



# Cross-Language IR

many slides courtesy

James Allan@umass

Jimmy Lin, University of Maryland

Paul Clough and Mark Stevenson, University of Sheffield, UK



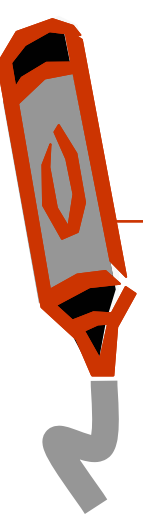
# What is Cross-Lingual Retrieval?

---

- Accepting questions in one language (English) and retrieving information in a variety of other languages
  - “questions” may be typical Web queries or full questions in across-lingual question answering (QA) system
  - “information” could be news articles, text fragments or passages, factual answers, audio broadcasts, written documents, images, etc.
- Searching distributed, unstructured, heterogeneous, multilingual data
- Often combined with summarization, translation, and discovery technology

# Current Approaches to CLIR

---



- Typical approach is to translate query, use monolingual search engines, then combine answers
  - other approaches use machine translation of documents
  - Or translation into an interlingua
- Translation ambiguity a major issue
  - multiple translations for each word
  - query expansion often used as part of solution
  - translation probabilities required for some approaches
- Requires significant language resources
  - bilingual dictionaries
  - parallel corpora
  - “comparable” corpora
  - MT systems

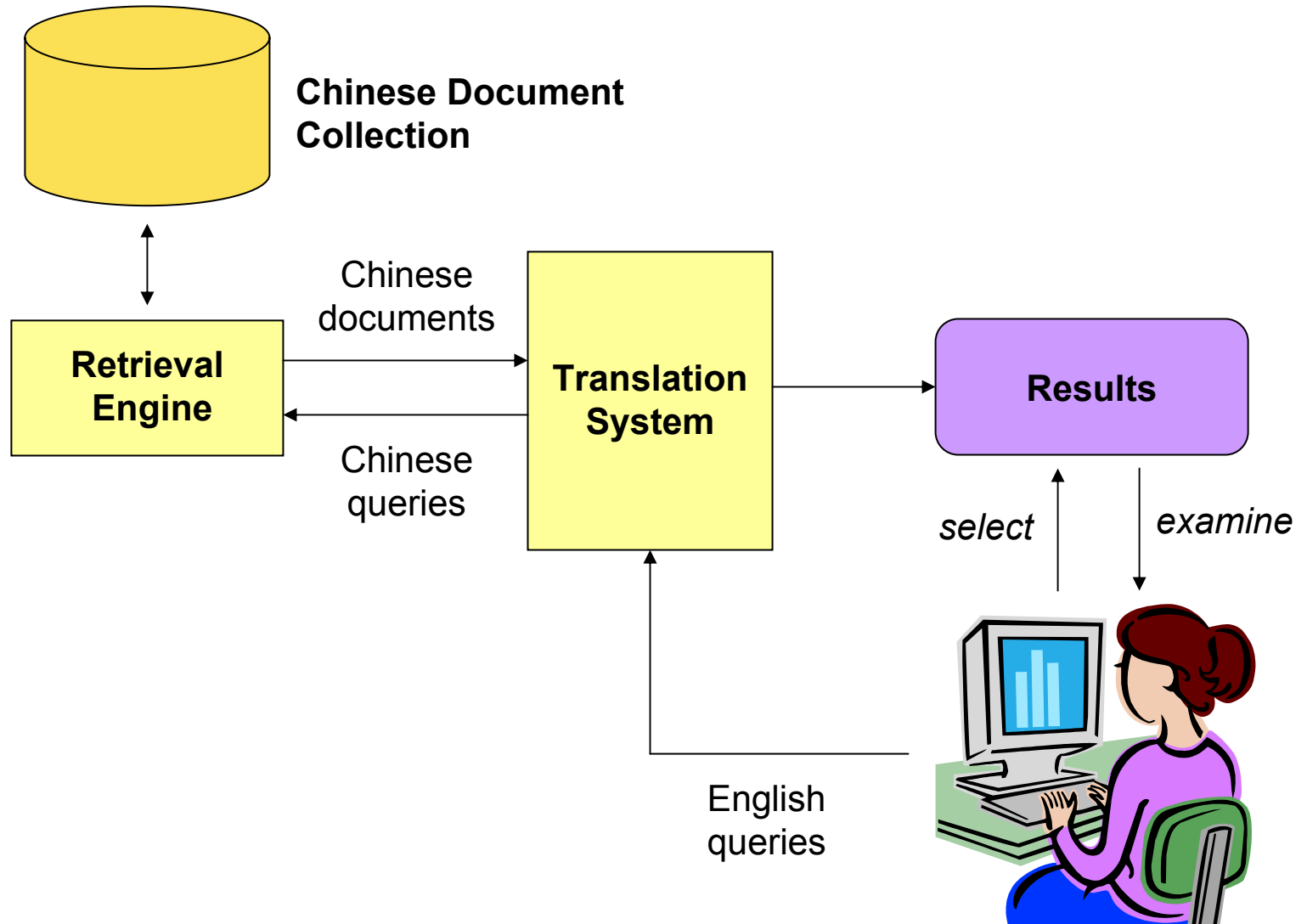


# Two Approaches

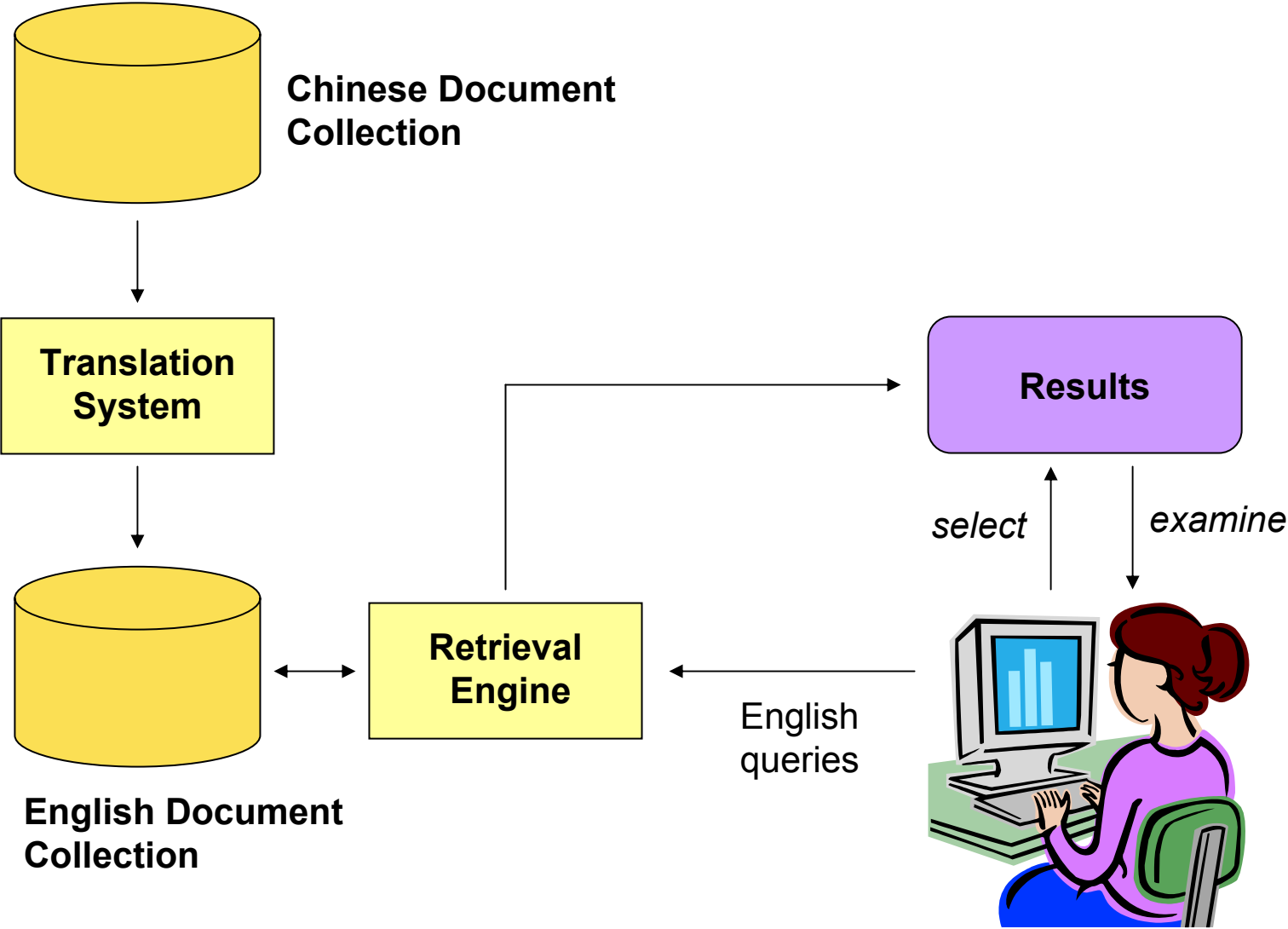
---

- Query translation
  - Translate English query into Chinese query
  - Search Chinese document collection
  - Translate retrieved results back into English
- Document translation
  - Translate entire document collection into English
  - Search collection in English
- Translate both?

# Query Translation

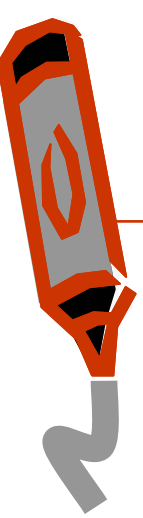


# Document Translation



# Tradeoffs

---



- Query Translation
  - Often easier
  - Disambiguation of query terms may be difficult with short queries
  - Translation of documents must be performed at query time
- Document Translation
  - Documents can be translate and stored offline
  - Automatic translation can be slow
- Which is better?
  - Often depends on the availability of language-specific resources (e.g., morphological analyzers)
  - Both approaches present challenges for interaction



# A non-statistical approach

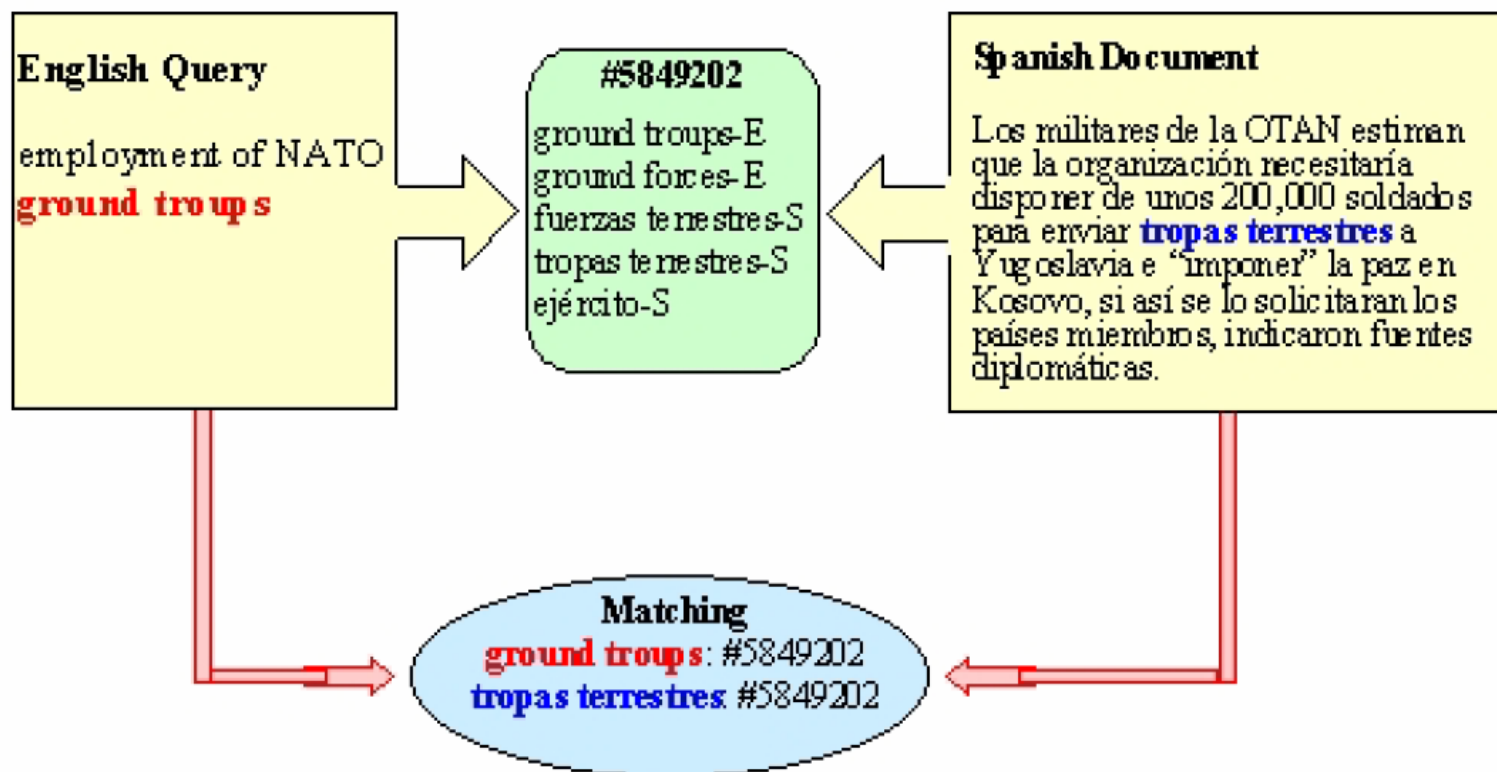
---

- A non-statistical approach
- Interlingua approaches
  - Translate query into special language
  - Translate all documents into same language
  - Compare directly
  - Cross-language retrieval becomes monolingual retrieval
- Choice of interlingua?
  - Could use an existing language (e.g., English)
  - Create own
- Textwise created a “conceptual interlingua”



# CINDOR

- Conceptual Interlingua for Document Retrieval



[Liddy, Infonortics 1999]

# CINDOR

## English Query

## Conceptual Interlingua

## Multilingual Documents

I would like information

about the **possible**

**employment** of

**NATO**

**ground troops**

in the **Kosovo**

**conflict.**

possible (E)  
conceivable (E)  
possible (F)  
conceivable (F)  
imaginable (F)  
factible (S)  
conceivable (S)  
posibilidad (S)

employment (E)  
engagement (E)  
commissioning (E)  
engagement (F)  
envoyé (F)  
empleo (S)  
uso (S)  
envío (S)

NATO (E)  
North Atlantic Treaty  
Organization (E)  
OTAN (F)  
Organisation du Traité  
de l'Atlantique Nord (F)  
OTAN (S)  
Organización del Tratado  
del Atlántico Norte (S)

Kosovo (E)  
Kosovo (F)  
Kosovo (S)

ground troops (E)  
ground forces (E)  
armées de terre (F)  
troupes (F)  
fuerzas terrestres (S)  
tropas terrestres (S)  
soldados (S)  
tropas (S)  
ejército (S)

conflict (E)  
discord (E)  
confit (F)  
désaccord (F)  
dissension (F)  
conflicto (S)  
discordia (S)  
enfrentamiento (S)  
crisis (S)

### English Document Excerpt:

WASHINGTON, March 29 (AFP) - The United States and Britain beefed up **NATO** forces as the bombing campaign against Yugoslavia entered a 24-hour phase and US officials warned **ground troops** in **Kosovo** were "no magic bullet."

### French Document Excerpt:

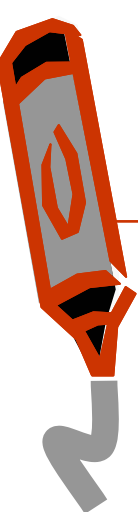
PARIS, 30 mars (AFP) - 25 MARS: Le président américain Bill Clinton déclare ne pas avoir l'intention "d'**envoyer** de **troupes**."

### Spanish document Excerpt:

BRUSELAS, Mar 28 (AFP) - De enviarse **tropas terrestres**, **posibilidad** que descartan actualmente todos los países de la organización, las pérdidas serían considerables, según los estrategas de la **OTAN**.

# Does it work?

---



- Some background research suggested large gains over word-by-word translation
- Fielded in TREC-7 cross-language task
- Performed poorly overall
  - System not completed at the time
  - Interlingua incomplete
  - Several small processing errors added up
  - On queries without problems, comparable to monolingual
- Statistical methods now dominate the field

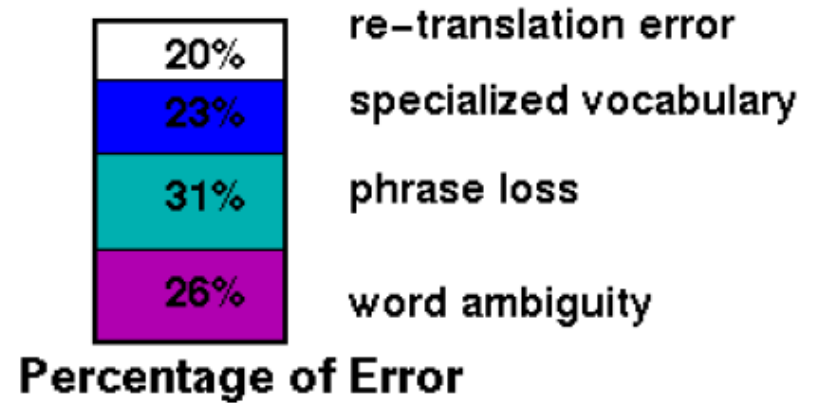
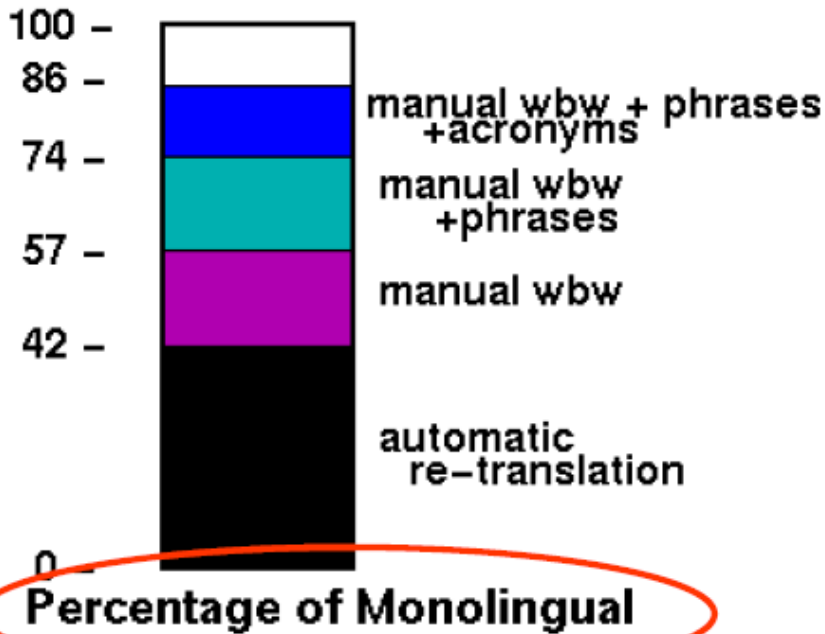
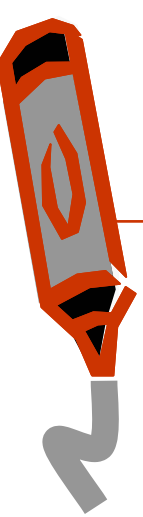


# Current Capabilities of CLIR

---

- Best performance obtained by
  - probabilistic approach using translation probabilities estimated from an aligned parallel corpus
  - “structured” query that treats translations from bilingual dictionary as synonyms and uses advanced search engine
  - Combination of techniques including MT
  - Most experiments done in Chinese, Spanish, French, German, and recently, Arabic
- Cross-lingual can achieve 80-90% effectiveness of monolingual
  - with sufficient language resources
  - sometimes does even better, but can also do worse

# CLIR errors





# But how good is “monolingual”?

---

- Not easy to summarize IR performance as a single number
  - We’ve considered average precision, Swet’s number, utility functions, expected search length, ...
- Based on measures of recall and precision...
  - Breakeven of 30% for “Web” queries, precision 40% in top 20, 20% in top 100
  - Breakeven of 45% for “analyst” queries, precision 65% in top 20, 45% in top 100
  - Recall can be improved through techniques such as query expansion and relevance feedback



# Adding New Languages

---

- Morphological processing
  - segmenting (what is a word?)
  - stemming (combining inflections and variants)
  - stopwords (words that can be ignored)
- Language resources
  - minimum is a bilingual dictionary
  - parallel or comparable corpora are even better
  - MT system is a luxury



# Problems with CLIR

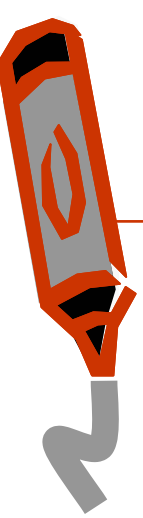
---

- Morphological processing difficult for some languages (e.g. Arabic)
  - Many different encodings for Arabic
    - Windows Arabic (e.g. dictionaries)
    - Unicode (UTF-8) (e.g. corpus)
    - Macintosh Arabic (e.g. queries)
  - Normalization
    - Remove diacritics  
**العَرَبِيَّة to العَرَبِيَّة** *Arabic (language)*
  - Standardize spellings for foreign names  
**كلينتون vs كلنتون** *"Kleentoon" vs "Klntoon" for Clinton*



# Problems with CLIR

---



- Morphological processing (contd.)
  - Arabic stemming
  - Root + patterns+suffixes+prefixes=word  
ktb+CiCaC=kitab
- All verbs and nouns derived from fewer than 2000 roots
- Roots too abstract for information retrieval
  - ktb → kitab *a book* kitabi *my book*
  - alkitab *the book* kitabuki *your book (f)*
  - kataba *to write* kitabuka *your book (m)*
  - maktab *office* kitabuhu *his book*
  - maktaba *library, bookstore ...*
  - Want stem=root+pattern+derivational affixes?
- No standard stemmers available, only morphological (root) analyzers

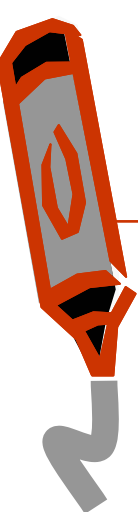


# Problems with CLIR

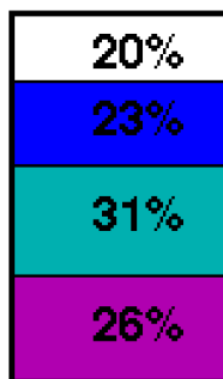
---

- Availability of resources
  - Names and phrases are very important, most lexicons do not have good coverage
- Difficult to get hold of bilingual dictionaries
  - can sometimes be found on the Web
- e.g. for recent Arabic cross-lingual evaluation we used 3 on-line Arabic- English dictionaries (including harvesting) and a small lexicon of country and city names
  - Parallel corpora are more difficult and require more formal arrangements

# Phrase translation



- Phrases are a major source of translation error
- How to get phrases translated properly?
- Assume that correct translations of words in phrase co-occur
  - Given two-word phrase “A B”
  - Look at all translations of A: A1 or A2 or ... or An (and B, similarly)
  - Look at all pairs “Ai Bj” and see which of them co-occur
- Probably in passages of the collection
  - Use the best pair as the phrase translation



**re-translation error**

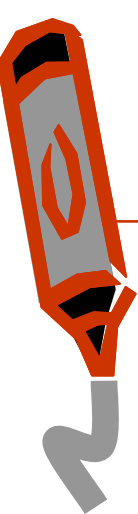
**specialized vocabulary**

**phrase loss**

**word ambiguity**

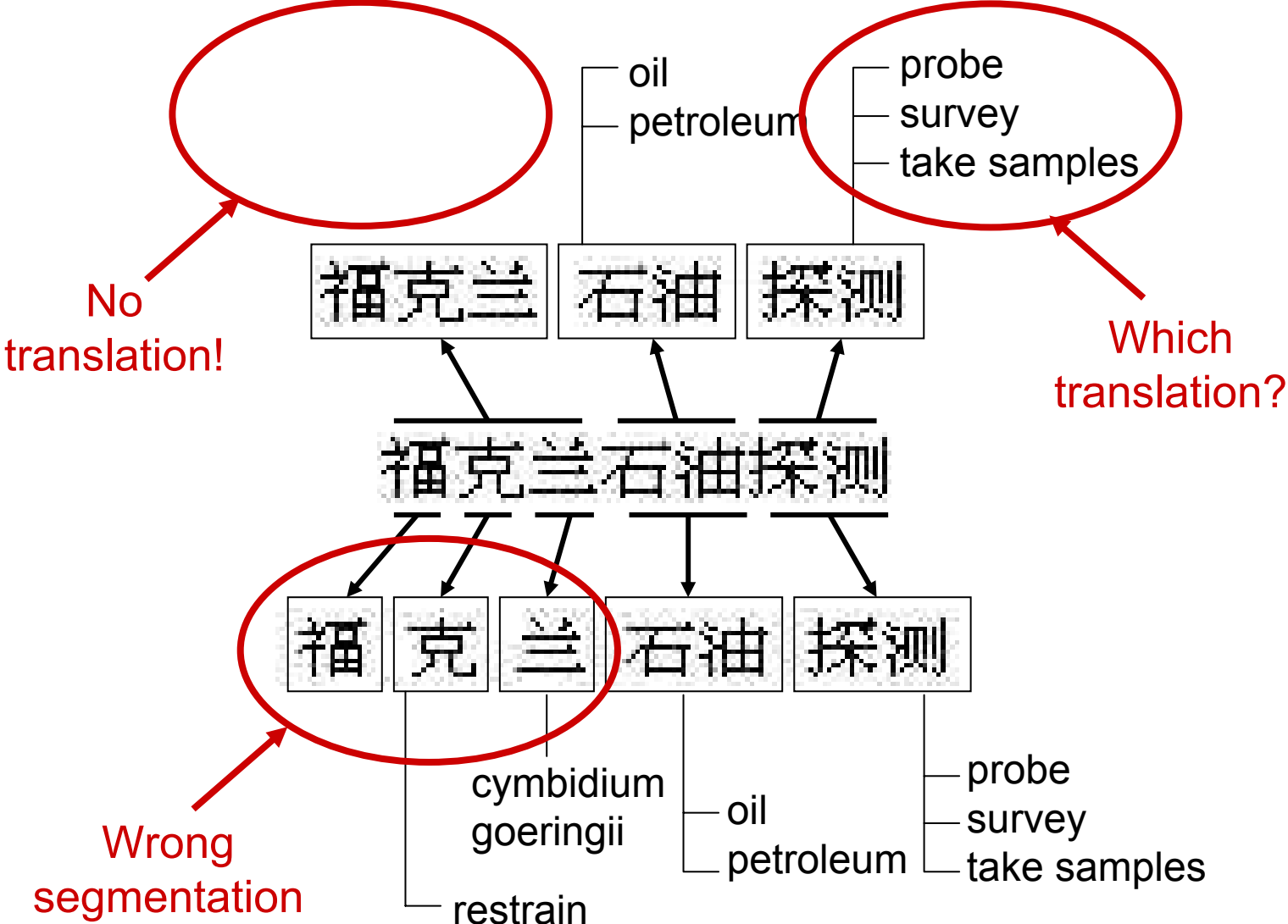
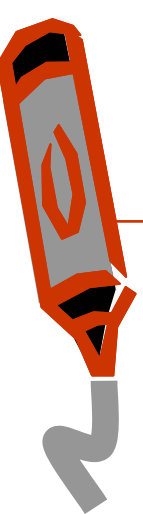
# Example Phrase

---



- Worked quite well in English-Spanish CLIR
- Consider Spanish phrase “Proceso Paz”
  - process, lapse of time, trial, prosecution, action, lawsuit, proceedings, processing
  - peace, peacefulness, tranquility, peace, peace treaty, kiss of peace, sign of peace
- Ranked possible translation pairs:
  - peace process
  - peacefulness process
  - tranquility process
  - ...

# CLIR Issues





# Learning to Translate

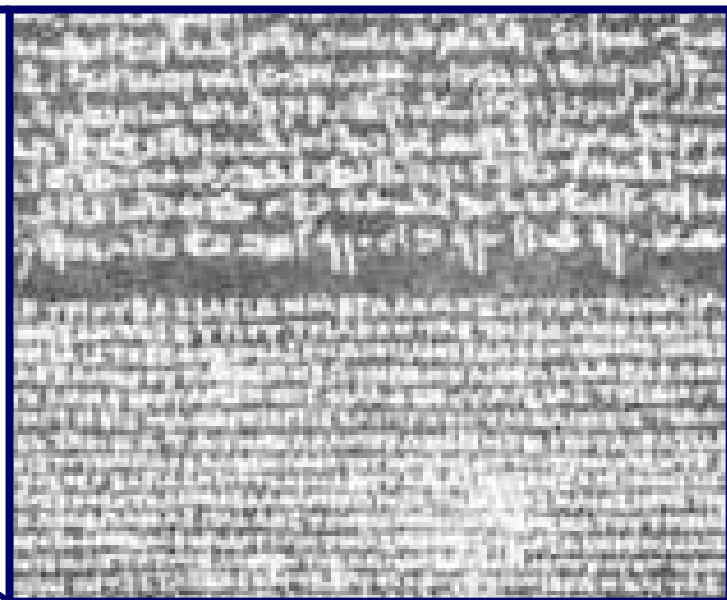
---

- Lexicons
  - Phrase books, bilingual dictionaries, ...
- Large text collections
  - Translations (“parallel”)
  - Similar topics (“comparable”)
- People

**Hieroglyphic**

**Demotic**

**Greek**





# Word-Level Alignment

---

**English**

Diverging opinions about planned tax reform

Unterschiedliche Meinungen zur geplanten Steuerreform

**German**

**English**

Madam President , I had asked the administration ...

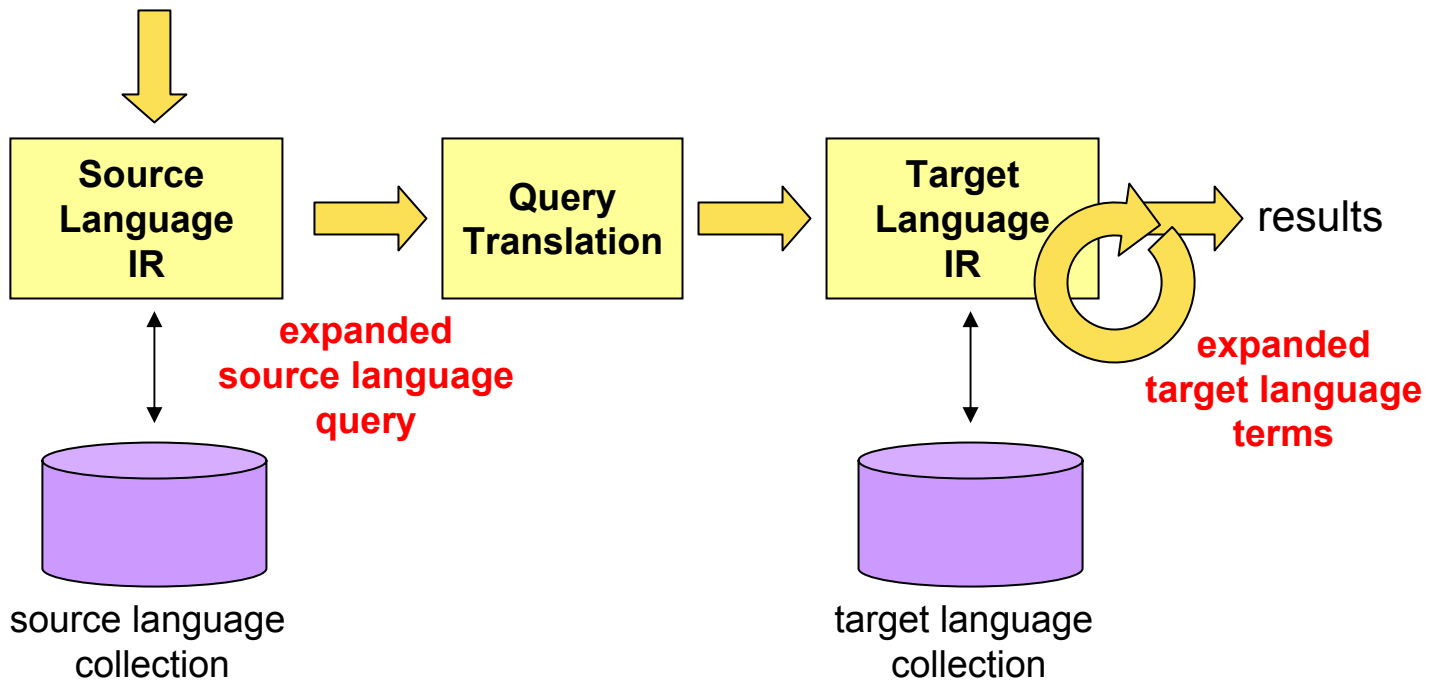
Señora Presidenta, había pedido a la administración del Parlamento ...

**Spanish**



# Query Expansion/Translation

source language query



**Pre-translation expansion**

**Post-translation expansion**



# TREC 2002 CLIR/Arabic

---

- Most recent (US-based) study in CLIR occurred at TREC
  - Results reported November 2002
- Problem was to retrieve Arabic documents in response to English queries
  - Translated Arabic queries provided for monolingual comparison
- Corpus of Arabic documents
  - 896Mb of news from Agence France Presse
  - May 13, 1994 through December 20, 2000
  - 383,872 articles
- Topics
  - 50 TREC topic statements in English
  - Average of 118.2 relevant docs/topic (min 3, max 523)
- Nine sites participated
  - 23 CL runs, 18 monolingual



# Sample topic

---

<top>

<num>Number: AR26</num>

<title>Kurdistan Independence</title>

<desc> Description:

How does the National Council of Resistance relate to the potential independence of Kurdistan?

</desc>

<narr> Narrative:

Articles reporting activities of the National Council of Resistance are considered on topic. Articles discussing Ocalan's leadership within the context of the Kurdish efforts toward independence are also considered on topic.

</narr>

</top>



# sample topic arabic document

<DOC>

<DOCNO>20000321\_AFP\_ARB.0001</DOCNO>

<HEADER>از01004 4 بش 8920 فير /اف-دزز03 اسرائيل/فلسطينيون</HEADER>

- <BODY>

<HEADLINE>جرح ثلاثة اسرئيليين اصابة اثنين منهم خطيرة في هجوم في الضفة الغربية</HEADLINE>

- <TEXT>

<P>القدس 12-3 (اف ب) - افادت حصيلة جديدة للجيش الاسرائيلي ان ثلاثة اسرئيليين جرحوا مساء امس الاثنين في هجوم جرى عندما اطلق</P>

<P>عليهم الرصاص من سيارة تجاوزت السيارة المدنية التي كانت تغلهم قرب ترقومية في محيط الخليل بالضفة الغربية</P>

<P>واوضح المتحدث باسم الجيش الاسرائيلي ان سائق السيارة التي كانت تغل الاسرئيليين، وهو من مستوطني الضفة الغربية اصيب بجروح</P>

<P>"طفيفة، ووصف حالة احد الجرحين الاخرين بانها "حرجة" وحالة الثاني بانها "خطيرة</P>

<P>وتشكل الخليل حيث يعيم 004 مستوطن يهودي بحماية الجيش الاسرائيلي ووسط 021 الف فلسطيني، بؤرة توتر بين الاسرئيليين والعرب، وقد</P>

<P>انسحبت اسرائيل في كانون الثاني/يناير 7991 من 08% من هذه المدينة وابقت على وجود عسكري كبير في الحي الذي يسكنه المستوطنون</P>

<P>وجرح الاسرئيليون الثلاثة عندما تعرضت السيارة التي كانوا فيها لاطلاق نار من سيارة اخرى تجاوزتها قرب بلدة ترقومية التي يؤدي اليها "الممر"</P>

<P>"الاهن" الذي يربط بين غزة وجنوب الضفة الغربية مروراً بالاراضي الاسرائيلية</P>

<P>وقد نقل الجرحان بسيارة اسعاف تم بمروحية الى مستشفى حداسا في القدس</P>

<P>وبدا الجيش عمليات بحث عن الفاعلين واقام حواجز على الطرقات</P>

<P>وابلغت السلطة الفلسطينية بملاسات الهجوم لتحاول العثور على مرتكبيه</P>

<P>واشاد المسؤولو الاسرائيليون في الفترة الاخيرة بالتعاون مع اجهزة الامن الفلسطينية في اطار مكافحة الارهاب</P>

<P>ونشرت بلدية مستوطنة كريات اريج الغربية من الخليل بيان احتجاج على سياسة السلام التي يتبناها رئيس الوزراء الاسرائيلي ايهود باراك الذي</P>

<P>"تهمة" نترك المستوطنين رهائن بايدي الفلسطينيين</P>

<P>وقال الجيش الاسرائيلي في تغذيرات اولية ان خلية تابعة لحركة المقاومة الاسلامية (حماس) قد تكون وراء الاعتداء</P>

<P>وتعارض حركة حماس بشدة اتفاقات اوسلو حول الحكم الذاتي الفلسطيني المبرمة عام 3991 وقد اعلنت مسؤوليتها عن غالبية الاعتداءات</P>

<P>التي استهدفت اسرائيل منذ ذلك الحين</P>

</TEXT>

<FOOTER>شفا/ 11 موا004 افب</FOOTER>

</BODY>

<TRAILER>جمت مار 00 405012</TRAILER>

</DOC>





# Stemming (Berkeley)

---

- Alternative way to build stem classes
- Trying to deal with complex morphology
- Use MT system to translate Arabic words
  - Now have (arabic, english) pairs
- Stop and stem all of the English words/phrases using favorite stemmer
  - (arabic, english-stem) pairs
  - If English stem is the same, then assume Arabic words should be in the same stem class
- (Also used a light stemmer)

# Stemming

Arabic word	English translation	Arabic word	English translation	Arabic word	English translation	Arabic word	English translation
أطفال	children	اطفالهم	their children	يطفل	by child	فما لطفلة	then the child
أطفالا	children	اطفالي	my children	يطفلة	by child	فطفل	then child
أطفالنا	our children	الاطفال	children	يطفالنا	by our child	كأطفال	as children
أطفاله	and his children	الاطفال	children	يطفلة	by his child	كالطفل	as the child
أطفاله	his children	الطفل	the child	يطفه	by his child	لأطفال	to children
أطفالها	her children	الطفلان	the children	يطفلها	by her child	لطفلها	to her child
أطفالهم	their children	الطفلة	the child	يطفلها	by their child	للطفلة	to the child
أطفالهن	their children	الطفلتان	the children	يطفلين	by children	وأطفالنا	and our children
أطفالي	my children	الطفلتين	the children	يطفلها	by her children	والأطفال	and the children
اطفال	children	الطفله	the child	طفل	child	ويطفل	and by child
اطفالا	children	الطفلين	the children	طفلا	child	ويطفلين	and by children
اطفالك	your children	بأطفال	by children	طفلان	children	وطفلة	and child
اطفالكم	your children	بأطفاله	by his children	طفلاها	her children	وطفتان	and children
اطفالكن	your children	بأطفالها	by her children	طفلة	child	وطقلنا	and our child
اطفالنا	our children	بالأطفال	by the children	طفلت	child	وطفلها	and her child
اطفاله	his children	بالطفل	by the child	طفلتان	children	وطقله	and his children
اطفالها	her children	بالطفلة	by the child	طفلة	his child	وطقلها	and her children
اطفالهم	their children	بالطفلتين	by the children	طفلتنا	our child	ولأطفالها	and to her children
اطفالها	their children	بالطفلين	by the children	طفلة	his child	وللطفل	and to the child



# UMass core approaches

---

- InQuery
  - For each English word, look up all translations in dictionary
- If not found as is, try its stem
  - Stem all Arabic translations
  - Apply operators
- Put Arabic phrases in `#filreq()` operator
- Use synonym operator, `#syn()`, for alternate translations
- Wrap all together in `#wsum()` operator
- Cross-language language modeling (after BBN)

$$P(Q_e|D_a) = \prod_{e \in Q_e} \left( \alpha \sum_{a \in Arabic} P(a|D_a)P(e|a) + (1 - \alpha)P(e|GE) \right)$$





# Breaking the LM approach apart

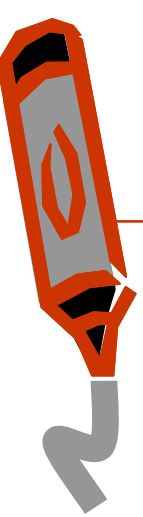
---

- Query likelihood model
- $P(a|D_a)$ 
  - Probability of Arabic word in the Arabic document
- $P(e|a)$ 
  - Translation probability (prob. of English word for Arabic word)
- $P(e|GE)$ 
  - Smoothing of the probabilities

$$P(Q_e|D_a) = \prod_{e \in Q_e} \left( \alpha \sum_{a \in Arabic} P(a|D_a)P(e|a) + (1 - \alpha)P(e|GE) \right)$$

# Calculating translation probabilities

---



- Dictionary or lexicon
  - Assume equal probabilities for all translations
  - Unless dictionary gives usage hints
- Parallel corpus
  - Assume sentence-aligned parallel corpora
- Know that sentence  $E$  is a translation of sentence  $A$ 
  - Estimate  $P(e|a)$  from those aligned sentences
  - Consider sentence pairs  $(E,A)$  where  $e$  is in  $E$  and  $a$  is in  $A$
  - To get  $P(e|a)$ , divide by number of Arabic sentences containing  $a$

$$P(e|a) = \frac{|\{(E, A) | e \in E \text{ and } a \in A\}|}{|\{A | a \in A\}|}$$



# Other techniques

---

- Query expansion
  - Useful to bring in additional related words
  - Same as in monolingual retrieval
- Expand query in English
  - Need comparable corpus (why comparable?)
  - Brings in synonyms and other words related to query
- Expand translated query in Arabic
  - Done on actual target corpus
  - Brings in Arabic synonyms not in dictionary
  - BBN in TREC 2002 was careful to expand only by translation of original query words
- Can do neither, either, or both
- UMass added 5 terms from English and 50 from Arabic
  - For LM runs, used “relevance modeling” instead in Arabic



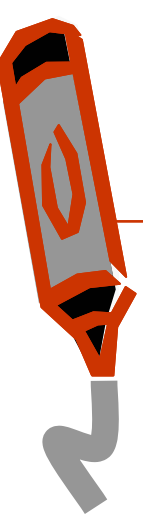
# CLIR better than IR?

---

- How can cross-language beat within-language?
  - We *know* there are translation errors
  - Surely those errors should *hurt* performance
- Hypothesis is that translation process may disambiguate some query terms
  - Words that are ambiguous in Arabic may not be ambiguous in English
  - Expansion during translation from English to Arabic prevents the ambiguity from re-appearing
- Has been proposed that CLIR is a model for IR
  - Translate query into one language and then back to original
  - Given hypothesis, should have an improved query
  - Should be reasonable to do this across many different languages

# International Research Programs

---



- Major ones are
  - TREC(US DARPA under TIDES program),
  - CLEF(EU) and
  - NTCIR (Japan)
- Programs were initially designed for ad-hoc cross-language text retrieval, then extended to multi-lingual, multimedia, domain specific and other dimensions.



# CL image retrieval, CLEF 2003

---

- A pilot experiment in CLEF 2003
- Called ImageCLEF
- Combination of image retrieval and CLIR
- An ad hoc retrieval task
- 4 entries
  - NTU (Taiwan)
  - Daedalus (Spain)
  - Surrey (UK)
  - Sheffield (UK)



# Why a new CLEF task?

---

- No existing TREC-style test collection
- Broadens the CLEF range of CLIR tasks
- Facilitates CL image retrieval research
- International forum for discussion



# CL image retrieval, CLEF 2003

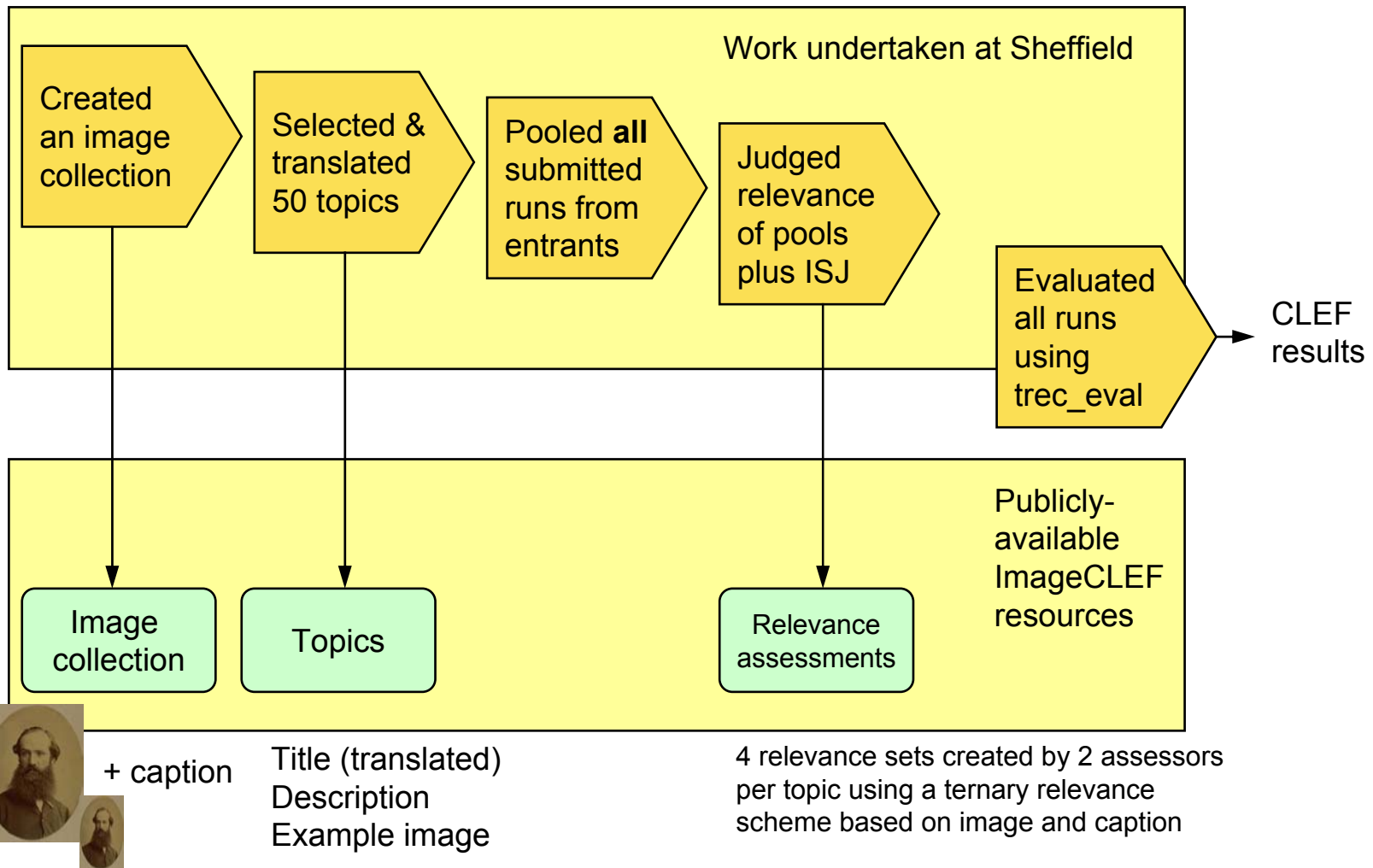
---

Given a user need expressed in a language different from the document collection, find as many relevant images as possible

- Fifty user needs (topics):
  - Expressed with a short (title) and longer (narrative) textual description
  - Also expressed with an example relevant image (QBE)
  - Titles translated into 5 European languages (by Sheffield) and Chinese (by NTU)
- Two retrieval challenges
  - Matching textual queries to visual documents (*use captions*)
  - Matching non-English queries to English captions (*use translation*)
- Essentially a bilingual CLIR task
- No retrieval constraints specified



# Creating the test collection





# Evaluation

---

- Evaluation based on most stringent relevance set (strict intersection)
- Compared systems using
  - MAP across all topics
  - Number of topics with no relevant image in the top 100
- 4 participants evaluated (used captions only):
  - NTU – Chinese->English, manual and automatic, Okapi and dictionary-based translation, focus on proper name translation
  - Daedalus – all->English (except Dutch and Chinese), Xapian and dictionary-based + on-line translation, Wordnet query expansion, focus on indexing query and ways of combining query terms
  - Surrey – all->English (except Chinese), SoCIS system and on-line translation, Wordnet expansion, focus on query expansion and analysis of topics
  - Sheffield – all->English, GLASS (BM25) and Systran translation, no language-specific processing, focus on translation quality



# Results

---

- Surrey had problems
- NTU obtained highest Chinese results
  - approx. 51% mono and 12 failed topics (NTUiaCoP)
- Sheffield obtained highest
  - Italian: 72% mono and 7 failed topics
  - German: 75% mono and 8 failed topics
  - Dutch: 69% mono and 7 failed topics
  - French: 78% mono and 3 failed topics
- Daedalus obtained highest
  - Spanish: 76% mono and 5 failed topics (QTdoc)
  - Monolingual: 0.5718 and 1 failed topic (Qor)
- For more information ... see the ImageCLEF working notes