

Key problem = translation (bis)

2. Query translation – Translate query into document language(s)
 - flexibility (translation on demand)
 - less text to translate
- Cons:
- less precise (2/3-word queries)
 - The retrieved documents need to be translated (gist) to be readable.

Needs for CLIR and MLIR

- Why is CLIR and MLIR useful?
 - An information searcher wants to retrieve relevant documents in whatever language.
 - Intelligence:
 - CIA,
 - companies (finding competing companies, finding calls for tenders, ...)
 - A user speaking several languages also may want an MLIR to avoid typing the same query several times in different languages.

Cross-lingual and Multilingual IR

Jian-Yun Nie
DIRO, University of Montreal
nie@iro.umontreal.ca
<http://www.iro.umontreal.ca/~nie>

How to translate

1. Machine Translation (MT)
 2. Bilingual dictionaries, thesauri, lexical resources, ...
 3. Parallel texts: translated texts
- Parallel texts encompass translation knowledge

Problems in CLIR and MLIR

- CLIR and MLIR are based upon monolingual IR; all the problems of monolingual IR.
Document representation v.s. query representation
- Problems due to the differences in languages.
 - Documents in E, F, I, ...
 - Query in E
- 1. Documents in F ↓ document representation in E
Query in E ↓ query representation in E
2. Query in E ↓ query representation in F
Documents in F ↓ document representation in F

Definitions

- Cross-lingual (cross-language) IR (CLIR):
Retrieval of documents in a language different from that of a query
- synonym: bilingual IR
- Multilingual IR (MLIR)
Retrieval of documents in several languages from a query

Outline

- Needs for CLIR and MLIR
 - Problems in CLIR and MLIR
 - Approaches to CLIR
 - MT
 - Bilingual dictionary
 - Parallel texts
 - Approaches to MLIR
 - Experiments and evaluation campaigns
 - A better integration of translation and retrieval?

Key problem = translation

1. Document translation - Translate documents into the query language
- Pros:
- translation may be (theoretically) more precise
 - documents become "readable" by the user
- Cons:
- huge volume to be translated
 - impossible to translate them in all the languages (translate English documents in F, I, ..., Chinese, Thai, ...)

Approach 1: Using MT

- Seems to be the ideal tool for CLIR and MLIR (if the translation quality is high)
Query in F → Translation in E
MT ↓
Documents in E
- Problems:
 - Quality
 - Availability
 - Development cost

Word-by-word translation

- Select the first translation word
 - Assumption: The first translation is the most frequently used translation
 - Depends on the organization of the dict.
 - Not the case for: Freedict:
access: attaque, accéder, intelligence, entrée, accès
 - Problems:
 - May select a wrong translation
 - context-independent
 - May miss synonyms

Word/term level

- Choose one translation word
 - E.g. organic – organique
 - Better to keep all the synonyms (organique, biologique)? – query expansion effect
- Sometimes, use context to guide the selection of translation words
 - The boy grows: grandir
 - ... grow potatoes: cultiver

Problems of MT

- Translation quality for CLIR and MLIR
 - Wrong choice of **translation word/term**
 - organic food – nourriture organique
 - Train skilled personnel - personnel habile de train (ambiguity)
 - Wrong syntax
 - human-assisted machine translation - traduction automatique humain-aidee
 - Unknown words
 - Personal names:
Bérégovoy ⇄ Береговой, Beregovoy
Deng Xiaoping, Deng Hsiao-ping

Word-by-word translation (bis)

- Concatenate all the translation words
 - access: attaque, accéder, intelligence, entrée, accès
- Covers all the possible translation
- Keeps ambiguity and incorrect translations
- Noisy query translation

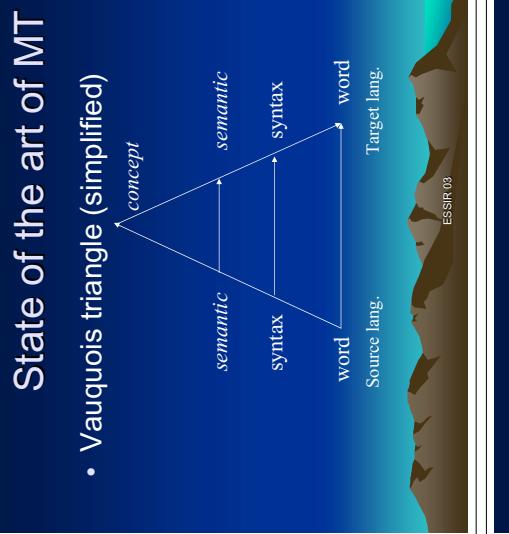
Syntax – unused effort for CLIR?

- Current IR approaches based on words (bag of words)
- Efforts on determining the correct syntax not used for IR
- However, useful for disambiguation (stem)
international terrorism v.s. tree stem

State of the art of MT

- Vauquois triangle (simplified)

```
graph TD; concept((concept)) --> semantic[semantic]; semantic --> syntax[syntax]; syntax --> word[word]; word --> targetLang[Target lang.]; word --> sourceLang[Source lang.]
```



State of the art of MT (cont'd)

- General approach:
 - Word / term: dictionary
 - Syntagm (phrase)
 - Syntax
 - “semantic”

Approach 2: Using bilingual dictionaries

- General form of dict. (e.g. Freedict)
 - access: attaque, accéder, intelligence, entrée, accès
 - academic: étudiant, académique
 - branch: filiale, succursale, spécialité, branche
 - data: données, matériel, data
- Approaches
 - For each word in a query
 - Select the first translation word
 - Select all the translation words
 - For all the query words
 - Select the set of translation words that produce the highest **cohesion**

Translate the query as a whole

- Best global translation for the whole query
- Candidates:
 - For each query word
 - Determine all the possible translations (through a dict.)
 - 2. Selection
 - select the set of translation words that produce the highest **cohesion**

Example of aligned sentences

Débat	Artificial intelligence	A Débat
L'intelligence artificielle		
Depuis 35 ans, les spécialistes d'intelligence artificielle cherchent à construire des machines pensantes.	Attempts to produce thinking machines have met during the past 35 years with a curious mix of progress and failure.	

Leurs avancées et leurs insuccès alternent curieusement.

Les symboles et les programmes sont des notions purement abstraites.

Two further points are important.

First, symbols and programs are purely abstract notions.

Parallel texts (cont'd)

- Training a translation model
- Principle:
 - train a statistical translation model from a set of parallel texts: $p(t_j|s_i)$
 - Principle: The more s_i appears in parallel texts of t_j , the higher $p(t_j|s_i)$.
- Given a query, use the translation words with the highest probabilities as its translation

Initial probability assignment

$$t(t_j|s_i, A)$$

même even
un a
cardinal cardinal
n' is
est not
pas safe
à from
l' drug
abri cartel
des cartels
de la
drogue .

even
a
cardinal
is
not
safe
from
drug
cartels

Application of EM: $p(t_j|s_i, A)$

même even
un a
cardinal cardinal
n' is
est not
pas safe
à from
l' drug
abri cartel
des cartels
de la
drogue .

Cohesion

- Cohesion ~ frequency of two translation words together

E.g.

- data: données, matériel, data access: attaque, accéder, intelligence, entrée, accès
- (accès, données) 152 *
- (accéder, données) 31
- (données, entrée) 21
- (entrée, matériel) 3

...
Freq. from a document collection or from the Web (Grefenstette 99)

Approach 3: parallel texts

- Parallel texts contain possible translations of query words
- First exploration: using IR methods
 - Given a query in F
 - Find relevant documents in the parallel corpus
 - Extract keywords from their parallel documents, and consider them as a query translation



Corresponding → Words in E

doc. F

Rel.

doc. F

ESSIR 03

- Second approach: LSI
 - In monolingual LSI, singular value decomposition (SVD) is able to group synonyms in the created structure
 - For CLIR, SVD is performed on a parallel corpus
 - The LSI encompasses translation relationships (special case of cross-language synonymy)
 - Query in F can match directly documents in E in LSI
- Problems:
 - Computational complexity
 - Number of singular value to choose (empirical setting)
 - Coverage of the parallel texts w.r.t. semantics

parallel texts (cont'd)

- Assumption:
 - The order of sentences in two parallel texts is similar
 - A sentence and its translation have similar length (length-based alignment, e.g. Gale & Church)
 - A translation contains some "known" translation words, or cognates (e.g. Simard et al 93)

Sentence alignment

Summary of the experimental results

- High-quality MT is still the best solution
- TM based on parallel texts can match MT
 - Dictionary
 - Simple utilization is not good
 - Complex approaches improve quality
 - The performance of CLIR usually lower than monolingual IR (between 50% and 90% of monolingual in general)

Summary of the existing approaches to CLIR

CLIR = Query Translation + IR

- Trend: Integrate QT with IR
 - QT is one step in the global IR process
 - E.g. Kraaij, Nie and Simard, 2003
 - Using language model



CLIR as a special case of query expansion

- Query expansion:
 - $Q \xrightarrow{\text{related terms}} Q' \xrightarrow{\text{trans. rel.}} D$
 - CLIR
 - $Q \xrightarrow{\text{inferred relation}} Q' \xrightarrow{\text{term relation}} D$
 - Translation relation ~ term relation
 - Inferential IR
 - Infer an expression Q' such that when Q' is satisfied, Q is too.

CLIR Results: C-E

- Test collections from TREC
 - Chinese: People's Daily, Xinhua news agency
 - English: AP

	C-E	E-C
Monolingual	0.3861	0.3976
Dictionary (EDict)	0.1530 (39.6%)	0.1427 (35.9%)
TM	0.2063 (53.4%)	0.2013 (50.6%)
TM + Dict	0.2811 (72.8%)	0.2601 (65.4%)

- MT system:
 - E-C: 0.2001 (50.3%)
 - C-E: (56 - 70%)

43

Problems of using parallel corpora

- Not strictly parallel (Web)
- Coverage
- In a different domain than the documents to be retrieved
- Not applicable to “minor” languages

44

Further verification of parallelism

- Download files (for verification with document contents)
- Compare file lengths
- Check file languages (by an automatic language detector – SIlC)
- Compare HTML structures
- (Sentence alignment)

45

Mining Results

- French-English
 - Exploration of 30% of 5,474 candidate sites
 - 14,198 pairs of parallel pages
 - 135 MB French texts and 118 MB English texts
- Chinese-English
 - 196 candidate sites
 - 14,820 pairs of parallel pages
 - 117.2M Chinese texts and 136.5M English texts
- Several other languages I-E, G-E, D-E, ...
 - ESSIR 03

46

Translation problems with TM

- Compound terms (e.g. *pomme de terre* – earth, apple, potato, soil, ...)
- Coverage (personal names, unknown words)
- Ambiguous words remain ambiguous (*drug* – médicament, drogue)
- Possible solutions
 - Recognize compounds before model training
 - treat named entities differently
 - Combine with dictionaries

47

	F-E / (Trec6)	F-E / (Trec7)	E-F / (Trec6)	E-F / (Trec7)
Monolingual	0.2865	0.3202	0.3686	0.2764
Systran	0.3098 (107.0%)	0.3293 (102.8)	0.2727 (74.0%)	0.2327 (84.2%)
Hansard TM	0.2166 (74.8%)	0.3124 (97.6%)	0.2501 (67.9%)	0.2587 (93.6%)
Web TM	0.2389 (82.5%)	0.3146 (98.3%)	0.2504 (67.9%)	0.2289 (82.8%)

- Web TM comparable to Hansard TM

48

ESSIR 03

49

ESSIR 03

Future problems

- Develop better translation tools for IR (e.g. for special types of data such as personal names)
- Integrating multiple translation results
- Translate non-English languages
- Integration of query translation and retrieval process
- Develop approaches to MLIR
- Make the retrieved documents readable

Some References

- Introduction and survey:
 - D. Oard, several survey papers <http://www.glue.umd.edu/~oard/research.html#litreview>
- MT-based approaches
 - Chen 02, A. Chen, Cross-language retrieval experiments at CLEF 2002, in QLIEF-2002 working notes pp. 5-20, 2002.
 - Savoy 02, J. Savoy, Report on CLEF-2002 experiments: Combining multiple sources of evidence, in QLIEF-2002 working notes, pp. 31-46, 2002.
 - Grefenstette 98, J. Grefenstette, A resource for example-based MT tasks, proc. ASIRB translating on the computer 21 conference London, 1999.
 - Grefenstette 99, J. Grefenstette, D. Evans, Resolving translation ambiguity using monolingual corpora - A report on CLEF-2002 experiments, in CLEF 2002 working notes 2002, pp. 115 - 202.
 - Gao et al. 02, J. Gao, L.-Y. Nie, H. He, W. Chen, M. Zhou, Resolving Query Translation Ambiguity using a Decaying Co-occurrence Model and Syntactic Dependencies Relations, 25th ACM SIGIR, 2002, Tampere, pp. 153-160.
 - Nie et al. 00, Y. Nie, M. Simard, P. Isabelle, R. Durand, "Cross-Language Information Retrieval Based on Parallel Texts and Automatic Mining of Parallel Texts in the Web," 22nd ACM SIGIR, Berkeley, 1999, pp. 72-81.
 - Nie et al. 03, Y. Nie, M. Simard, Embedding Web-based Statistical Translation Models in Cross-Language Information Retrieval, Computational Linguistics 2003, 2003.
 - Kraaij et al. 03, W. Kraaij, R. D. Brown, R. E. Frederick, Translating Information Retrieval Learning from Bilingual Corpora, Artificial Intelligence 103, 323-345, 1998.
 - Nie 02, J.-Y. Nie, F. Jin, A multilingual approach to multilingual information retrieval, Advances in Cross-language Information Retrieval, Third Workshop on the Cross-language Evaluation Forum, CLEF 2002, Rome, Sept. 2002, pp. 101-110.
- Based on parallel texts

ESSIR 03
52

- LSI
 - M. Litman and S. Dumais and T. Landauer, Automatic Cross-Linguistic Information Retrieval using Latent Semantic Indexing, in G. Grefenstette (Ed.), Cross-Language Information Retrieval, 1997, <http://tesser.nj.nec.com/damais/7automtic.html>.
 - Mori et al. 01, C. Mori et al., Cross-Language Information Retrieval based on LSI with multiple space, NTCIR-2, <http://research.nii.ac.jp/ntcir/workshop/onlineProceedings/2monlr.pdf>
 - CLIR & MLIR campaigns
 - TREC CLIR track (<http://trec.nist.gov>)
 - NTCIR (<http://ntcir.nict.go.jp/ntcir/ntcir2002/ntcir2002.html>)
 - Sentence alignment
 - Gale & Church 93, W. Gale, K.W. Church, A program for aligning sentences in bilingual corpora, Computational Linguistics, 19: 175-192, 1993.
 - Simard et al.92, M. Simard, G. Foster, P. Isabelle, Using Cognates to Align Sentences in Parallel Corpora, Proceedings of the 4th International Conference on Theoretical and Methodological Issues in Machine Translation, Montreal, 1992.
 - Translation models
 - (Brown et al. 93) P. F. Brown, S.A.D. Pietra, V. D. J. Pietra, and R. L. Mercer, The mathematics of machine translation: Parameter estimation, Computational Linguistics, vol. 19, pp. 263-312, 1992.
 - P. Resnik, Parallel Stands: A preliminary investigation into mining the Web for bilingual text, AMTA98, 1998.
 - J. Chen, I.Y. Nie, Parallel text mining for Chinese-English cross-language information retrieval, NAACL-ANLP, Seattle, May 2000.

ESSIR 03
53

ESSIR 03
54

Experiments

- CORI works well for monolingual distributed IR
- For MLIR, Raw score and Normalized score work better than CORI in CLEF01 and CLEF02
 - The effectiveness of MLIR is lower than CLIR (bilingual IR)

MLIR (cont'd)

- MLIR = mixed query for mixed doc. Collection (Chen 02, Nie 02)
 - Translate the query into all the languages
 - Concatenate them into a mixed query
 - IR using mixed query on mixed documents
- Avoiding merging
 - homograph in different languages (but, pour, ...)
 - Possible improvement: distinguishing language (add a tag to the indexes, e.g. but_f, pour_e)

ESSIR 03
55

Problems in MLIR

- Translation
- Merging different (incompatible) retrieval results
 - Is it necessary to produce a single mixed result list ?

Merging (cont'd)

- Normalized score
 - $S' = S / S_{\max}$
 - $S' = (S - S_{\min}) / (S_{\max} - S_{\min})$
 - CORI
 - $S' = S * (1 + (S - S_{\text{avg}}) / S_{\text{avg}})$
- Idea: modify the raw score according to the average score for a collection (language)

ESSIR 03
56

ESSIR 03
57

Multilingual IR

- **MLIR = CLIR + merging**
 - Translate the query into different languages
 - Retrieve doc. in each language
 - Merge the results into a single list

Merging – often used approaches

- Round-robin
 - Take the first from the list of F, E, I, ...
 - Take the second from the list of F, E, I, ...
 - ...
- Assumption: similar number of rel. doc., ranked similarly
- Raw score
 - Mix all the lists together
 - Sort according to the similarity score
- Assumption: similar IR method, collection statistics

ESSIR 03
58

Merging (cont'd)

- Normalized score
 - $S' = S / S_{\max}$
 - $S' = (S - S_{\min}) / (S_{\max} - S_{\min})$
 - CORI
 - $S' = S * (1 + (S - S_{\text{avg}}) / S_{\text{avg}})$
- Idea: modify the raw score according to the average score for a collection (language)

ESSIR 03
59

ESSIR 03
60