# Mining parallel corpora for multilingual machine translation system

*Dr. NGUYEN Viet Son*
*Dr. DO Thi Ngoc Diep*
*ASEAN IVO, Kuala Lumpur, 26/11/2015*

## International Research Institute MICA
**Multimedia, Information, Communication & Applications**
**UMI 2954**

**Hanoi University of Science and Technology**
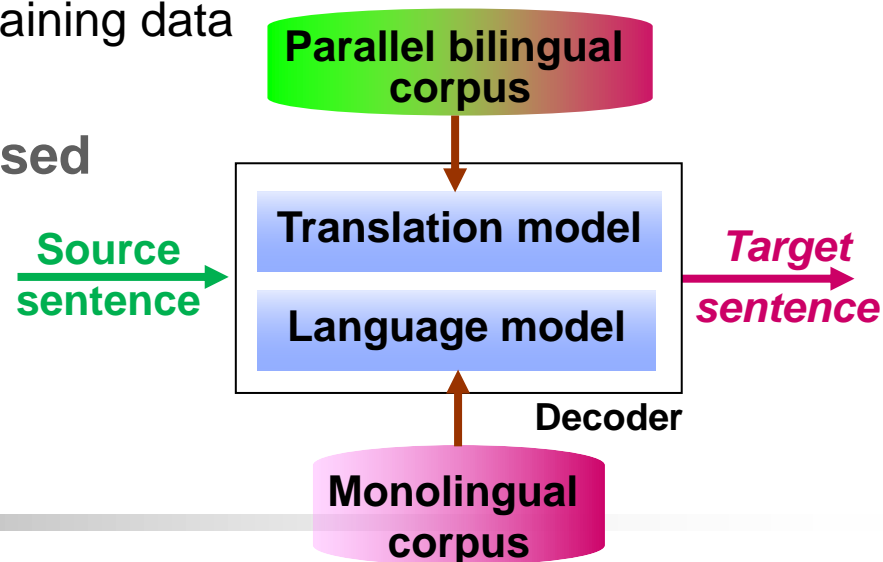**1 Dai Co Viet - Hanoi - Vietnam**

# Machine Translation Technology

- **Rules based (expert method)**
  - lexicons, morphological
  - syntactic, sematic

  → need **specialist knowledge**

  costly in **time** & **human resources**

- **Example-based**

- **Statistical**

  → **quickly build** a translation system

  based on **large parallel bilingual corpus**

  - Adapted with a lot of language pairs
  - Extract statistical information from learning databases
  - Depends: quantity & quality of training data

- **Hybrid = Statistical + Rules based**

- **For common language pairs:**
  - EN - FR
  - EN - CH
  - EN - JA
  - etc.

**Parallel bilingual corpus**

**Translation model**

**Language model**

**Source sentence** → → **Target sentence**

Decoder
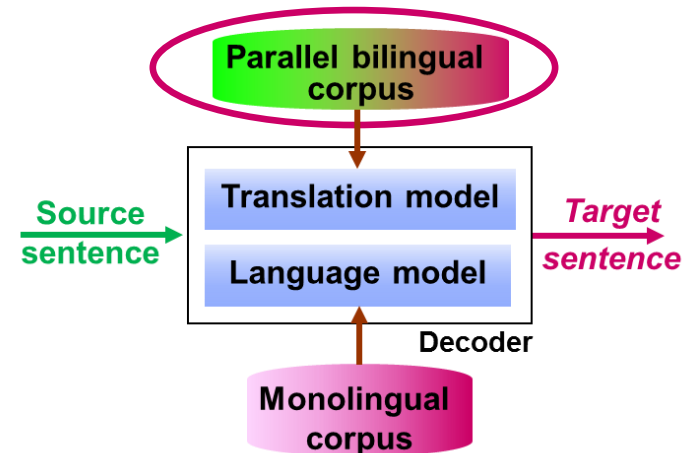
**Monolingual corpus**

# MT for ASEAN languages

- **For ASEAN (under-resourced) languages**

  - ◆ **DO NOT HAVE** any large parallel corpus for **ASEAN language pairs**

  - ◆ **FEW** available parallel corpus **even with common languages** & **LIMITTED** in some domains:

    - ★ BTEC (tourist): **NOT OPEN** for community

    - ★ OPENSubtitles (film subtitle): **SHORT** sentence, **DIALOG** format

    - ★ TATOEBA (human translation): starting in **TECHNICAL** domain

# Objectives

- **Methodology of retrieving a large parallel bilingual text corpus**

  - Firstly for ASEAN language pairs

  - Extracting parallel corpus from comparable corpora: Website, wiki, digital documents, etc.

  - Extendable for any under-resourced language pairs

- **Building a machine translation system for ASEAN languages in using the statistical approach.**

Parallel bilingual corpus → Translation model / Language model (Decoder) ← Monolingual corpus

Source sentence → Target sentence

# Our recent researches

## Extracting parallel corpus

- **Collect parallel text from multilingual websites**
- **Hypotheses:**
  - (One) multilingual website(s) of an under-resourced language
  - None of available parallel data
  - No supplementation data: lexical, morphological, dict. info.

# Our recent researches

- **Unsupervised method**
  - ◆ Based on Cross Language Information Retrieval
  - ◆ From scratch
    - ★ NO dictionary
    - ★ NO initial parallel data
  - ◆ Applicable to any language pairs
  - ◆ Verified: VI-FR, VI-EN



Multilingual web site

Cross-filtering module

Corpus to extract D

Corpus C

CLIR module

Source side

Target side

Translation module

Sentence in target lang.

Information retrieval module

Parallel sentence pairs

add

6

# Our recent researches

## Extracting parallel corpus

- **From website VnAgency**
  - ◆ **50K VI-FR** parallel sentence pairs
  - ◆ **40K VI-EN** parallel sentence pairs
- **Applications:**
  - ◆ **Webpage MT system**: VI-EN **(BLEU ~30%)**, VI-FR **(BLEU ~40%)**
  - ◆ **FR-VI Speech Translation** on smartphone **(ongoing)**

# Our recent researches

■ **Unsupervised method with triangulation**

# Our recent researches

- **2014: KH-EN language pair**

  - Try to apply our extracting methods

  - **Difficulties in Natural Language Processing** for KH language problems

  - Need **cooperation with Cambodia research groups**

- **2015: JA-VI language pair**

  - In cooperation with Multilingual Translation lab. (NICT)

  - Use available parallel text resources:

    - TED talks data: 100K pairs

    - Wiki data: 76K pairs

  - Hard task due to **different language families** (BLEU ~10%)

# Call for collaborations

- **Methodology of extracting parallel corpus:**
  - ◆ Multi-levels of extraction (sentence, phrase, fragment)
  - ◆ Independent language pairs: ASEAN, JA, EN, FR, etc.
  - ◆ From different resources: Website, Wiki, etc.
  - ◆ Multilingual extraction methods:
    - ★ **Unsupervised**: ASEAN - common language (EN/FR) pairs
    - ★ **Unsupervised with triangulation**: ASEAN languages pairs
- **Research on MT for different family language pairs**
- **Exchange researchers on MT**
- **Train students on MT**

EN / FR

high        high

ASEAN 1 - - - - - ▶ ASEAN 2

# References

[1]  Thi-Ngoc-Diep Do, Viet-Bac Le, Brigitte Bigi, Laurent Besacier, Eric Castelli. *Mining a Comparable Text Corpus for a Vietnamese - French Statistical Machine Translation System*. 4th Workshop on Statistical Machine Translation, Athens – Greece, 2009.

[2]  Thi-Ngoc-Diep Do, Laurent Besacier, Eric Castelli. *A Fully Unsupervised Approach for Mining Parallel Data from Comparable Corpora*. 14th Annual Conference of the European Association. for Machine Translation, Saint-Raphaël – France, 2010.

[3]  Thi-Ngoc-Diep Do, Eric Castelli, Laurent Besacier. *Mining Parallel Data from Comparable Corpora via Triangulation*. International Conference on Asian Language Processing, Penang – Malaysia, 2011.

[4]  Thi-Ngoc-Diep Do, Masao Utiyama, Eiichiro Sumita, *Machine Translation from Japanese and French to Vietnamese, the Difference among Language Families*, International Conference on Asian Language Processing, Oct 2015.

[5]  PhD thesis of DO Thi-Ngoc-Diep. *Extraction de corpus parallèle pour la traduction automatique depuis et vers une langue peu dotée*, 2011, Université de Grenoble, LIG - MICA

# Thank for your attention !

# Our recent researches

- **Fast collection method**
  - Lexical information, heuristic rules
  - Dictionary
  - Depend on language pairs

- **Unsupervised method**
  - Based on Cross Language Information Retrieval
  - From scratch
  - Applicable to any language pairs
  - Verified: VN-FR, VN-EN