# Speech Resources for Scientific Research and its Application

**NII**

**NII-SRC** Speech Resources Consortium

## Kimiko YAMAKAWA  Shuichi  ITAHASHI (NII)

## Objective?

NII has established Speech Resources Consortium (NII-SRC) so as to promote dissemination and distribution of speech resources.

NII-SRC conducts collection, distribution, investigation and research on speech resources (including speech data and software tools) necessary for developing science, education and industry related to speech.

## What are we doing?

We contribute to the development of various research including speech recognition and synthesis by collecting and distributing speech corpora or speech databases which are difficult to develop individually.

Another scientific contribution by supplying valuable material for phonetics and sociolinguistics by preserving dialects and minority languages.

## 1. What is "Speech Corpus/Corpora"

### What is corpus/corpora?

A corpus means a systematic collection of data for research with some additional information to be used for research.

(Ex.) Speech corpus, text corpus, multimedia corpus, image corpus, etc.

### Variety and use of speech corpora

【Use】Analysis, synthesis, recognition of speech; analysis of discourse and dialects; preservation of languages, etc.

【Variety】Isolated words, continuous speech, read speech, dialogues, dialects, multilingual speech; speech by non-native speakers, infants, aged people; speech in noisy or reverberant environments.

### Recording media of speech corpora

The major recording media are used though it varies according to the use or data size. DVD-R is the most common currently. On-line distribution will be available soon.

【Recording media for speech corpora】

CD-R, DVD-R, HDD, DAT, LD, etc.

### Contents of speech corpora

Analysis data

Video

Transcription data

## 2. What is SRC?

### Why speech corpora, now?

- Development of speech processing technology
- Quantitative research in linguistics-related areas
- Importance of preserving languages and dialects

Massive speech data of various kinds necessary
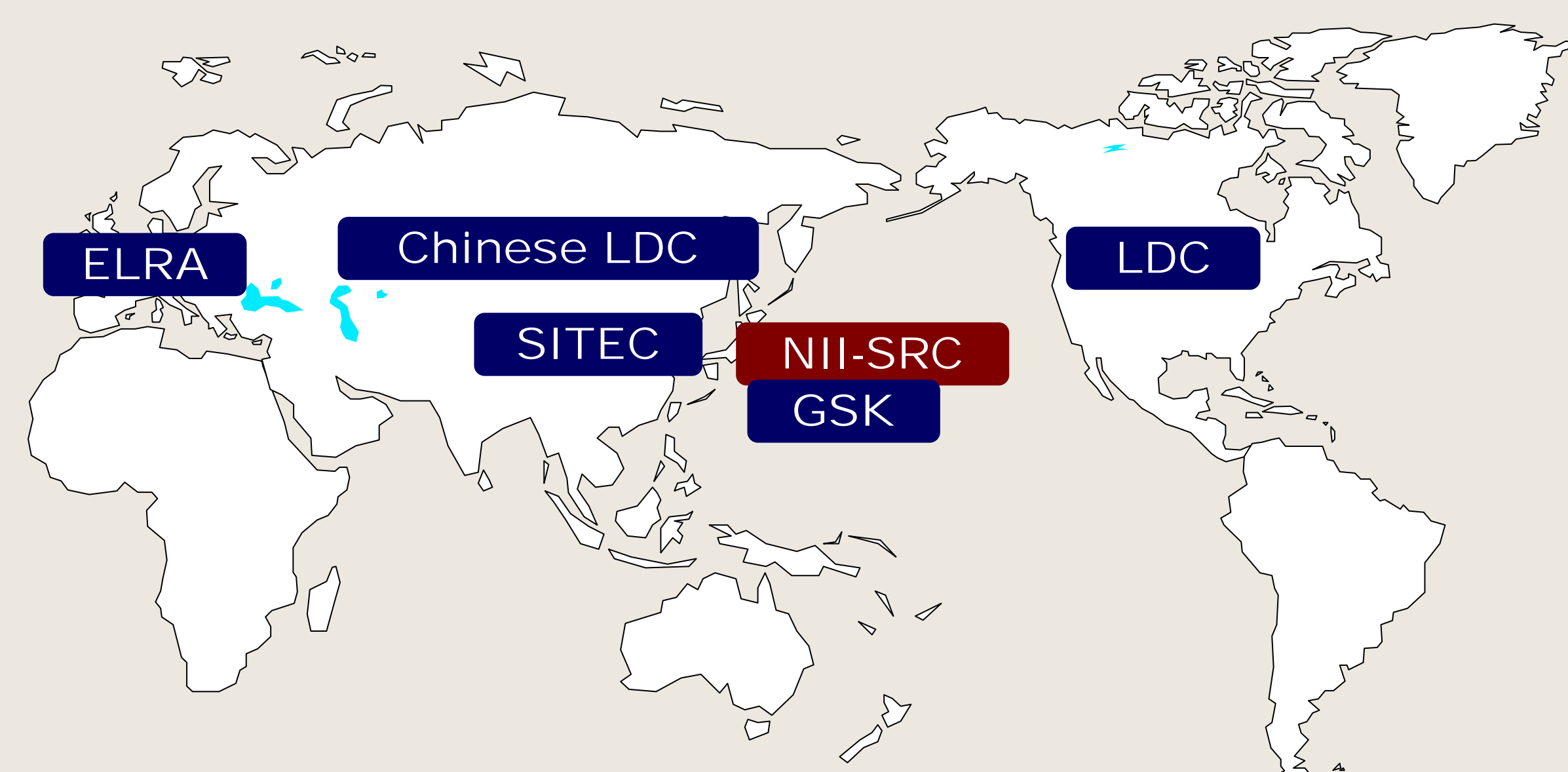
### Problems of speech corpora

- Most corpora are developed for a project.
- It requires cost, time and labor to create.
- Expensive.
- Not open to the public.

A common framework is required for creation, collection, accumulation, distribution and sharing

### Speech Resources Consortium (NII-SRC)

Launched by NII in 2006 in order to collect, manage, and distribute various speech corpora.

Currently, 31 corpora are available from NII-SRC.

ELRA    Chinese LDC    LDC    SITEC    NII-SRC    GSK

Speech-related Organizations in the World.

## 3. Categorization of speech corpora

Corpus attributes (8 attributes and 58 items)

| Attribute | | Item |
|---|---|---|
| Input device | 7 items | Type of input device (ex. Desk-top microphone) |
| Input environment | 5 items | Recording environment (ex. Soundproof room) |
| Number of speakers | 10 items | Number of speakers |
| Speaking style | 4 items | Style of speech (ex. Continuous speech) |
| Speech mode | 5 items | Speech mode (ex. dialog, read speech) |
| Data mode | 9 items | Other information (ex. Sampling frequency) |
| Language | 4 items | Type of language (ex. Monolingual) |
| Purpose | 14 items | Keyword for use or development (ex. Recognition) |

## 4. Corpus similarity visualization

Digit data corpora
Close-talking microphone
Read speech corpora
Digit data

Dialog corpora
Continuous speech
Dialog speech

Large read corpora
Read speech corpora
Monolingual
Over 100 Speakers

Small corpora
Fewer than 10 speakers

dim.2

dim.1

## Speech Resources Consortium, National Institute of Informatics

URI： http://research.nii.ac.jp/src/eng    E-mail： src@nii.ac.jp