# Rapid Deployment of Speech Processing Systems to New Languages and Domains

## Tanja Schultz

*University of Karlsruhe & Carnegie Mellon University*
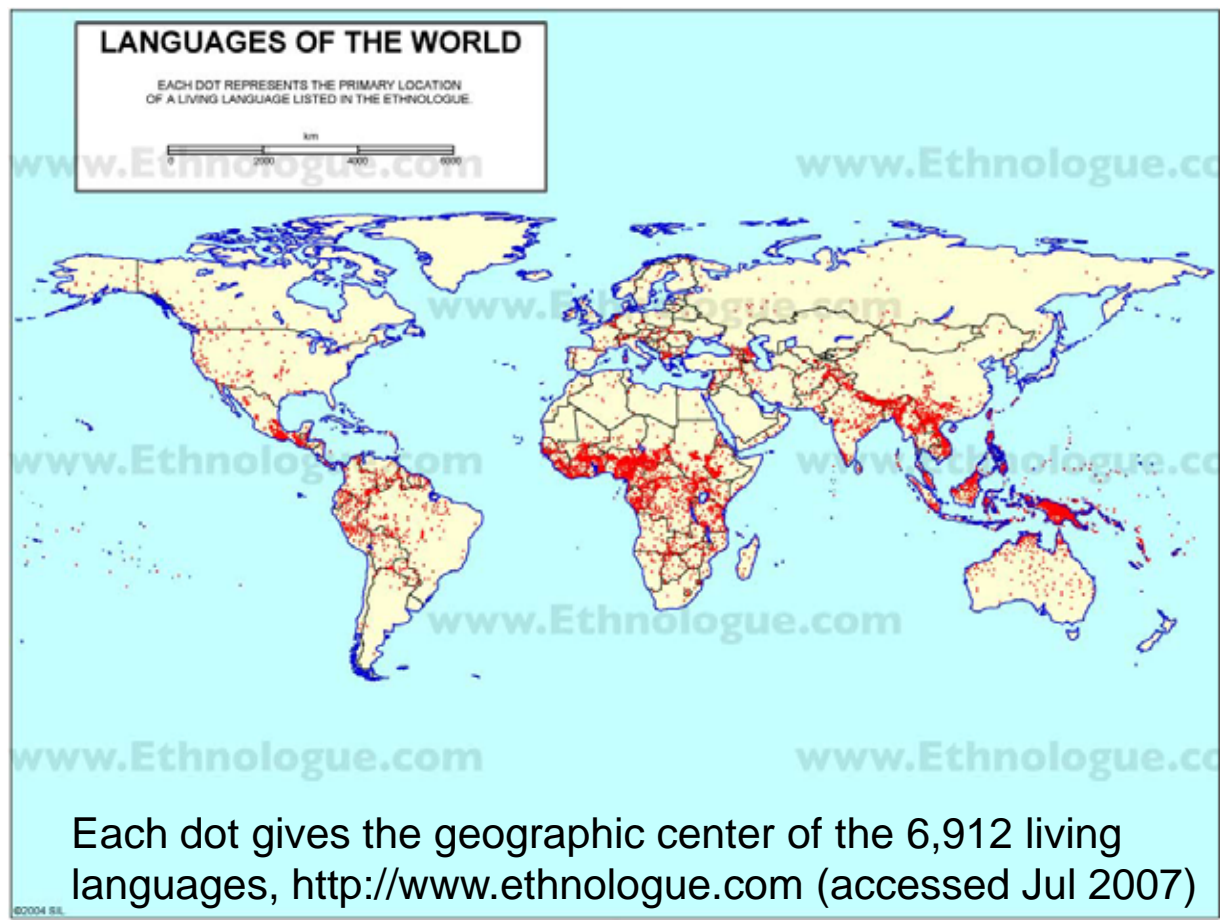Cognitive Systems Laboratory

http://csl.ira.uka.de

IS-LTC 2008, *Thursday, October 16th 2008,* Ljubljana

# Outline

o **The World's Languages**
  - o 6900 languages – So what?
  - o Language Extinction – What can the community do about it?
  - o Do we need Speech Processing for all of them?
  - o Is this really science – not just retraining on a new language?

o **Language Characteristics**
  - o Written form, scripts, letter-to-sound relationship
  - o Issues and Differences between languages

o **Challenges for Multilingual Speech Processing**
  - o Lack of Resources (Money, Data, Technical Support)
  - o Lack of Experts

o **Solutions**
  - o SPICE: A Rapid Language Adaptation Server
  - o Technologies: Leveraging off GlobalPhone & FestVox
  - o Experiments and Results

o **Conclusions and Future Work**

# Outline

o **The World's Languages**

  o 6900 languages – So what?

  o Language Extinction – What can the community do about it?

  o Do we need Speech Processing for all of them?

  o Is this really science – not just retraining on a new language?

o Language Characteristics

  o Written form, scripts, letter-to-sound relationship

  o Issues and Differences between languages

o Challenges for Multilingual Speech Processing

  o Lack of Resources (Money, Data, Technical Support)

  o Lack of Experts

o Solutions

  o SPICE: A Rapid Language Adaptation Server

  o Technologies: Leveraging off GlobalPhone & FestVox

  o Experiments and Results

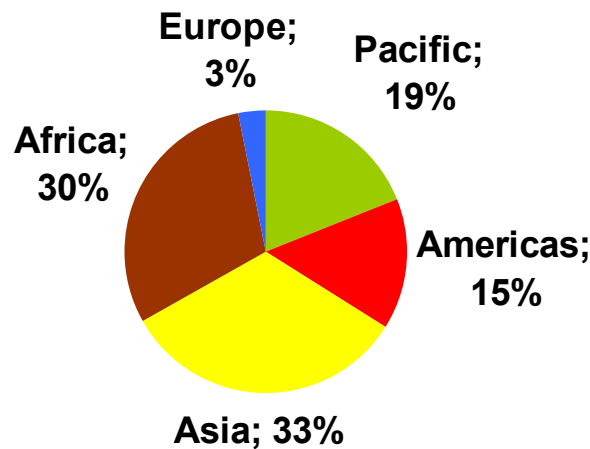o Conclusions and Future Work

# Everyone speaks English, why bother?

o Total number of Languages in the world: **6912**

o Language is not only a communication tool but fundamental to cultural identity and empowerment!

o Cultures, ideas, memories transmit *through language*

o Intellectual issues (e.g. world history) Practical issues (medical practices) Literature, …

Slovakian proverb: "with each newly learned language you acquire a new soul"

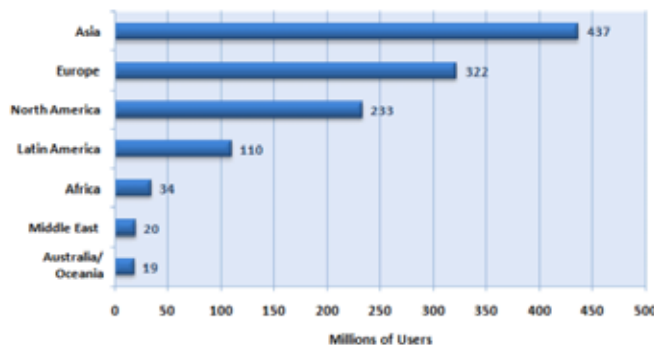o Preserve linguistic diversity! Similar to eco systems (David Crystal)

**LANGUAGES OF THE WORLD**

EACH DOT REPRESENTS THE PRIMARY LOCATION OF A LIVING LANGUAGE LISTED IN THE ETHNOLOGUE.

km

Each dot gives the geographic center of the 6,912 living languages, http://www.ethnologue.com (accessed Jul 2007)

# Increasing Language Diversity in Web

## Diversity of Languages in the Internet grows rapidly

o Top-10: 200%, All others: 440%

o Portuguese: 524%

o Arabic: 940%



Europe; 3%
Pacific; 19%
Africa; 30%
Americas; 15%
Asia; 33%

**Internet Usage by World Region**

Asia — 437
Europe — 322
North America — 233
Latin America — 110
Africa — 34
Middle East — 20
Australia/ Oceania — 19

Millions of Users
Copyright © 2007, www.internetworldstats.com

### Top Ten Languages Used in the Web
( Number of Internet Users by Language )

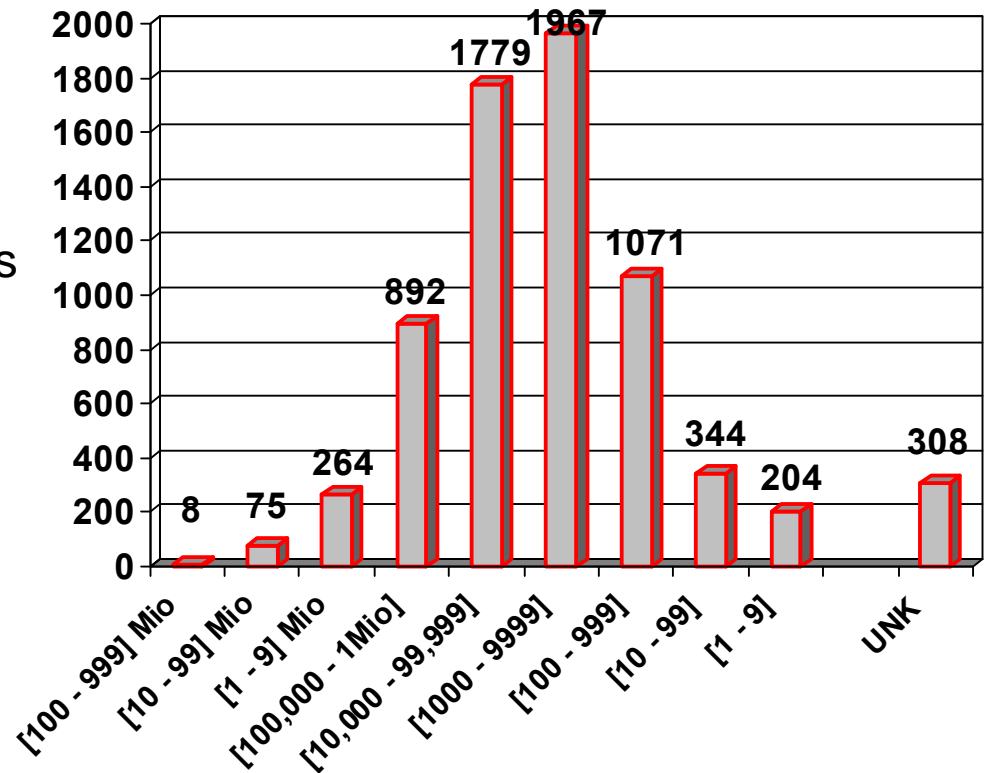| TOP TEN LANGUAGES IN THE INTERNET | % of all Internet Users | Internet Users by Language | Internet Penetration by Language | Internet Growth for Language ( 2000 - 2007 ) | 2007 Estimate World Population for the Language |
|---|---|---|---|---|---|
| English | 31.7 % | 365,893,996 | 17.9 % | 157.7 % | 2,042,963,129 |
| Chinese | 31.7 % | 166,001,513 | 12.3 % | 413.9 % | 1,351,737,925 |
| Spanish | 8.8 % | 101,539,204 | 22.9 % | 311.4 % | 442,525,601 |
| Japanese | 7.5 % | 86,300,000 | 67.1 % | 83.3 % | 128,646,345 |
| German | 5.1 % | 58,981,592 | 61.1 % | 112.9 % | 96,488,326 |
| French | 5.1 % | 58,456,702 | 15.1 % | 379.2 % | 387,820,873 |
| Portuguese | 4.1 % | 47,326,760 | 20.2 % | 524.7 % | 234,099,347 |
| Korean | 3.0 % | 34,120,000 | 45.6 % | 79.2 % | 74,811,368 |
| Italian | 2.7 % | 31,481,928 | 52.9 % | 138.5 % | 59,546,696 |
| Arabic | 2.5 % | 28,782,300 | 8.5 % | 940.5 % | 340,548,157 |
| TOP TEN LANGUAGES | 84.8 % | 978,883,995 | 19.0 % | 198.0 % | 5,159,187,766 |
| Rest of World Languages | 15.2 % | 175,474,783 | 12.4 % | 440.3 % | 1,415,478,651 |
| WORLD TOTAL | 100.0 % | 1,154,358,778 | 17.6 % | 219.8 % | 6,574,666,417 |

(*) NOTES: (1) Internet Top Ten Languages Usage Stats were updated for June 30, 2007. (2) Internet Penetration is the ratio between the sum of Internet users speaking a language and the total population estimate that speaks that specific language. (3) The most recent Internet usage information comes from data published by Nielsen//NetRatings, International Telecommunications Union, Computer Industry Almanac, and other reliable sources. (4) World population information comes from the world gazetteer web site. (5) For definitions and navigation help, see the Site Surfing Guide. (6) Stats may be cited, stating the source and establishing an active link back to Internet World Stats. Copyright © 2007, Miniwatts Marketing Group. All rights reserved.

# Currently 6900 Languages, but …

o **Extinction of languages on massive scale**
(David Crystal, Spotlight 3/2000)

o **Half of all existing languages die out over next century**
$\Rightarrow$ On Average: Every two weeks one language dies!

o **Survey Feb 1999 from Summer Institute of Linguistics**

51 languages with 1 speaker left

28 of those in Australia alone

500 languages with $<$ 500 spks

1500 languages with < 1000 spks

3000 languages with < 10.000

5000 languages with < 100.000

(not safe even for >100.000)

96% of world's languages are
spoken by only 4% of its people

Chart values:
- [100 - 999] Mio: 8
- [10 - 99] Mio: 75
- [1 - 9] Mio: 264
- [100,000 - 1Mio]: 892
- [10,000 - 99,999]: 1779
- [1000 - 9999]: 1967
- [100 - 999]: 1071
- [10 - 99]: 344
- [1 - 9]: 204
- UNK: 308

# How to safe Languages?

Prerequisites and Costs:

o   Community itself must want it, Surrounding culture must respect it

o   Funding for courses, materials, and teachers, support the community

o   Crystal estimates about $80.000 / year per language

o   3000 endangered languages is about $700Mio …

o   Foundation of endangered languages (FEL), UNESCO project

How could our community contribute:

o   Field Work and Community Outreach

   o   Get the tools to the people, i.e. flexible, portable, easy to handle

   o   Engage and actively involve native speakers

o   Lower the overall costs for data acquisition

   o   Automate the solicitation and data collection process

   o   Identify language specific aspects and focus

o   Reduce the data needs without sacrificing performance

   o   Streamline techniques & approaches to perform on small data

   o   Reuse language independent aspects of data and models

**Language support is good but why *Speech*?**

➢ Computerization: Speech is *the* key technology

➔ Ubiquitous Information Access: on the go, phone-based

➔ Mobile Devices: Too small and cumbersome for keyboards

➢ Globalization:

➔ Cross-cultural Human-Human Interaction

➔ Multilingual Communities: EU, South Africa, …

➔ Humanitarian needs, disaster, health care

➔ Military ops, communicate with local people

➔ Human-Machine Interfaces

➔ People expect speech-driven applications in their mother tongue

⇒ **Speech Processing in multiple Languages**

# ML Speech Processing – A Research Issue?

***Just retraining on foreign data? - No science!***

o  New language – new challenges

    o Writing system: different or no script, no vowelization, G-2-P

    o Word segmentation, morphology

    o Sound system: tonals, clicks

o  Different Cultures – social factors

    o Trust, access, exposure, background

o  Lack of Data and Resources

    o Audio, Transcripts, Pronunciations, Text, parallel bilingual data

o  Lack of Experts

    o Technology experts without language expertise

    o Native language experts without technology expertise

If we can solve the research issues for some languages, we might help the others along the way!

# Outline

# Language Characteristics

→ <u>Prosody, Tonality:</u>    Stress, Pitch, Lenght pattern, Tonal contours
<span style="color:red">(e.g. Mandarin 4, Cantonese 8, Thai & Vietnamese 5)</span>

→ <u>Sound system:</u>    simple vs very complex sound systems
<span style="color:red">(e.g. Hawaiian 5V+8C vs. German 17V+3D+22C)</span>

→ <u>Phonotactics:</u>    simple syllable structure vs complex consonant clusters
<span style="color:red">(e.g. Japanese Mora-syllables vs. German pf,st,ks)</span>

→ <u>Segmentation:</u>    Written form separate words by white space?
<span style="color:red">(NO: Chinese, Japanese, Thai, Vietnamese)</span>

→ <u>Morphology:</u>    short units, compounds, agglutination

English:    Natural segmentation into short units – great!

German:    Compounds – not quite so good

<span style="color:red">Donau-dampf-schiffahrt**s**-gesellschaft**s**-kapitän**s**-mütze …</span>

Turkish:    Agglutination – looooong phrases

<span style="color:red">Osman-lı-laç-tır-ama-yabil-ecek-ler-imiz-den-mi**ş**-siniz</span>

*behaving as if you were of those whom we might*

*consider not converting into Ottoman*

# Writing Systems

<u>Writing systems – basic unit is a Grapheme:</u>

Logographic: based on semantic units, grapheme represents meaning

Chinese: >10.000 *hanzi*; Japanese ~7000 *kanji,* Korean to some extend

Phonographic: based on sound units, grapheme represents sound

Segmental: grapheme roughly corresponds to phonemes
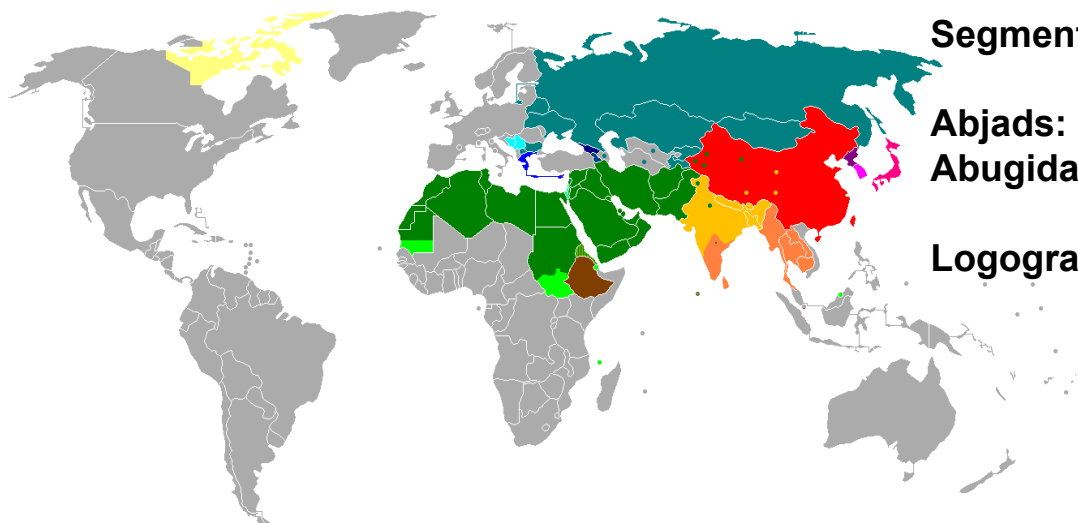
Latin (190), Cyrillic (65), Arabic (22) graphems

Abjads = consonantal segmental phonographic, e.g. Arabic

Syllabic: grapheme represents entire syllable, e.g. Japanese *kana*

Abugidas = mix of segmental and syllabic systems

Featural: elements smaller than phone, e.g. articulatory features

e.g. Korean: ~5600 *gulja*



**Segmental:** Latin  Cyrillic  Latin&Cyrillic  Greek
Georgian or Armenian

**Abjads:** Arabic, Arabic&Lat  Hebrew&Arabic
**Abugidas:** North Indic, South Indic, Ethiopic,
Thaana  Canadian Syllabic

**Logographic+syllabic:** Pure logographic,
Mixed logographic&syllabaries,
Featural syllabary+lmtd logographic
Featural-alphabetic syllabary

Wikipedia: August 2007

# Scripts – Some examples

العربى  булгарски  català  中国话  hrvatski  česky

english  ελληνικα  עברית  हिंदी  italiano  日本語

한국어  românește  русский  српски  ภาษาไทย

Scripts of some languages: Arabic, Bulgarian, Catalan, Chinese, Croatian, Czech, English, Greek, Hebrew, Hindi, Italian, Japanese, Korean, Romanian, Serbian, Thai

How many languages do have a written form?
- Omniglot lists about 780 languages that have scripts
- True number might be closer to 1000
  (Source Simon Ager, 2007, www.omniglot.com)

➔ Logographic scripts, mostly 2 representatives:
- Chinese: ~ 10.000 hanzi,
- Japanese: ~7000  kanji (+ 3 other scripts ☺)

➔ Phonographic:
- Korean: ~5600 gulja,
- Arabic, Devanagari, Cyrillic, Roman: ~100 characters

# Grapheme-to-Phoneme Relation

Grapheme-to-Phoneme (Letter-to-Sound) Relationship:

Logographic: NO relationship at all
   concern for Chinese, Japanese, Korean
Phonographic: segmental: close – far – complicated
   e.g. Finnish, Spanish: more or less 1:1, -- English: try „Phydough"
Phonographic: segmental – consonantal

   e.g. Arabic: no short vowels written
Phonographic: syllabic
   e.g. Thai, Devanagari: C-V flips



**Ratio Phonetic/Semantic Code**

DeFrancis/Unger

Finnish          Chinese
        French           Japanese
          English      Korean

Phonographic                    Logographic

➔ Automatic Generation of Pronunciations might get complicated

# Outline

# Challenges of MLSP

- o Lack of Resources: Stochastic approach needs **many** data
  - o Hundreds of hours audio recordings and corresponding transcriptions
    Audio data $\leq$ 40 languages; Transcriptions take up to 40x real time
  - o Pronunciation dictionaries for large vocabularies (>100.000 words)
    Large vocabulary pronunciation dictionaries $\leq$ 20 languages
  - o Mono- and bilingual text corpora: few language pairs, pivot mostly English
- o Algorithms are language independent – MLSP is not!
  - o Other Languages bring unseen challenges (segmentation, G2P, etc.)
  - o Have we already seen ALL or MOST of the language characteristics?
- o Social and Cultural Aspects
  - o Non-native speech and language, code switching
  - o Combinatorical explosion (domain, speaking style, accent, dialect, ...)
  - o Few native speakers at hand for minority (endangered) languages
- o Lack of Language Experts
  - o Bridge the gap between technology experts and language experts

# One Solution: Learning Systems

$\Rightarrow$ Intelligent systems that learn a language from the user

o Efficient learning algorithms for speech processing

- o <u>Learning:</u>
  - o Interactive learning with user in the loop
  - o Statistical modeling approaches
- o <u>Efficiency:</u>
  - o Reduce amount of data (save time and costs): at least by factor of 10
  - o Speed up development cycles: days rather than months

$\Rightarrow$ Rapid Language Adaptation from universal models

o Bridge the gap between language and technology experts

- o Technology experts do not speak all languages in question
- o Native users are not in control of the technology

# SPICE

## **S**peech **P**rocessing:

## **I**nteractive **C**reation and **E**valuation toolkit

- National Science Foundation, Grant 10/2004, 4 years

- Principle Investigator Tanja Schultz

- Bridge the gap between technology experts $\rightarrow$ language experts
  - Automatic Speech Recognition (ASR),
  - Machine Translation (MT),
  - Text-to-Speech (TTS)

- Develop web-based intelligent systems
  - Interactive Learning with user in the loop
  - Rapid Adaptation of universal models to unseen languages

- SPICE webpage *http://cmuspice.org*

# SPICE - Goals

Three main goals:

o Lower the overall costs for data acquisition

    o Automate the solicitation and data collection process

    o Identify language specific aspects and focus

o Reduce the data needs without sacrificing performance

    o Streamline techniques to perform on small data

    o Reuse language independent aspects of data/models

o Field Work and Community Outreach

    o Get the tools to the people, i.e.
       flexible, portable, easy to handle

    o Engage and actively involve native speakers

# CMU SPICE

# Welcome to SPICE

### Getting started

SPICE is a web-based system for building an end-to-end speech system (including Automatic Speech Recognition and Text-To-Speech) in your own language.

### Existing Users
Login with your account:

Login `tanja`

Password `***********`

Login

### New Users
Create a new account:

Login

Password

Re-type Password

Email

Create new account

# CMU SPICE

## Build Your System

- 🟢 Text and prompt selection (help)

- 🟢 Audio collection (help)

- 🟢 Phoneme selection (help)

- 🔴 Grapheme-to-phoneme rules (help)

  build language model first

- 🔴 Lexicon pronunciation creation (help)
  build language model first

- 🔴 Build acoustic model (help)

- 🔴 Build language model (help)

- 🔴 Create ASR system

- 🔴 Create speech synthesis voice

## SPICE Project

**You must do the following to build support for your language:**

- Text collection and selection
- Audio collection
- Phoneme set specification
- Lexicon pronunciation creation
- Speech recognition acoustic model creation
- Speech recognition language model creation
- Speech synthesis voice creation

# SPICE – System Functionalities

o SPICE Collects:
  o Appropriate text data
  o Appropriate audio data
o SPICE Defines and Refines:
  o Phoneme set
  o Rich prompt set
  o Lexical pronunciations
o SPICE Produces:
  o Vocabulary / Word lists (ASR, TTS, SMT)
  o Pronunciation model (ASR, TTS)
  o Acoustic model (ASR, TTS)
  o Language model (ASR, SMT)
  o Synthetic voices (TTS)
o SPICE Maintains:
  o Projects and users login
  o Data, Models, and Speech Processing Systems

# Building Process

User: Sameer Language: Hindi Project: Sameer_Hindi [Logout]

**Test acoustic model**

क्या तुम्हे अच्छा लगता है

Sessions Panel

Speech-to-Text | Text-to-Speech

Process Log

1. SUCCESS: Server path set to Sameer/Hindi/Sameer_Hindi
2. SUCCESS: Language set to Hindi
3. SUCCESS: Server address set to plan.is.cs.cmu.edu:7890
4. SUCCESS: File uploaded: 68204 Bytes transferred.
5. SUCCESS: क्या तुम्हे अच्छा लगता है

## **SPICE building process**

1. Collect a text corpus
2. Generate a 200-1000 ut
3. Record the prompt list from one or more native speakers
4. Define a phoneme set
5. Construct a lexicon and letter-to-sound rules
6. Build a language model from the text corpus
7. Build acoustic models for ASR
8. Build voice models for TTS
9. Evaluate ASR and TTS using "talk-back" function

# Outline

Phone set & Speech data

Hello
Input: Speech

AM   Lex   LM

hi /h//ai/
you /j/u/
we /w//i/

hi you
you are
I am

NLP / MT

TTS

สวัสดี ดรับ

Output:
Speech & Text

„adios"    → /a/ /d/ /i/ /o/ /s/
„Hallo"    → /h/ /a/ /l/ /o/
现在广播完了    → ???

**Pronunciation rules**

Hello

**Input: Speech**

**AM**

hi  /h//ai/
you /j/u/
we /w//i/

**Lex**

hi you
you are
I am

**LM**

**NLP / MT**

**TTS**

สวัสดี ดรับ

**Output: Speech & Text**

Resource rich languages $\leftrightarrow$ Resource low languages

Focused Re-crawling

Bridge Languages

Text data

Hello

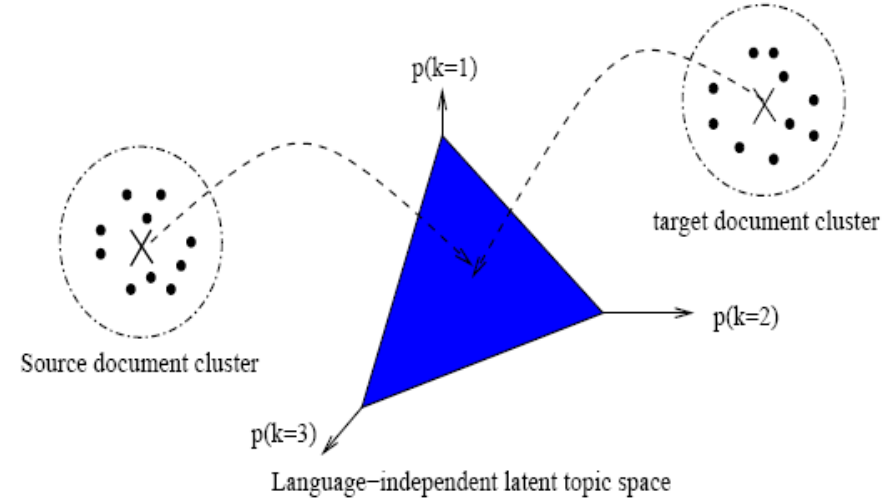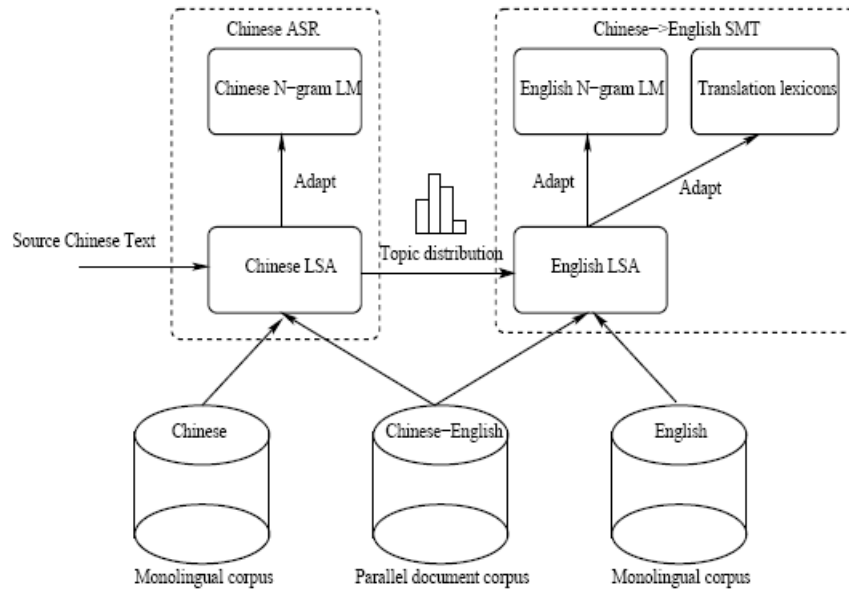Input: Speech

AM | Lex | LM | NLP | TTS

hi /h//ai/
you /j/u,
we /w//i/

hi you
you are
I am

สวัสดี ครับ

Output:
Speech & Text

# Bilingual LSA for Speech Translation



Yik-Cheung Tam, Tanja Schultz, Bilingual-LSA Based Translation Lexicon Adaptation for Spoken Language Translation, IS2007

# Multilingual Text and Speech Database



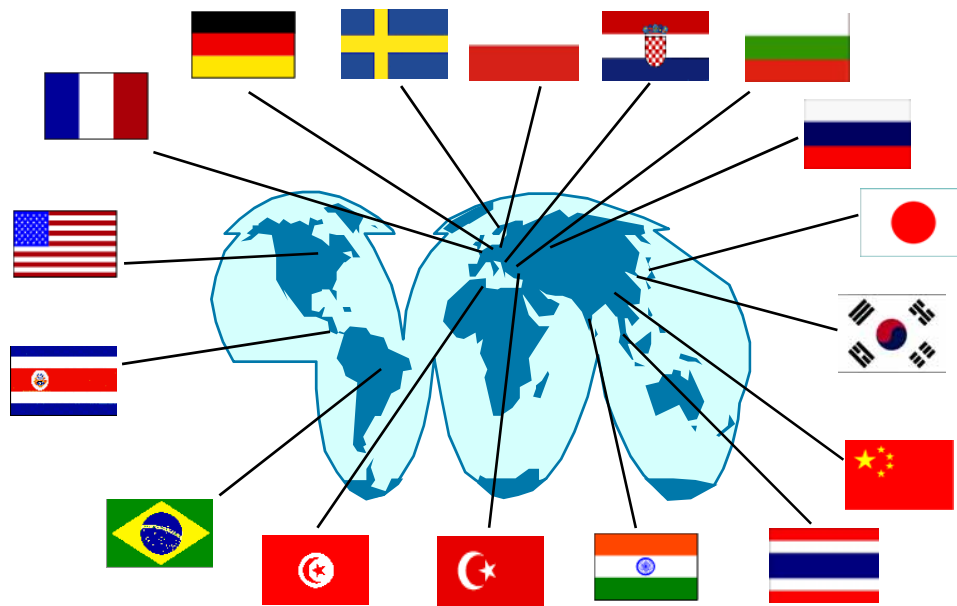Phone set & Speech data

First step for studies on Multilingual Speech Processing
and language dependencies:
Collect large amounts of data in many languages
Project GlobalPhone (since 1995)
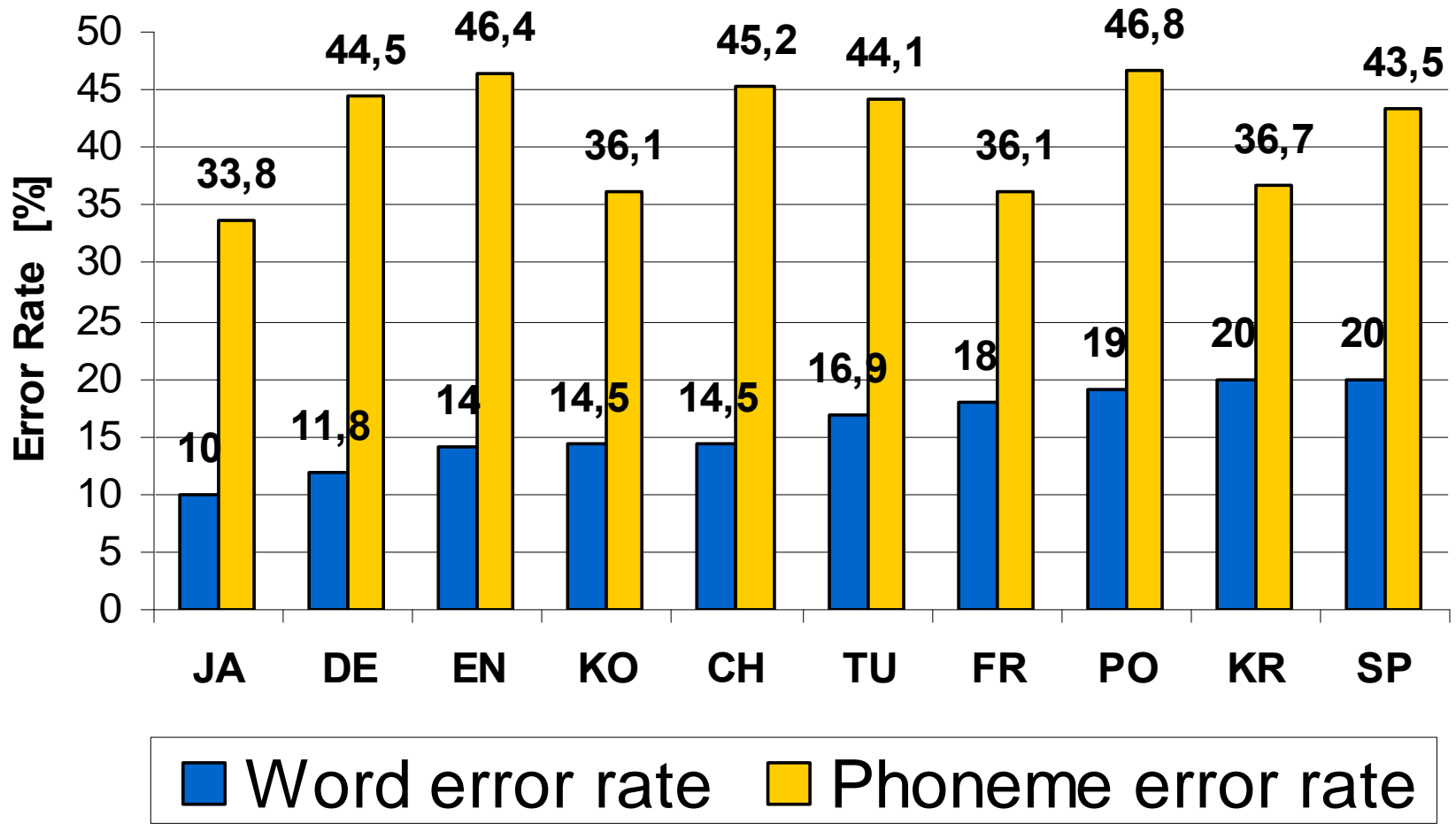
# GlobalPhone

## Multilingual Database

- Widespread languages
- Native Speakers
- Uniform Data
- Broad Domain
- Large Text Resources
    ➔ Internet, Newspaper

## Corpus

- 19 Languages … counting
- ≥ 1800 native speakers
- ≥ 400 hrs Audio data
- Read Speech
- Filled pauses annotated

Available from ELRA

| | | |
|---|---|---|
| Arabic | Croatian | Turkish |
| Ch-Mandarin | Czech | + Thai |
| Ch-Shanghai | Portuguese | + Creole |
| German | Russian | + Polish |
| French | Spanish | + Bulgarian |
| Japanese | Swedish | + ... ??? |
| Korean | Tamil | |

http://www.cs.cmu.edu/~tanja/GlobalPhone

Phone set & Speech data

Step 1:
- Uniform multilingual database (GlobalPhone)
- Build Monolingual acoustic models in many languages

## Step 2:
- Combine monolingual acoustic models to a set of multilingual "language independent" acoustic model

# Universal Sound Inventory

## Speech Production is independent from Language $\Rightarrow$ IPA

1) IPA-based Universal Sound Inventory

2) Each sound class is trained by data sharing

- Reduction from 485 to 162 sound classes
- *m,n,s,l* appear in all 12 languages
- *p,b,t,d,k,g,f* and *i,u,e,a,o* in almost all



ML-Sep                    ML-Mix                    ML-Tag

# Rapid Portability: Acoustic Models

Step 3:
- Define mapping between ML set and new language
- Bootstrap acoustic model of unseen language

# Polyphone Decision Tree Adaptation

**Problem:**

Context of sounds are language specific

How to train context dependent models for new languages?

**Solution:**

**1)** Multilingual Decision Context Trees

**2)** Specialize decision tree by Adaptation

# Rapid Portability: Acoustic Model

Pronunciation rules

„adios"       → /a/ /d/ /i/ /o/ /s/
„Hallo"       → /h/ /a/ /l/ /o/
现在广播完了    → ???

Hello

Input: Speech

AM    Lex    LM

| hi /h//ai/<br>you /j/u/<br>we /w//i/ | hi you<br>you are<br>I am |

NLP / MT

TTS

สวัสดี ดร้บ

Output:
Speech & Text

# Phoneme- vs Grapheme based ASR



Chart: Word Error Rate [%]

Legend: Phoneme, Grapheme, Grapheme (FTT)

| Language | Phoneme | Grapheme | Grapheme (FTT) |
|---|---|---|---|
| English | 11,5 | 19,2 | 18,4 |
| Spanish | 24,5 | 26,8 | |
| German | 15,6 | 14 | 12,7 |
| Russian | 33 | 36,4 | 32,8 |
| Thai | 16 | 26,4 | 18,3 |

Problem:
- 1 Grapheme ≠ 1 Phoneme

Flexible Tree Tying (FTT):
*One* decision tree
- Improved parameter tying
- Less over specification
- Fewer inconsistencies



AX-m
AX-b
IX-m

0=vowel?
0=obstruent?   0=begin-state?
-1=syllabic?   0=mid?   -1=obstruent?   0=end?

# Dictionary: Interactive Learning

```
                        Word list W
Delete wᵢ                                        Delete wᵢ
                    i:= best select
                       Word wᵢ

                      Generate           G-2-P
                   pronunciation P(wᵢ)

                        TTS                    Update G-2-P

        Yes                         No          Improve
                    P(wᵢ) okay?                  P(wᵢ)

     Lex
                        Skip
```

* Follow the work of
  Davel&Barnard

* Word list:
  extract from text

* G-2-P
  - explicit map rules
  - neural networks
  - decision trees
  - instance learning
    (grapheme context)

* Update after each $w_i$
  → effective training

User

**Build Your
System**

● Text and
prompt selection
(help)

● Audio
collection  (help)

● Phoneme
selection  (help)

●
Grapheme-to-phoneme

**Phoneme
labels for your
language:**

**P B T D K G M**

User: **awb** Language: **eng** Project: **aug19**   [Logout]
**Lexicon pronunciation creation**

**Rule entry**

3.0075187969925% Finished
new word:
**at**

system suggested pronunciation: AX T                              listen to
it  Accept Pronunciation

If you want to skip this word and work on it later, please click
Skip this word

If you don't think it's a valid word in your language, please click
Remove this word

# Lex Learner

**Build Your System**

- Text and prompt selection (help)

- Audio collection (help)

- Phoneme selection (help)

- Grapheme-to-phoneme rules (help)

**Phoneme labels for your language:**

**P B T D K G M**

User: **awb** Language: **eng** Project: **aug19**  [Logout]
**Lexicon pronunciation creation**

**Rule entry**

3.5087719298246% Finished
new word:
**Jeanne**

system suggested pronunciation: * AX N N    listen to
it  Accept Pronunciation

If you want to skip this word and work on it later, please click
Skip this word

If you don't think it's a valid word in your language, please click
Remove this word

# Issues and Challenges

o How to make best use of the human?

- o Definition of successful completion
- o Which words to present in what order
- o How to be robust against mistakes
- o Feedback that keeps users motivated to continue

o How many words to be solicited?

- o G2P complexity depends on the language (SP easy, EN hard)
- o 80% coverage hundred (SP) to thousands (EN)
- o G2P rule system perplexity

| Language | Perplexity |
|----------|------------|
| English | 50.11 |
| Dutch | 16.80 |
| German | 16.70 |
| Afrikaans | 11.48 |
| Italian | 3.52 |
| Spanish | 1.21 |

Phone set & Speech data

Hello

Input: Speech

AM    Lex    LM

hi  /h//ai/
you /j/u/
we /w//i/

hi you
you are
I am

NLP
/
MT

TTS

สวัสดี ดรับ

Output:
Speech & Text

# Statistical Parametric TTS

o Text-to-speech for Applications, Common technologies:
  o Diphone: too hard to record and label
  o Unit selection: too much to record and label

o Statistical Parametric Synthesis: "just right"
  o "HMM synthesis": ***clustergen*** trajectory synthesis
  o Clusters representing context-dependent allophones
  o Works robustly with as little as 10min speech data
  o But … Signal may sound "buzzy", can lack varied prosody

o Voice Building Process
  o Collect 300-500 utterances from single speaker, rich prompt set
  o Lexical coverage (from Lex Learner)
  o Automatic labeling from acoustic models
  o Automatic: spectral and prosodic models

o http://festvox.org  [Black and Lenzo 2000]

  o Documentation, Tools, Scripts, Examples

# TTS with very litte Data



Effect of Database Size on MCD - Multi-Lingual

Rule of Thumb for getting the best gain per amount of labor
  $\leq$ 30-60min speech: collect additional data
  \> 60min speech: improve lexicon

Kominek, J., Schultz, T., Black, A. *Synthesizer Voice Quality of New Languages Calibrated with Mean Mel Cepstral Distortion, SLTU-2008 Workshop, Hanoi, Vietnam.*

# Mono vs. Multilingual Models



Manual speaker selection

⇒ For all languages monolingual TTS performs best
⇒ Multilingual Models perform well …
    … only if knowledge about language is preserved (Multi+)
    (only small amount of sharing actually happens)

# Rapid Portability: Language Modeling

Focused Re-crawling

Text data

Hello

Input: Speech

**AM**

hi /h//ai/
you /j/u/
we /w//i/

**Lex**

hi you
you are
I am

**LM**

**NLP**

**TTS**

สวัสดี ดรับ

Output:
Speech & Text

# Language Model Building

Goal: Get as much relevant text data as possible

o Use the retrieved text data for

    o Generating recording prompts

    o Generating vocabulary lists

    o Build Language Models for ASR

Approach

1. User provides an URL or Text or Vocab list
2. Crawler retrieves N documents (web-pages)
3. Compute the statistics (TF-IDF) from the N documents
4. Terms with highest TF-IDF score form query terms
5. User may check terms for in/exclusion
6. Search engine (Google) gets URLs for the query terms
7. Crawl the top K URLs for the data

# Case Study with very small data - Hindi

o Targeted Domain in Hindi: Cooking recipes

o Data: 192 sentences, 1,523 words = 13min speech, 1 spk

o Use speech to adapt multi-lingual acoustic models

o Use transcripts to build bigram LM1

o LM2: Expanded by focused re-crawling to 159,995 words

o LM3: Expanded to 360,395 words

o Three evaluation sets (spoken by same speaker)

⇒ Focused recrawling significantly reduces the OOV rate and thus WER

| LM | word count | Word Error Rate (WER) (%) perplexity / OOV rate (%) | | | |
|---|---|---|---|---|---|
| | | split 1 | split 2 | split 3 | ave. |
| 1 | 1523 | 95.88 5.2/68.7 | 97.92 6.9/57.9 | 84.93 7.8/50.0 | 92.91 6.6/58.9 |
| 2 | 159995 | 55.15 177/16.8 | 56.25 93.4/27.4 | 51.81 165/13.4 | 54.41 145/19.2 |
| 3 | 360395 | 54.12 214/15.0 | 52.08 113/25.0 | 50.60 187/11.3 | 52.27 171/17.1 |

*John Kominek, Sameer Badaskar, Tanja Schultz, Alan W Black,* IMPROVING SPEECH SYSTEMS BUILT FROM VERY LITTLE DATA, Interspeech 2008, Brisbane

# SPICE 2005: Afrikaans – English

o Goal: Build Afrikaans – English Speech Translation System with SPICE

    o Cooperation with University Stellenbosch and ARMSCOR

    o Bilingual PhD visited CMU for 3 month

    o Afrikaans: Related to Dutch and English,
       g-2-p very close, regular grammar, simple morphology

o SPICE, all components apply statistical modeling paradigm

    o ASR: HMMs, N-gram LM (JRTk-ISL)

    o MT: Statistical MT (SMT-ISL)

    o TTS: Unit-Selection (Festival)

    o Dictionary: G-2-P rules using CART decision trees

o Text: 39 hansards; 680k words; 43k bilingual aligned sentence pairs;
Audio: 6 hours read speech; 10k utterances, telephone speech (AST)

# Time Effort

- o Good results: ASR 20% WER; MT A-E (E-A) Bleu 34.1 (34.7), Nist 7.6 (7.9)
- o Shared pronunciation dictionaries (for ASR+TTS) and LM (for ASR+MT)
- o Most time consuming process: data preparation $\rightarrow$ reduce amount of data!
- o Still too much expert knowledge required (e.g. ASR parameter tuning!)



Herman Engelbrecht, Tanja Schultz, Rapid Development of an Afrikaans-English Speech-to-Speech Translator, IWSLT 2005, Pittsburgh, PA, October 2005

# SPICE 2007: Field Experiments

o  Now targeting *more* languages in a *shorter* time frame

o  6-weeks Hands-on Course at CMU in Spring 2007
  - o  Adopt native languages of participating students as targets
  - o  Added up to 10 different languages: Bulgarian, English, French, German, Hindi, Konkani, Mandarin, Telugu, Turkish, Vietnamese

o  Teams of two students with different native language

o  Course goal was to build a simple S-2-S system and use this to communicate with each other in their mother tongue
  - o  Solely rely on SPICE tools
  - o  Build speech recognition components in two languages
  - o  Build simple SMT component in two directions
  - o  Build speech synthesis components in two languages
  - o  Report back on problems and system shortfalls

Schultz, T., Black, A., Badaskar, S., Hornyak, M., Kominek, J., *SPICE: Web-based Tools for Rapid Language Adaptation in Speech Processing Systems, Interspeech 2007, Antwerp.*

# Field Experiments (2)

o The 10 languages cover broad range of peculiarities

o <u>Writing system</u>:

  o Logographic Hanzi (Mandarin);

  o Cyrillic (Bulgarian);

  o Roman (German, French and English);

  o phonographic segmental (Telugu and Hindi);

  o phonographic featural (Vietnamese)

  o No script: Konkani

o <u>Segmentation:</u> No segmentation (Chinese); Segmentation white spaces do not necessarily indicate word (Vietnamese)

o <u>Morphology:</u> simple, low inflecting (English), compounding (German), agglutinating (Turkish) …

o <u>Sound System</u>: tonal (Mandarin and Vietnamese), stress (Bulgarian)

o <u>G-2-P</u>: straightforward (Turkish), challenging (Hindi), difficult (English), no relationship (Chinese), invented (Konkani)

# Lessons Learned

- o It is possible to create speech processing components for 10 languages in 6-weeks using SPICE
- o Each language brings new challenges
- o Many SPICE features turned out to be very helpful, e.g. only ONE speaker of Konkani in Pittsburgh, web recorder allowed remote collection of more speakers

- o Log: time spent in SPICE interface
- o Improve interface using breakdown
- o Use feedback
- o Interface allows for collaborative work

| Task | Time Spent [hh:mm] |
|------|--------------------|
| Text Collection | 8:35 |
| Audio Collection | 10:07 |
| Phoneme Selection | 4:05 |
| LM building | 1:25 |
| G-2-P specs | 1:30 |

# SPICE 2008: Cross-continental Course

o **SPICE-based course between CMU and UKA**
  - o Students at Carnegie Mellon University, PA
  - o Students at Karlsruhe University, Germany
  - o Linked by weekly meeting over VC

o **Similar to 2007 BUT distributed collaboration**
  - o Students create ASR & TTS in their native language
  - o Bonus for the ambitious: train SMT systems and create a speech-to-speech translation system

o **Evaluation includes**
  - o Time to complete
  - o Task difficulties
  - o ASR word error rate
  - o TTS voice quality

o **Fall 2008 course already in progress**

# Outline

o **The World's Languages**

    o 6900 languages – So what?

    o Language Extinction – What can the community do about it?

    o Do we need Speech Processing for all of them?

    o Is this really science – not just retraining on a new language?

o **Language Characteristics**

    o Written form, scripts, letter-to-sound relationship

    o Issues and Differences between languages

o **Challenges for Multilingual Speech Processing**

    o Lack of Resources (Money, Data, Technical Support)

    o Lack of Experts

o **Solutions**

    o SPICE: A Rapid Language Adaptation Server

    o Technologies: Leveraging off GlobalPhone & FestVox

    o Experiments and Results

o **Conclusions and Future Work**

# Conclusions

- **Challenges in Multilingual Speech Processing**

  - Well defined build processes: ASR, MT, TTS … BUT:

  - Every new language brings unseen challenges

  - Current (statistical) approaches require lots of data

  - … and native language expert and technology expertise

  - How to bridge the gap between language and tech expert?

- **Proposed solution: SPICE**

  - Learning by interaction from a cooperative (but naïve) user

  - Rapid adaptation from language universal models

  - Knowledge sharing across components

  - Development cycle: Days rather than weeks

# Next Steps

o **Continuous Server Support**

  o Improve Interface based on user feedback and lessons learned

  o Improve Language Robustness: font encoding, …

  o Software Engineering, Scaling

o **Collaboration**

  o Multiple people working on the same project

  o Leverage from archived projects

o **Cross-confirmation**

  o Multiple views for within and across project confirmation

  o Confidence measure to find appropriate combination

o **Error-blaming**

  o End-to-end system Evaluation vs Component Evaluation

  o Automatic Generation of Recommendations to improve systems

# Try This At Home

o System is online at http://cmuspice.org

o Use system for your own project
   o Create new login/passwd and project

o Preloaded Hindi Example
   o Login as
      o Login: demo
      o Passwd: demo
   o Chose project # (your birth day)

o Book on ML Speech Processing
  Elsevier, Academic Press, 2006