

Two-level Description of Turkish Morphology

KEMAL OFLAZER

Bilkent University, Turkey

Abstract

This paper describes a full two-level morphological description of Turkish word structures. The description has been implemented using the PC-KIMMO environment and is based on a root word lexicon of about 23,000 root words. The phonetic rules of contemporary Turkish (spoken in Turkey) have been encoded using 22 two-level rules while the morphotactics of the agglutinative word structures have been encoded as finite-state machines for verbal, nominal paradigms and other categories. Almost all the special cases of, and exceptions to phonological and morphological rules have been taken into account. In this paper, we describe the rules and the finite state machines along with examples and a discussion of how various special cases were handled. We also describe some known limitations and problems with this description

1. Introduction

Morphological analysis is an important component in any system for processing natural language. Spelling checkers, parsers, machine translation systems, dictionary tools and many other similar applications use a morphological analyzer to decompose words into their functional components. Morphological analysis is very important in languages such as Finnish or Turkish with agglutinative morphology where the concept of a word is much wider than what one finds in a dictionary. Such languages encode in a single word the equivalent of whole clauses and sentences in languages like English. Thus extensive morphological analysis is required to extract information that is to be used in any further syntactic and semantic analysis. This paper describes a full two-level morphological description (Karttunen, 1983; Koskeniemi, 1983) of Turkish word structures. The description has been implemented using the PC-KIMMO environment (Antworth, 1990) and is based on a root word lexicon of about 23,000 root words.

Turkish is an agglutinative language with word structures formed by productive affixations of derivational and inflectional suffixes to root words.¹ A popular (and rather exaggerated) example of a Turkish word formation is:

OSMANLILAŞTIRAMAYABİLECEKLERİMİZ-DENMİŞSİNİZCESİNE

which can be broken down into morphemes as follows:

OSMAN+LI+LAŞ+TIR+AMA+YABİL+ECEK
+LER+İMİZ+DEN+MİŞ+SİNİZ+CESİNE

where the +'s indicate morpheme boundaries. This adverb can be translated into English as '(behaving) as if you were of those whom we might consider not con-

verting into an Ottoman.' For the details of Turkish grammar and word formations rules one can refer to a number of books (Underhill, 1976; Lewis, 1991).

Turkish has finite-state but nevertheless rather complex morphotactics. Morphemes added to a root word or a stem can convert the word from a nominal to a verbal structure or vice-versa, or can create adverbial constructs as above. The surface realizations of morphological constructions are constrained and modified by a number of phonetic rules. Vowels in the affixed morpheme have to agree with the preceding vowel in certain aspects to achieve *vowel harmony*, although there are a small number of exceptions. Under certain circumstances vowels in the roots and morphemes are deleted. Similarly, consonants in the root words, or in the affixed morphemes undergo certain modifications, and may sometimes be deleted. We will see these when we discuss the two-level rules later. However, the assimilation of a large number of words into the language from various foreign languages—most notably Arabic and Persian—have resulted in word formations which behave as exceptions to many rules. Turkish morphology has been investigated from a computational point of view by Köksal (1975), by Hankamer (1986), and by Solak and Oflazer (1991, 1993).

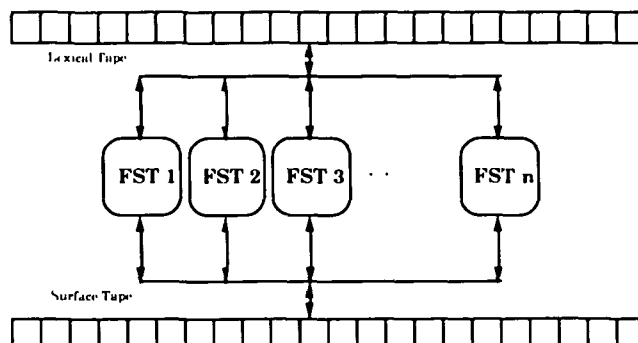


Fig. 1 Parallel acceptors for two-level morphology

2. Two-level morphology

Two-level morphology is a general purpose paradigm for morphological description of word structures (Karttunen, 1983; Koskeniemi, 1983; Antworth, 1990; Sproat, 1992). It has been used to describe the morphology of number of languages (Alam, 1983; Karttunen and Wittenberg, 1983; Khan, 1983; Lun, 1983; Koskeniemi, 1985).

A two-level description is based on a *lexical* and a *surface* representation of a word structure. The lexical level denotes the structure of the functional components of a word while the surface level denotes the standard orthographic realization of the word with the given lexical structure. The phonetic restrictions and

Correspondence: Kemal Oflazer, Department of Computer Engineering and Information Science, Bilkent University, Bilkent, Ankara 06533, Turkey. E-mail: ko@cs.bilkent.edu.tr

modifications are represented using four different rule types (Sproat, 1992):

- (1) $a:b \Rightarrow LC_RC$. This is the context restriction rule which states that a lexical **a** can be realized as a surface **b** only in the given left context (LC) and right context (RC), but not necessarily always.
- (2) $a:b \Leftarrow LC_RC$. This is the surface coercion rule which states that a lexical **a** has to be realized as a surface **b** always in the given context but not necessarily only in that context.
- (3) $a:b \Leftrightarrow LC_RC$. This is the composite rule which states that a lexical **a** has to correspond to a surface **b** in the given context and the correspondence is valid only in that context.
- (4) $a:b / \Leftarrow LC_RC$. This is the exclusion rule which states that a lexical **a** can never correspond to a surface **b** in the given context.

These rules are compiled into finite-state acceptors which check a lexical to surface correspondence in parallel. Any lexical to surface correspondence that is not rejected by any of these machines is assumed to be a valid correspondence and hence accepted. Fig. 1 shows the basic architecture of a two-level system.

The morphotactics which determine the proper sequencing of the morphemes are encoded as finite state machines using lexicons for root words and suffixes, and alternations for capturing the sequencing of suffixes. In the following sections we assume that the reader is familiar with the basic concepts of two-level morphology—for a good exposition to two-level morphology one can refer to Sproat (1992).

3. Two-level description of Turkish morphology

The Turkish language uses an alphabet of 29 letters in its current orthography using Latin characters. There are 8 vowels: **a, e, ı, i, o, ö, u, ü**, and 21 consonants: **b, c, ç, d, f, g, ğ, h, j, k, l, m, n, p, r, s, ş, t, v, y, z**. Tables 1 and 2 from van der Hulst and van der Weijer (1991) show the phonetic features that correspond to the sounds denoted by these letters.

There are however sounds not covered by these. Certain long vowels are mainly used in words borrowed from foreign languages most notably Arabic and Persian. Some vowels are sometimes distinguished in older

Table 1 Turkish vowels

	+Round		-Round	
	+Front	-Front	+Front	-Front
+High	u	ı	i	ı
-High	ö	o	e	a

Table 2 Turkish consonants

	Labial	Labio-dental	Dental	Palato-alveolar	Palatal	Velar	Glottal
Voiceless Stop	p		t	ç		k	
Voiced Stop	b		d	c		g	
Voiceless Fricative		f	s	ş			
Voiced Fricative		v	z	j			
Nasal	m		n				
Liquid			l,r				
Approximant					y		h

orthography by various means (such as with a $\hat{}$ on top of the vowel). In modern orthography such distinctions are not commonly used. There is also a certain phoneme known as 'yumuşak g' (soft g—denoted as ğ in orthography) which creates bisyllabic two-vowel sequences. At the end of a syllable, this phoneme causes the lengthening of the preceding vowel (van der Hulst and van der Weijer, 1991).

Consonants **k, g, and l** have palatal and non-palatal allophones. In certain cases the palatalization process has impact on the vowel harmony. Our purpose in this two-level model is not necessarily to account for all phonetic phenomena from a purely phonetic viewpoint, but to have a description which can deal with all observed surface forms for use in practical applications.

For the purposes of two-level specification, we assume an alphabet that includes the letters above and certain additional symbols (that will be described shortly) only used in the two-level description, and not in the orthography.³ First we define the following subsets:

- (1) Consonants: $C = \{b, c, ç, d, f, g, ğ, h, j, k, l, m, n, p, r, s, ş, t, v, y, z\}$
- (2) Surface vowels subject to ellipsis under certain cases: $V_s = \{ı, i, o, ö, u, ü\}$
- (3) Lexical vowels: $V = \{a, e, ı, i, o, ö, u, ü, H, A, â, ô, û\}$. These are the vowels used in the lexical level.
- (4) Back vowels: $V_b = \{a, ı, o, u\}$
- (5) Front vowels: $V_f = \{e, i, ö, ü\}$
- (6) Front unrounded vowels $V_{fu} = \{e, i\}$
- (7) Front rounded vowels $V_{fr} = \{ö, ü\}$
- (8) Back unrounded vowels $V_{bu} = \{a, ı\}$
- (9) Back rounded vowels $V_{br} = \{o, u\}$
- (10) Lexical consonants used as the first letter in a suffix but which may disappear on the surface under certain conditions $X = \{s, y, n\}$
- (11) Lexical consonants used as the first letter in a suffix but which are always realized on the surface $C_r = \{S, l, c, D\}$

Before going any further, it is necessary to point out a number of points about Turkish orthography and our two-level representation:

- (1) In the 3rd item above, **H** used at the lexical level only stands for the set of high vowels $\{ı, i, u, ü\}$, unresolved for the other features, **A** again used as the lexical level stands for nonhigh unrounded vowels $\{a, e\}$. **â, û** stand for long vowels which are present lexically but are not used in the orthography. These are realized as **a, u** on the surface. **ô** is the vowel in words of foreign origin such as *alkol* (alcohol) or *gol* (goal) where it is

always followed by an *l*. This vowel behaves like an *ö* in the vowel harmony process.

- (2) Proper nouns are separated from suffixes by an apostrophe ('). All vowel harmony rules and *some* of the consonant change rules are in effect in the orthography of proper nouns.
- (3) Some roots have vowels which are deleted when certain suffixes are affixed. Such vowels are prefixed with a \$ in the lexicon representations.
- (4) **D** is a lexical symbol denoting dental stop consonants which has default surface realization by the voiced **d** but may also be realized on the surface as the voiceless **t** under certain circumstances.
- (5) **S** represents a lexical **s** which never gets deleted during affixation. For example in the suffix +**sHz**, the **s** never gets deleted, while in the suffix +**sH**, the **s** may sometimes be deleted. We represent the **s** in the former by **S**, whose surface realization is always **s**.
- (6) **K** represents a root-final lexical **k** which never becomes a surface **ğ** during any affixation, e.g., in **erk** (power). Its surface realization is **k**.
- (7) The rule compiler KGEN heavily used in this development treats a regular expression of the sort **a*b*c*** as **(a*b*c*)***. This fact is used in the vowel harmony rules.

3.1 Two-level rules

The following are the two-level rules for the phonetic component of the description:

- (1) **A:a** ⇒ **V:V_b' : * C* @:0* + :0* _**
- (2) **A:e** ⇒ **[V:V_f|â:a|û:u|ô:o] ' : * @:0* + :0* _**
- (3) **H:u** ⇒ **[V:V_{br}|V_{br}:0 + :0] ' : * C* + :0* @:0* _**
- (4) **H:ü** ⇒ **[V:V_{fr}|V_{fr}:0 + :0|û:u|ô:o] ' : * C* + :0* @:0* _**
- (5) **H:i** ⇒ **[V:V_{bu}|V_{bu}:0 + :0] ' : * C* + :0* @:0* _**
- (6) **H:i** ⇒ **[V:V_{fu}|V_{fu}:0 + :0|â:a] ' : * C* + :0* @:0* _**

The first two rules force the agreement of an **A** vowel to a preceding vowel in the backness attribute. The last four rules force the agreement of an **H** vowel to a preceding one in backness and roundedness.

- (7) **H:0** ⇒ **V (' :) + :0 _**
An **H** vowel is deleted if the last phoneme of the stem it is being affixed to is a vowel. For example:

Lexical: masa+Hm N(table)+1PS-POSS
Surface: mas00m masam

- (8) **H:0 / ⇐ V:0 + :0_y o r**
However an **H** vowel is not deleted in the verbal suffix +**Hyor** denoting the progressive tense (present continuous) if the last phoneme of the stem corresponds to a vowel which gets deleted under this circumstance. For example:

Lexical: kapa+Hyor V(close)+PR-CON+3PS
Surface: kap00ıyor kapıyor

- (9) **A:0 ⇐ _ + :0 H:@ y o r**
The cases covered by this rule are the vowel ellipsis in certain verbal stems when +**Hyor** denoting the progressive tense is added to a stem ending in a vowel. An example is:

Lexical: ilaç+1A+Hyor N(drug)+
NtoV(1a)+PR-CON+3PS
Surface: ilaç0100ıyor ilaçlıyor

- (10) **V_s :0 ⇐ \$:0 _ C + :0 (X:0) [A:@|H:@] _ + :0 H:@ y o r**

A vowel in the lexical representation of certain roots will have to be deleted on the surface due to a vowel ellipsis phenomenon. Instead of adding these vowels to the alphabet we have indicated their context with a prefix \$ which eventually gets deleted on the surface. An example is:

Lexical: bur\$un+Hm H(nose)+1PS-POSS
Surface: bur00n0um burnum

This rule also deals with the vowel ellipsis for verbal roots ending in a vowel (cf. Rule 9 above). See the example for Rule 8 which also contains an example of the application of this rule.

- (11) **X:0 ⇐ C (' :) + :0 (C) V**

This rule deletes the beginning **s**, **n**, or a **y** of a suffix when it gets affixed to a stem ending in a consonant.

Lexical: ev+nHn N(house)+GEN
Surface: ev00ın evin

Lexical: kalem+sH N(pencil)+3PS-POSS
Surface: kalem00ı kalemi

Lexical: ağ \$1z+yH N(mouth)+ACC
Surface: ağ00z00ı ağızı

- (12) **D:t ⇐ [h|@:ç|ş|@:k|@:p|@:t|f|s] (' :) + :0 (@:0) _**

This rule realizes a **D** as a **t** whenever it is preceded by one of the consonants in the option list across a morpheme boundary. **D** is otherwise realized as a **d** by default. Some examples are:

Lexical: kitab+DA N(book)+LOC
Surface: kitab0ta kitapta

Lexical: yulaf+dAn N(oat)+ABL
Surface: yulaf0tan yulaftan

Lexical: aç+DHk V(open)+PERF
Surface: aç0tak açtık

- (13) **{b, d} : {p, t} ⇒ _ #|_ + :0 (X:0) [C|c:]**

This rule realizes voiced obstruents **b**, **d** as **p**, **t** respectively either when they end a word or when they are followed by a morpheme beginning with a consonant. Some examples are:

Lexical: kitab+1Ar N(book)+PLU
Surface: kitab0lar kitaplara

Lexical: kitab+cH N(book)+NtoN(c1)
Surface: kitab0ç1 kitapçı

Lexical: dolab+nHn N(closet)+GEN
Surface: dolab00ın dolabın

Lexical: tad+DHk V(taste)+PERF
Surface: tat0tak tattık

- (14) **c:ç ⇐ [@:ç|ş|@:k|@:p|@:t|f|s] + :0 [H:@|A:@]**

- (15) **ç:c ⇐ _ + :0 X:0 V**

c is another voiced obstruent like those above except that it also appears in certain suffixes as the first letter where it gets modified to a **ç**.⁴ Some examples are:

Lexical: haraç+cH N(tribute)+NtoN(ci)
 Surface: haraçOç1 haraçci
 Lexical: yaş+cA N(age)+NtoAdv(ca)
 Surface: yaşOça yaşça
 Lexical: haraç+yA N(tribute)+DAT
 Surface: haracOOa haraca

(16) **k:ğ** ⇒ **__ + :0 (X:0) V**

The last **k** sound at the end of a word becomes a **ğ** when a morpheme starting with a vowel is affixed. Some examples are:

Lexical: ayak+nHn N(foot)+GEN
 Surface: ayakOOin ayakın
 Lexical: tarak+Hm N(comb)+1PS-POSS
 Surface: tarakOım taracağım

(17) **k:g** ⇒ **n__ + :0 (X:0) V**

However under certain circumstances the same **k** changes to a **g** as in:

Lexical: renk+yH N(color)+ACC
 Surface: rengOOi rengi
 Lexical: ahenk+yA N(harmony)+DAT
 Surface: ahengOOe ahenge

(18) **g:ğ** ⇒ **__ + :0 (X:0) V**

The last **g** at the end of words (of foreign origin) also becomes a **ğ** when certain suffixes are added. Examples are:

Lexical: radyolog+yA N(radiologist)+DAT
 Surface: radyologOOa radyoloğa

(19) **g:ğ** / ⇐ **[n|r]__**

However under the circumstances above if the final **g** is preceded by another consonant (only **n** and **r** seem to be such consonants) then the **g** does not become a **ğ** when suffixes are added. Examples are:

Lexical: brifing+Hm N(briefing)+1PS-POSS
 Surface: brifingOım brifingim
 Lexical: aysberg+HnHz N(iceberg)+2PP-POSS
 Surface: aysbergOıniz aysberginiz

(20) **Y:y** ⇔ **__ + :0 [X:0|H:@]**

(21) **Y:0** ⇔ **__#|__ + :0 C,**

These two rules deal with a large number of nominal roots ending with *su* (water) e.g., *akarsu* (river). *Su*, along with *ne* (what) does not obey the standard inflection rules. For example **su+sH** (water +3PS-POSS) is **suyu** and not **susu** and **su+nHn** (water+GEN) is **suyun** and not **sunun**.⁵ We have added a lexical **Y** to such words so that further inflections can proceed properly but the **Y** is realized as **0** at the end of a word or if followed by a consonant which never drops in affixation. Here are some examples:

Lexical: akarsuY+yHnHz N(river)+2PP-POSS
 Surface: akarsuyOOunuz akarsuyunuz
 Lexical: akarsuY+lar N(river)+PLU
 Surface: akarsuOOlar akarsular

(22) **' :0** ⇔ **__ [#|1A:@]**

Finally we have a clean-up rule to remove the ' from a lexical proper noun under certain circumstances.

3.2 Root word lexicons

Our word list is based on the word list that we have compiled for our spelling checker (Solak and Oflazer, 1993). This list is pretty comprehensive and is annotated with categorical and exceptions information. It however contains a very small number of technical words and other domain specific jargon. This word list has been divided into a number of lexicons for:

- (1) *Nouns*,
- (2) *Adjectives*: This contains all nominal roots whose main usage is adjectival. They have the same morphotactics as the nouns. Current implementation has about 18,500 nominal (noun + adjective) roots.
- (3) *Verbs*: This lexicon contains about 2,450 verbal roots.
- (4) *Compound nouns*: This lexicon contains the two-compound nouns that are considered as a single lexical unit.
- (5) *Proper nouns*: This lexicon contains the proper nouns. All the roots in this lexicon end with a ' which gets deleted on the surface when necessary.
- (6) *Pronouns*,
- (7) *Adverbs*,
- (8) *Connectives*,
- (9) *Exclamations*,
- (10) *Postpositions*,
- (11) *Acronyms*: Most of our intended applications that will use this description as a kernel, have to deal with real text which contains many acronyms which nevertheless have their own suffixation schemes that we have to deal with.
- (12) *Special cases*: This contains nominal word structures which exhibit a number of special cases.

3.3 Finite state machines for morphotactics

The morphotactics of a language determine the structure and the ordering of morphemes. In agglutinative languages the morphemes are affixed to a root word like 'beads on a string' (Sprout, 1992). Thus words in such languages formed by such a scheme can correspond to whole sentences or clause in languages like English (see the example in the introduction). Turkish morphotactics allow productive formation of words whose part-of-speech categories may change a number of times during affixation. One can start with a nominal root, then form a verbal form with a suffix which can then take an aspect suffix and then become a nominal form again through for example a gerund suffix, and then take the standard nominal suffixes (plural, possessive case, etc.). It is also possible to have circular constructions (an example of which is given later). This however does not mean that there are no restrictions on such formations. In fact there are semantic restrictions on the formations, that is, one can technically form the word except it may not make much or any sense. It is possible to enforce such restrictions in morphotactics except the mechanisms one would need would have been much more sophisticated than the simple provisions provided by the most morphological analyzers. We will discuss this in a later section.

The two-level approach to computational morphology has concentrated mainly on the phonetic component

and providing very simple finite-state mechanisms for describing the ordering of morphemes. While this is sufficient and convenient for a morphotactically simple languages, descriptions of the morphotactics of languages like Turkish may become unwieldy unless the descriptive power (not necessarily the computational power) are significantly enhanced. Nevertheless, our description implements the morphotactics of the Turkish language with the provided mechanisms and we later indicate the limitations of our descriptions.

In PC-KIMMO, the morphotactics component is described by root and suffix lexicons which are linked to each other. In the figures describing our morphotactic component (such as Fig. 3), the boxes indicate suffixation states, the arrows indicate the next states which can be reached when a suffix matching one of the labels is found. The circles indicate the final states for complete and valid word formations with the labels in parentheses near these states labeled **End** indicate the class of the word construction when the machine ends up in that final state. The **O** on the transitions indicate that the transition can be taken with *null* input. The states drawn in bold correspond to references to states in other figures. For example, the state labeled **Possessive-3** indicates the state of a nominal construction which has been affixed a third person possessive suffix. Form that state one can go to a final state indicating a nominal in accusative case with suffix +**nH**, or to the states labeled **Case-1** or **Case-2** with the relevant case suffixes, or to another final state with the suffix +**cA**.

The morphotactic description for PC-KIMMO consists of lexicons of the sort shown in Fig. 2. Each lexicon entry consists of a pattern, a pointer (or a set of pointers) to the next lexicon to be processed and the gloss to be output in case of a successful match.

LEXICON NOUNS

ab	POST-NOUN	"N(ab)"
abad	POST-NOUN	"N(abat)"
abadan	POST-NOUN	"N(abadan)"
abadi	POST-NOUN	"N(abadi)"
abajur	POST-NOUN	"N(abajur)"
...		
...		
zürriyet	POST-NOUN	"N(zürriyet)"
züyuf	POST-NOUN	"N(züyuf)"

LEXICON POST-NOUN

+lHk	POST-NOUN	"NtoN(lik)"
+cH	POST-NOUN	"NtoN(ci)"
+cHk	POST-NOUN	"NtoN(cik)"
+lAş	POST-VERB	"NtoV(laş)"
+lAr	PLURAL	"PLU"
+lArH	POSSESSIVE-3	"3PP-POSS"
+lArH	POSSESSIVE-3	"PLU+3PP-POSS"
O	PLURAL	""

Fig. 2 Example of PC-KIMMO lexicons for Turkish morphotactics

3.3.1 Finite state machine for nominal morphotactics.

Figure 3 shows the finite state machine for the nominal paradigm. We would like to point out one caveat: These finite state machines have been designed with the assumption that the input (for recognition) is always legal and the recognizer will then chop it up into the proper sequence of morphemes. This assumption makes the finite state machines a bit simpler.

The morphotactics for the nominal paradigm is relatively simple. There are mainly two parts: The top part corresponds to nominal constructions with plural, possessive, case and relativization suffixes. It is technically possible to go around the loop through the state labeled **Relative** a number of times though in practice such constructions are rarely used. For example it is possible to have a word structure like:

MASA+LAR+IM+DA+Kİ+LER+İN+Kİ+NDE

which roughly means 'at those (things) which belong to those (other things) at my tables.'

The bottom part of the nominal morphotactics state diagram corresponds to the nominal verb and adverbial constructions like:

- (1) *evdeydi*—(S/he/it) was at the house.
- (2) *evdeyse*—If (s/he/it) is at the house.
- (3) *evdeymiş*—(s/he/it) was at the house. (Narrative)
- (4) *evdeyim*—I am at the house.
- (5) *evdedirler*—They are (definitely) at the house.
- (6) *evdeyken*—While (someone) is (was) at the house.
- (7) *evdeymişcesine*—(behaving) as if he is at the house

The nominal morphotactics are a bit different for compound nouns. The additional states required by these compound nouns are shown in Fig. 4. Compound nouns in Turkish which are treated as single lexical unit have almost always two component both of which are nominal roots. Thus Turkish does not have a productive compounding paradigm such as in German. The second component in such compound nouns is always affixed compound marker, which is the same as the third person possessive suffix, when the compound noun is used in the nominative case. For example *bitpazarı* (flea market) (Lexical **bit**+**pazar**+**sH**), is used as both the nominative form and the third person possessive form. However further affixation does not proceed as in other nominals. For example the plural of *bitpazarı* is *bitpazarları* where the plural suffix is now affixed to the nominative form of the second part of the compound and then the third person possessive is added. Similarly in *bitpazarım* (my flea market) or *bitpazarın* (your flea market) the affixation is onto the nominative form of the second component and not on to the nominative form of the compound noun. In our lexicon we represent such compound nouns lexically without the third person suffix and then add this suffix on the surface when only the nominal form is used. Otherwise, we treat it as a normal nominal root.

Within the noun lexicon there are a number of roots which are already in plural form. For those cases the plural affix and/or the possessive suffixes are skipped in the morphotactics. For example:

- (1) *amcamlar* (the family/home of my uncle):⁶ This is already in plural form and does not take any possessive suffix either. Hence the suffix lexicon that follows this is the **CASE-1** lexicon.
- (2) *bakliyat* (legumes), *baklagiller* (leguminous plants) are already in plural form.

For nouns already in plural form and ending in **+IAr**, the possessive suffix **+sH**, can be interpreted as both the third person singular possessive or third person plural possessive.

3.3.2 *Finite state machine for verbal morphotactics.* Figures 5 and 6 show the finite state machine for the verbal paradigm. The verbal morphotactics is significantly more complicated than the nominal morphotactics. Turkish verbal structures can take a sequence of reflexive⁷ or reciprocal, causative and passive suffixes which can then be followed by a compound verb, and then by aspect, tense and person suffixes. Verbal structures can also be made into nominal or adverbial structures with the addition of yet other suffixes. When a verbal root takes no reflexive or reciprocal suffix, the

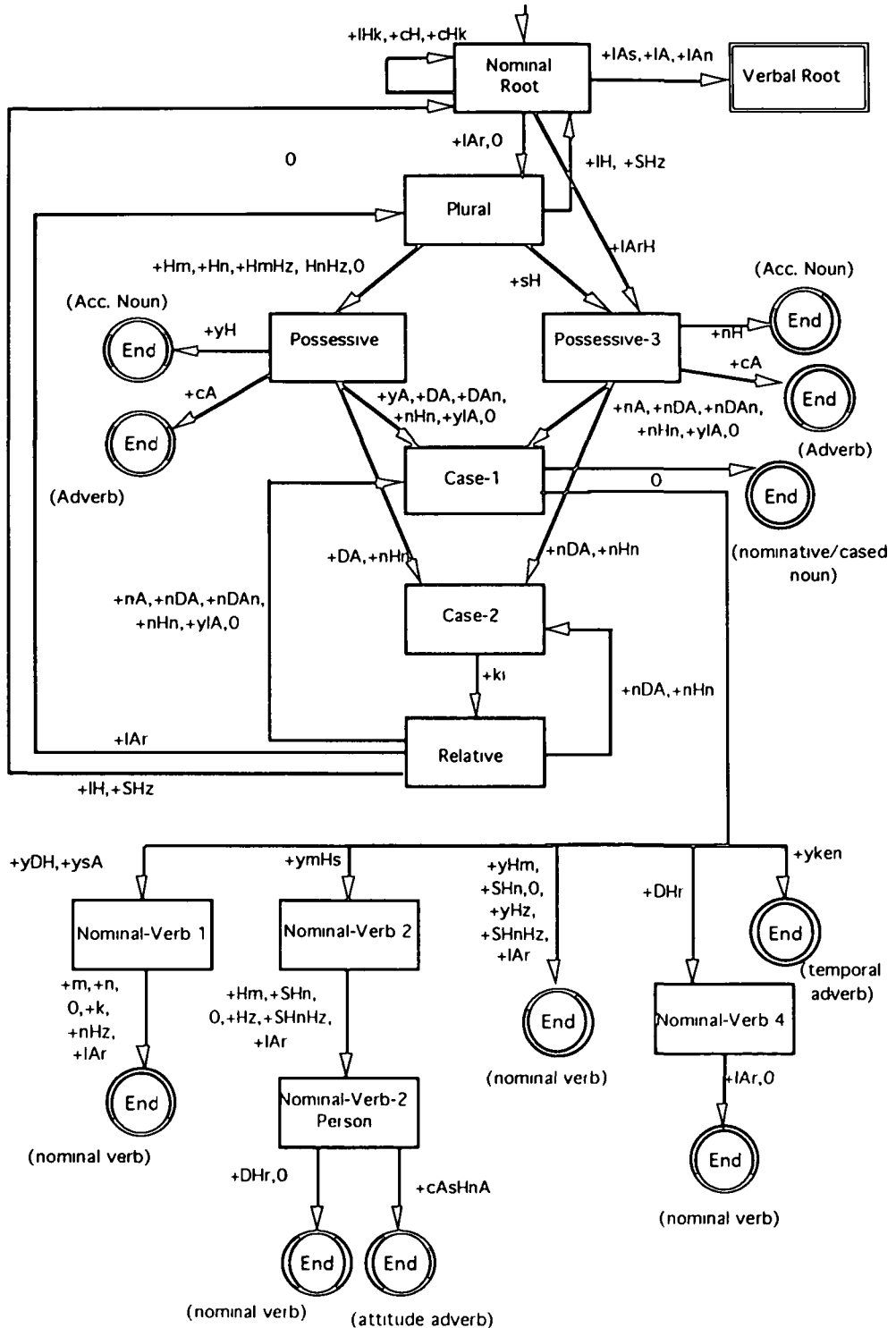


Fig. 3 Finite state machine for nominal morphotactics

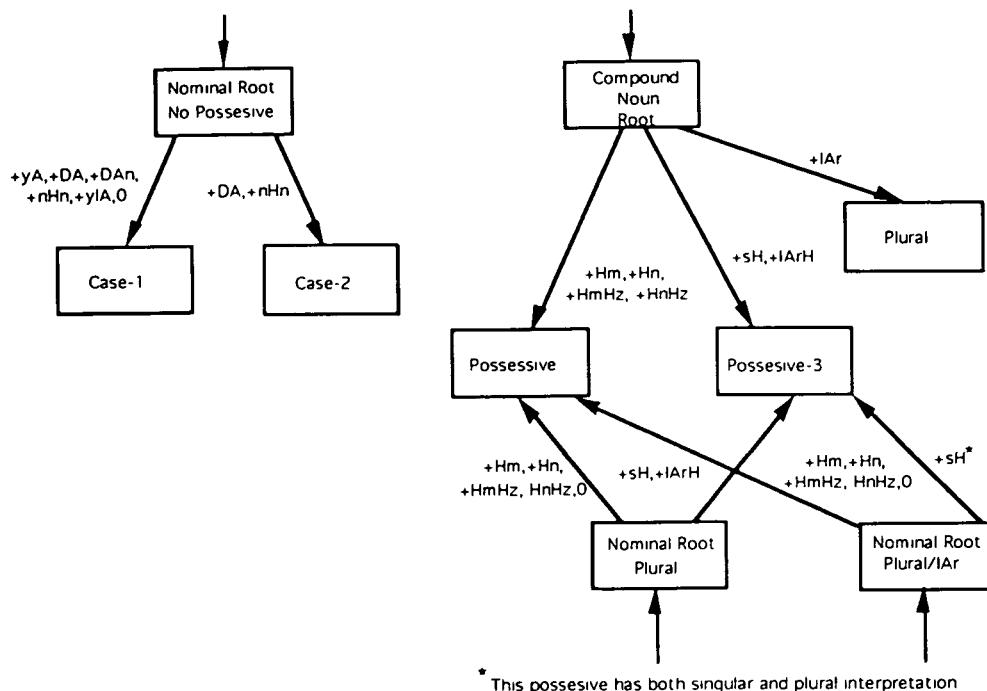


Fig. 4 Finite state machine for compound noun morphotactics

causative or the passive suffixes can take a variety of forms depending on a number of criteria on the roots.⁸ If, however, they take either of the reflexive or the reciprocal suffixes (which are mutually exclusive), then the causative and passive formations are very simple as shown on Fig. 5. After state labeled **Passive Hn** which corresponds to a verbal stem with all the reflexive/reciprocal, causative, and passive suffixes accounted for, we can construct a negative form by +mA and +yAmA or directly go into positive verb construction. In any case, we can possibly add from a small number of auxiliary (or compound) verbs (the most common being +yAbil indicating potentiality) to get a verbal stem to which we can now add tense and person suffixes, or suffixes which form nominal structures, infinitives and adverbs.

Turkish verbs can have at most two suffixes indicating aspect and tense. The first one can be one of *narrative, future, aorist, present continuous, necessitative, optative, imperative, perfect* and *conditional* suffixes. These can take possibly different sets of person suffixes to form a verbal structure, or take a second morpheme indicating *perfect, conditional* or *narrative*. As can be seen from the morphotactics diagram, not all possible combinations of the aspect and tense suffixes are possible. The second set of suffixes will only be allowed if the first suffix is one of narrative, future, aorist, present continuous and necessitative. There are a number of nonstandard cases especially involving the third plural person and these are accounted for in the state diagrams.

An example will clarify the general idea behind verbal constructions. Consider the verb: *görülemiyormuşum* which can be translated into English as ('it is said that) I was not able to be seen.' The morpheme structure is:

gör +Hl +yAmA +Hyor +ymHş +yHm
 gör +ül +OemO +tyor +Omuş +Oum
 see +PASS +NEG +PRES-CONT +NARR +IPS

This verbal root *gör* will generate the structure above by going through the states labeled.

- (1) **Verbal Root** (root)
- (2) **Passive Hl** with +Hl
- (3) **Passive Hn** with 0
- (4) **Negative yama** with +yAmA
- (5) **Verbal Stem** with 0
- (6) **Other Tense** with +Hyor
- (7) **Second Tense Other** with +ymHş
- (8) **End** with +yHM

Readers familiar with details of verb formation in Turkish will note that our morphotactic model does not deal with the three groups of a total of 13 verbal roots whose aorist forms are exceptions to the rules.⁹ Based on the caveat expressed above, our model can correctly identify the morphemes given correctly formed input but may not occasionally reject ill-formed aorist forms. We have opted to live with this as the alternative would have been to add another orthogonal categorization to the already large number of verb root groups and then replicate the remaining state machine three times to account for the aorist special cases. It is also possible to deal with special entries in the lexicon but this would not have prevented acceptance of ill-formed input. We expect to add this refinement in our model as we start using our two-level in applications like corpus tagging.

3.4 Special cases and exceptions

There are also a number of exceptions to most of the rules for word formation. We have opted to deal with these as special root lexicon entries. These exceptions are typically of the following sorts:

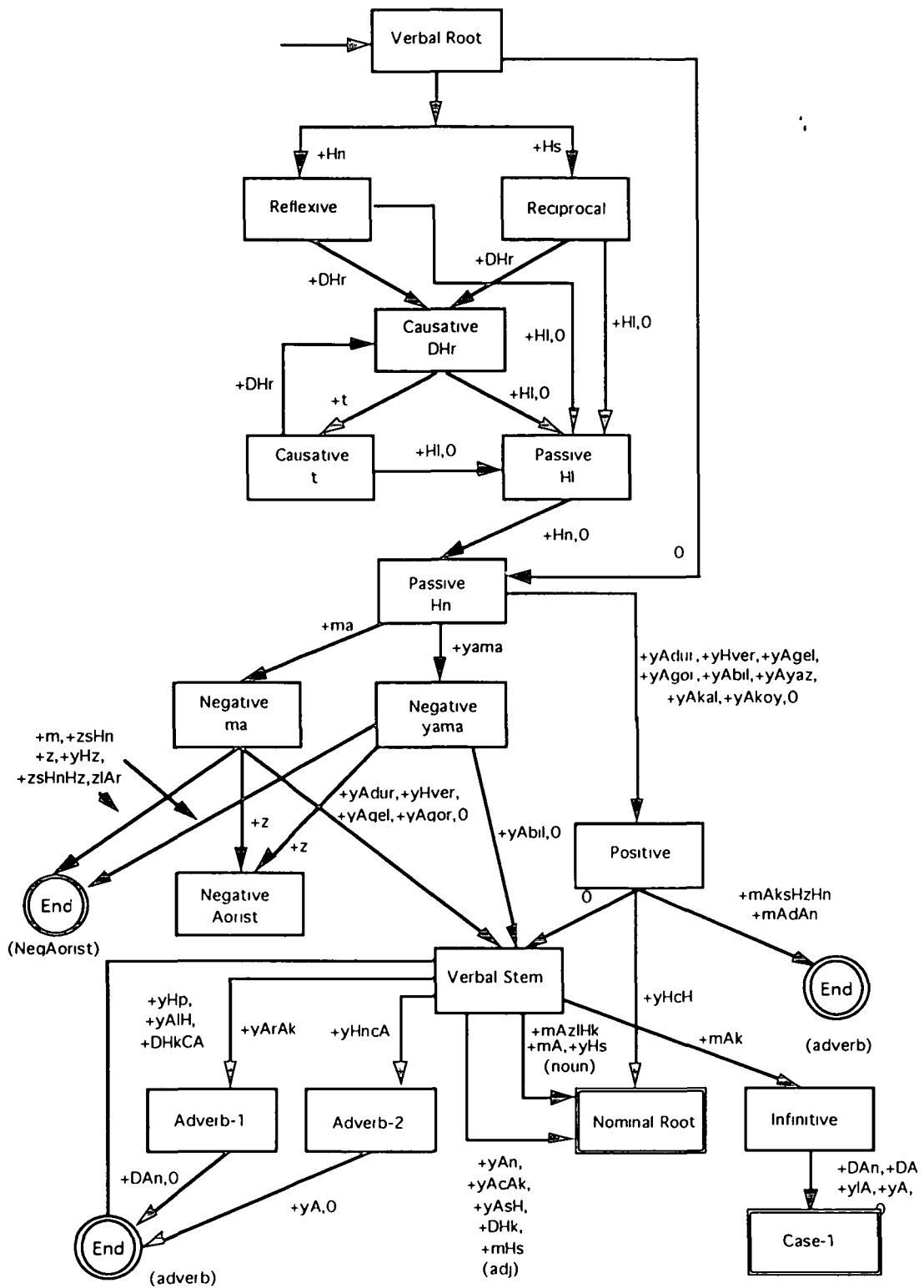


Fig. 5 Finite state machine for verbal morphotactics

(1) Pronominal or adjectival formations indicating temporal and spatial relative positions can be directly formed with the +**ki** suffix to a nominal root denoting a time or place. For example:

- (1) *önceki*: (the one) before,
- (2) *yarunki*: (the one) tomorrow,
- (3) *bugünkü*: (the one) today,¹⁰
- (4) *karşiki*: (the one) across from here.

(2) Third person singular forms of certain roots of

Arabic origin ending with a vowel do not have an intervening **s**. For example:

- (1) *sanayii* as opposed to *sanayisi*: industry+3PS-POSS
- (2) *mevkii* as opposed to *mevkisi*: place+3PS-POSS

(3) There are adverbs formed from roots indicating numbers by the suffixes +**şAr** or +**Ar**, which are only applicable to such roots. For example:

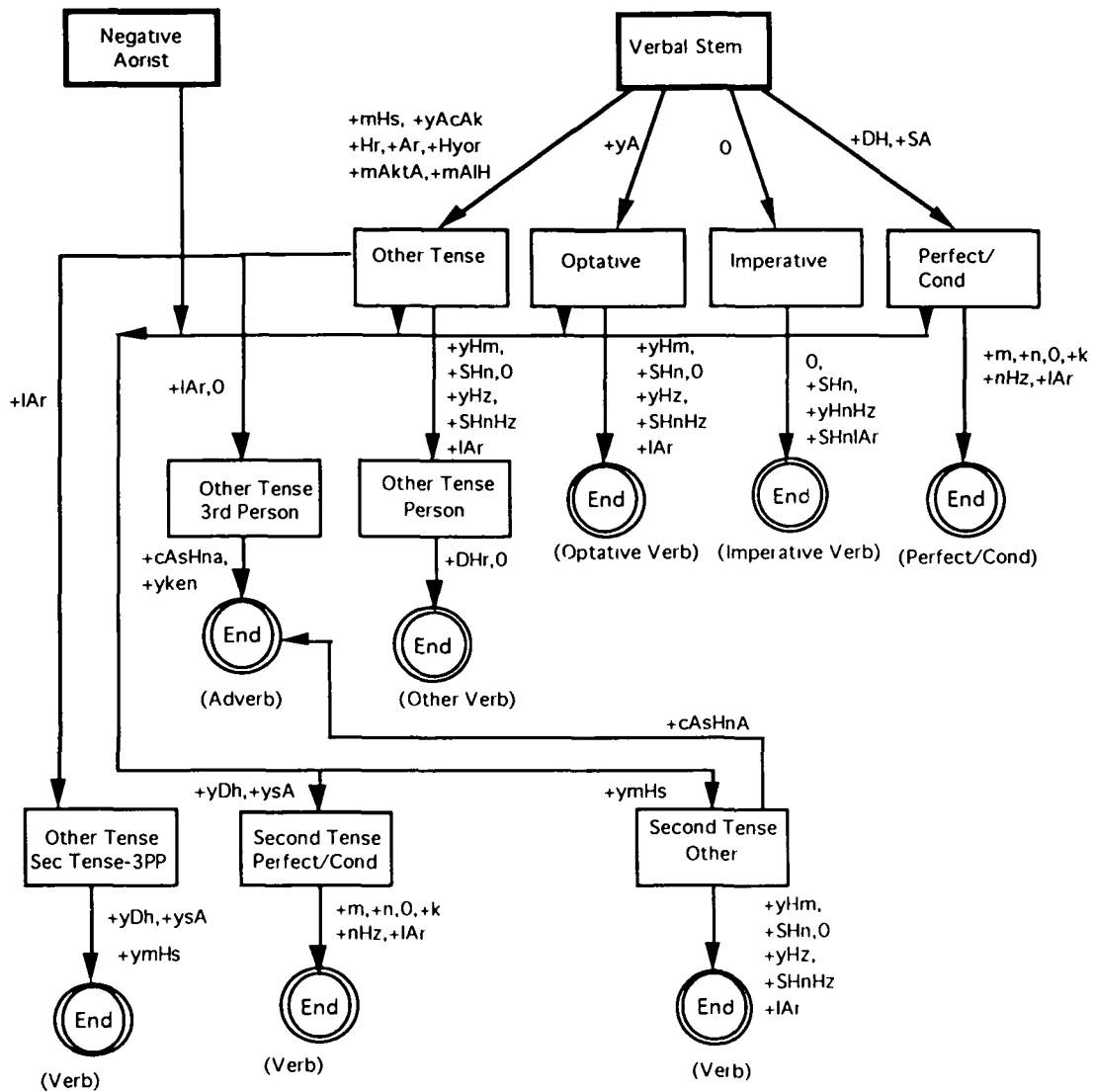


Fig. 6 Finite state machine for verbal morphotactics (cont.)

- (1) *birer*: one each,
- (2) *ikişer*: two each,
- (3) *ellişer*: fifty each.
- (4) There are nominal roots ending with consonant which is duplicated (possibly after going through a modification) before certain suffixes. For example:
 - (1) *hakki* (**hak**+**yH**): right+ACC (but *haklar* is **hak**+**PLU**)
 - (2) *tibba* (**tip**+**yA**): medicine+DAT
- (5) For certain polysemous nominal roots, the vowel ellipsis phenomenon (dealt with \$ above) does not occur when usage is for one of the meanings, but occurs for the other usage. For example:
 - (1) **metin**+**yH** (text+ACC) is *metni* but
 - (2) **metin**+**yH** (strong+ACC) is *metini*.
- (6) There are certain polysemous roots, where the last consonant is duplicated before certain suffixes when used with one of the meanings and is not duplicated with the other meaning.
 - (1) **şık**+**yH** (chic+ACC) is *şıkı* but
 - (2) **şık**+**yH** (item+ACC) is *şıkkı*.

In Turkish question suffixes starting with **mH** are written as a separate word, but the lexical **H** has to harmonize with the last vowel of the preceding word. For example:

geldiler mi?: Did they come?
görmüşsün you saw (it), **görmüş müsün?**: Did you see (it)? but not **görmüşsün mü?**

Since PC-KIMMO does not deal with sequences of words for recognition, this implementation does not deal with this harmonization nor with whether the prior morphotactics is valid if it is followed by the question part. It may be possible to deal with this by introducing a dummy space character which gets incorporated into the rules and gets inserted between the preceding word and the question part when it is detected.¹¹ In addition to this construct there are a large number of multi-word constructs in Turkish grammar where the grammatical role of the words involved do not have much relevance to the grammatical role of the multi-word construct. We are currently working on an intermediate level as part of a project of part-of-speech tagging of Turkish text (Ofazer, 1993) which sits upon the morphological processor and handles multi-word and idiomatic con-

structs. This level will both provide a more accurate analysis of text and reduce the load on syntactic and semantic analyzers (e.g. Güngördü and Oflazer, 1993) that now have to do this handling at their level. Examples of cases that can be handled at that stage are:

- (1) double optative and 3SG verbal constructions like *koşa koşa* which actually are used as adverbs,
- (2) question suffix formations, discussed above,
- (3) aorist verbal constructions like *yapar yapmaz* which function as temporal adverbs,
- (4) emphatic adjectival forms involving the question suffix, such as *güzel mi güzel*,
- (5) word sequences with idiomatic usage such as *yanı sıra*,
- (6) various multiple word proper names.

4. Examples

In this section we give a number of examples from our PC-KIMMO implementation.¹² The examples follow the following format:

Input	Morpheme Structure	Gloss English meaning
çalışmanın	çalış+mA+Hn+nHn	[V(çalış)+VtoN(ma)+2PS-POSS+GEN] <i>of your work(ing)</i>
çalış+mA+nHn		[V(çalış)+VtoN(ma)+GEN] <i>of the work(ing)</i>
çocukları	çocuk+IAr+sH	[N(çocuk)+PLU+3PS-POSS] <i>his/her children</i>
çocuk+IAr+yH		[N(çocuk)+PLU+ACC] <i>children (accusative)</i>
çocuk+IArH		[N(çocuk)+3PP-POSS] <i>their child</i>
çocuk+IArH		[N(çocuk)+PLU+3PP-POSS] <i>their children</i>
ödüllendirilmiş	ödül+IAAn+DHR+HI+mHş	[N(ödül)+NtoV(lan)+CAUS+PASS+VtoAdj(mis)] <i>rewarded (person)</i>
ödül+IAAn+DHR+HI+mHş		[N(ödül)+NtoV(lan)+CAUS+PASS+NARR+3PS] <i>(s/he has been) rewarded</i>
alınmış	al+Hn+ymHş	[N(al)+2PS-POSS+NtoV()+NARR+3PS] <i>(it) was your red (one)</i>
al+nHn+ymHş		[N(al)+GEN+NtoV()+NARR+3PS] <i>(it) belongs to the red (one)</i>
al\$ın+ymHş		[N(alın)+NtoV()+NARR+3PS] <i>(it) was a forehead</i>
al+Hn+mHş		[V(al)+PASS+VtoAdj(mis)] <i>(a) taken (object)</i>
al+Hn+mHş		[V(al)+PASS+NARR+3PS] <i>it was taken</i>

alın+mHş **[V(alın)+VtoAdj(mis)]**

(an) offended (person)

alın+mHş

[V(alın)+NARR+3PS]

s/he was offended

boynu

boy\$un+sH

[N(boyun)+3PS-POSS]

(his/her) neck

boy\$un+yH

[N(boyun)+ACC]

neck (accusative)

5. Using the two-level description of Turkish in natural language processing applications

Our main reason for developing this description was to use it as a the morphological analysis component in a number of applications in the context of a large scale project on natural language processing in Turkish. Some of these applications are:

- (1) Parsing Turkish with a lexical functional grammar parser for Turkish (Güngördü 1993, Güngördü and Oflazer, 1993),
- (2) Parsing Turkish with a augmented transition network grammar (ATN) (Demir, 1993),
- (3) Part-of-speech tagging of Turkish text (Oflazer, 1993).

6. Suggestion for a more powerful morphotactic descriptions for computational morphology

We have noted a certain difficulty in the coding of the morphotactic component of this description. As the mechanisms for describing state transitions are rather basic, the size of the finite state machines become unwieldy when portions of the machines have to be duplicated to deal with minor variations. For example, in our verbal morphotactics for the verbal structures, we have split the set of verbal roots into a number of groups depending on the kinds of causative and passive suffixes they can take and we have deliberately opted not to deal with the aorist suffix in this way. The distinctions among these groups are typically based on features like whether the root lexicon entry ends with a vowel or not, or ends with a certain consonant or not, or whether the root is polysyllabic, etc. For example, one of the groups corresponds to those verbal roots (taking the causative suffix +DHR) which end with a consonant other than *l*, so that when they do not take the causative suffix, the passive is formed by *HI*. If one could say *bind a variable* to the lexicon names and entry tokens matched as a transition is being taken, and refer to it in a later transition through an additional condition, structure of the finite state machines could be simpler. For example, in the state diagrams for the verbal structures, the aorist suffix is dealt in the transition coming out the state **Verbal Stem** with suffixes +*Hr* or +*Ar*. The +*Ar* should be taken only by those mono-syllabic roots ending with a consonant (except for a small finite set of roots) and those compound verbal roots formed with +*et*.¹³ The +*Hr* is taken by all multi-syllabic stems ending with a consonant except those formed by +*et*, and those mono-syllabic ones in that finite set above.¹⁴ It would be advanta-

geous to refer to certain properties of the variables bound, while taking a transition. Suppose we have a lexicon entry structure such as

```
token {COND C} alternation gloss {BIND V}
```

where {COND C} and {BIND V} are optional transition condition and the variable to bind to the matching token respectively (cf. Fig. 2). Then we could simply write the aorist transition as

LEXICON AORIST

```
+Ar COND Check-AR-Aorist POST-AORIST
      "+AOR" BIND Aorist-Suffix
+Hr COND Check-HR-Aorist POST-AORIST
      "+AOR" BIND Aorist-Suffix
```

For instance we could write the first condition **Check-AR-Aorist** as (expressed informally here):

```
if Causative-Suffix is NULL and Passive-Suffix is
   NULL and Negative-Suffix is NULL, and
   [Root has a single vowel (hence mono-syllabic) and
    Root ends with a consonant and
    Root Lexicon is not VerbExceptions1] or
    Root ends with +et
then return TRUE else return FALSE
```

with boldface indicating variable bound during recognition. Thus if this condition is checked when a suffix matching **+Ar** is found we would have a much more accurate morphotactics component with much simpler descriptions. The mechanism that we are suggesting are essentially the same as those provided by augmented transition networks with registers, minus the recursive network calls. Although we have not yet investigated the influence this would have on the overall formalism and performance of the two-level morphology approach, we have reason to believe that for languages like Turkish such additional functionality would simplify constructions of very accurate and maintainable morphotactic descriptions with simple machines.

7. Conclusions

In this paper we have presented a two-level description of Turkish morphology that we have implemented in the PC-KIMMO environment. We have used 22 two-level rules and the lexicon is based on a root word list of about 23,000 words. Being the first large scale two-level description of Turkish, this description can be used in applications like morphosemantic analysis, corpus tagging, sentence level parsing etc. From a computational speed viewpoint, PC-KIMMO is rather slow as it is a general purpose environment. Our PC-KIMMO implementation can process an input word in about one half second on a SparcStation ELC. On the other hand, a highly optimized parser that we have developed for spelling checking Turkish (Solak and Oflazer, 1993) is about 1000 times faster,¹⁵ but then ours is not a general purpose tool. This description of Turkish morphology is publicly available via ftp¹⁶ and also via the European Corpus Initiative collection.

8. Acknowledgements

I would like to thank Dr Lauri Karttunen of XEROX-PARC for encouraging the implementation of this description, and to Dr Evan Antworth of Summer Institute of Linguistics for providing the PC-KIMMO environment and numerous discussions. I would also like to thank Aysin Solak for constructing the very original word list on which this description is based, and to Zelal Güngördü and Coşkun Demir, who while struggling through their theses found many bugs in this description and helped me improve it tremendously.

Notes

1. Turkish is an exclusively suffixing language. There are however a few very unproductive prefixes of foreign origin, such as *na-* (*un-*).
2. Familiarity with Turkish morphology, and the two-level paradigm will certainly make this exposition clearer.
3. The PC-KIMMO implementation which is now in public domain uses different characters for the symbols used below as the ASCII character set does not support certain characters in the Turkish alphabet.
4. Normally words in Turkish do not end with **c**. Our original approach used this in the lexical representation. However we were not able to encode the reciprocal assimilation process that is observed when a suffix beginning with a **c** is affixed to a root ending in **c** in which both **c**'s change to **ç**'s by mutual influence.
5. For **ne** the normal inflections are also valid.
6. Note that this looks like it has a possessive suffix (**+Hm**) followed by the plural suffix (**+IAr**). However morphotactics puts the possessive after the plural, hence this can not be parsed as such within the normal paradigm.
7. The reflexive suffix is not very productive.
8. See (Solak, 1991) for the details of these criteria.
9. This corresponds to the whether **+Hr** or **+Ar** is to be selected as the aorist suffix from the state labeled **Verbal Stem**. See (Solak, 1991) for details.
10. Note that the relative suffix conforms to vowel harmony to his case.
11. Current PC-KIMMO interface would not let you do this but it is possible to embed the recognition capability in a stand alone program that does the question form checking and combining input strings for a recheck.
12. These outputs have been edited for proper character orthography, and then annotated.
13. These are explicitly listed in the verbal lexicon.
14. For stems ending in a vowel, the suffix is **+r** which is obtained when the vowel is either of the aorist suffixes drop.
15. Although due to the nature of the process, it will return after the first successful parse.
16. Site and Directory: <ftp.bilkent.edu.tr/pub/Turklang>

References

- Alam, Y. S. (1983). 'A two-level morphological analysis of Japanese', *Texas Linguistic Forum*, 22, 229–52.
- Antworth, E. L. (1990). 'PC-KIMMO: A two-level processor for Morphological Analysis', Summer Institute of Linguistics, Dallas, Texas.
- Demir, C. (1993). 'An ATN grammar for Turkish', Master's thesis, Department of Computer Engineering and Information Sciences, Bilkent University, Ankara, Turkey.
- Güngördü, Z. (1993). 'A lexical-functional grammar for

- Turkish', Master's thesis, Department of Computer Engineering and Information Sciences, Bilkent University, Ankara, Turkey
- Gungördü, Z. and Ofazer, K. (1993) 'Parsing Turkish with the Lexical Functional Grammar formalism' To appear in *Proceeding of COLING 94, the 15th International Conference on Computational Linguistics*.
- Hankamer, J. (1986). 'Finite state morphology and left to right phonology'. In *Proceedings of the West Coast Conference on Formal Linguistics*, 5. Stanford University.
- Karttunen, L. (1983). 'KIMMO: a general morphological processor', *Texas Linguistic Forum*, 22, 163-86.
- Karttunen, L. and Wittenburg, K. (1983). 'A two-level morphological analysis of English', *Texas Linguistics Forum*, 22, 217-28.
- Khan, R. (1983). 'A two-level morphological analysis of Rumanian', *Texas Linguistic Forum*, 22, 253-70.
- Köksal, A. (1973). 'Automatic Morphological Analysis of Turkish', PhD thesis, Hacettepe University, Ankara, Turkey.
- Koskeniemi, K. (1983). 'Two-level morphology: A general computational model for word form recognition and production', Publication no: 11, Department of General Linguistics, University of Helsinki, Finland.
- Koskeniemi, K. (1985). 'An application of the two-level model to Finnish' In F. Karlsson (ed.), *Computational morphosyntax: a report on research 1981-1984*, University of Helsinki, Department of General Linguistics, Helsinki, Finland
- Lewis, G. L. (1991). *Turkish Grammar*. Oxford: Oxford University Press.
- Lun, S. (1983) 'A two-level morphological analysis of French', *Texas Linguistic Forum*, 22, 271-8.
- Ofazer, K. (1993). 'A text tagger for Turkish' In *Proceedings of Second Turkish Symposium on Artificial Intelligence and Neural Networks*. Bosphorus University, Turkey.
- Solak, A. (1991). 'Design and implementation of a spelling checker for Turkish', Master's thesis, Bilkent University, Dept of Computer Engineering and Information Science, Ankara, Turkey.
- Solak, A. and Ofazer, K. (1992). 'Parsing agglutinative word structures and its application to spelling checking for Turkish'. In *Proceedings of the 15th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, 1, pp. 39-45
- Solak, A. and Ofazer, K. (1993). 'Design and implementation of a spelling checker for Turkish', *Linguistic and Literary Computing*.
- Sproat, R. (1992). *Morphology and Computation*. MIT Press.
- Underhill, R. (1976). *Turkish Grammar*. MIT Press.
- van der Hulst, H. and van der Weijer, J. (1991). 'Topics in Turkish phonology' In H. Boeschoten and L. Verhoeven (eds.), *Turkish Linguistics Today*. E. J. Brill