# The Challenges of Multilingual Search

**Paul Clough**

The Information School
University of Sheffield

# Outline

* What is multilingual search?
* MLIR and CLIR techniques
* Four challenges
    * Thinking beyond search
    * Translation and language resources
    * Providing effective user support
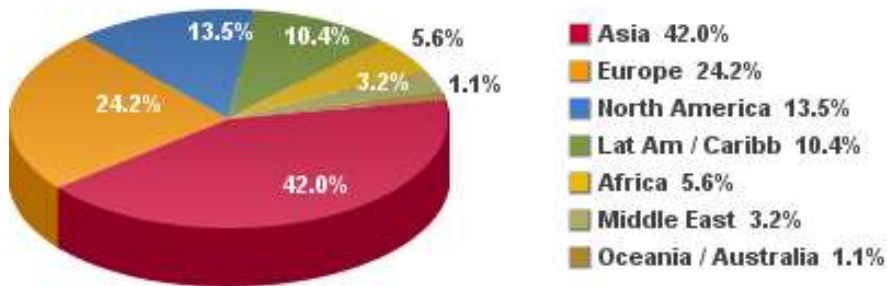    * Going from research to practice

# The "Grand [search] Challenge"

*"Given a query in any medium and **any language**, select relevant items from a **multilingual** multimedia collection which can be in any medium and **any language**, and present them in the style or order most likely to be useful to the querier, with identical or near identical objects in **different** media or **languages** appropriately identified."*

Douglas W. Oard and David Hull, AAAI Symposium on Cross-Language IR, Spring 1997, Stanford, USA
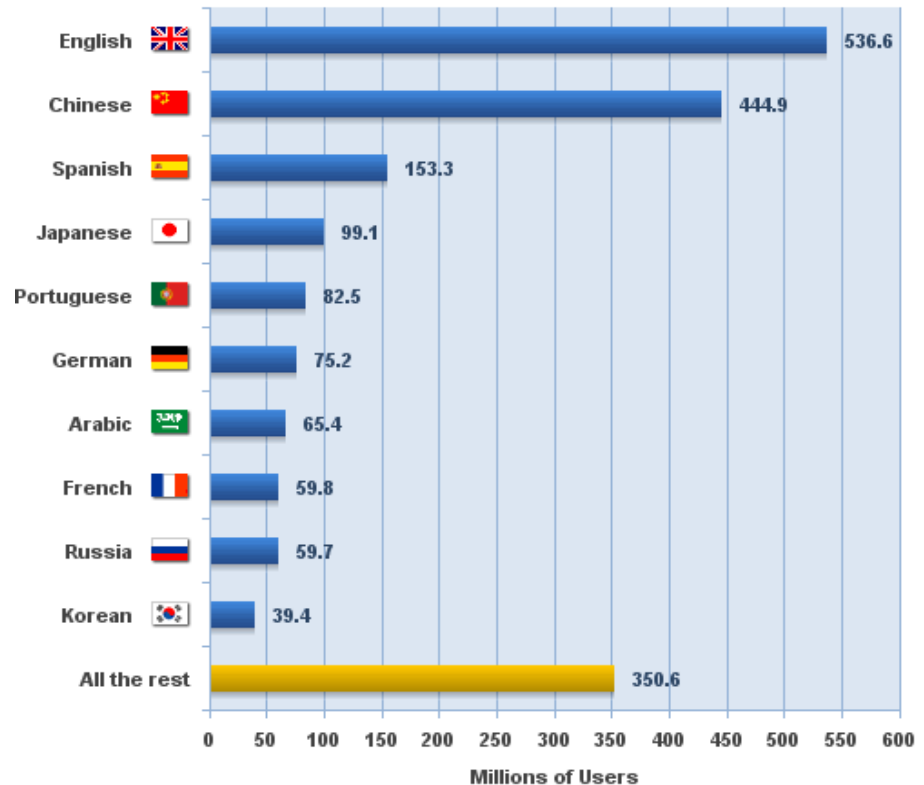
# The need for multilingual search

## Internet Users in the World
## Distribution by World Regions - 2010

- Asia 42.0%
- Europe 24.2%
- North America 13.5%
- Lat Am / Caribb 10.4%
- Africa 5.6%
- Middle East 3.2%
- Oceania / Australia 1.1%

Source: Internet World Stats - www.internetworldstats.com/stats.htm
Basis: 1,966,514,816 Internet users on June 30, 2010
Copyright © 2010, Miniwatts Marketing Group

http://www.internetworldstats.com/stats.htm

## Top Ten Languages in the Internet
## 2010 - in millions of users

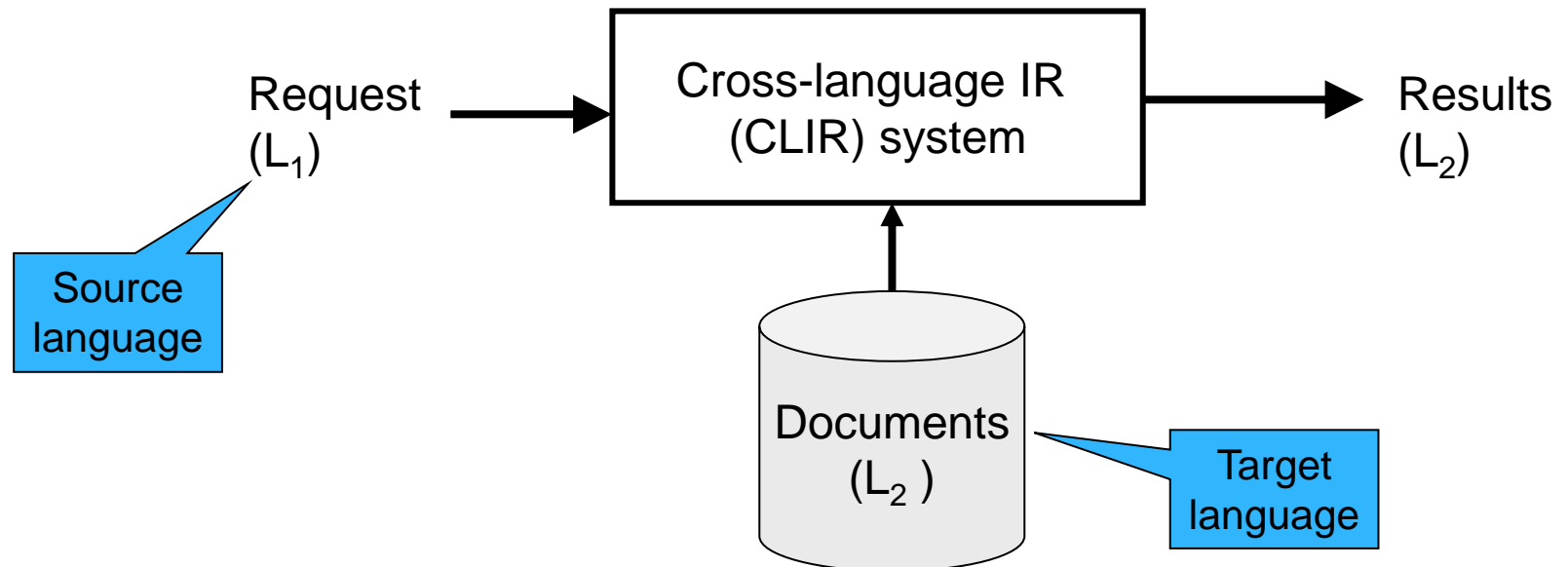| Language | Millions of Users |
|---|---|
| English | 536.6 |
| Chinese | 444.9 |
| Spanish | 153.3 |
| Japanese | 99.1 |
| Portuguese | 82.5 |
| German | 75.2 |
| Arabic | 65.4 |
| French | 59.8 |
| Russia | 59.7 |
| Korean | 39.4 |
| All the rest | 350.6 |

Source: Internet World Stats - www.internetworldstats.com/stats7.htm
Estimated Internet users are 1,966,514,816 on June 30, 2010
Copyright © 2000 - 2010, Miniwatts Marketing Group

ISKO UK conference 8-9 July 2013

# What is multilingual search?

* **Multilingual Information Access (MLIA)**
  * Accessing, querying and retrieving information from collections in any language
  * Includes search and browse functionalities
* **Cross-Language Information Retrieval (CLIR)**
  * Querying multilingual collections in one language in order to retrieve documents in other languages
  * Bilingual information retrieval
* **Multilingual Information Retrieval (MLIR)**
  * Process information (queries, documents, both) in multiple languages
  * Includes cross-language retrieval

# Cross-language search

* Documents and user requests are in different languages (**bilingual retrieval**)

Request $(L_1)$ → Cross-language IR (CLIR) system → Results $(L_2)$

Source language

Documents $(L_2)$

Target language

# Multilingual search

* Documents in collection in different languages, search requests in any language

Request
(L?)  →  **Multilingual IR (MLIR) system**  →  Results ($L_2$, $L_3$ and/or $L_4$)

e.g. the Web

Documents
($L_2$)

Documents
($L_3$)

Documents
($L_4$)

# Google web search

Google

still life paintings | Search

About 2,510,000 results (0.88 seconds)

**Everything**
Images
More

**Translated foreign pages** ☒

Translated results for **still life paintings** - My language: English ▼

| Language | Translated query | |
|---|---|---|
| French ☒ | **natures mortes** - Edit | 532,000 results |
| Spanish ☒ | **bodegones** - Edit | 380,000 results |
| Russian ☒ | **Натюрморты** - Edit | 1,410,000 results |
| Japanese ☒ | **静物画** - Edit | 188,000 results |

Add language ▼ - Automatically select languages to search

**Still Life** - Wikipedia
Translated from: French
The term refers to a **still life** subject matter, inanimate objects (fruit, flowers, vases, etc.). Or dead animals, then, by metonymy, a work (in...
➕ Show original text

**Open**

fr.wikipedia.org/wiki/Nature_morte

**Still Life** - Wikipedia
Translated from: Japanese
**Still Life** (The Animals blame) is one of the genres of Western **painting**, of a stationary nature (flowers, skulls, hunting prey, shells, vegetables, fruits, fish and kitchen) and artifacts (glass cups, ceramics,...
➕ Show original text
ja.wikipedia.org/wiki/静物画

The web
Pages from the UK

Any time
Latest
Past 24 hours
Past week
Past month
Past year
Custom range...

Standard view
Related searches
Wonder wheel
Timeline

Standard results
Sites with images
**Translated foreign pages**

# Obvious question …

*… "Why do users want to retrieve documents they presumably can't read?"*

* Some users are **multilingual**
  * Can formulate searches and judge relevance in many languages
  * Want the convenience of a single query
* Some users are **monolingual**
  * Want to query in their native language
  * Can judge relevance even if results not translated
  * Have access to document translation
  * Objects retrieved are language-independent (e.g. images)

# Matching queries and documents

* MLIR and CLIR involves matching queries in one language with documents in another
  * Translate the query (**query translation**)
  * Translate the documents (**document translation**)
  * Translate queries and documents (e.g. into an intermediate language)
  * No translation
* MLIR may also involve a **merging** step

# Translation resources

* The use of machine-readable, human-crafted **bilingual dictionaries** (and thesauri, word-lists and other similar resources)
* The use of translation resources generated by statistical approaches from suitable training data
  * Require parallel or comparable corpora
* The use of a 'full' **machine translation** system
  * Commercial or open source

# Does CLIR/MLIR work?

* In the lab CLIR/MLIR using various approaches and translation resources is effective

* Results from the Cross-Language Evaluation Forum (CLEF) have shown steady increases

  * **1997** 50-60% of monolingual baseline for various tasks
  * **2003** 80% onwards for most commonly used language pairs
  * **2009** up to 99% of monolingual baseline

* Main findings

  * Well-tuned MT systems are highly effective
  * Combinations of approaches work well

# Challenges for multilingual search

## 1. Thinking beyond search

# Search in the broader context

Information
seeking context
**(includes culture and language)**

Information
Need

Seeking
Strategy

Identify suitable sources
Decide search plan

IR system

Verbalisation

Formalise information need

IR
application

Decision-making

1. Reformulate query
2. Begin new search
3. Terminate search

Search

Browse

Evaluate
Outcomes

Examine results
Examine specific document

Use
Information

# Developing global applications

"Globalization is the process of making applications work seamlessly **utilizing the user's preferred language and culture**. It involves not just programming and deployment skills but also cultural, translation and language expertise as well."

http://www-01.ibm.com/software/globalization/topics/index.html

*Multilingual search will typically be part of a wider globalization process and concerned with developing multilingual applications (or websites) that include search functions (i.e. search is not the end goal)*

# The impact of culture

* Culture is the behaviour typical of a group of people
    * Members of the same culture are likely to have the same knowledge of certain things and would think and act similarly in certain situations
* Aspects to consider (for design) include
    * Religion, customs, colours, metaphors, icons and flags, and **language**
* Crossing the cultural boundary includes adapting to a given market's cultural conventions (**localisation**)

# Aspects to consider for localisation

* **Translation** of the product's interface and documentation
* **Colours, images, graphics and icons** adapting to cultural and legal requirements
* **Rendering** displaying text correctly (e.g. does the new text fit inside the allocated space?)
* **Fonts** making sure the correct fonts and characters for the target language
* **Bi-directional text** needed in Arabic and other languages
* **Locale data** how to display dates, time, number, currency and other regional data

# Formatting and page layout



For languages read right to left right alignment is predominantly used

# Differences in terminology

# Challenges for multilingual search

## 2. Languages and translation

# What is translation?

* Isn't translation just replacing words in one language (source) with words in another language (target)?

Diverging opinions about the planned tax reform

Unterschiedliche Meinungen zur geplanten Steuerreform

Simples!

* Not quite so simple …
  * Often includes changes in syntax
  * Often includes translation of semantics and use of culturally-specific terms or concepts

# Further translation challenges

* Out-Of-Vocabulary (OOV) terms
* Ambiguity (source and target languages)
    * Lexical ambiguity (e.g. bat: cricket or animal?)
    * Syntactic ambiguity (e.g. "I saw the boy with the telescope")
* Translation of phrases (e.g. "George Bush")
* Translation of proper names
* Word inflections (e.g. in Swedish)
* Compounds (e.g. in German and Dutch)
* Word segmentation (e.g. in Arabic and Chinese)
* …

# Handling different languages

* Character sets
* Punctuation and marks
* Word separation
* Digits
* Writing direction
* Formatting styles
* Character shapes
* Sort order and case
* Date and time formatting

รายงานข่าวจากธนาคารแห่งประเทศไทยเปิดเผยว่า การให้บริการ
บัตรเครดิตแยกตามประเภทบัตรเครดิตไตรมาส 2 ของปีนี้เทียบกับ
ไตรมาสแรก โดยปริมาณการใช้จ่ายผ่านบัตรเครดิตโดยรวมลดลง
จาก 8.3 หมื่นล้านบาทในสิ้นไตรมาสแรก เหลือ 7.8 หมื่นล้านบาท
ซึ่งทั้งปริมาณการใช้จ่ายในประเทศลดลงจาก 6.4 หมื่นล้านบาท
เหลือ 5.9 หมื่นล้านบาทในไตรมาสสองของปี ส่วนปริมาณการใช้จ่าย
ในต่างประเทศลดลงจาก 2.6 พันล้านบาท เหลือ 2.1 พันล้านบาท
ขณะที่ยอดสินเชื่อคงค้างปรับตัวเพิ่มขึ้นจากสิ้นไตรมาสแรกที่ระดับ
5.8 หมื่นล้านบาท เป็น 6.3 หมื่นล้านบาทในไตรมาสที่สอง

# Lexical resources for translation

* The use of machine-readable, human-crafted **bilingual dictionaries** (and thesauri, word-lists and other similar resources)

* The use of translation resources generated by statistical approaches from suitable training data
  * Require parallel or comparable corpora

* The use of a 'full' **machine translation** system
  * Commercial or open source

# The "resource bottleneck"

* There are approximately 6,800 known languages in the world and just over 2,000 have a writing system
  * **Around 300** have some kind of language processing tools
  * CLIR/MLIR performance depends on the availability of high-quality translation resources and language processing tools
* **"Resource bottleneck"**
  * Finding ways to acquire, maintain and update language tools and resources in economic manner
  * Building and sustaining resources is costly

# Example research: ACCURAT



* ACCURAT project investigated automatically gathering materials from various online sources to train MT systems
  * Focused on generating resources for "rare" language pairs and specific domains (e.g. automotive industry)
  * Used **comparable corpora** to derive translations
* Used various online resources
  * Wikipedia, news sites, parallel versions of websites, social media and crawls of general web pages
* Download Toolkit: http://www.accurat-project.eu/

# Challenges for multilingual search

## 3. Providing effective user interaction

# Adapting to language differences

* Individuals can have a **range of foreign language** abilities and knowledge, e.g. passive vs. active skills
    * Can use translation to deal with language differences
* Many people can use **English**
    * Often an interlingua (or default language) for many interfaces
* User's language skills will affect the design of interfaces, for example for search
    * e.g. monolingual users may need help formulating queries

# Supporting CLIR/MLIR



Help the user formulate their query

Help the user select translations

Help the user identify possible relevant result items

Formulate Query — query — Translate Query — translated query — IR application

Evaluation of resources

results (ranked or clustered)

reformulate query

reselect translation

Examine Results — selected document — Examine Document

reselect document

**Decision-making**

1. Reformulate query
2. Begin new search
3. Terminate search

Iterative search process

Help the user read/view a selected document

# Query translation (Mulinex)

# Translation of search results

# Inputting non-ASCII characters



Example of non-ASCII keyboard character input (Arabvista.com).
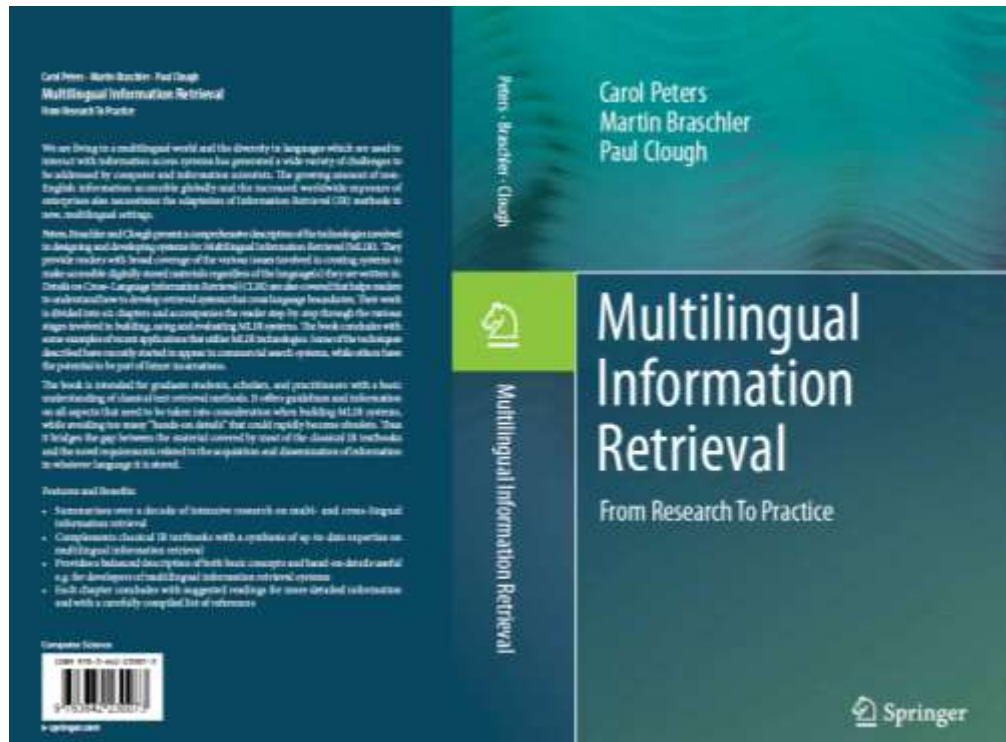
# Giving users a choice of interface language

# Challenges for multilingual search

## 4. Going from research to practice

# Are we stuck in the lab?

* The technologies for CLIR/MLIR are proven in the lab, but where are the applications in practice?
    * For example, consider the lack of CLIR/MLIR in large-scale information services, such as Amazon
* We need to go from research to practice
    * Guidelines for developing CLIR/MLIR applications
    * Collaborations between academics and business
    * Effective methods for knowledge transfer (in both directions)

# Research to practice

# Summary

* Search systems are used interactively so you must consider and design for **end users**
  * **Implication**: multilingual search is NOT just an engineering problem; you will also need to understand users' cultural background and language abilities
* Searching is part of wider information seeking activities supported by search **applications**
  * **Implication**: you may have to consider localisation of the whole service and not just querying support
* Research has shown we can do cross-language search in multiple languages well (**in the lab**)
  * **Implication**: the challenge is going from research to practice

# Questions?

**Paul Clough**

**Senior Lecturer in Information Retrieval**

p.d.clough@sheffield.ac.uk

http://ir.shef.ac.uk/cloughie/