Presentation on the 11th China-Japan Natural Language Processing Workshop

Multilingual Machine Translation Research in CAS-ICT

LIU QunLÜ Yajuan(刘群)(吕雅娟)

Miyazaki, Japan, 2011-10-30





Outline

Background

Challenges

Progress

Conclusion





Minority Languages in China

- China is a big country with 56 nationalities.
- Besides the Chinese language, there are more 80 minority spoken languages, and about 30 minority written languages in China.
- Some of the minority written languages have a very long history, such as Tibetan, Mongolian, Yi language, while some are new created written languages.





Minority Languages in China

- Some minority languages are the official languages in the autonomous regions.
- The minority languages with population more than one million:





Foreign Languages in China

- Besides English, the market demand of translation between Chinese and other foreign languages also increase very rapidly
- The demands in the translation market include:

English Japanese / Korean German / French /Russian Swedish / Dutch / Spain / Italian / Arabic Thai / Malay / Vietnamese / India





Multilingual Machine Translation Research in China

- In China, most machine translation research focus on between Chinese and English
- Only few research groups did research or development on machine translation between Chinese and Japanese or other languages
- We think it is time for the us to pay more attention to minority languages and other foreign languages besides English





Outline

Background

Challenges

Progress

Conclusion





Challenges in Multilingual MT

- Resource Scarcity
- Multiple Encoding and writing systems
- Spelling Check
- Sentence and Word Segmentation
- Morphological Analysis and Generation
- Syntactic Parsing
- Translation Modeling





Resource

- Bilingual Corpus
 - Collecting Existing Data: Limited Data
 - Purchasing
 - Web Crawling
 - Domain Limit
 - Source Bias
 - Translating: high cost





Resource

- Translation Dictionary
 - General
 - Proper Noun
 - Person names
 - Location names





Resource

Tools

- Article Alignment
- Sentence Alignment
- Word Alignment

Language Independent





Encoding and Writing System

- Some minority languages have many different encoding system or writing system.
- Encoding system:
 - Unicode

中斜抗计算好

- Many language has its own encoding system
- Writing system:
 - Traditional system
 - Romanization system
- The existing corpus usually use different encoding system.
- It is very important for the detection the encoding system and writing system and transfer between these systems.



Writing System for Mongolian

- For Mongolian, there are two different writing systems
 - Traditional writing system:
 - Traditional alphabets
 - Direction: Top Down, Left to Right
 - Being used in China
 - Cyrillic writing system

中斜於计复码

- Cyrillic alphabets (like Russian)
- Direction: Left to Right, Top Down
- Being used in Mongolia

بر ند ه ويلا



- For traditional Mongolian, there are many different encoding systems.
- Problem:
 - A character may have different glyph (shape) when it appear in different position in a word
 - Different character my have the same glyph (shape)
- There are two different kinds of encoding systems:
 - Glyph-based Encoding system
 - Pronunciation-based Encoding System





- Glyph-based Encoding system
 - One glyph, one code
 - Exact description of the shape
 - Linguistically incorrect
 - Broad used in electronic publication system
- Pronunciation-based Encoding System
 - One alphabet, one code
 - Linguistically correct
 - Adopt by Unicode and Romanized writing systems





- In machine translation, we will adopt the pronunciation-based encoding systems, since if we use the glyph-based encoding system, it will be impossible for us to do morphological analysis
- Mongolian is agglutinative language, where each word may have more 2000 different variations.
- Without morphological, there will be too many OOVs, which will bring great difficulties in MT





- It is a very difficult problem to transform between the glyph-based encoding systems and the pronunciation-based encoding system.
- Such a transform problem is a many-to-many mapping problem





Spelling Check

- Spelling check is a very big problem in some languages, such as Uygur, and Cyrillic Mongolian.
- Even for the famous writers, there are also many spelling mistakes in there works.
- For human readers, such kind of spelling mistake will make no uncomfortable.
- However, for a computer, spelling mistake is a great challenge since it will make a word unrecognized or wrongly recognized.
- So we will have to do spelling check for these kinds of languages for real application.





Sentence Segmentation

- For some languages, such as Tibetan and Thai, there are no explicit boundary between sentences.
- So sentence segmentation is a big problem for the beginning of many NLP applications in these languages.





Word Segmentation

- Word segmentation are also problems for the languages such as Tibetan and Thai.
- The technologies for Chinese word segmentation can also be used for these languages.
- However, there are also some difference.
- For example, Thai is a alphabet-based language rather than a character-based language like Chinese, so the character-labeling-based word segmentation approach can not be used directly

character-labeling ==> alphabet-labeling





Morphological Analysis

- For morphologically rich languages, such as Mongolian, Uygur, Japanese and Korean, morphological analysis is very important.
- For Mongolian and Uygur, there are no existing morphological analysis software.
- To develop a morphological system for a new language, we need help from the linguistics for the specific language, rather than a normal native speaker of that language, or we need a large stemming corpus.





Morphological Generation

- Morphological generation is also very important for machine translation which takes a morphologically rich language as target language.
- Very few researchers pay attention on this issue currently.





Syntactic Parsing

- Syntactic parsing will benefit machine translation between the language pairs with big differences.
- For the minority languages, there are almost no researches on parsing technology, and almost no treebank is available.
- For the morphologically rich languages, we wonder if the popular PCFG-like approach will get satisfied result. The problem is how to use the rich morphological information.





Translation Modeling

- Current statistical translation model regards each word as a symbol.
- It is not nature to use the rich morphological information for such languages.





Challenges in Multilingual MT

		Resou.	Spelli. Check	Sent. Segm.	Word Segm.	Encod.	Morph. Analy.	Morph. Gener.	Synt. Parsing	Trans. Modeli.
	Mongoli an	*	*			*	*	*	*	*
	Tibetan	*		*	*	*			*	
	Uygur	*	*			*	*	*	*	*
	Korean	*			*		*	*	*	*
	Vietnam ese	*	*						*	
	Thai	*		*	*				*	
INSTIT	NSTITUTE OF COMPUTING TECHNOLOGY									

Outline

Background

Challenges

Progress

Conclusion





Multilingual MT Research in ICT

- Minority Language MT
 - Mongolian ---- Inner Mongolian University
 - Tibetan ---- Qinghai Normal University
 - Uygur ----Xinjiang University
- Other Language MT
 - Korean
 - Russian
 - Thai
 - Vietnamese
 - Japanese

中斜柱计算好 INSTITUTE OF COMPUTING TECHNOLOGY



Resources

- Bilingual Corpus & Translation Dictionary
 - Uygur-Chinese
 - Mongolian-Chinese
 - Tibetan-Chinese
 - Russian-Chinese
 - Thai-Chinese
 - Vietnamese-Chinese
- Annotation Tools





🖗 align



30

INSTITUTE OF COMPUTING TECHNOLOGY

Encoding System & Writing System

Mongolian

- Romanization <=> Unicode (Traditional Mongolian)
- Traditional Mongolian <=> Cyrillic Mongolian (in development)
- Uygur
 - Traditional Uygur <=> Romanization Uygur
- Tibetan

中斜院计算所

- Banzhida code <=> Unicode
- Other 7 Encoding system <=> Unicode



Spelling Check

Uygur Spelling Check System





Mongolian Morphological Analysis

A sample: HUURNILDU/HU-DU



$$W_{1} \qquad W_{2} \\ - S_{1}A_{11}A_{12}...A_{1m} \qquad S_{2}A_{21}A_{22}...A_{2n} \\ W_{1} \qquad W_{2} \\ - S_{1} \rightarrow A_{11} \rightarrow A_{12} \rightarrow A_{1m} \rightarrow S_{2} \rightarrow A_{21} \rightarrow A_{22} \rightarrow A_{2n} \\ \end{array}$$





Mongolian Morphological Analysis

 A Directed-Graph-based Model for Joint Stemming and tagging



Mongolian Morphological Analysis

- Experiment Result
 - Corpus: 200K Words: Stemming + Annotation
 - Training: 90%, Development: 5% Test: 5%
 - Performance:
 - Stemming: 98%
 - Stemming+Annotation: 96%

设置	Pw+t	Pw-t	Fsa+t	Fsa-t
原有方法	93.0	95.1	93.0	94.7
+判别式词干词缀切分	95.2	97.8	96.1	98.0





35

Uygur Morphological Analysis

- Experiment Results
 - Corpus: 1M Words: Stemming + Annotation
 - Training: 90%, Development: 5% Test: 5%
 - Performance:
 - Stemming: 94.7%
 - Stemming+Annotation: 92.6%
- Chinese Person Name Recognition in Uygur Text
 - Chinese Person Names may have inflection (in development)

中斜院计算好 INSTITUTE OF COMPUTING TECHNOLOGY



Tibetan Sentence Segmentation and Word Segmentation

- Tibetan Sentence Segmentation: 97%
- Tibetan Word Segmentation: 96%





Mongolian Dependency Parsing

 Mongolian Dependency Parsing based on Bilingual Restriction



38



Multigrain Machine Translation

Uygur → Chinese
Mongolian → Chinese

Grain	BLEU	Grain	BLEU
Word	0.3962	Word	0.1373
Stem	0.4024	Stem	0.1610
Morpheme	0.3989	Morpheme	0.1615
Multigrain	0.4165	Multigrain	0.1818





Outline

Background

Challenges

Progress

Conclusion





Conclusion

- It's time for us to pay more attention to languages other than English
- Many new challenges will be encountered when we work on multilingual machine translation
- We have made some progress in multilingual machine translation
- However, there are still a long way to go ...







ر مخمەت ७९॥ इसम्बें नगपन्त्रेमके॥

Welcome to our website: <u>http://nlp.ict.ac.cn/new/</u>

http://nlp.ict.ac.cn/new/demo/demo-MT-index.php





1

£

STATI