

CORPUS INFORMATION FORM

Check (x) or fill in the blank (_____) if appropriate.

1. Name of the Corpus: BIT-MonoSyllable
(Mandarin Mono-Syllable Corpus)

2. IPR Holder:

Modern Communication Lab Of Beijing Institute Technology

3. Corpus Type:

- Speech Corpus (x)
- Text Corpus ()

4. If it is a speech corpus:

- Purpose:

- ASR (x)
- TTS ()
- Other, please specify acoustic study of monosyllable

- Language:

- Putonghua (SC) (x)
- Mandarin in Taiwan (TC) ()
- Cantonese in HK (TC) ()
- Other, please specify _____

- Style:

- Read speech (x)
- Spontaneous speech ()
- Conversational speech ()
- Other, please specify _____

- Channel:

- Close-talk Microphone (x)
- Telephone ()
- Mobile phone ()
- Other, please specify _____

- Sampling Rate: 16 k Hz

- Sampling Precision:

- PCM (x) 16 bits per sample
- A-law ()
- Miu-law ()
- Other, please specify _____

-
- Corpus size: _____ 3 _____ hours _____ 10 _____ speakers _____ ~500MB
 - SNR level: _____ 30 - 50 _____ dB
 - Transcriptions:
 - Character tier (SC) (x)
 - Character tier (TC) ()
 - Canonical Pinyin tier (x)
 - Other canonical pronunciation tier, please specify _____
 - Surface form IF tier ()
 - Surface form IPA tier ()
 - Surface form SAMPA-C tier ()
 - Other surface form tier, please specify _____
 - Other transcription, please specify _____
 - Other transcription, please specify _____
 - Other transcription, please specify _____

5. If it is a text corpus:

- Language:
 - SC ()
 - TS ()
 - Other, please specify _____
- Domain:
 - Culture ()
 - Economy ()
 - Military ()
 - News ()
 - Politics ()
 - Sciences ()
 - Sports ()
 - Other, please specify _____
- Corpus size: _____ Mega characters
- Tag Information:
 - Word segmentation: ()
 - Part-of-Speech ()
 - Other, please specify _____

6. A brief Description of the Corpus:

This is a single syllable corpus of Chinese mandarin by 10 speakers, including 6 male and 4 female speakers, every speakers reads 1~2 times of all syllables. The corpus covers the most-frequently-used 1259 tonal syllables in mandarin (corresponding to 406 basic syllables without tone). The corpus was recorded in quiet environment at BIT audio lab. This corpus can be used for acoustic and phonetic study of mandarin. For speech recognition, it can be used to train the acoustic models of single mandarin syllables. It can also be used for initialization of HMM models for LVCSR task.