

# Language Specific Peculiarities Document for Cantonese as Spoken in the Guangdong and Guangxi Provinces of China

## 1. Dialects

The name "Cantonese" is used either for all of the language varieties spoken in specific regions in the Guangdong and Guangxi Provinces of China and Hong Kong (i.e., the Yue dialects of Chinese), or as one particular variety referred to as the "Guangfu group" (Bauer & Benedict 1997). In instances where Cantonese is described as 'Cantonese "proper"' (i.e. used in the narrower sense), it refers to a variety of Cantonese that is spoken in the capital cities Guangzhou and Nanning, as well as in Hong Kong and Macau.

This database includes Cantonese as spoken in the Guangdong and Guangxi Provinces of China only (i.e. not in Hong Kong); five dialect groups have been defined for Cantonese (see the following table)<sup>1</sup>. Three general principles have been used in defining these dialect groupings: (i) phonological variation, (ii) geographical variation, and (iii) lexical variation. With relation to phonological variation, although Cantonese is spoken in all of the regions listed in the table, there are differences in pronunciation. Differences in geographic locations also correlate with variations in lexical choice. Cultural differences are also correlated with linguistic differences, particularly in lexical choices.

<b>Area</b>	<b>Cities (examples)</b>
Central Guangdong Group	Guangzhou, Conghua, Fogang (Shijiao), Longmen, Zengcheng, Huaxian
Northern Guangdong Group	Shaoguan, Qijiang, Lian Xian, Liannan, Yangshan, Yingde, Taiping
Northern Pearl River Delta Group	Nanhai, Gaoya, Sanshui, Shunde, Gaoming, Foshan, Wuchan, Huazhou
Southern Pearl River Delta Group	Shenzhen, Dongguan, Baoan, Zhongshan
Guangxi and Western Guangdong	Guangxi cities, Zhanjiang, Maoming, Lianjiang, Guazhou

<sup>1</sup> Although Chinese dialectologists (such as Yue-Hashimoto (1991) and Bauer and Benedict (1997)) classify the Siyi region, also known as Taishan, as one of the dialect groups for Cantonese, Bauer and Benedict note (1997: xxxvi – xxxvii) that the dialects spoken in this region are markedly different phonologically from Cantonese spoken in the other regions. The extent of these phonological differences has been considered significant enough to exclude the Siyi-Liangyang group of dialects from the current speech collection of Cantonese.

## 2. Deviation from native-speaker principle

No special deviation: only native speakers of Cantonese born in mainland China will be recruited for voice recording in this project. No speakers from Hong Kong will be included.

## 3. Special handling of spelling

Spelling is not commonly used for Cantonese words. Chinese characters can be decomposed into radicals, which are sequences of individual strokes.

When spelling foreign words, the Roman letters are pronounced in the English way; Cantonese syllables are used to mimic the English pronunciation (see Spelling alphabet table in Section 9). For Cantonese recordings, only English words and sequences of isolated letters from the English alphabet are used in the prompt files. The elicited spelling item is also replaced by an English word letter sequence.

## 4. Description of character set used for orthographic transcription

Unicode UTF8 is used for the orthographic transcription of Cantonese.

## 5. Description of Romanization scheme

For Cantonese, there is no recognized standard comparable to the Pinyin system used to Romanize Mandarin. A number of Romanization systems have been developed to represent Cantonese sounds, the most well known are Meyer-Wempe, Sidney Lau, LSHK (Jyut-ping) and Yale. The Hong Kong and Guangdong/Guangxi Province governments do not follow any particular Romanization system. Matthews and Yip (1994:7) admit that “[r]epresenting Cantonese in alphabetic form is an intrinsically challenging problem, and none of the current systems is ideal.”

The Romanization system adopted for transcribing the current Cantonese database is based on the Yale system. This system was developed by Parker Huang and Gerald Kok and has the following advantages:

- its conventions resemble those of the IPA and Pinyin systems (Matthews and Yip 1994);
- most current textbooks and dictionaries adopt the Yale system or some modified version of it;
- it is relatively easy to learn to use and read. This is beneficial not only for those learning Cantonese as a foreign language, but also for native speakers of Cantonese<sup>2</sup>, since they tend not to Romanize Chinese in practice; in typing email, for instance, Cantonese speakers either use Chinese characters or type in English.

Cantonese is a tonal language in which most morphemes are one syllable long. As in all languages, each syllable is composed of one initial and one final (or ‘rime’) and an inherent tone.

The following table provides the Romanization symbol for each of the initials and rimes in Cantonese. The first column provides the phones and the second column provides the (distinctive) phonemes.

Phone (IPA)	Phoneme		Romanization	Example word
			<b>Initials</b>	
p	/b/ <sup>1</sup>	1	b	ba0 ( <i>father</i> )
p <sup>h</sup>	/p/	2	p	pa3 ( <i>fear</i> )
t	/d/ <sup>1</sup>	3	d	da2 ( <i>hit</i> )

<sup>2</sup> All transcribers are native speakers of Cantonese.

t <sup>h</sup>	/t/	4	t	tong4 ( <i>sugar</i> )
k	/g/ <sup>1</sup>	5	g	gam0 ( <i>gold</i> )
k <sup>h</sup>	/k/	6	k	kam2 ( <i>cover</i> )
k <sup>w</sup>	/gw/ <sup>1</sup>	7	gw	gwa0 ( <i>melon</i> )
k <sup>hw</sup>	/kw/	8	kw	kwa0 ( <i>boast</i> )
m	/m/	9	m	ma0 ( <i>mother</i> )
n	/n/	10	n	nei5 ( <i>you</i> )
ŋ	/ŋ/	11	ng	nga4 ( <i>tooth</i> )
f	/f/	12	f	fa0 ( <i>flower</i> )
s	/s/	13	s	sa0 ( <i>sand</i> )
ɕ			s <sup>2</sup>	syu0 ( <i>book</i> )
h	/h/	14	h	ha0 ( <i>shrimp</i> )
ts	/dz/ <sup>1,3</sup>	15	j	jaak3 ( <i>narrow</i> )
tʃ			j	jyu0 ( <i>pig</i> )
ts <sup>h</sup>	/ts/ <sup>3</sup>	16	ch	chaang2 ( <i>orange</i> )
tʃ <sup>h</sup>			ch	cheung4 ( <i>long</i> )
w	/w/	17	w	waak6 ( <i>draw</i> )
l	/l/	18	l	laang5 ( <i>cold</i> )
j	/j/	19	y	yut6 ( <i>month</i> )
			<b>Rimes</b>	
i:	/i:/ <sup>4</sup>	1	i	ji6 ( <i>Chinese character</i> )
		2	ip	dip6 ( <i>plate</i> )
		3	im	lim4 ( <i>blinds, curtains</i> )
		4	it	chit3 ( <i>cut</i> )
		5	in	nin4 ( <i>year</i> )
e <sup>ɿ</sup>		6	ik	sik6 ( <i>eat</i> )
		7	ing	sing1 ( <i>star</i> )
ɛ:	/ɛ:/	8	e	se2 ( <i>write</i> )
		9	ek	sek6 ( <i>rock</i> )
		10	eng	peng4 ( <i>cheap</i> )
a:	/a:/	11	a	fa0 ( <i>flower</i> )
		12	aap	laap6-saap6 ( <i>garbage</i> )
		13	aam	saam1 ( <i>clothes</i> )
		14	aat	baat3 ( <i>eight</i> )
		15	aan	ngaan5 ( <i>eye</i> )
		16	aak	jaak3 ( <i>narrow</i> )
		17	aang	chaang0 ( <i>support</i> )
ɐ	/ɐ/	18	ap	yap6 ( <i>enter</i> )
		19	am	gam0 ( <i>gold</i> )
		20	at	gwat1 ( <i>bone</i> )
		21	an	yan4 ( <i>person</i> )
		22	ak	hak1 ( <i>black</i> )

		23	ang	dang1 ( <i>lamp</i> )
u:	/u:/ <sup>5</sup>	24	u	gu0 ( <i>aunt</i> )
		25	ut	fut3 ( <i>wide</i> )
		26	un	wun2 ( <i>bowl</i> )
o <sup>w</sup>		27	ung	fung0 ( <i>wind</i> )
		28	uk	luk6 ( <i>six</i> )
y:	/y:/	29	yu	syu0 ( <i>book</i> )
		30	yu	hyut3 ( <i>blood</i> )
		31	yu	dyun2 ( <i>short</i> )
ɔ:	/ɔ:/	32	o	ngo5 ( <i>fl</i> )
		33	on	gon0 ( <i>dry</i> )
		34	ong	bong0 ( <i>help</i> )
		35	ot	hot3 ( <i>thirsty</i> )
		36	ok	wok6 ( <i>frying pan</i> )
œ:	/œ:/ <sup>6</sup>	37	eu	heu0 ( <i>boots</i> )
		38	eung	geung0 ( <i>ginger</i> )
		39	euk	geuk3 ( <i>foot</i> )
ə		40	eung	cheun0 ( <i>spring</i> )
		41	eut	cheut1 ( <i>go out</i> )
			<b>Diphthongs</b>	
i:w	/i:w/	42	iu	tiu3 ( <i>jump</i> )
ej	/ej/	43	ei	nei5 ( <i>you</i> )
ew	/ew/	44	au	sau2 ( <i>hand</i> )
ej	/ej/	45	ai	sai3 ( <i>small</i> )
a:j	/a:j/	46	aai	daai6 ( <i>big</i> )
a:w	/a:w/	47	aau	maau1 ( <i>cat</i> )
ɔ:j	/ɔ:j/	48	oi	choi3 ( <i>vegetable</i> )
ow	/ow/	49	ou	dou1 ( <i>knife</i> )
u:j	/u:j/	50	ui	bui1 ( <i>cup</i> )
ey	/ey/	51	eui	seui2 ( <i>water</i> )
			<b>Colloquial and loanwords only</b>	
ɛ:m	/ɛ:m/	52	em <sup>7</sup>	kem1 ( <i>game</i> )
ɛ:p	/ɛ:p/	53	ep <sup>7</sup>	tep1 ( <i>taste, sip</i> )
ɛw	/ɛw/	54	ew <sup>7</sup>	tew6 ( <i>throw away</i> )

<sup>1</sup> Non aspiration has been represented as voiced

<sup>2</sup> [s] and [ɕ] : allophones, [ɕ] for GZ speakers

<sup>3</sup> [ts] and [tʃ]: allophones; [ts<sup>h</sup>] and [tʃ<sup>h</sup>]: allophones

<sup>4</sup> [i:] and [e] are treated as allophonic

<sup>5</sup> [u:] and [o<sup>w</sup>] are treated as allophonic (NB no –p or –m)

<sup>6</sup> [œ:] and [ə] are treated as allophonic (NB no –p or –m)

<sup>7</sup> No official Yale Romanization for these colloquial rimes.

To mark tones of words, the Yale system uses diacritics. Traditionally, Cantonese has been considered to have nine distinctive tones, with three tones belonging to a category called “entering tone”, “stopped

tones”, “dead syllables” or “checked syllables”; these syllables end abruptly with an unreleased  $-p$ ,  $-t$  or  $-k$  stop. However, these three tones are now usually considered to be allotones of the three level tones for the live syllables. Accordingly, the Yale system provides for seven tones: three high, three low and one mid. However, Matthews and Yip (1994) recognise only six tones, arguing that the contrast between the high falling tone and the high level tone is not contrastive (p8). In acknowledgement of the difference between the well-known Yale notation and Matthews and Yip’s description, our numbering system provides for the high falling tone (tone 0), but we have normalized virtually all instances of the high falling tone in transcripts to the high level tone (tone 1). (The one exception is the syllable 拜 *baai0*, which we have corrected to the mid level tone *baai3*). Consequently, no syllables in the high falling tone actually appear in the database.

There are also two (sometimes three are recognized) other tones called “changed tones” (Ramsay 1987, Matthews and Yip 1994) or pi:n-jɛm (Bauer and Benedict 1997), which are associated with colloquial Cantonese speech and are generally not used in formal speech contexts. The changed tone of certain words modifies the original meaning either grammatically or lexically, and can mean “that familiar thing one often speaks of” or can indicate the speaker’s attitude towards it. While some Romanization systems do have diacritics for these changed tones, the Yale system does not.

The numbering system is as follows:

TONE	Yale diacritic	Numbering
high falling	à	a0
high level	ā	a1
high rising	á	a2
mid level	a	a3
low falling	àh	a4
low rising	áh	a5
low level	ah	a6

Note then that the Romanization system adopted for this database is a modified version of Yale.

## 6. Description of method for word boundary detection

Word boundaries in Chinese character orthography are not marked in the text. Word boundaries in the Romanized orthography are determined by localization of white spaces (blank, tab, etc.). They have been added manually in accordance with principles developed by Appen. Compounds and reduplicated word groups are treated as one unit; furthermore, affixes and verbal complements are joined to their noun/verb ‘stems’. Tense-aspect markers, however, are not considered affixes by Appen and are not joined to their ‘stems’. The word boundaries have been verified manually for the phonetically rich material.

## 7. Table containing all phonemes in the stipulated notation

The phonemic transcription of the words in this database uses the official SAMPA symbol set for Cantonese, which can be found at <http://www.phon.ucl.ac.uk/home/sampa/cantonese.htm>. The total number of phonemes is 27 (19 consonants and 8 vowels). There are also 11 diphthongs.

Symbol	Word	Transcription	Remark
<b>Consonants</b>			
<b>Plosives</b>			
b	ba0 ( <i>father</i> )	ba:_0	Non aspiration has been represented as voiced.
p	pa3 ( <i>fear</i> )	pa:_3	

d	da2 ( <i>hit</i> )	da:_2	
t	tong4 ( <i>sugar</i> )	tO:N_4	
g	gam0 ( <i>gold</i> )	g6m_0	
k	kam2 ( <i>cover</i> )	k6m_2	
gw	gwa0 ( <i>melon</i> )	gwa:_0	
kw	kwa0 ( <i>boast</i> )	kwa:_0	
<b>Fricatives</b>			
f	fa0 ( <i>flower</i> )	fa:_0	[s] and [ç] are allophones. GZ speakers use [ç]
s	sa0 ( <i>sand</i> )	sa:_0	
h	ha0 ( <i>shrimp</i> )	ha:_0	
<b>Affricates</b>			
dz	jaak3 ( <i>narrow</i> )	dza:k_3	[ts] and [tʃ]: allophones; [tsʰ] and [tʃʰ]: allophones
ts	chaang2 ( <i>orange</i> )	tsa:N_2	
<b>Nasals</b>			
m	ma0 ( <i>mother</i> )	ma:_0	
n	nei5 ( <i>you</i> )	nej_5	
N	nga4 ( <i>tooth</i> )	Na:_4	
<b>Lateral</b>			
l	laang5 ( <i>cold</i> )	la:N_5	
<b>Semi-vowels</b>			
w	waak6 ( <i>draw</i> )	wa:k_6	
j	yut6 ( <i>month</i> )	ju:t_6	
<b>Vowels</b>			
i:	dip6 ( <i>plate</i> )	di:p_6	[i:] and [eʰ] allophones of this phoneme
E:	sek6 ( <i>rock</i> )	sE:k_6	
a:	baat3 ( <i>eight</i> )	ba:t_3	
ɔ:	gam0 ( <i>gold</i> )	g6m_0	
u:	wun2 ( <i>bow</i> )	wu:n_2	[u:] and [oʷ] allophones of this phoneme
y:	hyut3 ( <i>blood</i> )	hy:t_3	
O:	ngo5 ( <i>fl</i> )	NO:_5	
ɔ:	geung0 ( <i>ginger</i> )	g9:N_0	[œ:] and [ə] allophones of this phoneme
<b>Diphthongs</b>			
iw	tiu3 ( <i>jump</i> )	tiw_3	
ej	nei5 ( <i>you</i> )	nej_5	
Ew	[no Yale Romanization] ( <i>throw away</i> )	tEw_6	colloquial only
ɔw	sau2 ( <i>hand</i> )	s6w_2	
ɔj	sai0 ( <i>west</i> )	s6j_0	
a:j	daai6 ( <i>big</i> )	da:j_6	
a:w	maau1 ( <i>cat</i> )	ma:w_1	
O:j	choi3 ( <i>vegetable</i> )	tsO:j_3	
ow	dou1 ( <i>knife</i> )	dow_1	
u:j	bui1 ( <i>cup</i> )	bu:j_1	
9y	seui2 ( <i>water</i> )	s9y_2	

## 8. Complete list of all rare phonemes

/kw/ is a relatively rare phoneme in Cantonese.

According to Bauer and Benedict (1997:486-487) some rimes (finals) occur in loan/colloquial words only, but they are combinations of phonemes, rather than individual phonemes.

## 9. Other language specific items

### 9.1. Spelling alphabet

A	ei1
B	bi1
C	si1
D	di1
E	yi1
F	e1fu4
G	ji1
H	ek1chu3
I	aai1
J	chei1
K	kei1
L	e1lou4
M	em1
N	en1
O	ou1
P	pi1
Q	kiu1
R	a1-lou4
S	e1-si4
T	ti1
U	yu1
V	wi1
W	dak1-bi1-yu1
X	ek1si1
Y	waai1
Z	yi1-set1

### 9.2. Chinese written script

There are two versions of Chinese written script: the Simplified-characters version, which is used in Mainland China (including the Guangdong Province) and the Traditional-characters version, which is used in Hong Kong.

Because recordings will take place in mainland China only, the simplified character set is to be used for both display of the prompts and orthographic transcription.

The lexicon for the Cantonese database thus contains three fields: simplified characters, Romanization, and SAMPA representation.

### 9.3. Spoken and written Cantonese

It is generally agreed (e.g. Bauer and Benedict 1997, Matthews and Yip 1994, Ramsey 1987, Yue-Hashimoto 1991) that Cantonese is different from Mandarin not only in the pronunciation of written words, but also in the extent to which the spoken register diverges from the standard written register. The differences between the spoken and written varieties of Cantonese apply not only to the syntax but also to the lexicon, phonology and graphology. In other words, Cantonese can be regarded as having two “realization systems”: one typically used in spoken or informal contexts<sup>3</sup> and the other in written or formal contexts. Note, however, that the former also has a written realization and the latter also has a Cantonese spoken realization! For simplicity’s sake, we use the term “spoken form” for the informal variety, and “written form” for the formal variety.

The written form is represented graphically by standard Chinese characters, those which are readable by members of other Chinese dialects, with one caveat: Cantonese speakers in Hong Kong use the traditional characters, whereas those on the mainland use the simplified characters. When read aloud, these characters have Cantonese pronunciation, unintelligible to speakers of other Chinese dialects.

Words in the spoken form may or may not be represented by written characters. Bauer and Benedict (1997:xiv) call such characters “colloquial” or “dialectal” characters. These are often borrowed or adapted homophonous standard Chinese characters. These colloquial characters have not been officially standardized and are not taught to Cantonese speakers at school. Cantonese speakers are exposed to them in popular print media texts, e.g. magazines, certain sections of newspapers, song lyrics, personal letters. The use of colloquial characters in newspapers and magazines has become much more common over the last decade or so, particularly in Hong Kong publications.

There are also a number of commonly used colloquial words that do not have a written form. Since it is not unlikely that these will be used when speaking to a computer, they need to be included in the reading scripts. For such words, the Romanization will be used for identification, and a homophonous character will be used to represent the vernacular morpheme.

### 9.4. Spelling Chinese characters

When explaining (orally) how to write a new word, Cantonese speakers have a number of options:

- (i) Give a homograph of the word:  
For instance, this method is used for describing proper names.
- (ii) Use a distinguishing radical:  
This method is the most common, and is based on the assumption that Cantonese speakers are familiar with the written characters associated with each set of homophones. To use this method, the speaker will refer to the radical that distinguishes the character to be written from its homophones.  
Some radicals are themselves characters. For instance, the character for heart “sam1” 心 (‘heart’) can function as a radical in many other characters, such as the word “seung2” 想 (‘to think’).  
Other radicals may not be words themselves but have names. For instance, one of the radicals in the character 河 “ho4” (‘river’), is made up of three strokes, and is commonly referred to as

---

<sup>3</sup> Note that it is not uncommon to find the word “Cantonese” confusingly being used to refer to this spoken, informal variety only.



“saam1 dim2 seui2” 三点水 (‘the three water strokes’); it is found in many words associated with water and the sea.

If the speaker does not know the name of a particular radical of a character, s/he may use the specific characteristics of a character to indicate which character s/he is referring to. For instance, there are two different characters for the surname “Wong”, a speaker may say “saam1 waak6 Wong” 三画王 (three-strokes Wong) to indicate his/her surname is this character王, and may say “daai6 tou5 Wong” 大肚黄 (big-stomach Wong) to refer to the other Wong character, i.e., 黄.

(iii) Describe the component strokes:

NB: This method is used as a last resort. The speaker says the names of individual strokes in the order in which they are written to form the character. For instance, this character 谭 “taam4” can be broken down into three parts: 言 “yin4” (‘speech’), 西 “sai0” (‘west’), and 早 “jou2” (‘morning’).

There are six principles which determine the stroke order of Chinese characters: (i) left to right, (ii) top to bottom, (iii) horizontal before vertical, (iv) from outside to inside, (v) middle before two sides, and (vi) inside before closing. A radical can be situated on the top-centre, or left, or right, or at the bottom part of a character. Most (literate) Cantonese speakers will know whether the stroke or radical “spelled” first is situated on the top left-hand-side or the top centre of the word.

## 11. References

Bauer, Robert S. & Benedict, Paul, K. 1997. *Modern Cantonese Phonology*. Mouton de Gruyter, Berlin/NY

Matthews, Stephen & Yip, Virginia. 1994. *Cantonese: a comprehensive grammar*. Routledge, London/New York

Ramsey, S. Robert. 1987. *The Languages of China*. Princeton University Press, Princeton, New Jersey

Yue-Hashimoto, Anne. 1991. “The Yue dialect.” In *Languages and Dialects of China*, [Journal of Chinese Linguistics Monograph Series Number 3], ed. William S-Y Wang. Berkeley, CA, University of California, pp. 294-324.