Mohammed Ali Al-Najjar, editors. 1972. *Al-Moajam Al-Waseet*. Cairo Arab Language Academy, Cairo, Egypt.
- B. Thalouth and A. Al-Dannan. 1990. A Comprehensive Arabic Morphological Analyzer Generator. In Pierre Mackay, editor, *Computers and the Arabic Language*. Hemisphere Publishing, New York.
- Jim Yaghi. 2004. Computational Arabic Verb Morphology: Analysis and Generation. Master's thesis, University of Auckland.
- Sane M. Yagi and Saad Harous. 2003. Arabic Morphology: An Algorithm and Statistics. In *Proceedings of the 2003 International Conference on Artificial Intelligence (IC-AI 2003)*, Las Vegas, Nevada.