#### **CROSS-LANGUAGE INFORMATION RETRIEVAL** AND BEYOND

Jian-Yun Nie

University of Montreal http://www.iro.umontreal.ca/~nie

#### Problem of CLIR

 Cross-language IR (CLIR) • Use a query in a language (e.g. English) to retrieve documents in another language (Chinese)

3

4

- Multilingual IR (MLIR)
- Use a query in one language to retrieve documents in several languages

#### Outline

- What are the problems in CLIR?
- Recall: General approaches to IR
- The approaches to CLIR proposed in the literature

2

- Their effectiveness
- Remaining problems
- Applications

### • In 1970s, first papers on CLIR

•	TREC-3 (1994)	Spanish (monolingual): El Norte Newspaper	SP 1-25
*	TREC-4 (1995)	Spanish (monolingual): El Norte Newspaper	SP 26-50
*	TREC-5 (1996)	Spanish (monolingual): El Norte newspaper and Agence France Presse	SP 51-75
		Chinese (monolingual): Xinhua News agency, People's Daily	CH 1-28
*	TREC-6 (1997)	Chinese (monolingual), The same documents as TREC-6	CH 29-54
		CLIR:	
		English: Associated Press	CL 1-25
		French, German: Schweizerische Depeschenagentur (SDA)	
	TREC-7 (1998)	CLIR:	
		English, French, German, Italian (SDA)	CL 26-53
	TREC-8 (1999)	+ German: New Zurich Newspaper (NZZ) CLIR (English, French, German, Italian): as inTREC-7	CL 54-81
•	TREC-9 (2000)	English-Chinese: Chinese newswire articles from Hong Kong	CH 55-79
•	TREC 2001	English-Arabic: Arabic newswire from Agence France Presse	1-25
•	TREC2002	English-Arabic: Arabic newswire from Agence France Presse	26-75



#### Why is it necessary to do CLIR?

- When the relevant information only exists in another language
- Local information
- · Information about local companies, events, ...
- Patents
- When there is not enough relevant information in the given language
  - Recall-oriented search
- · When one can read several languages
- · Avoid submitting multiple queries in different languages
- When the information is language-independent
  - Image
  - · Programs, formal specification, ...

#### CLIR problem

- English query → Chinese document
- Does the Chinese document contain relevant information?
- Readability: In many cases, the translation of retrieved documents into the language of the query is still necessary (goal of machine translation)

#### Strategies for CLIR

- Translate the query
- Translate the documents
- Translate both query and documents into a third language (pivot language)
- (Related) Transitive translation English-> Chinese

- · Less effective than direct translation
- Used only when direct translation is impossible



- Query is short: less contextual information available for translation
- Query translation is more flexible: the user indicates the language(s) of interest, and there is less to be translated
- Document translation: richer context
- More to translate
- Impossible to predetermine in which languages the document should be translated
- · Generally, comparable effectiveness

#### How to translate

- 1. Machine Translation (MT)
- 2. Bilingual dictionaries, thesauri, lexical resources, ...
- 3. Parallel or comparable corpora

#### Recall: Basics on IR

- What is IR problem
- Basic operations in IR systems
- Models

#### What is IR?

- IR aims at retrieving relevant documents from a large document collection for a user's information need
- Information need  $\rightarrow$  query
- · Document collection: a static set of documents
- Relevance: the document contains desired information of the user





# Traditional document and query representation

- Using keywords (terms)
- A term is independent from any other term
- Term  $\leftarrow \rightarrow$  Meaning
- In NL, a term (word) may mean the same thing as another term (synonymy)

15

 It may mean different things in different contexts (polysemy)









#### **Common Preprocessing Steps**

- Strip unwanted characters (e.g. punctuation, numbers, etc.).
- · Break into tokens (keywords) on whitespace.
- Analyse structure
- Stem tokens into "root" words
   computational → comput
- Remove common stopwords (e.g. a, the, it, etc.).
- Detect common phrases (possibly using a domain specific dictionary).
- Build inverted index (keyword  $\rightarrow$  list of docs containing it + positions).

#### **Retrieval Models**

- A retrieval model specifies the details of:
- Document representation
- Query representation
- Retrieval function
- Determines a notion of relevance.

#### **Classes of Retrieval Models**

- · Boolean models (set theoretic)
  - Extended Boolean
- Vector space models (statistical/algebraic)
- Generalized VSM
- Latent Semantic Indexing
- Probabilistic models
- Language models

#### **Boolean Model**

- A document is represented as a set of keywords.
- Queries are Boolean expressions of keywords, connected by AND, OR, and NOT, including the use of brackets to indicate scope.
  - [[Rio & Brazil] | [Hilo & Hawaii]] & hotel & !Hilton]
- Output: Document is relevant or not. No partial matches or ranking.

#### • Assumption: Each term corresponds to a dimension in a vector space • Vector space = all the keywords encountered in the collection $< t_1, t_2, t_3, ..., t_n >$ • Document $D = < a_1, a_2, a_3, ..., a_n >$ $a_i = weight of t_i in D$ • Query $Q = < b_1, b_2, b_3, ..., b_n >$ $b_i = weight of t_i in Q$ • R(D,Q) = Sim(D,Q)







Probabilistic model
• Given D, estimate P(R D,Q) and P(NR D,Q) $P(R Q,D) = \frac{P(D RQ) * P(R Q)}{P(D Q)}$
• P(D Q), P(R Q) constant $\rightarrow$ P(R Q,D) $\propto$ P(D R) • D = {t <sub>1</sub> =x <sub>1</sub> , t <sub>2</sub> =x <sub>2</sub> ,} $X_i = \begin{cases} 1 \text{ if present} \\ 0 \text{ if absent} \end{cases}$
$P(D   R Q) = \bigcap_{\substack{(t_i = x_i) \ D}} P(t_i = x_i   R_0) $ $P(D   NR Q) = \bigcap_{\substack{(t_i = x_i) \ D}} P(t_i = x_i   NR_0) $ $Odd(D) = \log \frac{P(D   R Q)}{P(D   NR Q)}$



Some additional factors
<ul> <li>Smoothing (Robertson-Sparck-Jones formula)</li> </ul>
$Odd(D) = \sum_{t_i} x_i \log \frac{(r_i + 0.5)(N - R_i - n_i + r_i + 0.5)}{(R_i - r_i + 0.5)(n_i - r_i + 0.5)} = \sum_{t_i \in D} x_i w_i$ • When no sample is available:
P(t=1 R)=0.5, P(t=1 NR)=(n <sub>i</sub> +0.5)/(N+0.5)≈n <sub>i</sub> /N
<ul> <li>Use relevance feedback to get more samples for more precise estimation (usually unavailable)</li> </ul>



32	
Language models	
Question: Is the document likelihood increased when a query is submitted?	
$LR(D,Q) = \frac{P(D \mid Q)}{P(D)} = \frac{P(Q \mid D)}{P(Q)}$ (Is the query likelihood increased when D is retrieved?) - P(Q D) calculated with P(Q M_{D}) - P(Q) estimated as P(Q M_{C})	
$Score(Q, D) = \log \frac{P(Q \mid M_D)}{P(Q \mid M_C)}$	











	38
<ul> <li>Approach 1: Using MT</li> <li>Seems to be the ideal tool for CLIR and I the translation quality is high)</li> </ul>	MLIR (if
Query in F ──────────────────────────────────	
<ul> <li>Typical effectiveness: 80-100% of the monolingual effectiveness</li> <li>Problems: <ul> <li>Quality</li> <li>Availability</li> </ul> </li> </ul>	

Google	state of the art in information retrieval
Search	About 6,560,000 results (0.29 seconds)
Everything	[PDF] Content-based Multimedia Information Retrieval: State of the Art
	and
mages	www.liacs.nl/home/mlew/mir.survey16b.pdf
aps	File Format: PDF/Adobe Acrobat - Quick View
lideos	b) MS LEW - Cited by 119 - Related ancies In ACM Transactions on Multimedia Computing, Communications, and Applications, Feb. 2006. Content-based Multimedia Information Retrieval: State of the Art
News	
Shopping	[PDF] State of the Art in Cross-Lingual Information Retrieval 202.141.152.9/clir/papers/clir.pdf
More	File Format: PDF/Adobe Acrobat - Quick View by A Ramanathan - Cited by 1 - Related articles
Any time	<ul> <li>State of the Art in Cross-Lingual Information Retrieval. Ananthakrishnan Ramanathan.</li> <li>National Centre for Software Technology. Rain Tree Marg, Sector 7, CBD</li> </ul>
Past hour	
Past 24 hours	[PDF] State of the Art in Web Information Retrieval
Past week	paginas.fe.up.pt/~ssn/2005/prodel/soa-webir.pdf
Past month	hy S Nunes - 2006 - Cited by 3 - Related articles
Past year	the concept of temporal analysis of the web for information retrieval. 1 Introduction.
ustom range	This is a state of the art report on the broad subject of Information Retrieval on
All results	INFORMATION RETRIEVAL; STATE OF THE ART Don R. Swanson
Sites with images	dl.acm.org/ft_gateway.cfm?id=1460717&type=pdf
telated searches	by DR Swanson - 1961 - Cited by 16 - Related articles
fisited pages	problems. A proper perspective for assessing the state of the information retrieval art can best be achieved through considering first the broader problem context.
ot yet visited	can been be consider an eage considering that the product problem context.
Dictionary	IPDEL State of the Art on Systems for Data Analysis Information
Reading level	Retrieval and
learby	www.infobiomed.org/paginas_en/INFOBIOMED_D13_final.pdf
Sustom Location	File Format: PDF/Adobe Acrobat
translated lofeign pages	State of the Art on Systems for Data. Analysis, Information Retrieval and Decision

	40	
Web	Translated foreign pages	×
Images	Translated results for state of the art in information retrieval	
Maps	- My language: English v	
Videos	French V état de l'art dans la recherche d'information - Edit	
News	Add language V - Automatically select languages to search	
Books	(PDF) 20 - state of the art	
Places	www.innovatech.be/files/espace_innovation//20etat_de_lart.pdf	
Blogs	Translated from: French Summary Making a state of the art is to gather as much information about a subject	
Flights	or aFind the relevant information sources to consultA literature search is often	
Discussions	<ul> <li>done by using keywords.</li> <li>+ Show original text</li> </ul>	
Applications	IPPE Personalization of information: an overview of the state of the ar	
Patents	apmd.prism.uvsq.fr//Articles/	
Fewer	Translated from: French	
	Personalization of information: an overview of the state of the art andaddressed in the community Information Retrieval (IR) community Databases	
The web	+ Show original text	
Pages from Canada	state of the art definition state of the art and synonymous with state	
	dictionnaire.sensagent.com/etat+de+i+ <b>aru</b> /tr-tr/ Translated from: French	
Any time	Definitions of state of the art, synonyms, antonyms, derivatives of state of the art, A	
Past hour	window (pop-into) of information (main content of Sensagent) is invokedMake the state of the art is to seek all information, publications	

	41					
Everything	Translated foreign pages	×				
Images	Translated results for state of the art in information					
Maps	retrieval - My language: English V					
Videos	Language Translated query					
videos	Chinese (Simplified) V 最先进的国家在信息检索 - Edit					
News	Add language v - Automatically select languages to search					
Shopping	National Library of China					
More	www.ccnt.gov.cn/xxfb/igsz/zsdw/gt/					
	Translated from: Chinese (Simplified)					
	, Wenjin Street premises completed (now the National Library of Ancient Museum), to					
Any time	become the largest, most state-of-the-art libraryNational Library with a wealth of					
Past nour	service functionsUnified portal through the national digital library resources.					
Past 24 nours	information retrieval services, push the one-stop service; open + Show original text					
Past week						
Past voar						
Custom range	State Intellectual Property Office Patent Search Consulting Center					
ouston range	Baidu Baike					
	baike.baidu.com/view/551522.htm					
All results	Translated from: Chinese (Simplified)					
Sites with images	State Intellectual Property Office Patent Search Advisory Center: Founded in April 1995,					
Related searches	the Department of State Intellectual Property Office (formerlyAuthority to the field of					
Visited pages	domestic science, technology and intellectual property to provide advanced information					
Not yet visited	retrieval and consulting services The most advanced search technology to provide					
Dictionary	analysis and other high-and services					

Query translation using MT • Query translation: to make query compara documents	ble to
<ul> <li>Similarities with MT</li> <li>Similar translation problems</li> <li>Similar methods can be used</li> </ul>	
- A good MT $\rightarrow$ good CLIR effectiveness	

	42
Neb	Translated foreign pages
mages	Translated results for state of the art in information retrieval
Maps	- My language: English V
lidoos	Language Translated query
videos	Add language V - Automatically select languages to search
News	
Books	[PDF] DV Lande BASIS OF INFORMATION INTEGRATION
Places	poiskbook.kiev.ua/art/monogr-osnov/spusk3.pdf
Bloge	Translated from: Russian
bioga	The third chapter covers the current state of integration of information flows, attempts to solve the problems of search and navigation within
lights	+ Show original text
Discussions	the second se
Applications	Information technology to find information
Patente	Translated from: Russian
atomo	They have developed a means of information retrieval, calledIn practice, it is largely determined by the art of achievinginformation needs are determined by the type and
ewer	condition
	+ Show original text
The web	GotAI.NET - Materials - Publications
Canada	www.gotai.net/documents/doc-art-005.aspx - Russia
	I ranslated from: Russian In the first chapter held analysis of the current the state of of information-search
Any time	systems, of the modern the state of research in the the field of search
Past hour	+ Snow original text



















Approac	ch 2: Using bilingual dictionaries	
General access: academic: branch: data:	form of dict. (e.g. Freedict) attaque, accéder, intelligence, entrée, accès étudiant, académique filiale, succursale, spécialité, branche données, matériau, data	
<ul> <li>LDC English-C</li> <li>AIDS</li> <li>data</li> <li>prevention</li> <li>problem</li> <li>structure</li> </ul>	Chinese / 艾滋病/爱滋病/ /材料/资料/事实/数据/基准/ /恒碍/防止/妨碍/预防/预防法/ /问题/准愿/旗向/到题/作图题/将军/课题/困难/难/题是/ /构造/构成/结构/组织/化学构造/石理/纹路/构造物/ 建筑物/建造/物/	









#### Simple utilization

• Directly use the probability of a word translation

$$P(f \mid Q_E) = t(f \mid e) \quad P(e \mid Q_E)$$

 One should also take into account the discriminant power of the translation (IDF) (better in VSM)

$$w(f, Q_E) = t(f \mid e) \quad P(e \mid Q_E) \quad \log \frac{\mid C_F}{n_f}$$



 Multiple translations, but ambiguity is kept • Difficult to distinguish usual translation words from coincidental translations (e.g. pouvoir) • Unknown word in target language stupéfiant=0.005265 produit=0.004789

#### 61

## IBM1 + dictionary: a simple combination in VSM

 Increase the weight of usual translation words in the dictionary (TREC-6)

		Number of translation words							
Defau	lt 1	10	20	30	40	50	100	Number N of translation words	MAP (% monolingual IR)
probabi	lity							10	0.2546 (68.24%)
0.005	5 0.2	2671	0.2787	0.2812	0.2813	0.2829	0.2671	20	0.2635 (70.62%)
0.01	0.2	2755	0.2873	0.2891	0.2896	0.2906	0.2742	20	0.2055 (70.0270)
0.02	0.3	2972	0.2050	0.2062	0 2067	0 2085	0.2825	30	0.2660 (71.30%)
0.02	0.2	2013	0.2939	0.2902	0.2907	0.2765	0.2825	40	0.2664 (71.40%)
0.03	0.2	2811	0.2906	0.2898	0.2897	0.2904	0.2744	50	0.2671 (71.50%)
0.04	0.2	2751	0.2842	0.2827	0.2826	0.2831	0.2683	.30	0.20/1 (71.39%)
0.05	0.2	2687	0.2761	0.2729	0.2729	0.2730	0.2578	100	0.2506 (67.14%)

MAP-mono = 0.3731

• MAP-LOGOS = 0.2866 (76.8%), MAP-Systran = 0.2763 (74.1%)

	62
Principle of model training	
<ul> <li>t(t<sub>j</sub> s<sub>i</sub>) is estimated from a parallel training corpus, aligned into parallel sentences</li> </ul>	
• IBM models 1, 2, 3,	
Process:	
<ul> <li>Input = two sets of parallel texts</li> </ul>	
• Sentence alignment A: $S_k \Leftrightarrow T_l$ (bitext)	
• Initial probability assignment: $t(t_j s_j, A)$	
• Expectation Maximization (EM): $t(t_j s_i,A)$	
• Final result. $l(l_j S_i) = l(l_j S_i, A)$	

# Integrating translation in language model (Kraaij et al. 2003) • The problem of CLIR: $\begin{aligned} & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & &$

#### Results (CLEF 2000-2002)

Run	EN-FR	FR-EN	EN-IT	IT-EN
Mono	0.4233	0.4705	0.4542	0.4705
MT	0.3478	0.4043	0.3060	0.3249
QT	0.3878	0.4194	0.3519	0.3678
DT	0.3909	0.4073	0.3728	0.3547

64

- Translation model (IBM 1) trained on a web collection

- TM can outperform MT (Systran)



#### Sentence alignment

- Align a sentence in the source language to its translation(s) in the target language
- Translation model
- Extract translation relationships
- Various models (assumptions)

# Example of aligned sentences (Canadian Hansards)

Débat L'intelligence artificielle	Artificial intelligence A ebate	2-2
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.	2-1
Leurs avancées et leurs insuccès alternent curieusement.		
	Two further points are important.	0-1
Les symboles et les programmes sont des notions purement abstraites.	First, symbols and programs are purely abstract notions.	1-1







#### Word alignment for one sentence pair

Source sentence in training:  $\mathbf{e} = e_1, \dots e_l (+NULL)$ Target sentence in training:  $\mathbf{f} = f_1, \dots f_m$ Only consider alignments in which each target word (or position *j*) is aligned to a source word (of position *aj*)



# IBM models (Brown et al.) IBM 1: does not consider positional information and sentence length

- · IBM 2: considers sentence length and word position
- IBM 3, 4, 5: fertility in translation
- For CLIR, IBM 1 seems to correspond to the current (bagof-words) approaches to IR.









Sum up all the alignments
$p(\mathbf{f}   \mathbf{e}) = \frac{\mathcal{E}}{(l+1)^m} \sum_{a_1=0}^{l} \dots \sum_{a_m=0}^{l} \prod_{j=1}^{m} t(f_j   e_{a_j})$
$= \frac{\mathcal{E}}{\left(l+1\right)^m} \prod_{j=1}^m \sum_{i=0}^l t(f_j \mid e_i)$
<ul> <li>Problem: We want to optimize t(f<sub>j</sub>   e<sub>i</sub>) so as to maximize the likelihood of the given sentence alignments</li> <li>Solution: Using EM</li> </ul>











#### Candidate Site Selection

By sending queries to AltaVista, find the Web sites that may contain parallel text.

#### File Name Fetching

For each site, fetching all the file names that are indexed by search engines. Use host crawler to thoroughly retrieve file names from each site.

#### Pair Scanning

From the file names fetched, scan for pairs that satisfy the common naming rules.



#### File Name Fetching

- Initial set of files (seeds) from a candidate site: host:www.info.gov.hk
- Breadth-first exploration from the seeds to discover other documents from the sites

#### Pair Scanning

• Naming examples:

index.html v.s. index\_f.html

/english/index.html v.s. /french/index.html

• General idea:

parallel Web pages = Similar URLs at the difference of a tag identifying a language

#### Mining Results (several years ago)

- 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, ...

#### CLIR results: F-E

	F-E	F-E	E-F	E-F
	(Trec6)	(Trec7)	(Trec6)	(Trec7)
Monolingual	0.2865	0.3202	0.3686	0.2764
Hansard TM	0.2166	0.3124	0.2501	0.2587
	(74.8%)	(97.6%)	(67.9%)	(93.6%)
Web TM	0.2389	0.3146	0.2504	0.2289
	(82.5%)	(98.3%)	(67.9%)	(82.8%)

- Web TM comparable to Hansard TM
- Parallel texts for CLIR can tolerate more noise



# Alternative methods – LSI [Dumais et al. 97] Monolingual LSI : Create a latent semantic space (using SVD) Each dimension represents a combination of initial dimensions (terms, documents) Comparison of document-query in the new space Blingual LSI : Create a latent semantic space for both languages on a parallel corpus Concatenate two parallel texts together Convert terms in both languages into the semantic space Problems: The dimensions in the latent space are determined to minimize some representational error – may be different from translational error Coverage of terms by the parallel corpus Complexity in creating the semantic space Effectiveness – usually lower than using a translation model

#### Explicit Semantic Analysis (ESA) [Gabrilovich et al. 07, Spitkovsky et al. 12]

- Assume that each Wikipedia article corresponds to one explicit representation dimension
- A term → association with an article (tf\*idf or conditional probability) based on text body or anchor text
- Document ranking: compare the document and the query representations in the ESA space
- · CLIR:
- · Using cross-lingual links between Wikipedia articles
- Possible problems:
- Completeness of ESA space
- Representation granularity

## Summary on using document-level term relations

- Pseudo-RF, LSI, ESA
- Assumption: Terms occurring in the same text (or parallel text) are related
- (Translation) relations between terms are coarse-grained
   Related to the same topics
- · Query "translation" by topic-related terms
- Usually lower retrieval effectiveness than explicit translation relations
- May be suitable for comparable texts
- Do not exploit fully parallel texts

#### Using a comparable corpus

- Comparable: articles about the same topic
   E.g. News articles on the same day about an event
- Impossible to train a translation model
- Estimate cross-lingual similarity (less precise than translation)
- Similar methods to co-occurrence analysis
- $\ \cdot \$  Conditional probability, mutual information, ...
- Less effective than using a parallel corpus
- To be used only when there is no parallel corpus, or to complement a dictionary or parallel corpus
  - Helpful to further expand the translation by a dictionary, MT or TM

91

92

#### Other problems – unknown words

- Proper names ('Vladmir Ivanov' in Chinese?)
- New technical terms ('Latent Semantic Analysis' in Chinese?)
- Possible solutions
- Transliteration
- Mining the web







# Mining the Web (3) Wikipedia – a rich source for translations Links between articles in different languages Translation of concepts (Wiki titles) and named entities (e.g. proper names) through cross-language links Linked articles as comparable texts → term similarity



#### Current state

#### · Effectiveness of CLIR

- Between European languages ~90-100% monolingual
- Between European and Asian languages ~ 80-100%
- A usable quality
- One usually needs translation of the retrieved documents

99

100

- The use of CLIR by Web users is still limited / Tools for CLIR are limited
- To be increased (integration in the main search engines)

#### Structured query

- Traditional method: all the translation terms in a bag of words
  - data /材料/资料/事实/数据/基准/
  - structure /构造/构成/结构/组织/化学构造/石理/纹路/构造物/ 建筑物/建造/物/
- Consider all translations for the same source word as synonyms
- •#syn(材料 资料 事实 数据 基准) #syn(构造 构成 ...)
- sum up the occurrences of all the synonyms in a document (v.s. sum up the log probabilities without #syn)
- Pirkola: structured query > bag-of-words query
- Probabilistic structured query (#wsyn): add weights to synonyms

#### Remaining problems

#### • Current approaches :

- · CLIR= translation + monolingual IR
- E.g. Using MT as a black box + monolingual IR
- $\,\cdot\,$  The resources and tools are usually developed for MT, not for CLIR
- MT: create a readable sentence
- · CLIR: retrieving relevant documents
- Problems of translation selection
- MT: Select one best translation  $\rightarrow$  multiple translations
- Phrases in MT: consecutive words
- But dependent words do not always form a phrase
- "Mixing drug cocktails for mental illness is still more art than science"
- $\cdot \ { \rightarrow } {\sf Take}$  into account more flexible dependencies (even proximity)
- How to train a translation model in such a context?
- These are not only a problem in CLIR but also in general IR.



#### What have we learnt from CLIR?

- The original query is not always the best expression of information need.
- We can generate better or alternative query expressions
   Query expansion, reformulation, rewriting



#### Beyond

- Current CLIR approaches can succeed in crossing the language barrier
- MT, Dictionary, STM
- Can one use CLIR methods for monolingual (general) IR?
   Basic idea

102

- IBM1 model
- Phrase translation model and some adaptations

#### Basic idea for general IR

- Assumption: Lexical gap between queries and documents
  - Written by searchers and by authors, using different vocabularies
- Translate between query language and document language
- How?
  - · Dictionary, thesaurus (~dictionary-based translation)
  - Statistical "translation" models (\*)



Query logs		
msn web	0 6675749	
Webmensseger	0.6621253	
msn online	0.6403270	
windows web messanger	0.6321526	
talking to friends on msn	0.6130790	
school msn	0.5994550	
msn anywhere	0.5667575	
web message msn com	0.5476839	Title: msn web messenger
msn messager	0.5313351	nael nee neeelige
hotmail web chat	0.5231608	
messenger web version	0.5013624	
instant messager msn	0.4550409	
browser based messenger	0.3814714	
im messenger sign in	0.2997275	
igure 1: A fragment of the query cli	ck field for the	page
ttp://webmessenger.msn.com [16].		







					110
	Results				
#	# Models	NDCG@1	NDCG@3	NDCG@10	
1	BM25	0.3181	0.3413	0.4045	
1	2 WTM_M1 (β=1)	0.3202	0.3445	0.4076	No TM
1	3 WTM_M1	0.3310	0.3566	0.4232	
4	4 WTM_M1 (β=0)	0.3210	0.3512	0.4211	Only TM



### Phrase translation in IR (Riezler et al. 2008)

Phrase = consecutive words

Determining phrases:

- Word translation model
- Word alignments between parallel sentences
- Phrase = intersection of alignments in both directions

	QE Models	NDCG@1	NDCG@3	NDCG@10
1	NoQE	0.2803	0.3430	0.4199
2	PhraseMT	0.3002	0.3617	0.4363
3	WordModel	0.3117	0.3717	0.4434
<ul> <li>Note:</li> </ul>	on a different colle	ection		

# Why isn't phrase-model better than word-model?

 Sentence (Riezler et al 2008) herbs for chronic constipation

Words and phrases

herbs

- herbs for
- herbs for chronic

for

- for chronic
- for chronic constipation
- chronic constipation
- constipation

# More flexible "phrases" – n-grams for query

Word generation process:

- Select a segmentation S of Q according to P(S|Q),
- Select a query phrase **q** according to  $P(\mathbf{q}/S)$ , and
- Select a translation (i.e., expansion term) *w* according to  $P(w|\mathbf{q})$ .





instant messenger msn instant-messenger messenger-msn instant-messenger-msn

instant messenger msn

#### Comparison

	QE Models	NDCG@1	NDCG@3	NDCG@10
1	NoQE	0.2803	0.3430	0.4199
2	WM	0.3117	0.3717	0.4434
3	PM <sub>2</sub>	0.3261	0.3832	0.4522
4	$\mathbf{PM}_8$	0.3263	0.3836	0.4523

115

WM: Word translation model PM<sub>n</sub>: n-gram Phrase translation model

#### "Phrase" in IR

- Not necessarily consecutive words ("information ...retrieval")
- Words at distance (context words) may be useful to determine the correct translation of a word
  - the new drug developed by X was approved by Y

- drug ... smuggle ...

- Question
  - How to use context words (within a text window) to help translation?
  - $\rightarrow$  Some approaches in SMT but will need to be further extended



Comparison of different translation models				
	QE Models	NDCG@1	NDCG@3	NDCG@1 0
1	NoQE	0.2803	0.3430	0.4199
2	WM	0.3117	0.3717	0.4434
3	PM <sub>2</sub>	0.3261	0.3832	0.4522
4	$PM_8$	0.3263	0.3836	0.4523
5	CM <sub>T-B</sub>	0.3208	0.3786	0.4472
6	CM <sub>T-P2</sub>	0.3204	0.3771	0.4469
7	CM <sub>T-B-P3</sub>	0.3219	0.3790	0.4480
8	CM <sub>T-B-P5</sub>	0.3271	0.3842	0.4534
9	CM <sub>T-B-P8</sub>	0.3270	0.3844	0.4534
CM: Co	ncept translation me	odel (T-term, B-t	pigram, Pn-prox	ximity within n





	121
Markov Random Field Score(D,Q) = $P_{\Lambda}(Q,D) = \frac{1}{Z_{c C(Q)}}$	$ \psi(c,\Lambda) $ $ \lambda_o f_o(c) + \lambda_U f_U(c) $
• Three components: - Term – T - Ordered term clique – O - Unordered term clique – U • $\lambda_{\mathcal{P}} \lambda_{\mathcal{O}} \lambda_{\mathcal{U}}$ : fixed parameters - Typical good values (0.8, 0.1, 0.1)	$f_{T}(c) = \log P(q_{i}   D)$ $f_{O}(c) = \log P(\#1(q_{i},,q_{i+k})   D)$ $f_{U}(c) = \log P(\#UW(q_{i},,q_{i+k})   D)$

				122
Res	sults			
	QE Models	NDCG@1	NDCG@3	NDCG@10
1	MRF (NoQE)	0.2802	0.3434	0.4201
2	$1 + M-CM_{T-B-P8}$	0.3293	0.3869	0.4548
3	M-CM <sub>T-B-P8</sub>	0.3271	0.3843	0.4533

- Using MRF to evaluate "phrase" queries helps slightly.
- 1 + M-CM<sub>T-B-P8</sub>: Generate and use phrase translations
- ${}^{\bullet}$   $M\text{-}CM_{\text{T-B-P8}}{:}$  Generate phrase translations but break them into bag of words





#### Some applications of CLIR (1)

#### CLIR in specialized areas

- Patent retrieval
- Patents are written in local languages (Chinese, Korean, ...)
- · Companies doing international business are supposed to be aware of the existing patents
- · CL patent retrieval in practical uses: increasing (patent searhc on Google)
- Medical IR
- · Relevant information in another language may be useful (Chinese medicine)
- Specific problems?
- Patent/medical language?
- Patent/medical document structure?
- Term mismatch (standardization/expansion is important)
- →Talk on domain-specific IR

#### Some applications of CLIR (2)

CL Question-Answering

#### Yahoo! QA. Baidu knows. ...

- Some answers may only be available in a foreign language
- Translate the question
- · Several possible formulations to select/combine
- → Talk on community QA

#### Some applications of CLIR (3)

- Image retrieval
- Mostly language-independent
- but the annotations and surrounding texts are languagedependent
- Query-Annotation search ~ CLIR problem
- Specific problem: annotations are limited → Expansion is required

#### Summary

- High-quality MT usually offers a good solution
- Well-trained TM based on parallel texts can match or outperform MT (Kraaij et al. 03)
- Dictionary
  - Simple utilization is not good
- Translation selection is important → can match monolingual IR
   The performance of CLIR usually lower than monolingual IR (between 80% and 100% for advanced approaches)
- Better translation means → better CLIR effectiveness
- Better to combine translation means
- CLIR is now integrated in some search engines (wide use to be expected in the future)
- Remaining problems:
- Improve translation quality (select good translation terms)
- IR-oriented translation
- Beyond: CLIR is not separated from general IR

#### Some references

Jian-Yun Nie, Cross-language information retrieval, Morgan-Claypool, 2010. (a survey book on CLIR)

survey book on CLIR Adriani, M. van Rijsbergen, C.J. (2000), Phrase identification in cross-language inflammation, International Conference on Conference on Clifford International Conference on Conference on Clifford International Internation International Internation International International International International International Internation International Internationa

Kraaij, W., Nie, J.Y., and Simard, M. (2003). Embedding Web-Based Statistical Translation Models in Cross-Language Information Retrieval. Computational Linguistics, 29(3): 381-420. (CLIR using language modeling

129

Computational Linguistics, 25(3): 381–420. (CLIH using language modeling and Wab-mindparallel corpora) Liu, Y., Jin R. and Chai, Jayoe Y. (2005). A maximum colerence model for SIGH cort. (pp. 585–561 (transition selection) J. Scott McCatrley, Should were Translate the Documents or the Queries in Cross-singuage SIG-561 (transition selection) (consel-inguage SIGH Cort. (pp. 243–254. (Document translations, euery translation) Neuron McCatrley, Should were translated the Automatic Mining of Parallel Tests in the Weik, Indiversity SIGH Cort. (pp. 243–214. (Document Feats in the Weik, Indiversity SIGH Cort. (pp. 243–214. (Document Feats in the Weik, Indiversity SIGH Cort. (pp. 243–214. (Document SIGH Cort. (pp. 243–244. (Document))

Tests in the Web, in Proceedings of SIGIR Cort, pp. 74-81 (CLIR tailing Web-mined parallel pages). Pixola, A. (1989) The effects of query structure and detionary setups in SIGIR Cort, pp. 55-63. (Structured paral translation) Pixola, A., Toivonen, J., Keskustalo, H., Visala K., Jarvein K., (2003) Fuzzy translation of cores impairs genity availants, In Proceedings of SIGIR Cort, pp. 45-352. (Piuzzy makin) Pixola, A., Toivonen, J., Keskustalo, H., Visala K., Jarvein K., (2003) Fuzzy translation of cores impairs genity availants, In Proceedings of SIGIR Cort, pp. 45-352. (Piuzzy makin) Pixola, M., Visala, M., Visala, M., Visala, M., Savein, K., (2003) Fuzzy translation of cores impairs genity available, the Sidir Sidir Pixola and Sidir Core and Sidir Sidir Sidir Sidir Sidir Sidir Sidir Retire Core (TECC-6). (Document translation) N. Query translation) Reter S., and U.V., Y 2010. Query reveniting using monolingui statistical machine translation. Computational Anguistes, 39(5): 569-582 (dmaxe

translation for CLR) Seq. H-G. Km, S. A. Brin, H-G. and Mayeng S-H. (2005) Imposing quary translation in English-Korean cross-language information retrieval, Information Processing and Management, 41: 552-522. (translation assellation) Sheridan, P. and Ballerini, J. P. (1996). Experiments in multilingual information retrieval using the SIDER system. In Proceedings of SIG Cont., pp. 58-56. (CLIR using comparable corpona) Valentin 1. Splicovsky and Anga X. Chang. 2012. A Cross-Lingual Dictionary for English Wikipeda Concepts. In Proceedings of the Eighth International Conference on Language Resources and Evaluation (LPCC 2012) (ESA For Conference on Language Resources and Evaluation (LPCC 2012) (ESA For Conference on Language Resources and Evaluation (LPCC 2012) (ESA For Conference on Language Resources and Evaluation (LPCC) 2012) (ESA For Conference on Language Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Language Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Language Resources and Evaluation (LPCC) 2012) (ESA For Section (LPCC) 2012) (ESA For Conference on Language Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Language Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Lenguage Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Lenguage Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Lenguage Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Lenguage Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Lenguage Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Lenguage Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Lenguage Resources and Evaluation (LPCC) 2012) (ESA For Section 2012) (ESA For Conference on Lenguage Resources and Evaluation (LPCC) 2012) (ESA For

CLIR)

CLIR) Wen, J., Ne, J-Y., and Zhang, H. 2002. Query clustering using user logs. ACM 70/S, 20(1): 59-81. (Query clustering based on query logs) Yiming Yang, Jaime G. Carbonell, Rati D. Brown, and Robert E. Frederking, "Translingual information retrieval: Learning from bilingual corpora", Artificial Intelligence, vol. 103, pp. 523–455, 1998. (Pseudo-Frad or Clust)