# ASEAN Language Speech Translation thru' U-STAR

Ms. Ai Ti Aw
Prof. Haizhou Li

# Agenda

- U-STAR Consortium

- Collaboration Projects
  - Language Analysis
  - Technology Core
  - Application & Localization

- I$^2$R Speech Translation Technologies

- Conclusions

# U-STAR Consortium



Universal
Speech
Translation
Advanced
Research

http://www.ustar-consortium.com/members.html

➡ A consortium for research and collaboration of speech translation technologies

➡ A global network for enabling real-time, location-free, multi-party communication among different language speakers

➡ An open platform based on ITU-T standards for speech translation

➡ Comprise of 32 institutes from 27 countries/regions

# ASEAN Collaborations

1. Agency for the Assessment and Application of Technology, Indonesia (BPPT)
2. Institute for Infocomm Research, Singapore ($I^2R$)
3. Institute of Information Technology, Vietnam Academy of Science and Technology, Vietnam (IOIT)
4. National Electronics and Computer Technology Center, National Science and Technology Development Agency, Thailand (NECTEC)
5. National Institute of Post Telecommunication Information Communication Technology, Cambodia (NIPTICT)
6. University of Computer Studies, Yangon, Myanmar (UCSY)
7. University of the Philippines Diliman, Philippines (UPD)

To leverage on this infrastructure and call for collaboration to develop language resources and technologies within the ASEAN community.

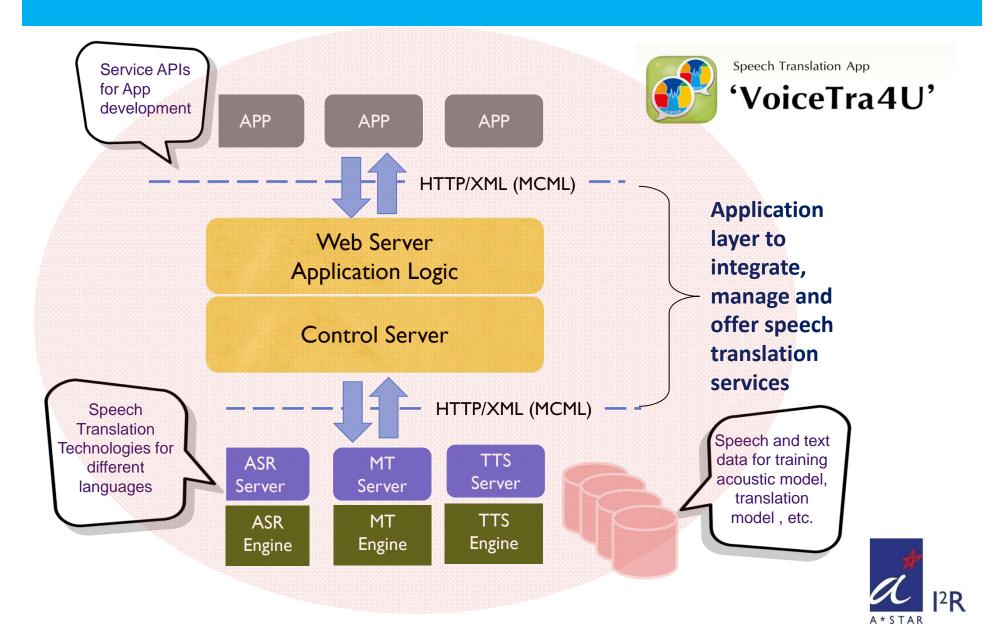To work with more institutes and promote the use of U-STAR within the ASEAN community

# Global Network for Speech Translation

- Infrastructure
  - ✓ Network and cloud with sufficient computational power and speed
  - ✓ Server system to integrate, manage and offer the technology services
- Technology Core
  - ✓ Automatic Speech Recognition (ASR), Machine Translation (MT) and Text To Speech (TTS)
  - ✓ Speech and language resources
- App Developer and Market
  - ✓ User to develop and use the app
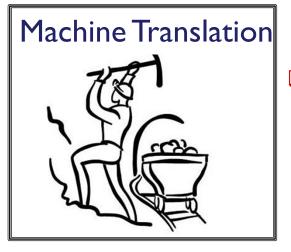
*All components are equally important !!!*

# U-STAR Architecture

# Language Analysis Project

➡ Language resources
- ❏ Linguistically tagged data

➡ Language Analysis Tools
- ❏ Morphological Analyzer
- ❏ Part-of-Speech Tagging Tool
- ❏ Named-Entity Tagger
- ❏ Word/Phrase/Sentence Segmentation Tool
- ❏ Syntactic Parser

➡ Languages
- ❏ Bahasa Indonesian, Bahasa Melayu, Chinese, English, Khmer, LAO, Myanmar, Filipino, Thai, Vietnamese
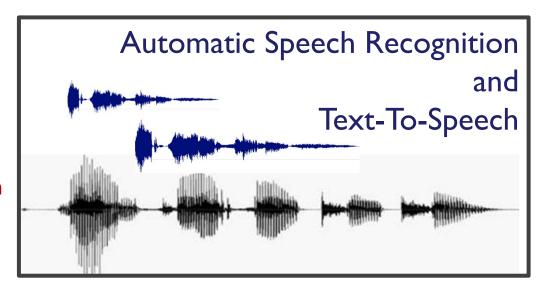
# Technology Core Project

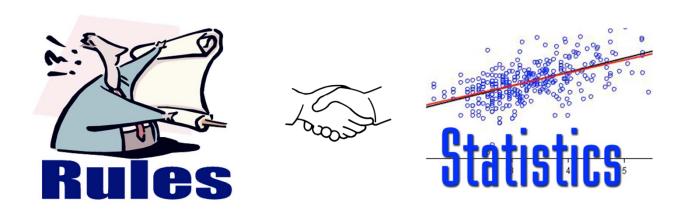## Machine Translation

- Mining of parallel data from open source through semi-automated means

- Construction of speech transcribed data for training acoustic model
- Construction of pronunciation dictionary

## Automatic Speech Recognition and Text-To-Speech

# Challenge: Overcome Low Resources

1.  How to acquire language resources ?
    - Cheaper, faster and better
2.  How to build system with limited language resources?
3.  How to leverage on human translation knowledge for SMT?

# Application Project

Tourism

eLearning

Disaster Management

Speech Recognition

Machine Translation

Text to Speech

Content Services

Medical Consultation

Customer Service

Tele-Conferencing

# Localization Project

Terjemahan Pertuturan

语音翻译



**Meruntuhkan rintangan antara bahasa-bahasa di dunia**

建立各国语言沟通的桥梁

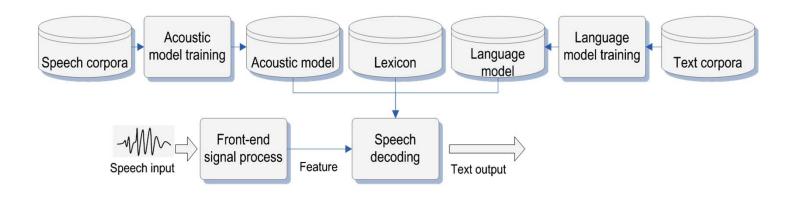**BREAKING LANGUAGE BARRIERS IN THE WORLD**

**To promote and encourage the use of U-STAR App in the local community**

# I²R ASEAN Language Technology

➡ Bahasa Melayu, Bahasa Indonesian, English Transducer based Morphological Analyzer and POS tagger
- ❑ 13,000 POS tagged sentences (Bahasa Melayu, NEWS)

➡ Thai word segmentation, POS and NER tagger (NECTEC, Thailand)
- ❑ 3 million words (NEWS), 1 million words (NE&L) annotated with WS, POS and NE
- ❑ 12,000 marked with SB in 33 paragraphs

➡ Vietnamese word segmentation, POS and NER tagger (HCMUS, Vietnam)
- ❑ 200K sentences (NEWS), 100K sentences (NE&L) annotated with WS, POS and NE

➡ Chinese word segmentation engine
- ❑ 650K sentences on informal text (SMS/Chat, CTS)

# I²R Technology – ASR



- Multilingual, Speaker-Independent Engine
  - Support continuous & naturally spoken sentences or phrases
  - Chinese, English, Malay

- Speech Training Data
  - More than 2000 hours speech data for Mandarin
  - More than 200 hours speech data for Malay

- Language Model
  - 5-gram trained with more than 100G Chinese text extracted from websites
  - 3-gram trained with more than 1G Malay text extracted from newspapers

# I²R Technology– MT

➡️ **Technology**
- ❑ Rule-based, Statistical-based, Hybrid

➡️ **Features**
- ❑ Sparse features
- ❑ Operation Sequence Model
- ❑ Minimum-Perplexity Translation Model Combination
- ❑ Source-to-Target and Target-to-Source bidirectional NNJM(Neural Network Joint Model)
- ❑ Forward and Backward RNNLM(Recurrent Neural Network Language Model)
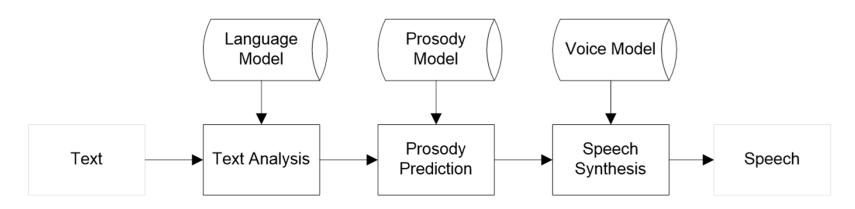
# I²R Technology – MT Resources

➡️ **Languages**

❑ Bahasa Indonesian, Bahasa Melayu, Chinese, Thai, Vietnamese

➡️ **Resources**

✓ 300K Bahasa Indonesia-English
✓ 500K Bahasa Melayu-English
✓ 400K Thai-English
✓ 500K Vietnamese-English

# I²R Technology – TTS



❑ **Multilingual Text Support**
- ✓ Understand free text, correctly read non-alphabet text and new words.
- ✓ Chinese, English, Malay

❑ **Natural prosody**
- ✓ Intelligently understand sentences and read with correct phrasing
- ✓ Natural tone, intonation, and rhythm

❑ **High quality speech**
- ✓ Apply both concatenative and statistical approaches for best quality

# Conclusions

➡ U-STAR Membership Application
- Annual meeting during Interspeech conference

➡ Research & Collaboration Project

## Contact US

**Prof. Li Haizhou**
hli@i2r.a-star.edu.sg

**Ms. Aw Ai Ti**
aaiti@i2r.a-star.edu.sg

**U-STAR Secretariat**
u-star-sec@i2r.a-star.edu.sg

Confidential