# A Country Report – COCOSDA  Activities in China

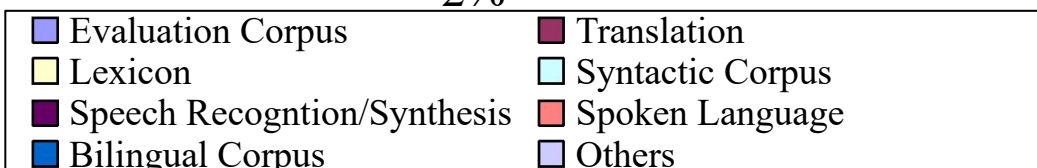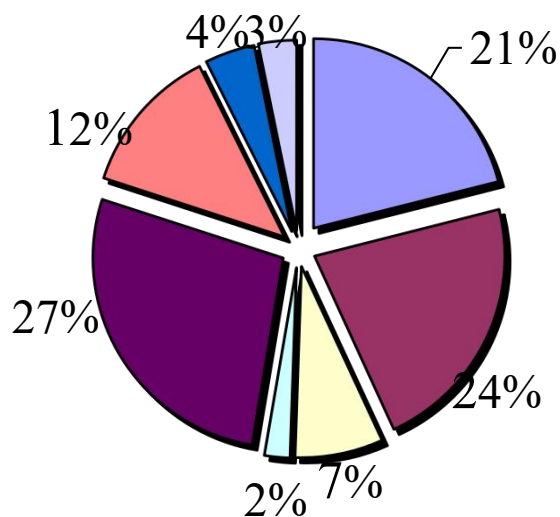O-COCOSDA 2014，Phuket, Thailand

Aijun LI , *Dong WANG

Institute of Linguistics, Chinese Academy of Social Sciences

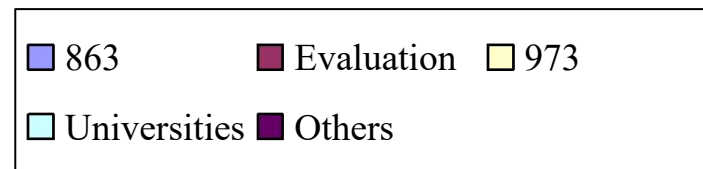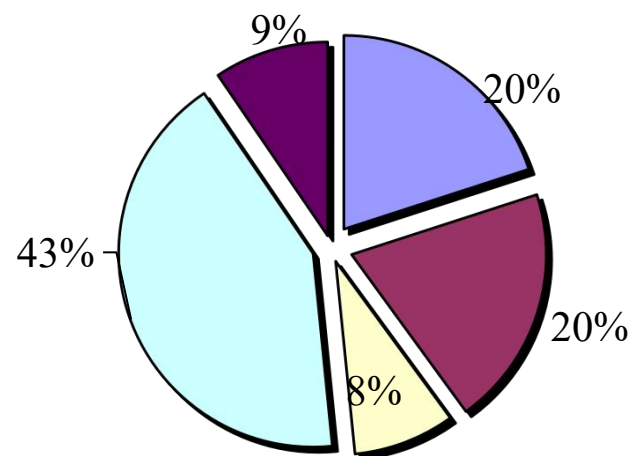*Research Institute of Information Technology, Tsinghua University

# 1 Activities in Chinese LDC

- Till now, there are 97 corpora, including speech synthesis/recognition corpora, corpora for machine translation, lexicon and other natural language processing corpora.

- In 2014, there are 2 new corpora added;14 corpora have been distributed to 8 institutes and companies.



Types of the corpora

Legend (Types of the corpora):
- Evaluation Corpus
- Translation
- Lexicon
- Syntactic Corpus
- Speech Recogntion/Synthesis
- Spoken Language
- Bilingual Corpus
- Others

Providers of the corpora

Legend (Providers of the corpora):
- 863
- Evaluation
- 973
- Universities
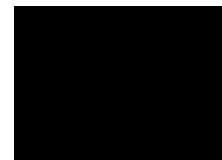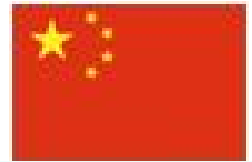- Others

# 2 Activities in Institute of Linguistics, Chinese Academy of Social Sciences-1

- **Phonetic Research on AESOP-CASS**
  - Question of learners from Northern Guan Dialect
  - Yes-no question of learners from Beijing
  - Yes-no question of learners from Wu dialect
  - Segmental errors of English learners from Shandong dialect

- **3D Articulatory speech database**
  - spoken by a female professional speaker
  - 3D articulatory data are recorded by NDI Waves
  - 2206 syllables,11204 words and 222 sentences in Standard Chinese
  - 3D IPA training database for 4 phoneticians and modeling

□ ... n of

内容检索按钮

下拉菜单按钮

九州语言学网

192.168.1.101

中国社会科学院语言研究所
九州网：语言资源网络服务平台

分类导航按钮

语法　　　词汇　　　字音

项目内容显示区

1. 巴
|帮系|帮组|帮| |假|开二|麻| |平|全清|
《现代汉语方言音库：常用字》 ☆43

1. 鹅卵石
|词汇|方言词汇|贰：地理|石沙土块矿物|
《方言调查词汇表》

1. 张三跟李四一样高。
|语法|句法|比较|等比|
《现代汉语方言语法语料库调查方案》 ☆13

2. 沙子
|词汇|方言词汇|贰：地理|石沙土块矿物|
《方言调查词汇表》

项目查询基本
信息显示区

条数:6855个；页面:1/1714；分类:全部

向前翻页按钮

hildren's
r each child
All speech
road phonetic

✓
✓                                                           r aged between
                                                            : 2- or 3-month,

□ **Jiuzhou linguistic resources (九州语言网）,** http://9zhou.phonetics.org.cn/.
✓  an open  platform sharing resources for linguistics, phonetics and dialect research.
✓ Chinese  dialect data topologically organized  by pronunciation, lexicon/morphology and syntax.

- **Speech corpora**
  - Broadcast Media Language Corpus
    - More than 15000h audio & video recordings of TV and Radio programs
    - Over 1000h Sentence-matched speech data from TV or radio news programs
  - Mandarin Chinese Digital Sampling Corpus (MCDSC)
    - 3543h video or audio files with metadata tagging since 1933
  - An Online TV Programs Transcription Corpus (http://ling.cuc.edu.cn/RawPub/)
- **Tools**
  - The parallel corpus retrieval tool (http://ling.cuc.edu.cn/chs/download/CUC_ParaConc_V0.3.rar)
  - The dialect homophonic character generation tool TYZH_v0.4.rar (http://ling.cuc.edu.cn/chs/download/TYZH_v0.4.rar)

# 4 Activities in Linguistics Lab, Department of Chinese Language and Literature, Peking University （Jiangping Kong）

- **Speech Multimodal corpora**
  - X-ray vocal tract database of Mandarin
  - 2D MRI static and 3D MRI vocal tract of Mandarin
  - 2D high-speech imaging database of dynamic glottis
  - Breathing band rhythm database of Mandarin
  - electropalatography of Mandarin Tibetan and Mongolian
- **Oral culture Multimodal corpora**
  - Liao songs, Epic of Buluotuo and <span style="color:red">wizard religious rituals</span> of Zhuang minority
  - King Pan grand song of Yao minority
  - Grand song of Dong minority
  - Bimo scriptures and Suni wizard religious rituals of Yi minority
  - Dongba <span style="color:red">scriptures and folk-song</span> of Naxi minority
  - Tibetan Buddhism recite Buddhist
  - Humai and long tuned of Mongolian nationality

# 5 Activities in Institute of Ethnology & Anthropology, Chinese Academy of Social Sciences

**Built Acoustic Parameter Database**
– Acoustic Parameter Databases of 7 languages (multi-syllabic, reading): Mongolia, Tu, East YuGu, Uighur, Hasak, Ewenki, Xibo. The databases cover most acoustic parameters in time and frequency domains.

**Developed Auto-Annotation & Auto-Extraction Software**
– 8 levels annotation file format
– Auto-Annotation Software: version 1.0 – 3.3
– Auto-Extraction Software: version 1.0 – 3.9

**Developed Acoustic Parameter Management Software**
– implement acoustic parameter inquiry & retrieval

**Published Related Research Papers**
– 20 papers in academic journals
– 39 papers in academic conferences
– 7 thesis papers

# 6 Activities in Minority Languages and literature College of the Minzu University of China

- Publishing a "Series of Endangered Chinese Minority Languages Folklore Audio Book"
- This Series audio book are part of the Speech Corpus of Endangered Chinese Minority Languages, which are established by Minority Languages and literature College of the Minzu University of China, and which are financially supported by Chinese government.
- "Tujia Language Folklore" and "De'ang language Folklore" were published in 2013 and 2014.
-
- CD in the book
  - to retrieve the Folklore