

A Pronunciation Lexicon for Turkish Based on Two-Level Morphology

Kemal Oflazer^a, Sharon Inkelas^b

^aFaculty of Engineering and Natural Sciences
Sabancı University,
Istanbul, Turkey
oflazer@sabanciuniv.edu

^bDepartment of Linguistics
University of California,
Berkeley, CA, USA
inkelas@socrates.berkeley.edu

Abstract

This paper describes the implementation of a full-scale pronunciation lexicon for Turkish based on a two-level morphological analyzer. The system produces at its output, a parallel representation of the pronunciation and the morphological analysis of the word form so that morphological disambiguation can be used to disambiguate pronunciation when necessary. The pronunciation representation is based on the SAMPA standard and also encodes the position of the primary stress. The computation of the position of the primary stress depends on an interplay of any exceptional stress in root words and stress properties of certain morphemes, and requires that a full morphological analysis be done. The system has been implemented using XRCE Finite State Toolkit.

1. Introduction

Pronunciation lexicons are computational devices that map the graphemic representation of words to a representation of their pronunciation. They are a valuable resource in annotating speech data used in training automatic speech recognizers, and in generating accurate speech in text-to-speech systems. Finite state transducers are commonly used in building pronunciation lexicons for mapping between orthographic strings and phonological strings, either as an efficient encoding of direct mapping, or a mapping involving some kind of morphological processing [1]. In this paper, we present the design and implementation of a full scale finite state pronunciation lexicon of Turkish, an agglutinating language with extremely productive inflectional and derivational morphology and an essentially infinite lexicon. The agglutinating nature of the language implies that any corpus based compilation of a word list for use in speech applications will be rather inadequate [2]. Noun roots typically have a few hundred inflected forms – many more if one considers productive derivations; verbs have a few million forms. This necessitates that one employ a generative model that is able to recognize all possible words in the language, and base the pronunciation lexicon on this to avoid a significant out-of-vocabulary word problem. The lexicon implemented as a finite state transducer takes as input a word form and produces all possible pronunciations of the word, paired with the corresponding morphological analyses. The pronunciations are encoded using the SAMPA encoding, and also include the marking of the position of the primary stress.¹ The dependence of the stress computation on the proper identification of morphemes and their morphotactical function requires that the pronunciation lexicon

be built on top of a morphological analyzer. Although Turkish superficially seems to have an almost one-to-one mapping between graphemics and pronunciation, there are quite a number of subtle phenomena especially in loan words. Furthermore there are a number of additional processes in morphophonology such as suffixation induced exceptional vowel lengthening and palatalizations in the surface realizations of certain suffixes depending on vowel harmony resolution. Such phenomena are not distinguished in orthography but have to be handled while representing pronunciation.

2. Turkish

Turkish is an Ural-Altaic language, having agglutinative word structures with productive inflectional and derivational processes. Turkish word forms consist of morphemes concatenated to a root morpheme or to other morphemes, much like “beads on a string.” The morphotactics of word forms can be quite complex when multiple derivations are involved.

Turkish has a 8-vowel inventory which is symmetrical around the axes of backness, roundness, and height: /i, y, e, ɛ, a, o, ɯ, u/ which correspond to *i, ü, e, ö, a, o, ı, u* in Turkish orthography. Suffix vowels typically harmonize in backness, and (if high) in roundness to the preceding stem vowel (compare e.g., *ev+ler /evler/* ‘house-plural’ to *at+lar /atlar/* ‘horse-plural’), although there are several suffixes whose vowels do not harmonize.²

Turkish has 26 consonants: /p, t, tʃ, k, c, b, d, dʒ, g, ɟ, f, s, ʃ, v, w, z, ʒ, m, n, N, l, ʃ, r, j, h, ʁ/. /G/ represents the velar fricative or glide corresponding to the historical voiced velar fricative that was lost in Standard Turkish; but we treat it as a consonant for the purposes of this work. On the other hand, orthography uses only 21 letters for consonants: /g/ and its palatal counterpart /ɟ/ are written as *g*, while /k/ and its palatal counterpart /c/ are written as *k, ç* and /l/ are written as *l, ll, w* are written as *v* and /n/ and its nasal counterpart /N/ are written as *n*. Palatalized segments (/ɟ, c, l/) contrast with their nonpalatalized counterparts only in the vicinity of back vowels (thus *sol* is pronounced /soʃ/ when used to mean ‘left’ vs. /sol/ when used to mean ‘note in scale’). In the neighborhood of front vowels, palatality is predictable (*lig /lig/* ‘league’).

Turkish has lexical stress: each word has exactly one primary-stressed syllable.³ Some roots are lexically stressed⁴ and certain suffixes are prestressing, meaning if not overridden, they will stress the preceding *syllable*. A word composed of

¹See <http://www.phon.ucl.ac.uk/home/sampa/turkish.htm>. We use the SAMPA notation to show pronunciations in the text, where necessary.

²We use – to denote syllable boundaries and + to denote morpheme boundaries wherever appropriate.

³The existence of secondary stress is controversial.

⁴We call this *exceptional stress*.

only unstressed free and bound morphemes exhibits the default stress pattern where the last syllable is stressed. In a word with stressed root and/or prestressing suffixes, the stress of the leftmost such morpheme will prevail (see e.g., Inkelas [3] for an overview.)

In place names and foreign names used in Turkish, a different default pattern is used, termed here the “Sezer” stress pattern, in tribute to its description in Sezer [4]. For such words the antepenultimate syllable is stressed when it is heavy (meaning containing a long vowel or ending in a consonant) and the penultimate is light (meaning ending in a short vowel), e.g., /‘an-ka-ra/; otherwise stress is penultimate (e.g., /is-‘tan-bul/, /a-‘da-na/. The Sezer pattern can be imposed on any word if used as a place name (thus *kandil+li* /kan-dil-‘li/ ‘oil lamp-with’, but *Kandilli* /kan-‘dil-li/ (same word used as place name).

3. Computational Considerations

The problem of grapheme-to-phoneme mapping for Turkish is simpler than for languages such as English or French. Orthography more or less maps one-to-one to pronunciation yet there are quite a number of cases where orthography is ambiguous. These cases usually stem from the fact that a loan word (usually from Arabic, Persian or French) is a homograph of another Turkish word. The once used accent marks to mark the distinctions are not longer consistently used, if at all, and one is left to rely on the context for inferring the correct pronunciation. The other major source of pronunciation ambiguity is the location of primary stress in the word. Certain morphemes which are homographs (but marking different morphosyntactic features) may appear in the same inflection paradigm in morphotactics but have different stress marking properties. For instance, the *+ma/+me* suffix in the verb paradigm will (depending its position in the suffix sequence) either mark negative polarity or an infinitive derivation. So, a word form like *okuma* would either mean the imperative ‘don’t read’ or the infinitive ‘to read/reading’.⁵ In the imperative reading, the stress will be on the *syllable* preceding the *+ma* suffix, while in the other reading the suffix is neutral and stress is on the last syllable. Thus, morphological analysis is necessary to determine the morpheme structure which along with any stress markers in the root morpheme, then determines the location of the stress. On the other hand, in an application context such as text-to-speech, the appropriate pronunciation of the word has to be determined by a morphological disambiguation and/or word sense disambiguation process. For instance, morphological disambiguation of the readings of the word *okuma* would be necessary to select the appropriate pronunciation in a context, while a process akin to word sense disambiguation would be necessary to disambiguate the appropriate pronunciation of the word *sol* mentioned earlier.⁶

It is however necessary to generate the representations of pronunciation and morphological analysis in a paired and parallel fashion. Generating them separately and independently would not be of much use. It would not be possible to associate a given pronunciation with an analysis as this in general is an *n-to-m* mapping. These considerations have prompted us to build the pronunciation lexicon on the scaffolding provided by the wide coverage morphological analyzer for Turkish [5] that we have built using XRCE finite state tools [6, 7]

⁵In addition to a nominal reading *ok+um+a*, ‘arrow-1sg possessive-dative’ meaning ‘to my arrow’ which has the same pronunciation as the infinitive reading.

⁶Though the two are not senses of a word in the lexicographic sense.

4. The architecture of the pronunciation lexicon

The word pronunciation lexicon transducer is the composition of a series of transducers that transform a surface form into all possible and ambiguous parallel representations of its pronunciation and morphological features. The overall internal structure of the transducer is shown in Figure 1. All the boxes in this figure are finite state transducers, and in implementation, they are all composed at compile-time to give one (very large) transducer. Let us now describe the function each of these transducers in detail.

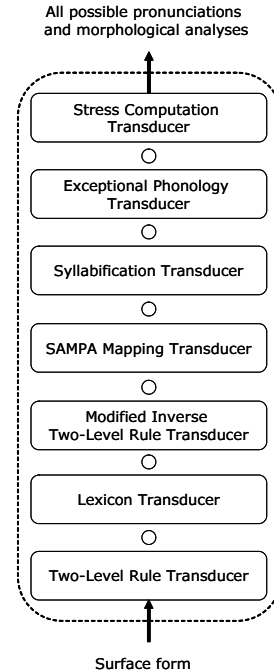


Figure 1: The structure of the word pronunciation lexicon transducer

The *Two-level Rule Transducer* at the bottom (Figure 1) implements the morphographemic mapping described by a set of parallel two-level rules [8]. It is the intersection of the transducers for about 35 morphographemic rules that capture the morphographemics of Turkish [5]. It maps the surface representation of a word into possible lexical representations.⁷ For example the surface form *karın* would map to five lexical forms:⁸

Lexical Form	Gloss
kar+Hn	your snow/profit
kar+[n]Hn	of the snow/profit
kar+[y]Hn	mix (it)!
karı+[H]n	your wife
karın	belly

The next transducer is the *Lexicon Transducer* comprising the root and the affix lexicons. In addition to the standard ordering of the suffix lexicons in the inflectional and deriva-

⁷Though such lexical representations do not necessarily make distinctions not required by morphographemic processes.

⁸H represents a lexical meta-phoneme that denotes a high vowel unresolved for frontness and roundness. [. .] denotes segments of morphemes that are deleted on the surface

tional paradigms, it also comprises a couple hundred finite state constraints motivated by morphosyntactic and semantic concerns. These constraints impose fine grained tuning on word structures and significantly tame the overgeneration of the paradigm-based morpheme lexicon ordering. In the context of the pronunciation lexicon the *Lexicon Transducer* has two main functions: (i) it produces all possible morphosyntactic feature representations of the free and bound morphemes, (ii) it replicates the lexical form at its output possibly augmented with any stress markers induced by *specific prestressing bound morphemes*. Continuing the example earlier, the five lexical forms above would map to the following at the output of this transducer.⁹

Input Lexical Form → Output of the Lexicon Transducer

kar+Hn → (kar) kar+Noun+A3sg (+Hn) +P2sg+Nom
 kar+nHn →
 (kar) kar+Noun+A3sg+Pnon (+nHn) +Gen
 kar+yHn →
 (kar) kar+Verb+Pos (+#p#yHn) +Imp+A2sg
 kar1+Hn →
 (kar1) kar1+Noun+A3sg (+Hn) +P2sg+Nom
 kar1n → (kar1n) kar1n+Noun+A3sg+Pnon+Nom

A couple of remarks are in order here: (i) the morphological analysis process is essentially complete. We have an interleaved representation of lexical morphemes (between (...)’s) and the morphosyntactic features they map to. The concatenation of the contents in pairs of parentheses comprise the original lexical form with possible addition of certain stress markers, while the rest when concatenated gives the morphological analysis; (ii) the morpheme +yHn in morphotactics is a prestressing morpheme (hence the marker #p#), that is, it may eventually cause the primary stress to appear in the *syllable* before this marker.

From this point upwards in the structure, we carry the morphological features around, manipulating the lexical representation sandwiched between the (...) to generate the representation of the corresponding pronunciation.

The *Inverse Two-level Transducer* is essentially (but not exactly) the inverse of the *Two-level Transducer*. We have the same set of rules and a slightly different set of (inverse) feasible pairs (different for a variety of technical reasons.) The only difference in the rules is that the context regular expressions of the two-level rules are modified to ignore the boundary symbols (,), the stress markers, and the feature symbols outside the parentheses. The function of this transducer is to map the lexical form back to the *surface morphemes*, the concatenation of which will give the original surface form (plus any stress markers.) So, for the five outputs above we will get the following:

Output of the Inverse Two-level Rule Transducer

(kar) kar+Noun+A3sg (1n) +P2sg+Nom
 (kar) kar+Noun+A3sg+Pnon (1n) +Gen
 (kar) kar+Verb+Pos (#p#1n) +Imp+A2sg
 (kar1) kar1+Noun+A3sg (n) +P2sg+Nom
 (kar1n) kar1n+Noun+A3sg+Pnon+Nom

Now that we have the original surface form of the word, the

⁹Note that some default feature are produced with zero morphemes. The morphological features used in the paper, other than the obvious POSs are: Imp: imperative mood, +P2sg: 2nd person possessive agreement, +A1sg: 1st person singular agreement, +A2sg: 2nd person singular agreement, +A3sg: 3rd person singular agreement, +Pnon: No possessive agreement, +Nom: Nominative case, +Gen: Genitive case, +Pos: Positive Polarity, +Neg: Negative Polarity, +Become: Become verb, +Caus: Causative verb, +Prog1: Progressive aspect, +Past: Past tense ^DB denotes a derivation boundary.

surface morphemes can now be mapped to the representation of their pronunciations by the *SAMPA Mapping Transducer*. This transducer employs a separate lexicon which maps from the root words and a subset of their morphosyntactic features into a representation of their pronunciation in the SAMPA standard of phonetic encoding. The root words explicitly encoded here are those whose pronunciation (including stress) can not be handled by a default mapping. All other root words and all bound morphemes are handled by a grapheme level default mapping lexicon. The mapping of the roots involves a number of issues: (i) any Sezer/exceptional stress in root words are generated during this mapping; (ii) homograph roots with different pronunciations and/or different stress locations have to be handled by this transducer (e.g., the word kar when used to mean *profit* would be pronounced as /car/ but when used to mean *snow* (or when used as a verb), would be pronounced as /kar/); (iii) certain roots (such as compounds) which go through morphographemic changes during affixation have to be listed here by explicitly modeling any pronunciation changes if they can not be handled by the default mapping. Thus the outputs above now become:

Output of the SAMPA Mapping Transducer (Gloss)

(car) kar+Noun+A3sg (1n) +P2sg+Nom (your profit)
 (kar) kar+Noun+A3sg (1n) +P2sg+Nom (your snow)
 (car) kar+Noun+A3sg+Pnon (1n) +Gen (of the profit)
 (kar) kar+Noun+A3sg+Pnon (1n) +Gen (of the snow)
 (kar) kar+Verb+Pos (#p#1n) +Imp+A2sg (mix (it)!)
 (kar1) kar1+Noun+A3sg (n) +P2sg+Nom (your wife)
 (kar1n) kar1n+Noun+A3sg+Pnon+Nom (belly)

The concatenation of the segments within (...) comprise the preliminary SAMPA encoding of the word’s encoding. A number of minor phenomena have to be fixed in the SAMPA representation. These are handled by the exceptional phonology transducer later.

The *Syllabification Transducer* performs the syllabification of the complete SAMPA encoding. This involves breaking up (by inserting a -) the consonant cluster between any two vowels into the coda of the left syllable and the onset of the right syllable. The implementation of this consists of the composition of a series of finite state transducers for inserting the syllable boundary marker in the right context. The outputs of the *Syllabification Transducer* for the cases listed earlier would be follows.

Output of the Syllabification Transducer

(ca-r) kar+Noun+A3sg (1n) +P2sg+Nom
 (ka-r) kar+Noun+A3sg (1n) +P2sg+Nom
 (ca-r) kar+Noun+A3sg+Pnon (1n) +Gen
 (ka-r) kar+Noun+A3sg+Pnon (1n) +Gen
 (ka-r) kar+Verb+Pos (+#p#1n) +Imp+A2sg
 (ka-r1) kar1+Noun+A3sg (n) +P2sg+Nom
 (ka-r1n) kar1+Noun+A3sg+Pnon+Nom

The *Exceptional Phonology Transducer* handles a set of phenomena for certain exceptional roots and morphemes. The most important of these is the handling of the lexically long vowel in selected roots. For such roots, the last vowel in the root will be coded as a long vowel when it is part of an open syllable and as a short vowel when it is part of a closed syllable. It turns out that *kar* (/car/) is one of those root words (but not *kar* (/kar/)), and in the examples above, the root is followed by the surface morpheme /1n/ which forces the last vowel in the root to an open syllable /ca-/. So /a/ is lengthened and for those we get.¹⁰

¹⁰Since the original morphological analyzer did not make vowel length distinctions (as they were not needed while handling the mor-

(ca:-r)kar+Noun+A3sg(1n)+P2sg+Nom
 (ca:-r)kar+Noun+A3sg+Pnon(1n)+Gen

A second common phenomenon handled by this transducer is the palatalization of /5/ to /l/ and of /k/ to /c/ in certain suffixes when they are followed or preceded by front vowels /e, i, y/.¹¹ Since this is not a morphographemic process, the underlying morphological analyzer is oblivious to this in earlier stages. The output of this transducer has the correct SAMPA encoding except that the stress location is not yet determined.

The *Stress Computation Transducer* computes the location of the primary stress by examining any Sezer/exceptional stress markings in the roots and prestressing markers in the morphemes. The stress is determined in a series of steps:

1. All prestressing markers (#p#) interspersed in the SAMPA representation, that have a prestressing marker or a Sezer/exceptional stress marker somewhere on their left, are removed, that is, the left-most marker “wins.” For example, the representation at this point of the surface word *taşlaştıramıyorduk* (‘we were not being able to petrify (them)’) would be

(taS)taş+Noun+A3sg+Pnon+Nom
 (-laS)ˆDB+Verb+Become
 (-t1-r)ˆDB+Verb+Caus
 (a-#p#m)ˆDB+Verb+Able+Neg
 (1-j#p#or)+Prog1
 (-#p#du)+Past(k)+A1p1

where three morphemes have prestressing markers. This step deletes all such prestressing marker except the first one (in the surface morpheme (a-#p#m)).

2. If there is a Sezer/exceptional root stress marker left then that is also the location of the final stress. So for example the word *pencerede* (on the window) would have the output (pen-dZˆe-re)pencerede+Noun+A3sg+Pnon(-de)+Loc
3. If there is a prestressing marker (which should be the only marker left at this point), then the vowel of the preceding syllable receives the stress. Thus the word above in item 1 would have the final representation of its pronunciation taS-laS-t1-rˆa-m1-jor-duk.
4. If there are no stress markers at this point, then the final stress mark is inserted just before the vowel of the last syllable (all except the second from the last in the examples below).

For representational purposes the final stress mark is then moved to the preceding syllable boundary.¹² For the examples that we have been tracing all along, the final outputs will be:

Output of the Stress Computation Transducer

(ca:-"r)kar+Noun+A3sg(1n)+P2sg+Nom
 (ka-"r)kar+Noun+A3sg(1n)+P2sg+Nom
 (ca:-"r)kar+Noun+A3sg+Pnon(1n)+Gen
 (ka-"r)kar+Noun+A3sg+Pnon(1n)+Gen
 ("ka-r)kar+Verb+Pos(1n)+Imp+A2sg
 (ka-"r1)kar1+Noun+A3sg(n)+P2sg+Nom
 (ka-"r1n)kar1n+Noun+A3sg+Pnon+Nom

phographemic process), the process has to be handled here as a vowel length adjustment. The end result from a computational point of view does not change.

¹¹/2/ (ö in orthography) does not appear in any suffixes.

¹²The only exceptions to these rules are forms of the question clitic and the emphasis clitic which do not bear any primary stress.

The complete pronunciation model is implemented using the XRCE finite state tools, *xfst*, *lexc*, and *twolc*. Apart from various lexicons making the morphological analyzer, the complete system is described by about 500 regular expressions. The resulting transducer has about 6.5 million states and about 9 million transitions. A complete build will require about 30 minutes on 1.7 Ghz Pentium 4 running Windows 2000.

The underlying morphological analyzer is very high coverage and implements all word formation processes of Turkish. It has a noun root lexicon of about 25K entries, a verb root lexicon of about 5K entries and a proper noun lexicon of about 40K entries.

5. Conclusions

We have presented the design and implementation of a finite state pronunciation lexicon for Turkish and agglutinating language with an infinite lexicon. The lexicon produces a representation that encodes in parallel, the SAMPA representation of the all possible pronunciations of the word along with the corresponding morphological analyses. The correct computation of the pronunciation and the location of the stress requires a full morphological analysis be done so the pronunciation has been built on top of a two-level morphological analyzer. The system has been implemented using the XRCE finite state tools.

6. Acknowledgements

This work was supported in part by a joint National Science Foundation and TÜBİTAK (Turkish Scientific and Technological Research Foundation) project “A Unified Electronic Lexicon of Turkish”. We also thank XRCE for making the finite state tools available.

7. References

- [1] D. Jurafsky and J. H. Martin, *Speech and Language Processing*. Upper Saddle River, New Jersey: Prentice Hall, 2000.
- [2] D. Hakkani-Tür, K. Oflazer, and G. Tür, “Statistical morphological disambiguation for agglutinative languages,” *Computers and the Humanities*, vol. 36, no. 4, 2002.
- [3] S. Inkelas, “Exceptional stress-attracting suffixes in Turkish: Representations vs. the grammar,” in *The Prosody-Morphology Interface*, R. K. abd Harry van der Hulst and W. Zonneveld, Eds. Cambridge: Cambridge University Press, 1999, pp. 134 – 187.
- [4] E. Sezer, “On non-final stress in Turkish,” *Journal of Turkish Studies*, vol. 5, pp. 61–69, 1981.
- [5] K. Oflazer, “Two-level description of Turkish morphology,” *Literary and Linguistic Computing*, vol. 9, no. 2, 1994.
- [6] L. Karttunen and K. R. Beesley, “Two-level rule compiler,” Technical Report, XEROX Palo Alto Research Center, 1992.
- [7] L. Karttunen, “Finite-state lexicon compiler,” XEROX, Palo Alto Research Center– Technical Report, 1993.
- [8] K. Koskenniemi, “Two-level morphology: A general computational model for word form recognition and production.” Publication No: 11, Department of General Linguistics, University of Helsinki, 1983.