



ISCA Archive  
<http://www.isca-speech.org/archive>

## AN APPLICATION OF SAMPA-C FOR STANDARD CHINESE

Chen Xiaoxia\*, Li Aijun, Sun Guohua, Hua Wu and Yin Zhigang

Institute of Linguistics  
Chinese Academy of Social Sciences  
5 Jianguomennei Rd., Beijing 100732, China

6<sup>th</sup> International Conference on Spoken  
Language Processing (ICSLP 2000)  
Beijing, China  
October 16-20, 2000

### ABSTRACT

Labeling segment is an important work in database building. This paper presents a labeling system for Standard Chinese named SAMPA-C. We give some charts: consonant chart, vowel chart, tone chart, retroflex final chart, sound variation chart and non-speech symbol chart. Then this labeling system is used in two corpora labeling. The result shows that the labeling system is suitable for Standard Chinese.

### I. INTRODUCTION

It is significant work to label segment in speech corpora for Standard Chinese. With the labeled material we can do many research work. Hanyupinyin is an effective way to transcribe Standard Chinese. But it is not entirely corresponding to IPA. For example, “i” representing [i], [i̥], [ī]. It is not easy to be a machine-readable symbol system. According to international machine readable symbol system SAMPA[1], Zhu Weibin and Zhang Jialu have transcribed a symbol system with SAMPA for labeling syllable.[2,3]. They give Chinese SAMPA symbols including consonant, vowel and tone charts according to Xu Shirong’s view. They label isolated syllable in a database. This is an important work for transcribing Standard Chinese. But it is not enough for us. We hope <sup>1</sup>to label continuous speech. The representation of continuous speech is more complex than

isolated syllable. There are sound variation phenomena in continuous speech such as centralization, reduction, insertion etc.. The detailed labeling must include them. So we must formulate symbols to label them. Based on these, we design SAMPA-C labeling system for Standard Chinese. We have made a labeling system in syllable tier last year[4]. Now we make it in a continuous speech tier. What we refer to is Luo Changpei’s view[5] for consonant and vowel. For retroflex final, we refer to Wang Lijia’s result[6]. Then, we give diacritics for sound variation and give non-speech symbols.

We have two-speech corpora, which are read speech corpora and spontaneous corpora in CASS. The first one includes 18 articles and 10 speakers. The materials are read in recording room with normal rate. The second is originated from 19 cassettes provided by the Broadcast Station of Tsinghua University (BSTHU), Beijing, China. Most of the speech in the cassettes is causally given without paper preparation. Thus it is natural and covers a lot of valuable spontaneous phenomena. Those cassettes are then digitized into mono waveform at 16-bit precision and 16-kHz sampling rate. Through a standard Sound Blaster card on the PC, resulting in the 1.5 GB raw speech database totally[7]. With Pinyin and SAMPA-C, we label the two corpora.

### II. LABELING SYSTEM

2.1 The principles of labeling system are as follows:

- (1) Accurate: It need to transcript each segment precisely.

<sup>1</sup> \*The author is studying in Peking University for Ph.D.

This paper is supported by the National Social Sciences Foundation of China.

- (2) Systematic : For one phenomenon, we use a consistent manner to transcribe it. For example, there are not voiced stop and voiced affricate consonants in isolated syllable. But, for continuous speech, there are many voiced stops and voiced affricates. We just give voiced symbol “\_v” to represent those consonants becoming voiced in SAMPA-C not give voiced stop or voiced affricate.
- (3) Consistency: For most segments, we give the precise corresponding from segment’s

IPA to SAMPA. But there are seldom symbols we must redefine them in Standard Chinese. For example, retroflex final is a special phonetic characteristic. We give retroflex symbol “r” in SAMPA to represent it. For voiceless, we don’t use “\_0” but “\_u” because “\_0” is used in neutral tone.

### 2.2 Labeling system SAMPA-C

We give SAMPA-C as follows: consonant chart, vowel chart, retroflex final chart sound variation chart and also non-speech chart.

**Table 1 Consonant Chart for Standard Chinese**

| PinYin | IPA | SAMPA-C | PinYin | IPA | SAMPA-C |
|--------|-----|---------|--------|-----|---------|
| b      | p   | p       | z      | ʈʂ  | ʈʂ      |
| p      | pʰ  | p_h     | c      | ʈʂʰ | ʈʂ_h    |
| m      | m   | m       | s      | s   | s       |
| f      | f   | f       | zh     | tʂ  | ʈʂ`     |
| d      | t   | t       | ch     | tʂʰ | ʈʂ_h`   |
| t      | tʰ  | t_h     | sh     | ʂ   | s`      |
| n      | n   | n       | r      | ʀ   | ʂ`      |
| (a)n   | n   | _n      |        |     |         |
| l      | l   | l       | j      | tʃ  | ʈʂ\     |
| g      | k   | k       | q      | tʃʰ | ʈʂ_h    |
| k      | kʰ  | k_h     | x      | ç   | s\      |
| h      | x   | x       |        | ?   | ?       |
| ng     | ŋ   | ŋ       |        |     |         |

**Table 2 Vowel Chart For Standard Chinese**

| PinYin | IPA | SAMPA-C |
|--------|-----|---------|
| a      | ɑ   | A       |
| o      | o   | o       |
| e      | ɛ   | 7       |
| i      | i   | I       |
| u      | u   | u       |
| ü      | y   | y       |
| (zh)i  | ʅ   | i`      |
| (z)i   | ɿ   | i\      |
| er     | ɚ   | @`      |

**Table 3 Tone Chart For Standard Chinese**

| TONE   | IPA | SAMPA-C | Example |
|--------|-----|---------|---------|
| Tone 0 | 0   | ba_0    | ba0 吧   |
| Tone 1 | 1   | ba_1    | ba1 巴   |
| Tone 2 | 2   | ba_2    | ba2 拔   |
| Tone 3 | 3   | ba_3    | ba3 把   |
| Tone 4 | 4   | ba_4    | ba4 罢   |

**Table 4 Retroflex Final For Standard Chinese**

| NAME | PINYIN | IPA | SAMPA-C | EXAMPLE |
|------|--------|-----|---------|---------|
|      | ar     | ar  | a`      | par     |

|           |       |        |       |           |
|-----------|-------|--------|-------|-----------|
| opened    | or    | or     | o`    | mor       |
|           | er    | °r     | 7`    | ger       |
|           | (zh)i | Èr     | i@`   | zhir,shir |
|           | (z)i  | Èr     | i@`   | zir       |
|           | air   | ar     | a`    | bair      |
|           | eir   | Èr     | @`    | leir      |
|           | aor   | šor    | Ao`   | daor      |
|           | our   | our    | ou`   | gour      |
|           | anr   | ar     | a`    | ganr      |
|           | enr   | Èr     | @`    | genr      |
|           | angr  | a!!<r  | a~`   | gangr     |
|           | engr  | È!<r   | @~`   | dengr     |
| stretched | ir    | iÈr    | i@`   | jir       |
|           | iar   | iar    | ia`   | iar       |
|           | ier   | iEr    | ie_r` | jier      |
|           | iaor  | išor   | iAo`  | jiaor     |
|           | iour  | iour   | iou`  | qiur      |
|           | ianr  | iar    | ia`   | jianr     |
|           | inr   | iÈr    | i@`   | jinr      |
|           | iangr | ia!!<r | ia~`  | liangr    |
|           | ingr  | iÈ!<r  | i@~`  | ingr      |
|           | iongr | iu!<r  | iu~`  | xiongr    |
| rounded   | ur    | ur     | u`    | gur       |
|           | uar   | uar    | ua`   | guar      |
|           | uair  | uar    | ua`   | guair     |
|           | ueir  | uÈr    | u@`   | gueir     |
|           | uanr  | uar    | ua`   | tuanr     |
|           | uenr  | uÈr    | u@`   | lunr      |
|           | uor   | uor    | uo`   | luor      |
|           | uangr | ua!!<r | ua~`  | kuangr    |
|           | uengr | uÈ!<r  | u@`   | uengr     |
|           | ongr  | u!<r   | u~`   | kongr     |
| protruded | ŭr    | yÈr    | y@`   | yur       |
|           | ŭer   | yEr    | yE_r` | yuer      |
|           | ŭanr  | yar    | ya`   | yuanr     |
|           | ŭnr   | yÈr    | y@`   | qunr      |

Table 5 Diacritics Chart For Standard Chinese

| NAME           | IPA  | SAMPA-C | EXAMPLE |
|----------------|------|---------|---------|
| nasalized      | a<   | ~       | e~      |
| centralized    | e!f  | _”      | e_”     |
| voiceless      | n!%  | _u      | n_u     |
| voiced         | !dœ  | _v      | t_v     |
| rounded        | t!!f | _O      | O_O     |
| syllabic       | \    | =       | M=      |
| pharyngealized | øt/  | _?\     | A_?\    |
| silence        |      | sil     | sil     |
| silence voiced |      | silv    | silv    |

Table 6 Non-speech Chart for Standard Chinese

| PHENOMENA    | SAMPA-C            |
|--------------|--------------------|
| repairs      | repair <...repair> |
| disfluencies | disfl <...disfl>   |
| silences     | silen <...silen>   |

|                    |                                |
|--------------------|--------------------------------|
| <b>laughing</b>    | <b>laugh&lt;...laugh&gt;</b>   |
| <b>coughing</b>    | <b>cough&lt;...cough&gt;</b>   |
| <b>breathing</b>   | <b>breath&lt;...breath&gt;</b> |
| <b>crying</b>      | <b>cry&lt;...cry&gt;</b>       |
| <b>noise</b>       | <b>noise&lt;...noise&gt;</b>   |
| <b>lengthening</b> | <b>leng&lt;...leng&gt;</b>     |
| <b>modal</b>       | <b>mod&lt;...mod&gt;</b>       |
| <b>murmur</b>      | <b>mum&lt;... mum&gt;</b>      |
| <b>smack</b>       | <b>smack&lt;... smack&gt;</b>  |

### III. LABELING RESULT

Consistence is high: Using the labeling system, we segment and label the two corpora with Pinyin and SAMPA-C. Three tiers are given. The first tier is pinyin, the second is semi-syllable and the third is sound variation or other speaking phenomena. With manual work, we give a consistence test for labeler for nearly 15 minutes. It is about from 82.39% to 88.25%. The consistency is high. It shows that the labeling system is feasible. Most symbols are used during the labeling. The other result will be showed in another paper [7].

### IV. DISCUSSION

We change some symbols in Standard Chinese. Next, we explain them as follows:

- (1) Retroflex final is an important phonetic representation. We give final plus r as retroflex final.
- (2) There is not voiced consonant in isolated syllable in Standard Chinese. But it is common that stop, affricate or fricative can be voiced. So, we give voiced symbol to represent the phenomenon in continuous speech, but not give voiced stop, voiced affricate or voiced fricative.
- (3) The silence before stop and affricate often becomes voiced. It can be a long time. We just give a symbol “silv” to represent that duration.
- (4) The neutral tone is a special tone in Standard Chines. We use “\_0” as the symbol to consist with the other tones. So, for voiceless, we use “\_u”. It is not consistence with SAMPA.

- (5) For apical nasal “n”, we give two varieties according to their place in a syllable. As initial, it is showed with “n”. But as final, it is showed with “\_n”.

### REFERENCES

- [1] J. Wells, “Computer-coding the IPA: a proposed extension of SAMPA”, 2000, <http://www.phon.ucl.ac.uk/home/sampa/>
- [2] Zhu, Weibin & Zhang, Jialu(1997), Manual segmentation & labeling in Chinese speech database, The first China-Japan Workshop on Spoken Language Processing (CJSPL’97)
- [3] Zhang, Jialu(1999), A SAMPA system for PUTONGHUA(Standard Chinese) Oriental COCOSA’99 PROCEEDINGS
- [4] Chen, Xiaoxia Zu, Yiqing & Li Aijun(1999), A cardinal labeling system for Standard Chinese, The fourth phonetics conference in China
- [5] Luo, Changpei & Wang, Jun(1957), An outline of general phonetics, Science press
- [6] Wang, Lijia(1992), The principle of phonology, YUWEN Press
- [7] Li, Aijun etc.(2000), A phonetic labeling on read and spontaneous discourse corpora, ICSLP2000