

# SEMĤE: A Generalised Two-Level System

George Anton Kiraz\*

Computer Laboratory

University of Cambridge (St John's College)

Email: [George.Kiraz@cl.cam.ac.uk](mailto:George.Kiraz@cl.cam.ac.uk)

URL: <http://www.cl.cam.ac.uk/users/gk105>

## Abstract

This paper presents a generalised two-level implementation which can handle linear and non-linear morphological operations. An algorithm for the interpretation of multi-tape two-level rules is described. In addition, a number of issues which arise when developing non-linear grammars are discussed with examples from Syriac.

## 1 Introduction

The introduction of two-level morphology (Koskeniemi, 1983) and subsequent developments has made implementing computational-morphology models a feasible task. Yet, two-level formalisms fell short from providing elegant means for the description of non-linear operations such as infixation, circumfixation and root-and-pattern morphology.<sup>1</sup> As a result, two-level implementations – e.g. (Antworth, 1990; Karttunen, 1983; Karttunen and Beesley, 1992; Ritchie et al., 1992) – have always been biased towards linear morphology.

The past decade has seen a number of proposals for handling non-linear morphology;<sup>2</sup> however, none

(apart from Beesley's work) seem to have been implemented over large descriptions, nor have they provided means by which the grammarian can develop non-linear descriptions using higher level notation.

To test the validity of one's proposal or formalism, minimally a medium-scale description is a desideratum. SemHe<sup>3</sup> fulfils this requirement. It is a generalised multi-tape two-level system which is being used in developing non-linear grammars.

This paper (1) presents the algorithms behind SemHe; (2) discusses the issues involved in compiling non-linear descriptions; and (3) proposes extension/solutions to make writing non-linear rules easier and more elegant. The paper assumes knowledge of multi-tape two-level morphology (Kay, 1987; Kiraz, 1994c).

## 2 Linguistic Descriptions

The linguist provides SemHe with three pieces of data: a lexicon, two-level rules and word formation grammar. All entries take the form of Prolog terms.<sup>4</sup> (Identifiers starting with an uppercase letter denote variables, otherwise they are instantiated symbols.) A lexical entry is described by the term

`synword(<morpheme>, <category>).`

Categories are of the form

`<category-symbol> : [<feature_attr1 = value1>,  
...,  
<feature_attrn = valuen>]`

a notational variant of the PATR-II category formalism (Shieber, 1986).

\* Supported by a Benefactor Studentship from St John's College. This research was done under the supervision of Dr Stephen G. Pulman. Thanks to the anonymous reviewers for their comments. All mistakes remain mine.

<sup>1</sup>Although it is possible to express some classes of non-linear rules using standard two-level formalisms by means of *ad hoc* diacritics, e.g., infixation in (Antworth, 1990, p. 156), there are no means for expressing other classes as root-and-pattern phenomena.

<sup>2</sup>(Kay, 1987), (Kataja and Koskeniemi, 1988), (Beesley et al., 1989), (Lavie et al., 1990), (Beesley, 1990), (Beesley, 1991), (Kornai, 1991), (Wiebe, 1992), (Pulman and Hepple, 1993), (Narayanan and Hashem, 1993), and (Bird and Ellison, 1994). See (Kiraz, 1996) for a review.

<sup>3</sup>The name SemHe (Syriac *šemĥê* 'rays') is not an acronym, but the title of a grammatical treatise written by the Syriac polymath (*inter alia* mathematician and grammarian) Bar 'Ebrōyō (1225-1286), viz. *ktōbō dšēmĥê* 'The Book of Rays'.

<sup>4</sup>We describe here the terms which are relevant to this paper. For a full description, see (Kiraz, 1996).

```

tl_alphabet(0, [k,t,b,a,e]). % surface alphabet
tl_alphabet(1, [c1,c2,c3,v,b]). tl_alphabet(2, [k,t,b,b]). tl_alphabet(3, [a,e,b]). % lexical alphabets
tl_set(radical, [k,t,b]). tl_set(vowel, [a,e]). tl_set(c1c3, [c1,c3]). % variable sets
tl_rule(R1, [[], [], []], [[b], [b], [b]], [[], [], []], =>, [], [], [],
[], [[], [], []]).
tl_rule(R2, [[], [], []], [[P], [C], []], [[], [], []], =>, [], [C], [],
[c1c3(P), radical(C)], [[], [], []]).
tl_rule(R3, [[], [], []], [[v], [], [v]], [[], [], []], =>, [], [V], [],
[vowel(V)], [[], [], []]).
tl_rule(R4, [[], [], []], [[v], [], [v]], [[c2,v], [], []], <=>, [], [], [],
[vowel(V)], [[], [], []]).
tl_rule(R5, [[], [], []], [[c2], [C], []], [[], [], []], <=>, [], [C], [],
[radical(C)], [[], [root:[measure=p'al]], []]).
tl_rule(R6, [[], [], []], [[c2], [C], []], [[], [], []], <=>, [], [C,C], [],
[radical(C)], [[], [root:[measure=pa'el]], []]).

```

Listing 1

A two-level rule is described using a syntactic variant of the formalism described by (Ruessink, 1989; Pulman and Hepple, 1993), including the extensions by (Kiraz, 1994c),

```

tl_rule(<id>, <LLC>, <Lex>, <RLC>, <Op>,
<LSC>, <Surf>, <RSC>,
<variables>, <features>).

```

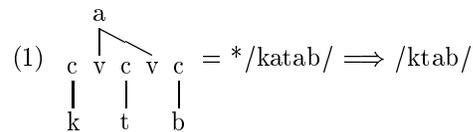
The arguments are: (1) a rule identifier, *id*; (2) the left-lexical-context, *LLC*, the lexical center, *Lex*, and the right-lexical-context, *RLC*, each in the form of a list-of-lists, where the *i*th list represents the *i*th lexical tape; (3) an operator, => for optional rules or <=> for obligatory rules; (4) the left-surface-context, *LSC*, the surface center, *Surf*, and the right-surface-context, *RSC*, each in the form of a list; (5) a list of the *variables* used in the lexical and surface expressions, each member in the form of a predicate indicating the set identifier (see *infra*) and an argument indicating the variable in question; and (6) a set of *features* (i.e. category forms) in the form of a list-of-lists, where the *i*th item must unify with the feature-structure of the morpheme affected by the rule on the *i*th lexical tape.

A lexical string maps to a surface string iff (1) they can be partitioned into pairs of lexical-surface subsequences, where each pair is licenced by a rule, and (2) no partition violates an obligatory rule.

Alphabet declarations take the form `tl_alphabet(<tape>, <symbol_list>)`, and variable sets are described by the predicate `tl_set(<id>, <symbol_list>)`. Word formation rules take the form of unification-based CFG rules, `synrule(<identifier>, <mother>, [[<daughter1>, ..., <daughtern>]])`.

The following example illustrates the derivation

of Syriac /ktab/<sup>5</sup> ‘he wrote’ (in the simple *p’al* measure)<sup>6</sup> from the pattern morpheme {cvcvc} ‘verbal pattern’, root {ktb} ‘notion of writing’, and vocalism {a}. The three morphemes produce the underlying form \*/katab/, which surfaces as /ktab/ since short vowels in open unstressed syllables are deleted. The process is illustrated in (1).<sup>7</sup>



The *pa’el* measure of the same verb, viz. /katteb/, is derived by the gemination of the middle consonant (i.e. t) and applying the appropriate vocalism {ae}.

The two-level grammar (Listing 1) assumes three lexical tapes. Uninstantiated contexts are denoted by an empty list. R1 is the morpheme boundary (= b) rule. R2 and R3 sanction stem consonants and vowels, respectively. R4 is the obligatory vowel deletion rule. R5 and R6 map the second radical, [t], for *p’al* and *pa’el* forms, respectively. In this example, the lexicon contains the entries in (2).<sup>8</sup>

```

(2) synword(c1vc2vc3, pattern: []).
    synword(ktb, root: [measure = M]).
    synword(aa, vocalism: [measure = p'al]).
    synword(ae, vocalism: [measure = pa'el]).

```

Note that the value of ‘measure’ in the root entry is

<sup>5</sup>Spirantization is ignored here; for a discussion on Syriac spirantization, see (Kiraz, 1995).

<sup>6</sup>Syriac verbs are classified under various measures (forms). The basic ones are: *p’al*, *pa’el* and *af’el*.

<sup>7</sup>This analysis is along the lines of (McCarthy, 1981) – based on autosegmental phonology (Goldsmith, 1976).

<sup>8</sup>Spreading is ignored here; for a discussion, see (Kiraz, 1994c).

uninstantiated; it is determined from the feature values in R5, R6 and/or the word grammar (see *infra*, §4.3).

### 3 Implementation

There are two current methods for implementing two-level rules (both implemented in **SemHe**): (1) compiling rules into finite-state automata (multi-tape transducers in our case), and (2) interpreting rules directly. The former provides better performance, while the latter facilitates the debugging of grammars (by tracing and by providing debugging utilities along the lines of (Carter, 1995)). Additionally, the interpreter facilitates the incremental compilation of rules by simply allowing the user to toggle rules on and off.

The compilation of the above formalism into automata is described by (Grimley-Evans et al., 1996). The following is a description of the interpreter.

#### 3.1 Internal Representation

The word grammar is compiled into a shift-reduce parser. In addition, a first-and-follow algorithm, based on (Aho and Ullman, 1977), is applied to compute the feasible follow categories for each category type. The set of feasible follow categories, *NextCats*, of a particular category *Cat* is returned by the predicate FOLLOW(+*Cat*, -*NextCats*). Additionally, FOLLOW(bos, *NextCats*) returns the set of category symbols at the beginning of strings, and eos ∈ *NextCats* indicates that *Cat* may occur at the end of strings.

The lexical component is implemented as character tries (Knuth, 1973), one per tape. Given a list of lexical strings, *Lex*, and a list of lexical pointers, *LexPtrs*, the predicate

LEXICAL-TRANSITIONS(+*Lex*, +*LexPtrs*,  
-*NewLexPtrs*, -*LexCats*)

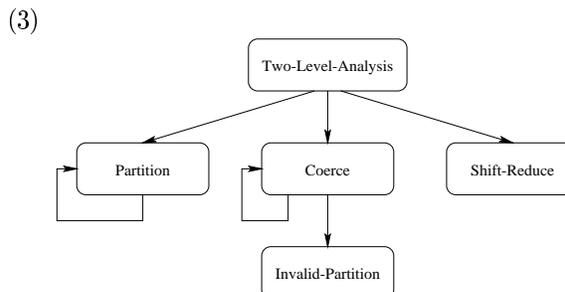
succeeds iff there are transitions on *Lex* from *LexPtrs*; it returns *NewLexPtrs*, and the categories, *LexCats*, at the end of morphemes, if any.

Two-level predicates are converted into an internal representation: (1) every left-context expression is reversed and appended to an uninstantiated tail; (2) every right-context expression is appended to an uninstantiated tail; and (3) each rule is assigned a 6-bit ‘precedence value’ where every bit represents one of the six lexical and surface expressions. If an expression is *not* an empty list (i.e. context is specified), the relevant bit is set. In analysis, surface expressions are assigned the most significant bits, while lexical expressions are assigned the least significant ones. In generation, the opposite state of

affairs holds. Rules are then reasserted in the order of their precedence value. This ensures that rules which contain the most specified expressions are tested first resulting in better performance.

#### 3.2 The Interpreter Algorithm

The algorithms presented below are given in terms of prolog-like non-deterministic operations. A clause is satisfied iff all the conditions under it are satisfied. The predicates are depicted top-down in (3). (**SemHe** makes use of an earlier implementation by (Pulman and Hepple, 1993).)



In order to minimise accumulator-passing arguments, we assume the following initially-empty stacks: *ParseStack* accumulates the category structures of the morphemes identified, and *FeatureStack* maintains the rule features encountered so far. (‘+’ indicates concatenation.)

PARTITION partitions a two-level analysis into sequences of lexical-surface pairs, each licenced by a rule. The base case of the predicate is given in Listing 2,<sup>9</sup> and the recursive case in Listing 3.

The recursive COERCE predicate ensures that no partition is violated by an obligatory rule. It takes three arguments: *Result* is the output of PARTITION (usually reversed by the calling predicate, hence, COERCE deals with the last partition first), *PrevCats* is a register which keeps track of the last morpheme category encountered, and *Partition* returns selected elements from *Result*. The base case of the predicate is simply COERCE([], -, []) – i.e., no more partitions. The recursive case is shown in Listing 4. *CurrentCats* keeps track of the category of the morpheme which occurs in the current partition. The invalidity of a partition is determined by INVALID-PARTITION (Listing 5).

TWO-LEVEL-ANALYSIS (Listing 6) is the main predicate. It takes a surface string or lexical string(s) and returns a list of partitions and a

<sup>9</sup>For efficiency, variables appearing in left-context and centre expressions are evaluated after LEXICAL-TRANSITIONS since they will be fully instantiated then; only right-contexts are evaluated after the recursion.

```

PARTITION(SurfDone, SurfToDo, LexDone, LexToDo, LexPtrs, NextCats, Result)
SurfToDo = [] & % surface string exhausted
LexToDo = [[], [], ..., []] & % all lexical strings exhausted
LexPtrs = [rt, rt, ..., rt] & % all lexical pointers are at the root node
eos ∈ NextCats & % end-of-string
Result = []. % output: no more results

```

Listing 2

```

PARTITION(SurfDone, SurfToDo, LexDone, LexToDo, LexPtrs, NextCats,
[ResultHead | ResultTail])
there is tl_rule(Id, LLC, Lex, RLC, Op, LSC, Surf, RSC, Variables, Features) such that
(Op = (=> or <=>), LexDone = LLC, SurfDone = LSC,
SurfToDo = Surf + RSC and LexToDo = Lex + RLC) &
LEXICAL-TRANSITIONS(Lex, LexPtrs, NewLexPtrs, LexCats) &
push Features onto FeatureStack & % keep track of rule features
if LexCats ≠ nil then % found a morpheme boundary?
  while FeatureStack is not empty % unify rule and lexical features
    unify LexCats with (pop FeatureStack) &
  push LexCats onto ParseStack & % update the parse stack
  if LexCats ∈ NextCats then % get next category
    FOLLOW(LexCats, NewNextCats)
end if &
ResultHead = Id/SurfDone/Surf/RSC/
  LexDone/Lex/RLC/LexCats &
NewSurfDone = SurfDone + reverse Surf & % make new arguments ...
NewSurfToDo = RSC & % ... and recurse
NewLexDone = LexDone + reverse Lex &
NewLexToDo = RLC &
PARTITION(NewSurfDone, NewSurfToDo,
NewLexDone, NewLexToDo,
NewLexPtrs, NewNextCats, ResultTail) &
for all SetId(Var) ∈ Variables % check variables
  there is tl_set(SetId, Set) such that Var ∈ Set.

```

Listing 3

```

COERCE([Id/LSC/Surf/RSC/LLC/Lex/RLC/LexCats | ResultTail], PrevCats,
[Id/Surf/Lex | PartitionTail])
if LexCats ≠ nil then
  CurrentCats = LexCats
else
  CurrentCats = PrevCats &
not INVALID-PARTITION(LSC, Surf, RSC, LLC, Lex, RLC, CurrentCats) &
COERCE(ResultTail, CurrentCats, PartitionTail).

```

Listing 4

```

INVALID-PARTITION(LSC, Surf, RSC, LLC, Lex, RLC, Cats)
there is tl_rule(Id, LLC, Lex, RLC, <=>, LSC, NotSurf, RSC, Variables, Features) such that
NotSurf ≠ Surf &
for all SetId(Var) ∈ Variables % check variables
  there is tl_set(SetId, Set) such that Var ∈ Set &
unify Cats with Features &
fail.

```

Listing 5

TWO-LEVEL-ANALYSIS(?*Surf*, ?*Lex*, -*Partition*, -*Parse*)  
 FOLLOW(**bos**, *NextCats*) &  
 PARTITION([], *Surf*, [[], [], ..., []], *Lex*, [**rt**, **rt**, ..., **rt**], *NextCats*, *Result*) &  
 COERCE(**reverse** *Result*, **nil**, *Partition*) &  
 SHIFT-REDUCE(*ParseStack*, *Parse*).

Listing 6

morphosyntactic parse tree. To analyse a surface form, one calls TWO-LEVEL-ANALYSIS(+*Surf*, -*Lex*, -*Partition*, -*Parse*). To generate a surface form, one calls TWO-LEVEL-ANALYSIS(-*Surf*, +*Lex*, -*Partition*, -*Parse*).

## 4 Developing Non-Linear Grammars

When developing Semitic grammars, one comes across various issues and problems which normally do not arise with linear grammars. Some can be solved by known methods or ‘tricks’; others require extensions in order to make developing grammars easier and more elegant. This section discusses issues which normally do not arise when compiling linear grammars.

### 4.1 Linearity vs. Non-Linearity

In Semitic languages, non-linearity occurs only in stems. Hence, lexical descriptions of stems make use of three lexical tapes (pattern, root & vocalism), while those of prefixes and suffixes use the first lexical tape. This requires duplicating rules when stating lexical constraints. Consider rule R4 (Listing 1). It allows the deletion of the first stem vowel by the virtue of *RLC* (even if  $c_2$  was not indexed); hence /katab/ → /ktab/. Now consider adding the suffix {eh} ‘him/it’: /katab/+{eh} → /katbeh/, where the second stem vowel is deleted since deletion applies right-to-left; however, *RLC* can only cope with stem vowels. Rule R7 (Listing 7) is required. One might suggest placing constraints on surface expressions instead. However, doing so causes surface expressions to be dependent on other rules.

Additionally, *Lex* in R4 and R7 deletes stem vowels. Consider adding the prefix {wa} ‘and’: {wa} + /katab/ + {eh} → /wkatbeh/, where the prefix vowel is also deleted. To cope with this, two additional rules like R4 and R7 are required, but with *Lex* = [[V], [], []].

We resolve this by allowing the user to write expansion rules of the form

`expand(<symbol>, <expansion>, <variables>).`

In our example, the expansion rules in (4) are needed.

- (4) `expand(C, [[C], [], []], [radical(C)]).`  
`expand(C, [[c], [C], []], [radical(C)]).`  
`expand(V, [[V], [], []], [vowel(V)]).`  
`expand(V, [[v], [], [V]], [vowel(V)]).`

The linguist can then rewrite R4 as R8 (Listing 7), and expand it with the command `expand(R8)`. This produces four rules of the form of R4, but with the following expressions for *Lex* and *RLC*:<sup>10</sup>

<i>Lex</i>	<i>RLC</i>
[[V1], [], []]	[[C, V2], [], []]
[[V1], [], []]	[[c, v], [C], [V2]]
[[v], [], [V1]]	[[C, V2], [], []]
[[v], [], [V1]]	[[c, v], [C], [V2]]

### 4.2 Vocalisation

Orthographically, Semitic texts are written without short vowels. It was suggested by (Beesley et al., 1989, et. seq.) and (Kiraz, 1994c) to allow short vowels to be optionally deleted. This, however, puts a constraint on the grammar: no surface expression can contain a vowel, lest the vowel is optionally deleted.

We assume full vocalisation in writing rules. A second set of rules can allow the deletion of vowels. The whole grammar can be taken as the composition of the two grammars: e.g. {cvcvc}, {ktb}, {aa} → /ktab/ → [ktab, ktb].

### 4.3 Morphosyntactic Issues

Finite-state models of two-level morphology implement morphotactics in two ways: using ‘continuation patterns/classes’ (Koskeniemi, 1983; Antworth, 1990; Karttunen, 1993) or unification-based grammars (Bear, 1986; Ritchie et al., 1992). The former fails to provide elegant morphosyntactic parsing for Semitic languages, as will be illustrated in this section.

#### 4.3.1 Stems and $\bar{X}$ -Theory

A pattern, a root and a vocalism do not always produce a free stem which can stand on its own. In Syriac, for example, some verbal forms are bound: they require a **stem morpheme** which indicates the measure in question, e.g. the prefix {ʔa} for *afʔel*

<sup>10</sup>Note, however, that the `expand` command does not insert  $\bar{b}$  randomly in context expressions.

```

tlrule(R7, [[], [], []], [[v], [], [V]], [[c3,b,e], [], []], <=>, [], [], [],
[vowel(V)], [[], [], []]).
tlrule(R8, [], [V1], [C,V2], <=>, [], [], [],
[vowel(V1),vowel(V2),radical(C)], [[], [], []]).

```

Listing 7

```

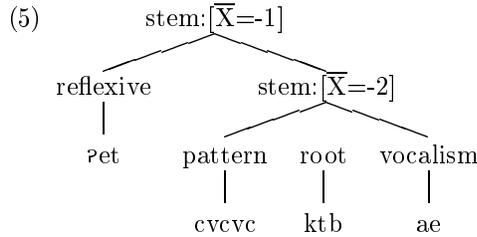
synrule(rule1, stem: [X̄=-2,measure=M,measure=p'al|pa'e1],
[pattern: [], root: [measure=M,measure=p'al|pa'e1],
vocalism: [measure=M,measure=p'al|pa'e1]]).
synrule(rule2, stem: [X̄=-2,measure=M],
[stem_affix: [measure=M],
pattern: [], root: [measure=M], vocalism: [measure=M]]).
synrule(rule3, stem: [X̄=-1,measure=M,mood=act],
[stem: [bar=-2,measure=M,mood=act]]).
synrule(rule4, stem: [X̄=-1,measure=M,mood=pass],
[reflexive: [], stem: [X̄=-2,measure=M,mood=pass]]).
synrule(rule5, stem: [X̄=0,measure=M,mood=MD,npq=s&3&m],
[stem: [X̄=-1,measure=M,mood=MD]]).
synrule(rule6, stem: [X̄=0,measure=M,mood=MD,npq=NPG],
[stem: [X̄=-1,measure=M,mood=MD], vim: [type=suff,circum=no,npq=NPG]]).
synrule(rule7, stem: [X̄=0,measure=M,mood=MD,npq=NPG],
[vim: [type=pref,circum=no,npq=NPG], stem: [X̄=-1,measure=M,mood=MD]]).
synrule(rule8, stem: [X̄=0,measure=M,mood=MD,npq=NPG],
[vim: [type=pref,circum=yes,npq=NPG], stem: [X̄=-1,measure=M,mood=MD],
vim: [type=suff,circum=yes,npq=NPG]]).

```

Listing 8

stems. Additionally, passive forms are marked by the **reflexive morpheme** {ʔet}, while active forms are not marked at all.

This structure of stems can be handled hierarchically using  $\bar{X}$ -theory. A stem whose stem morpheme is known is assigned  $\bar{X}=-2$  (Rules 1-2 in Listing 8). Rules which indicate mood can apply only to stems whose measure has been identified (i.e. they have  $\bar{X}=-2$ ). The resulting stems are assigned  $\bar{X}=-1$  (Rules 3-4 in Listing 8). The parsing of Syriac /ʔetkteb/ (from {ʔet}+/kateb/ after the deletion of /a/ by R4) appears in (5).<sup>11</sup>



Now free stems which may stand on their own can be assigned  $\bar{X}=0$ . However, some stems require

<sup>11</sup>In the remaining examples, it is assumed that the lexicon and two-level rules are expanded to cater for the new material.

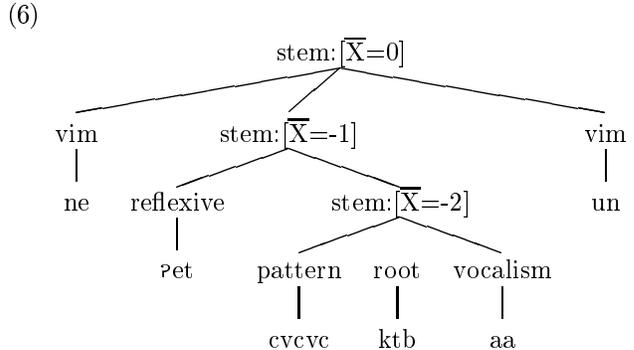
verbal inflectional markers.

### 4.3.2 Verbal Inflectional Markers

With respect to verbal inflectional markers (VIMs), there are various types of Semitic verbs: those which do not require a VIM (e.g. sing. 3rd masc.), and those which require a VIM in the form of a prefix (e.g. perfect), suffix (e.g. some imperfect forms), or circumfix (e.g. other imperfect forms).

Each VIM is lexically marked *inter alia* with two features: 'type' which states whether it is a prefix or a suffix, and 'circum' which denotes whether it is a circumfix. Rules 5-8 (Listing 8) handle this.

The parsing of Syriac /netkatbun/ (from {ne}+{ʔet}+/katab/{un}) appears in (6).



Verb Class	Inflections Analysed	1st Analysis (sec/word)	Subsequent Analysis (sec/word)	Mean (sec/word)
Strong	78	5.053	0.028	2.539
Initial <i>nūn</i>	52	6.756	0.048	3.404
Initial <i>ālah</i>	57	4.379	0.077	2.228
Middle <i>ālah</i>	67	5.107	0.061	2.584
Overall mean	63.5	5.324	0.054	<b>2.689</b>

Table 1

(Beesley et al., 1989) handle this problem by finding a logical expression for the prefix and suffix portions of circumfix morphemes, and use unification to generate only the correct forms – see (Sproat, 1992, p. 158). This approach, however, cannot be used here since, unlike Arabic, not all Syriac VIMs are in the form of circumfixes.

### 4.3.3 Interfacing with a Syntactic Parser

A Semitic ‘word’ (string separated by word boundary) may in fact be a clause or a sentence. Therefore, a morphosyntactic parsing of a ‘word’ may be a (partial) syntactic parsing of a sentence in the form of a (partial) tree. The output of a morphological analyser can be structured in a manner suitable for syntactic processing. Using tree-adjointing grammars (Joshi, 1985) might be a possibility.

## 5 Performance

To test the integrity, robustness and performance of the implementation, a two-level grammar of the most frequent words in the Syriac New Testament was compiled based on the data in (Kiraz, 1994b). The grammar covers most classes of verbal and nominal forms, in addition to prepositions, proper nouns and words of Greek origin. A wider coverage would involve enlarging the lexicon (currently there are 165 entries) and might triple the number of two-level rules (currently there are *c.* 50 rules).

Table 1 provides the results of analysing verbal classes. The test for each class represents analysing most of its inflexions. The test was executed on a Sparc ELC computer.

By constructing a corpus which consists only of the most frequent words, one can estimate the performance of analysing the corpus as follows,

$$P = \frac{5.324n + \sum_{i=1}^n 0.054(f_i - 1)}{\sum_{i=1}^n f_i} \text{ sec/word}$$

where  $n$  is the number of distinct words in the corpus and  $f_i$  is the frequency of occurrence of the  $i$ th word. The SEDRA database (Kiraz, 1994a) provides such data. All occurrences of the 100 most frequent lexemes in their various inflections (a total of 72,240

occurrences) can be analysed at the rate of 16.35 words/sec. (Performance will be less if additional rules are added for larger coverage.)

The results may not seem satisfactory when compared with other prolog implementations of the same formalism (cf. 50 words/sec, in (Carter, 1995)). One should, however, keep in mind the complexity of Syriac morphology. In addition to morphological non-linearity, phonological conditional changes – consonantal and vocalic – occur in all stems, and it is not unusual to have more than five such changes per word. Once developed, a grammar is usually compiled into automata which provides better performance.

## 6 Conclusion

This paper has presented a computational morphology system which is adequate for handling non-linear grammars. We are currently expanding the grammar to cover the whole of New Testament Syriac. One of our future goals is to optimise the prolog implementation for speedy processing and to add debugging facilities along the lines of (Carter, 1995).

For useful results, a Semitic morphological analyser needs to interact with a syntactic parser in order to resolve ambiguities. Most non-vocalised strings give more than one solution, and some inflectional forms are homographs even if fully vocalised (e.g. in Syriac imperfect verbs: sing. 3rd masc. = plural 1st common, and sing. 3rd fem. = sing. 2nd masc.). We mentioned earlier the possibility of using TAGs.

## References

- Aho, A. and Ullman, J. (1977). *Principles of Compiler Design*. Addison-Wesley.
- Antworth, E. (1990). *PC-KIMMO: A two-Level Processor for Morphological Analysis*. Occasional Publications in Academic Computing 16. Summer Institute of Linguistics, Dallas.
- Bear, J. (1986). A morphological recognizer with syntactic and phonological rules. In *COLING-86*, pages 272–6.

- Beesley, K. (1990). Finite-state description of Arabic morphology. In *Proceedings of the Second Cambridge Conference: Bilingual Computing in Arabic and English*.
- Beesley, K. (1991). Computer analysis of Arabic morphology. In Comrie, B. and Eid, M., editors, *Perspectives on Arabic Linguistics III: Papers from the Third Annual Symposium on Arabic Linguistics*. Benjamins, Amsterdam.
- Beesley, K., Buckwalter, T., and Newton, S. (1989). Two-level finite-state analysis of Arabic morphology. In *Proceedings of the Seminar on Bilingual Computing in Arabic and English*. The Literary and Linguistic Computing Centre, Cambridge.
- Bird, S. and Ellison, T. (1994). One-level phonology. *Computational Linguistics*, 20(1):55–90.
- Carter, D. (1995). Rapid development of morphological descriptions for full language processing systems. In *EACL-95*, pages 202–9.
- Goldsmith, J. (1976). *Autosegmental Phonology*. PhD thesis, MIT. Published as *Autosegmental and Metrical Phonology*, Oxford 1990.
- Grimley-Evans, E., Kiraz, G., and Pulman, S. (1996). Compiling a partition-based two-level formalism. In *COLING-96*. Forthcoming.
- Joshi, A. (1985). Tree-adjointing grammars: How much context sensitivity is required to provide reasonable structural descriptions. In Dowty, D., Karttunen, L., and Zwicky, A., editors, *Natural Language Parsing*. Cambridge University Press.
- Karttunen, L. (1983). Kimmo: A general morphological processor. *Texas Linguistic Forum*, 22:165–86.
- Karttunen, L. (1993). Finite-state lexicon compiler. Technical report, Palo Alto Research Center, Xerox Corporation.
- Karttunen, L. and Beesley, K. (1992). Two-level rule compiler. Technical report, Palo Alto Research Center, Xerox Corporation.
- Kataja, L. and Koskeniemi, K. (1988). Finite state description of Semitic morphology. In *COLING-88*, volume 1, pages 313–15.
- Kay, M. (1987). Nonconcatenative finite-state morphology. In *EACL-87*, pages 2–10.
- Kiraz, G. (1994a). Automatic concordance generation of Syriac texts. In Lavenant, R., editor, *VI Symposium Syriacum 1992*, *Orientalia Christiana Analecta* 247, pages 461–75. Pontificio Institutum Studiorum Orientalium.
- Kiraz, G. (1994b). *Lexical Tools to the Syriac New Testament*. JSOT Manuals 7. Sheffield Academic Press.
- Kiraz, G. (1994c). Multi-tape two-level morphology: a case study in Semitic non-linear morphology. In *COLING-94*, volume 1, pages 180–6.
- Kiraz, G. (1995). *Introduction to Syriac Spirantization*. Bar Hebraeus Verlag, The Netherlands.
- Kiraz, G. (1996). *Computational Approach to Non-Linear Morphology*. PhD thesis, University of Cambridge.
- Knuth, D. (1973). *The Art of Computer Programming*, volume 3. Addison-Wesley.
- Kornai, A. (1991). *Formal Phonology*. PhD thesis, Stanford University.
- Koskeniemi, K. (1983). *Two-Level Morphology*. PhD thesis, University of Helsinki.
- Lavie, A., Itai, A., and Ornan, U. (1990). On the applicability of two level morphology to the inflection of Hebrew verbs. In Choueka, Y., editor, *Literary and Linguistic Computing 1988: Proceedings of the 15th International Conference*, pages 246–60.
- McCarthy, J. (1981). A prosodic theory of non-concatenative morphology. *Linguistic Inquiry*, 12(3):373–418.
- Narayanan, A. and Hashem, L. (1993). On abstract finite-state morphology. In *EACL-93*, pages 297–304.
- Pulman, S. and Hepple, M. (1993). A feature-based formalism for two-level phonology: a description and implementation. *Computer Speech and Language*, 7:333–58.
- Ritchie, G., Black, A., Russell, G., and Pulman, S. (1992). *Computational Morphology: Practical Mechanisms for the English Lexicon*. MIT Press, Cambridge Mass.
- Ruessink, H. (1989). Two level formalisms. Technical Report 5, Utrecht Working Papers in NLP.
- Shieber, S. (1986). *An Introduction to Unification-Based Approaches to Grammar*. CSLI Lecture Notes Number 4. Center for the Study of Language and Information, Stanford.
- Sproat, R. (1992). *Morphology and Computation*. MIT Press, Cambridge Mass.
- Wiebe, B. (1992). Modelling autosegmental phonology with multi-tape finite state transducers. Master's thesis, Simon Fraser University.