

The multilingual facet of the CACAO project

Luca Dini, CELI, Italy
Frédérique Segond, Xerox Research Centre Europe, France

eContentplus Program



<http://www.cacao-project.eu>

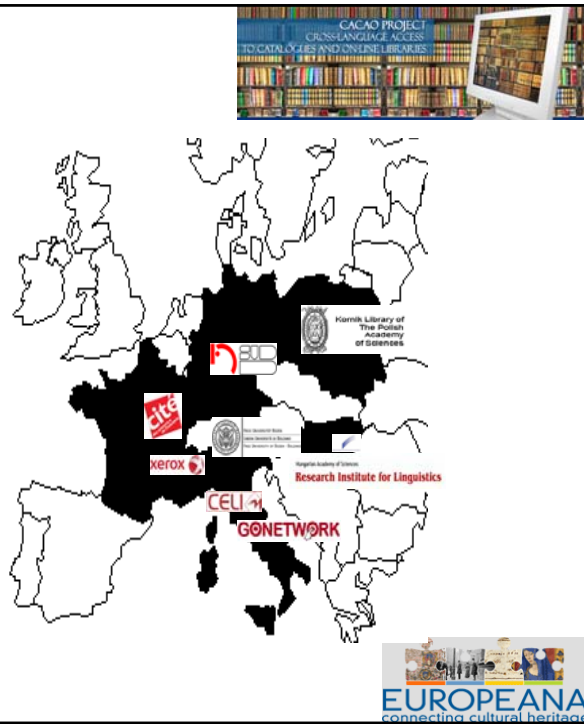


CACAO enables international users to better access and exploit the available European multilingual electronic A.O.C. content.



Who are we?

- ❖ Xerox Research Centre
Europe (France)
- ❖ Bibliothèque Cité des
Sciences (France)
- ❖ CELI (Italy)
- ❖ Free University of Bozen-
Bolzano (Italy)
- ❖ Gonetnetwork (Italy)
- ❖ University of Goettingen
(Germany)
- ❖ Kórnik Library (Poland)
- ❖ Hungarian Academy of
Sciences (Hungary)
- ❖ National Széchényi Library

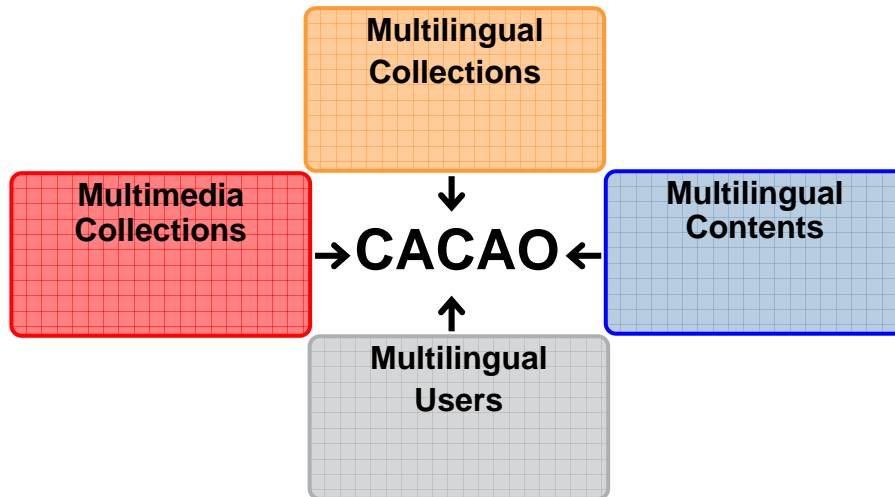


What CACAO
does for whom?

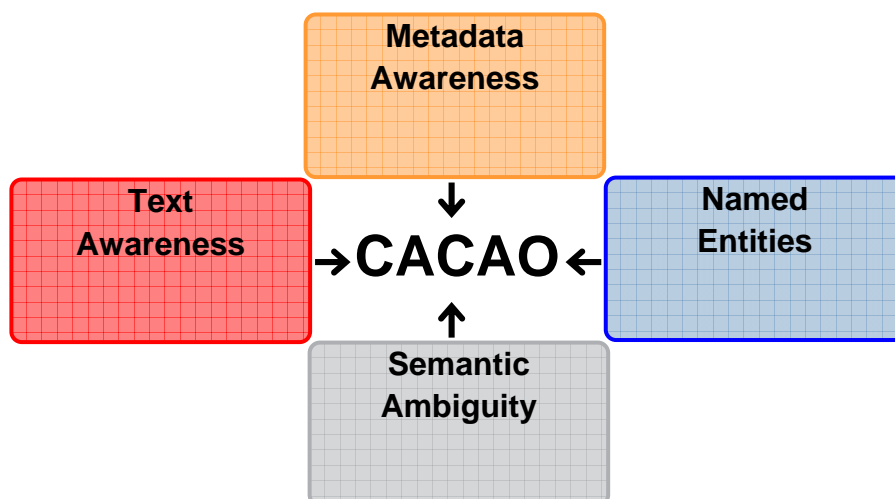
CACAO helps individual citizens
and librarians *access, understand and
navigate multilingual* textual digital
libraries and Online Public Access
Catalogue content.

CACAO
PROJECT

Facets of Multilinguality



Facets of Search in DL



What are the problems?

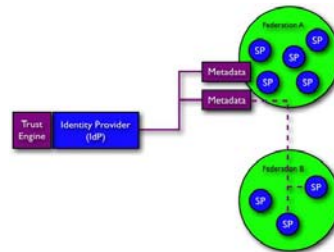
Information search and retrieval

Metadata search (text)

Multilingual search



What are the basic issues of multilingual search?



Structured documents

While traditionally a document is just a text, structured documents might contained metadata that are completely structured and expressed in different languages (titles, body, footnotes, etc.). These textual metadata go (might go) to a different *index field*.

Documents and fields

Metadata:

- ❖ Title: The Tempest
- ❖ Genre: Drama
- ❖ Abstract: The sorcerer Prospero, rightful Duke of Milan...

“prospero shakespeare”?

1. Document level : The Tempest, drama, shakespeare. The sorcerer Prospero, rightful Duke of Milan.
2. Metadata level: “+abstract:prospero +author:shakespeare”

As for text indexing it is also possible to give some more weight to certain metadata (e’g to boost generic fields). Metadata often contain free text that can be used for better indexing.

CACAO Indexing metadata content

Immagini Roma

~~> I papi, Roma e Dante: l'idea e le immagini di Roma nella
Commedia dantesca Di Nicola Longo Pubblicato da Bulzoni, 2004~~



Dante Alighieri

~~> All'angolo di questa via con **corso Dante Alighieri**...~~

- ~~> Books of *Dante Alighieri*?~~
- ~~> Books about *Dante Alighieri*?~~

CACAO: Indexing textual content

❖ CACAO offers extended indexing when free and/or full text is available. This can include an abstract, description or even the whole document.

❖ CACAO also offers thorough analyses of texts, thanks to automatic detection of personal names, locations, dates etc. which are usually difficult to determine in an approach that is not NLP-based.

CACAO: Indexing textual content

Title: *La nascita del teatro moderno*

Abstract: *La "Storia del teatro moderno e contemporaneo" è articolata in tre volumi, ai quali si aggiunge un volume finale che raccoglie le trame dei mille testi teatrali fondamentali dal Cinquecento a oggi. Il primo volume è dedicato alla nascita del teatro moderno, di cui si analizzano i molteplici aspetti della "rottura" rispetto alla tradizione medievale, dove l'esperienza teatrale era ampiamente diffusa ma non veniva riconosciuta come tale a causa del rifiuto ideologico di cui era fatta oggetto da parte della dominante cultura cristiana.*



A query such as: "*teatro storia*" would not return this document, as the only elements which might match are in the abstract and not in the title.

CACAO: Indexing textual content

~~> 0 Hits~~

pterodattilo

> "Manuale di anatomia comparata dei vertebrati" di Emanuele Padoa ('pterodattilo' pag.235)

~~> "La certosa di Parma" di Stendhal~~

Parma

> "Storia della città di Parma" di Ireneo Affò

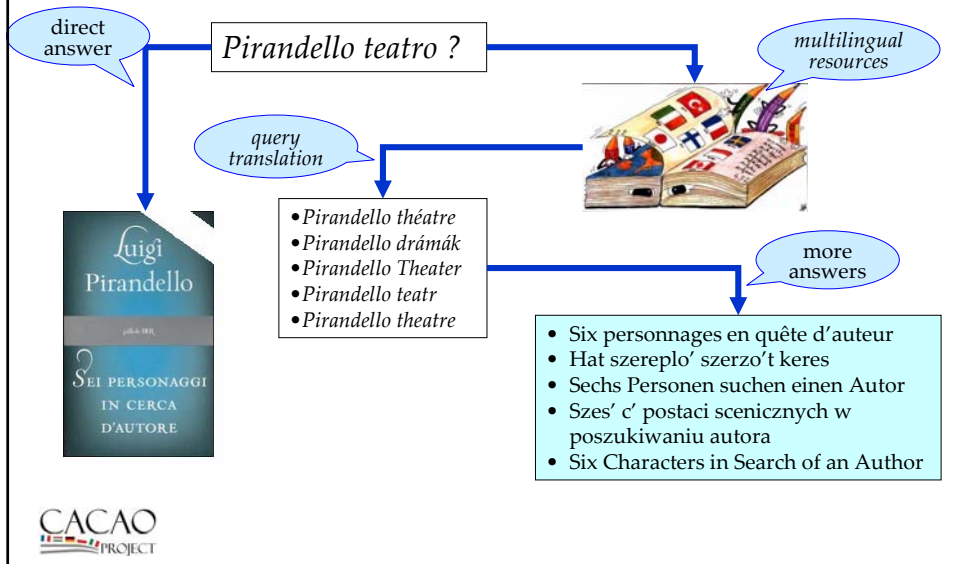
What are the modules for Multilingual and cross-lingual search?



The basic components to perform multilingual and cross lingual search

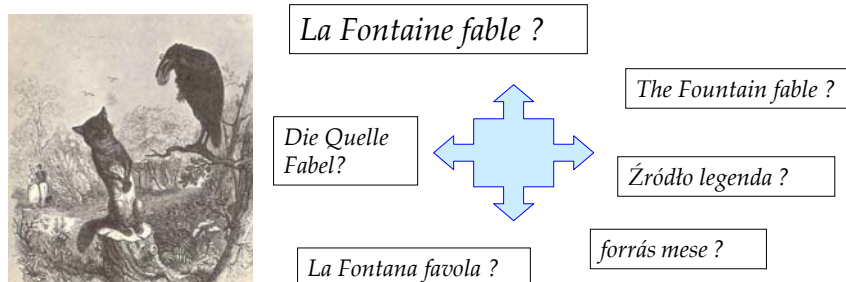
- ❖ All the basic components of Information retrieval and search plus a *query translator*:
 - ❖ It intersects with all other parameters (lemmatization/expansion/stopwording/NER * 2)
 - ❖ Additional issue:
 - ❖ Compounding
 - ❖ NER
 - ❖ Lemmatization level
 - ❖ Harmony of resources

A query in one language, answers in different languages



Proper names can be ambiguous

A query is not just a sequence of words that can be simply translated



Proper names and locations should be identified to avoid wrong translations.

Words are ambiguous

avocat

- ▶ ~~"Profession d'avocat" de Armand-Gaston Camus~~
- ▶ ~~"Le grand livre des fruits tropicaux" par F. Le Bellec.~~

- ▶ "avocat" as a profession?
- ▶ "avocat" as a fruit?

lawyer

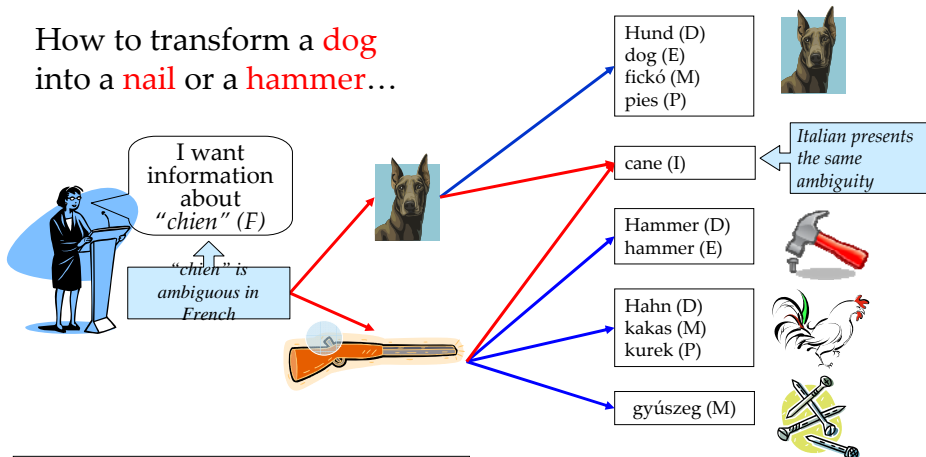
- ▶ ~~"Le grand livre des fruits tropicaux" par F. Le Bellec.~~

- ▶ "Profession d'avocat" de Armand-Gaston Camus



Words are ambiguous

How to transform a **dog**
into a **nail** or a **hammer**...



Translation does not reduce to a simple dictionary search...

D=Deutsch M=Magyar F=Français
E=English P=Polski I=Italiano



