

Cross-Language IR at
University of Tsukuba
Automatic Transliteration for
Japanese, English, and Korean

Atsushi Fujii and Tetsuya Ishikawa
University of Tsukuba

C26

Motivation

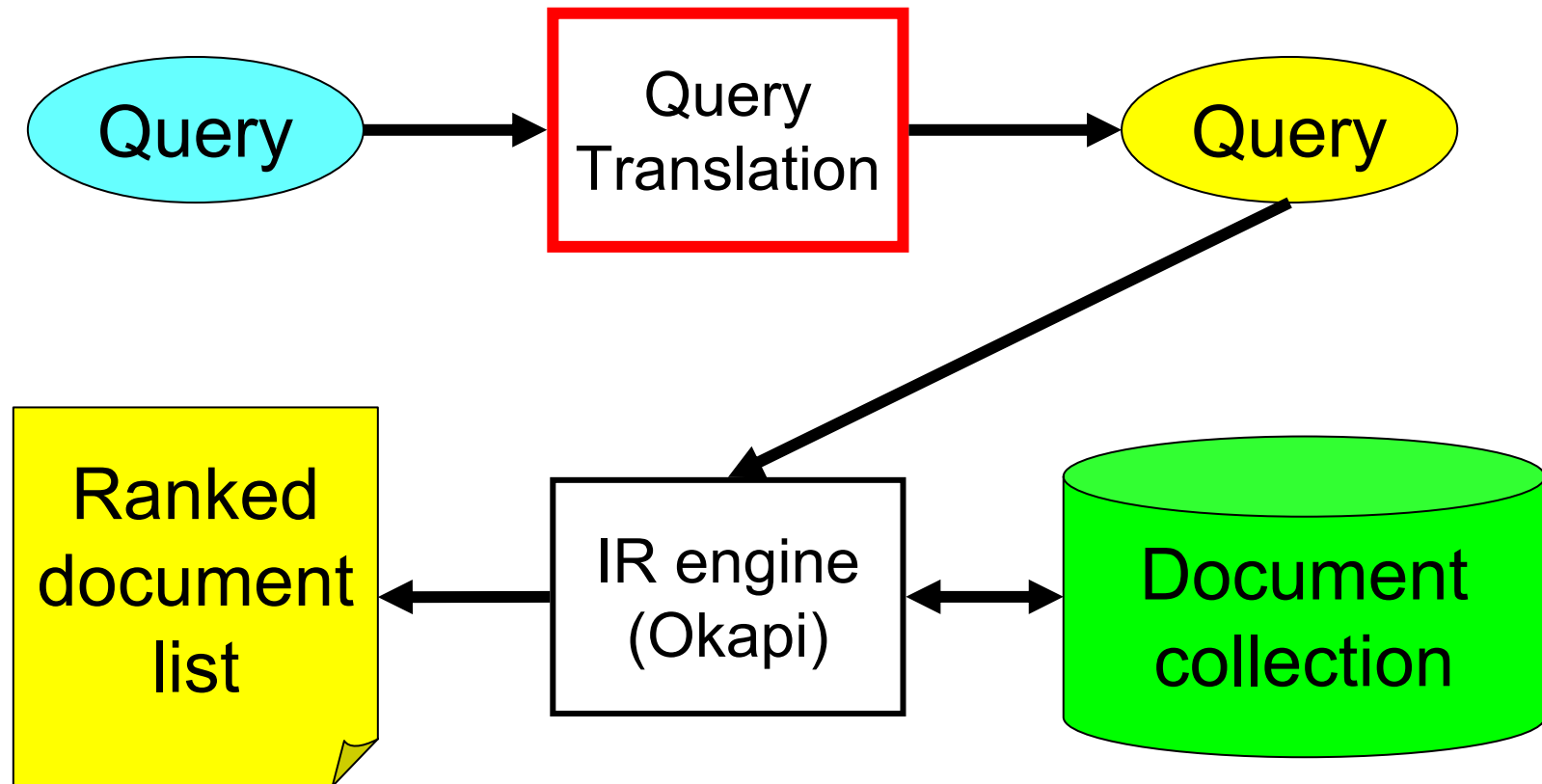
- We developed an automatic **transliteration** method for Japanese and English CLIR
 - effective in translating foreign words spelled out by phonetic alphabet (e.g., Katakana)
 - evaluation since NTCIR-1
 - the method has been used in commercial cross-language patent IR service
- In NTCIR-4 CLIR, we applied our method to Korean and **realized JEK transliteration in a single framework**

Basis of transliteration

- spelling out foreign words (loanwords) by phonetic alphabet
 - technical terms and proper names
 - often out-of-dictionary words
- examples
 - dioxin → ダイオキシン, 다이옥신
 - Yugoslavia → ユーゴスラビア, 유고슬라비아
- back-transliteration
 - process to identify the source English word

Overview of our CLIR system

Focus of
today's talk



Example of J-E Query Translation

レジスタ転送言語

consulting dictionary
lexical segmentation

レジスタ

転送

言語

transliteration

resister

transfer

language

resistor

transmission

register

transport

disambiguation

register transfer language

Query Translation (cont.)

- compound query term S and a translation candidate T

$$S = s_1, s_2, \dots, s_N$$

s_i and t_i are base words

$$T = t_1, t_2, \dots, t_M$$

- compute $P(T|S) = P(S|T) \cdot P(T)$

translation model

language model

- select the candidate with $\max P(T|S)$

Translation model

- $P(S|T) = \prod P(s_i | t_i)$
si and ti are base words comprising S and T
- heuristics and EM algorithm to correspond dictionary entries on a word-by-word basis

情報	検索	システム	Information	retrieval	system
検索	モデル		retrieval	model	
情報	抽出	システム	Information	extraction	system
特許	情報	処理	patent	information	processing

- estimate $P(s_i | t_i)$

Language model

- word-based trigram model
- 100K vocabulary in a target document collection
- Palmkit was used
 - compatible with CMU-LM toolkit

Transliteration method

- out-of-dictionary word S and a transliteration candidate T

$$S = s_1, s_2, \dots, s_N$$

$$T = t_1, t_2, \dots, t_M$$

s_i and t_i are letters (substrings of words)

- compute $P(T|S) = \underbrace{P(S|T)}_{\text{transliteration model}} \cdot \underbrace{P(T)}_{\text{language model (word unigram)}}$

transliteration model

language model
(word unigram)

- select the candidate with $\max P(T|S)$

Transliteration dictionary

- dictionary for transliteration includes correspondence b/w source and target words on a phonogram-by-phonogram basis
- we use Roman representation as a pivot

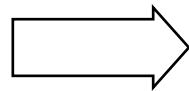
Producing J/E dictionary

1. extract Japanese **Katakana** words and English translations from J-E dictionary
2. romanize Katakana words
3. correspond romanized Katakana and English words on a letter-by-letter basis
4. find the best correspondence

Example matrix

テキスト (te-ki-su-to) text

	テ	キ	ス	ト	\$
t	3	1	2	3	0
e	0	0	0	0	0
x	1	2	1	1	0
t	3	1	2	3	0
\$	0	0	0	0	3



テ	te
キス	x
ト	t

By performing the same process for all Katakana entries, we produce transliteration dictionary

Extension to other languages

- our transliteration method can be applied to any language **if represented by Roman characters**
- no existing method has been used and evaluated in CLIR for more than two languages
 - our experiment was the first effort to explore this issue

Problems in Korean

- romanization of Korean words is more difficult than that of Katakana words
 - # of Hangeul characters is approx. 11,000
 - one-to-one mapping b/w Hangeul and Roman characters is not easy
- both conventional Korean words and foreign words are written by Hangeul characters
 - detection of foreign words in Korean dictionary is crucial

Romanizing Korean words

- Hangul character consists of three types of consonants
 - first consonant (19)
 - vowel (21)
 - last consonant (27 + 1)
- # of possible combinations is 11,172
(# of common characters is approx. 2,000)
- We used Unicode, in which characters are coded according to consonants

last consonant is optional



Detecting foreign words in Korean

- compute the phonetic similarity b/w romanized Hangeul words and their translations (either English or Japanese)
- discard translation pairs whose similarity is below a threshold
 - conventional Korean words are discarded
- foreign word entries remained

Experiments (J/E)

<TITLE>, mean average precision (rigid)

Languages	#Entries	w/o transliteration		w/ transliteration
J-E	1M	0.2174	<	0.2182
E-J	1M	0.1250	=	0.1250
J-E (EDICT)	108K	0.1147	<	0.1383
E-J (EDICT)	108K	0.0612	<	0.0857

transliteration was effective for small dictionaries

Experiments (Korean)

<TITLE>, mean average precision (rigid)

Languages	w/o transliteration		w/ transliteration
J-K	0.2177	<	0.2457
K-J	0.1486	<	0.1746
E-K	0.2026	<	0.2153
K-E	0.1017	<	0.1231

transliteration was also effective for Korean

Conclusion

- realized transliteration for Japanese, English, and Korean in a single framework
- evaluated its effectiveness in NTCIR-4 CLIR task