

## CORPUS INFORMATION FORM

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

**1. Name of the Corpus:** 500-People TRSC (Telephone Read Speech  
Corpus)

**2. IPR Holder:** d-Ear

**3. Corpus Type:**

- Speech Corpus (x )
- Text Corpus ( )

**4. If it is a speech corpus:**

- Purpose:
  - ASR (x )
  - TTS ( )
  - Other, please specify \_\_\_\_\_
- 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 ( )
  - Telephone (x )
  - Mobile phone ( )
  - Other, please specify \_\_\_\_\_
- Sampling Rate: 8 k Hz
- Sampling Precision:
  - PCM ( ), \_\_\_\_\_ bits per sample
  - A-law ( )
  - Miu-law (x )
  - Other, please specify \_\_\_\_\_
- Corpus size: 72 hours 500 speakers ~2GB
- SNR level: \_\_\_\_\_ 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 \_\_\_\_\_
  - Other, please specify \_\_\_\_\_
  - Other, please specify \_\_\_\_\_
  - Other, please specify \_\_\_\_\_
  - Other, please specify \_\_\_\_\_
  - Other, please specify \_\_\_\_\_
- Corpus size: \_\_\_\_\_ Mega characters
- Tag Information:
  - Word segmentation: ( )
  - Part-of-Speech ( )
  - Other, please specify \_\_\_\_\_
  - Other, please specify \_\_\_\_\_
  - Other, please specify \_\_\_\_\_
  - Other, please specify \_\_\_\_\_

## 6. A brief Description of the Corpus:

This is a large telephone read speech corpus by about 500 speakers, half male and half female. Every speaker reads about 110 sentences most of which are taken from *the China*

*People's Daily* newspapers. It also contains some sentences to cover the connected digits and English characters. This corpus has some characteristics:

- 1) Roughly di-IF (di-Initial/Final) balanced;
- 2) Digital and English character balanced;
- 3) Large number of speakers;

This corpus can be used to train the acoustic model on telephone channel for LVCSR. It can also be used to do many other research works, for example, speaker identification and verification, and so on.