

GILBERT PUECH et PIERRE BANCEL

Centre de Recherches Linguistiques et Sémiologiques
 Université Lumière-Lyon 2
 69500 Bron, France

ABSTRACT

The International Phonetic Alphabet (IPA) is the standard reference as a transcription system. With only minor variants, it is commonly used by linguists to record the pronunciation of languages whether they are supported by an orthographic tradition or not. The scope of this paper is to transpose the IPA to a computer-oriented coding system in order to use phonetic records in data bases and expert systems.

INTRODUCTION

A computer-oriented coding system for the representation of sounds should be viewed as an interface between linguists faced with the representation of a wide range of sounds and a Data Base Management System.

First the code corresponding to each sound must be a key to its major characteristics and, consequently, to the way it patterns with other sounds. The binary distinctive features theory seems to be the natural interface between phonetic analysis and the binary logic of computers. It turns out, however, that there is no clear agreement on how a number of complex or rare sounds should be treated in this approach; furthermore the built-in definition of some features is costly since it precludes some combinations - for instance [+High] is exclusive of [-Low] - or hardly satisfactory to account for some sounds - such as flaps and trills. On the other hand an IPA based classification presents several advantages: it is independent of any particular theory; it associates phonetic interpretation and a graphic representation in the same table; it allows a more compact code. This code can be easily converted into a matrix of distinctive features so that the exploitation of the data can be independent of the

coding system.

Secondly, the coding system must fit one of the standard formats for computer words. It should also be used to classify phonetically recorded words in the data bases in the same manner as the ASCII code is used to classify orthographically recorded words. If the data base is organized in n-ary trees, the algorithm will find all the relevant information necessary for the equilibration of the trees in the set of codes forming each word.

GENERAL ORGANIZATION

For maximal efficiency, each segment is coded in a short integer (16 bits word) noted by 4 hexadecimal figures. Consonants and vowels are coded independently of each other, thus it is necessary to know if one given code refers to a consonant or to a vowel before being interpreted. For languages - such as Bantu - in which words are built after a strict syllabic pattern, the data base may determine the fields composing the records as corresponding either to a consonant or to a vowel; in languages where no such syllabic regularity prevails, the first field of the record (a long integer) will in the first byte determine the number of segments included in the record and, in the three following bytes, select the V/C choice (bits 8 up to 31 set to 1 when the segment should be interpreted as a vowel and left at 0 if it is a consonant). Suprasegmental information - stress and pitch - is normally associated with vowels; provision is made however for consonants bearing a tone. A set of diacritics is used to give maximal versatility to this coding system which was designed both for narrow and broad transcriptions. Coding of morpheme boundaries for morphophonemic representations was not examined but could be accommodated.

CONSONANTS

A - Basic consonants are coded in the least significant byte of the short integer. Table 1 yields the phonetic interpretation of the coding and illustrates some of the realizations. The 4 most significant bits correspond to the lines (manner of articulation) and the 4 remaining bits to the columns (place of articulation):

Phonetic symbol	Code	Phonetic interpretation
b	0041	bilabial voiced stop
m	00C1	bilabial nasal stop
kp	001C	labiovelar unvoiced stop

Sonorants (lines B to F) are assumed to be voiced; implosives and ejectives are respectively voiced and unvoiced. For clicks, which may be voiced, aspirated, murmured etc., further qualification is needed. In order not to have more than 15 places of articulation, some choices had to be made; thus, apico-labial sounds, which are to be found in Umutina[1], are not included in the set of basic consonants but could be handled as a special case (see section F). To facilitate the editing on the lineprinter, it is convenient to have each basic symbol occupy one space only even if it is commonly transcribed as a sequence of two consonants (such as kp or ts).

B - Double consonants, geminates as well as complex segments, are coded in two morae and occupy two spaces:

bb	4141	geminate bilabial voiced stop
mb	C141	bilabial prenasalized voiced stop
nt	C414	alveolar prenasal. unvoiced stop
nts	C474	alveolar prenasalized unvoiced affricate

C - A release, transcribed by a right-adjacent diacritic occupying half a space, is coded in the least significant byte: the most significant bits refer to Table 2; the final hexadecimal zero is a flag indicating that the basic consonant (coded in the first byte) is followed by a release, the interpretation of which is given in Table 2:

kʸ	1B90	velar stop/palatal release
bʷ	41C0	bilabial stop/labiovelar release
dʳ	44F0	alveolar stop/alveolar trill release

Codes which are left free may be defined as necessary.

D - A segment synchronic property, transcribed by a subscribed diacritic, is coded in the most significant byte. The initial hexadecimal zero is a flag indicating that the first byte is to be interpreted as shown in Table 3:

y	0CB9	nasalized palatal approximant
z̥	0DA4	lateralized alveolar fricative
m̥	04C1	unvoiced bilabial nasal stop

Provision was made to code the lenis quality on a par with the fortis. However, the lenis quality is assumed to be the unmarked case and it is not associated with a graphic diacritic:

t	0114	lenis t
t̥	0214	fortis t

E - Consonants may be syllabic and bear tones. The syllabicity is coded by the least significant byte set to zero:

m	C100	syllabic bilabial nasal stop
t̥	9400	syllabic alveolar unvoiced fricative

Tones on consonants are coded as they are on vowels (see VOWELS, B); tone bearing consonants are assumed to be syllabic.

m̄	C104	syllabic nasal stop/high tone
ṁ	C102	syllabic nasal stop/low tone

F - The overwhelming majority of known consonants may be coded according to the preceding conventions. However it may be crucial in some languages to handle difficult cases as accurately as possible. We shall resort to the following system: the most significant byte is used as a pointer to a specific filter corresponding to the primary consonant coded in the second byte. One has access, through this filter, to a complementary code, so that the resulting code is extended to 3 bytes; the flag set to detect this situation is the zero corresponding to the least significant bits of the first byte:

ndʳ	10C4	Prenasalized stop/trill release	filter C4/1	: 44F0	extended code : C444F0
ŋ̥	10CB	Murmured prenasalized click	filter CB/1	: 0567	extended code : CB0567
ŋ̥	20CB	Voiced prenasalized click	filter CB/2	: 0467	extended code : CB0467

VOWELS

A - A short vowel - one mora - is coded on a short integer. A long vowel or a diphthong is coded as two morae. The most significant byte corresponds to segmental information. Vowels are plotted on an articulatory space defined by two axes: height (5 degrees) and tongue position in the oral cavity (front, central, back):

	Front	Central	Back
height	1	6	B
	2	7	C
	3	8	D
	4	9	E
	5	A	F

The most significant bits are interpreted as follows:

- bit 0 - approximant-like vowel
- 1 - marked tongue root
- 2 - nasal
- 3 - round

The bit 0 is used to mark superclosed vowels (like reconstructed proto-bantu ī/ȳ) or, more generally,

the non syllabic part of a diphthong:

- a_i 0A00 8100 diphthong with gliding i
- i_a 0100 8800 diphthong with gliding a

The bit 1 is used to interpret marked tongue root position (emphatic vowels in the Berber-Arabic domain or the harmonic set of vowels characterized by Advanced Tongue Root in a number of sub-Saharan languages). Nasality and roundness may combine with this feature:

- i 0100 (unrounded) i
- u 1B00 (round) u
- i̇ 2100 nasalized i
- u̇ 3B00 nasalized u
- ĩ 4200 ATR I

Basic symbols corresponding to the set of unrounded vowels and of rounded vowels are shown in Tables 4 and 5 respectively.

B - Suprasegmental information is coded in the second byte. Tonal languages use up to 5 levels of pitch, represented henceforth as accents. The code 06 is reserved for a downstepped High:

- 0101 Falling low ı̂
- 0102 Level low ı̄
- 0103 Mid ı̄
- 0104 High ı̇
- 0105 Suprahigh ı̇
- 0106 Downstepped High ı̂

Contour tones are coded by reference to their source/target pitch:

- 0142 Falling high-low ı̂ı̇
- 0124 Rising low-high ı̄ı̇

The bit 4 is set to 1 if the corresponding tone is floating:

- 014A High + Floating low ı̂̇
- 012C Low + Floating high ı̄̇

Double contours require two morae; we propose the convention that the first mora bear a level tone and the second a contour tone:

- 0104 0124 Falling-rising long i ı̂ı̇ı̂ı̇
- 0102 0142 Rising-falling long i ı̄ı̇ı̂ı̇

C - In order to maximally compact suprasegmental information the bit 0 is reserved for stress:

- 0180 stressed i ı̂

If the stressed vowel bears a tone, the code is modified accordingly:

- 0182 stressed i/low tone ı̂
- 01C2 stressed i/falling tone ı̂ı̇

The code A0 is assigned to pitch accent as required by some languages:

- 01A0 i associated with pitch accent ı̂

D - Hexadecimal codes 7 and F are left free in our system. Corresponding combinations will be used to account for marked voice quality:

- unvoicing 0107 unvoiced i ı̥̂
- 0147 unvoiced i ı̥̂
- high tone retained ı̥̂
- creaky voice 012F creaky i/low tone ı̥̂
- breathy voice 0172 breathy i/low tone ı̥̂

Special cases may be treated with an extended code as proposed for consonants: a flag (hexadecimal F) indicates that one has to go through a filter table, access to which is given by the code of the vowel mora and a pointer:

- 01F2 : go to case 2 of the filter table corresponding to vowel i.

Rhotacized vowels, for instance, could be conveniently dealt with in this way.

CONCLUSION

It is indeed possible to rely on the International Phonetic Alphabet to propose a comprehensive and versatile computer oriented coding system. The fact that the code is phonetically motivated makes it particularly attractive for expert systems aiming at comparing data or reconstructing proto-languages.

Reference

[1] P. Ladefoged, "Preliminaries to Linguistic Phonetics", The Univ. of Chicago Press, 1971.

		bilabial 1	labiodental 2	dental 3	alveolar 4	labiodental 5	retroflex 6	postalveolar 7	prepalatal 8	palatal 9	labiodental A	velar B	labiovelar C	uvular D	pharyngeal E	glottal F
unvoiced consonants	1	p		t	pt	ṭ				c		k	kp	q		ʔ
aspirated	2	ph		th								kh				
ejectives	3	p'		t'								k'				
voiced consonants	4	b		d	bd	ḍ				ɟ		g	gb	g		
implosives	5	ɓ		ɗ								ɠ				
clicks	6	ǀ	ǃ	ǂ			ǁ	ǁ								
unvoiced affricates	7		pf		ts			tʃ					kf			
voiced affricates	8		bv		dz			dʒ					gv			
unvoiced fricatives	9	ɸ	f	θ	s		ʃ	ʃ	ç	ç		x		χ	ħ	h
voiced fricatives	A	β	v	ð	z		ʒ	ʒ	ʒ	j		ɣ		ʁ	ʕ	ʕ
approximants	B	u		ɹ						y	ɥ		w			
nasals	C	m		n		ɳ				ɲ		ŋ		ɴ		
laterals	D			l		ɭ				ʎ		ʟ				
flaps / taps	E			r												
trills	F	ʙ		r			ʀ							ʀ		

Table 1

Symbol	Example	Code	Phonetic interpretation	Symbol	Example	Code	Phonetic interpretation
ʔ	tʔ	1410	unreleased	t		0114	lenis
h	tsʰ	7420	aspirated release	ṭ		0214	fortis
ʔ	tsʔ	7430	glottal release	̥		03C1	unvoicing
v	kʷ	1B90	palatal release	̣		0494	voicing
ɥ	kʷ	1BA0	labiodental release	̣		0541	murmur
n	ṭ̃	14B0	nasal release	̣		0A94	rounding
w	ḅ̃	41C0	labiovelar release	ṭ		0B14	velarization
l	ṭ̃	14D0	lateral release	w		0CBC	nasalization
ɹ	ṭ̃	14E0	pharyngeal release	z̃		0DA4	lateralization
r	ḍ̃	44F0	trill release	t̃		0E	pharyngalization
				̣̣		0F41	laryngalization

Table 2

Unrounded vowels

- i ı̂ u
- ɪ ı̄ u
- e ɛ y
- ɛ ɛ ʌ
- æ a ɑ

Table 4

Table 3 -

Round vowels

- y ɥ u
- y ɥ u
- ø ɘ o
- œ ɚ ɔ
- œ ɚ ɔ

Table 5