

VOWEL INSERTION GRAMMAR

MUHAMMAD KHURRAM RIAZ, MUHAMMAD MUSTAFA RAFIQUE, SYED RAZA SHAHID

ABSTRACT

In Urdu language the position of vowels within the syllable is highly considerable issue. This paper contains vowel insertion grammar for Urdu language. The grammar is presented in the form of an automaton. The presented automaton was tested successfully on many words of Urdu language. The tested words were supposed to have all the diacritics correctly marked on them. The following vowel insertion grammar is extremely useful for building speech synthesis for Urdu language.

1. INTRODUCTION

In Urdu text, diacritics are very important as they indicate the existence of vowels between two consonants. In an Urdu text-to-speech converter, it is extremely important to judge where to produce a vocalic sound, as there is no particular letter in Urdu language, which shows the existence of a vowel in Urdu text. Vowel insertion grammar serves the purpose in determining the existence of vowels due to the presence of diacritics in Urdu text.

2. LITERATURE REVIEW

Seventeen vowels have been identified during the analysis of Urdu vocalic sounds (Urdu Consonantal and Vocalic Sounds, 2002). According to Kachru (1990), there are seven long oral vowels, and three short oral vowels. Khan (1997) also agrees with the long and short vowel distribution of Kachru (1990).

Hussain (1997, p. 153) suggests that Urdu vowels are written with the help of three letters (phonemes) and three diacritics—namely, letters are 'ا', 'و' and 'ی', and diacritics are paish (ـِ), zair (ـِ) and zabar (ـَ). The long vowels are written using one of the three letters (phonemes) mentioned above with one of the diacritics on the preceding consonant. The short vowels are written with only a diacritic on the preceding

consonant. However, these diacritics are not necessarily written in Urdu script (they can be ignored if the writer thinks no ambiguity can occur).

There are some languages that expand their vowel inventory by opening the velar port during some of their vowels' utterance.

Such a language uses nasalized vowels to distinguish within pairs of words that are otherwise same. It uses nasal vowels that are produced with shapes of oral configuration similar to the oral shapes of some of its non-nasal vowels (Pickett, 1999, p. 71). Urdu is also one such language, which enhances its vowel inventory by adding nasal vowels in it. Each long vowel in Urdu language has its nasalized version, but there is no existence of any nasalized short vowel in Urdu language. Therefore, overall, there are seventeen vowels (including oral and nasal vowels) in Urdu language.

3. RESULTS

Vowel insertion grammar is represented in the form of an automaton shown in Figure 1. Since the automaton is represented in the form of a *Mealy* machine, so there is a corresponding output associated with each input in transition from one state to another. Circles in the figure represent states. Start and final states are mentioned explicitly in the given automaton. Each intermediate state remembers all the previous inputs, which lead the system to that particular state.

The grammar is also explained with examples in Table 1. First column in Table 1 shows the pattern for each of the possible input. The second column shows the next input character for each possible input. First and second column form the composite key for Table 1. Output vowel is shown in the column Output. Examples for each case are shown in IPA (International Phonetics Alphabet) symbols.

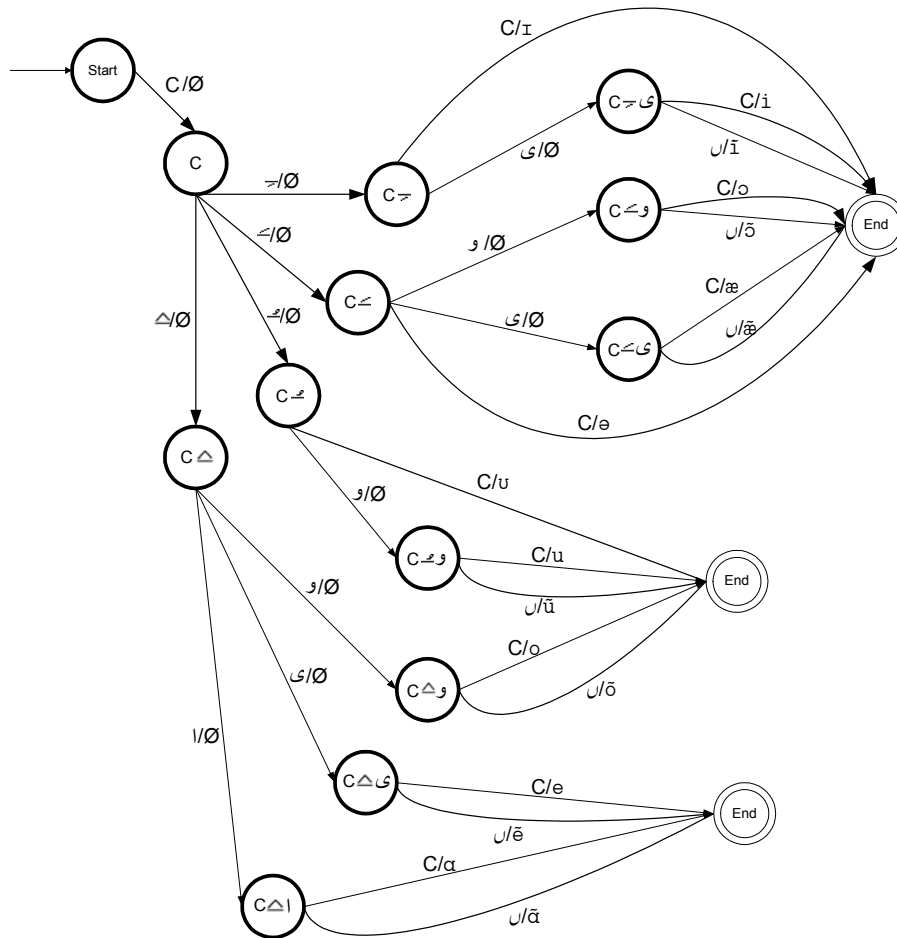


FIGURE 1 Automaton for vowel insertion grammar

4. DISCUSSION

The existence of vowels in Urdu language is partially similar to English language. But unlike English language, vowels in Urdu language are very much predictable because of the presence of punctuation marks (diacritics) in Urdu script. There are three punctuation marks in Urdu that help in the identification of vowels. These punctuation marks are paish (<math>\bar{Δ</math>), zair (<math>\bar{Δ</math>) and zabar (<math>\bar{Δ</math>). These punctuation marks are paish (<math>\bar{Δ</math>), zair (<math>\bar{Δ</math>) and zabar (<math>\bar{Δ</math>).

The results reveal that long vowels in Urdu language are inserted if a letter has any of the diacritics paish (<math>\bar{Δ</math>), zair (<math>\bar{Δ</math>) and zabar (<math>\bar{Δ</math>), followed by any of these three Urdu letters ' <math>\bar{Δ</math>', ' <math>\bar{Δ</math>' or ' <math>\bar{Δ</math>'. If any other letter apart from the above-mentioned letters follows

these diacritics, then a short vowel is inserted. Similarly if any of the Urdu letter ' <math>\bar{Δ</math>', ' <math>\bar{Δ</math>' or ' <math>\bar{Δ</math>' has a diacritic preceding and letter ' <math>\bar{Δ</math>' following it, then a nasalized long vowel is inserted.

Apart from the punctuation marks—paish (<math>\bar{Δ</math>), zair (<math>\bar{Δ</math>) and zabar (<math>\bar{Δ</math>), there is another punctuation mark jazm (<math>\bar{Δ</math>) that causes no insertion of vowel in Urdu language if the letter following the jazm (<math>\bar{Δ</math>) is any other than the letters ' <math>\bar{Δ</math>', ' <math>\bar{Δ</math>' or ' <math>\bar{Δ</math>'. So, only long vowels can exist if there is a jazm (<math>\bar{Δ</math>) on any letter. Jazm (<math>\bar{Δ</math>) is considered as default punctuation mark if a letter has no punctuation mark on it.

It was also observed that the letter ‘i’ can exist in Urdu script only if there is a jazm

(‘ Δ ’) on the previous letter.

TABLE 1 Possible diacritics combination. ‘C’ indicates any consonant

<i>Pattern</i>	<i>Next Level</i>	<i>Output</i>	<i>Example</i>
C $\bar{}$	C	ɪ	bɪtʃ ^h na
C $\bar{}$ ی	C	i	t̪in, amin
C $\bar{}$ ی	∪	ĩ	vəhĩ
C $\bar{}$	C	ə	t̪ərʃ, t̪ərɑdʒɪm
C $\bar{}$ و	C	ɔ	pɔdɑ, sɔdɑ
C $\bar{}$ ی	C	æ	tʃæn, bæɪ
C $\bar{}$ ی	∪	æ̃	hæ̃
C $\bar{}$	C	ʊ	ɖʊlɑhn, sʊkh
C $\bar{}$ و	C	u	χun, t̪u
C $\bar{}$ و	∪	ũ	hũ
C Δ ی	C	e	t̪evr, bekar
C Δ ی	∪	ẽ	rəhẽ
C Δ و	C	o	t̪oɾ, sonɑ
C Δ و	∪	õ	hõ
C Δ ا	C	ɑ	t̪ɑɪ
C Δ ا	∪	ã	hã

The automaton shown in Figure 1 has been tested on many words and successful results were noted. Since the automaton/grammar requires diacritics to decide between the vowels, so it becomes extremely important to mention all the diacritics to get the correct vowel as output. Generally, in Urdu text these diacritics are ignored and are considered to be understood by default. So, for using the automaton shown in Figure 1, extra effort is required to mention all the diacritics on the input text.

5. REFERENCES

Hussain, S. 1997. “Phonetic correlates of lexical stress in Urdu”; Unpublished Ph.D. dissertation, North Western University, IL, USA.

Kachru, Yamuna 1990. Hindi-Urdu in The Major Languages of South Asia, The Middle East and Africa, edited by Bernard Comrie.

Khan, Mahboob Alam 1997. Urdu ka Soti Nizaam.

Pickett, J. M. 1999. The Acoustics Of Speech Communication Fundamentals, Speech Perception Theory, And Technology. USA

Urdu Consonantal and Vocalic Sounds; Speech Synthesis Group, Center for Research in Urdu Language Processing (CRULP), National University of Computer and Emerging Sciences, Lahore.