

University of Tunis II
Ecole Nationale des Sciences de l'Informatique

Hannibal - A Search Engine
With a plug-in in Arabic

by Imed CHIHI

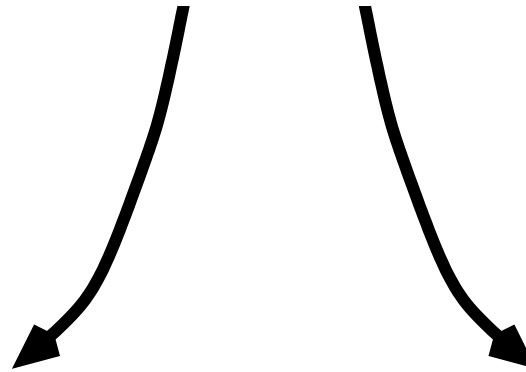
Supervised by

Mohamed Said OUERGHI, ENSI

Moez SOUABNI, MRS

The problem

**Web growth
in the Arab World**



Accessing documents
in Arabic

Finding
desired resources

Part 1: Problems of the Arabic language,
A more general issue: internationalization,
Possible approaches,
Technology used.

Part 2: A search service,
Functions.

Conclusion & perspectives

Problem categories

1. character representation,
2. inconsistencies due to multilingual typesetting,
3. satisfy the needs without breaking with the standards.

Multiple character sets,

Typesetting methods not well established,

Complex morphology,

Complex physiology,

Bidirectionnality.

US-ASCII used to be *the* character set

New tendencies:

- ISO/IEC 10646, then Unicode,
- Enhancements to HTML,
- Using MIME with HTTP,
- ISO 10646 as the privileged character set,
- RFC 1556 for bi-directionnality.

➔ Needs *intelligent* user agents.

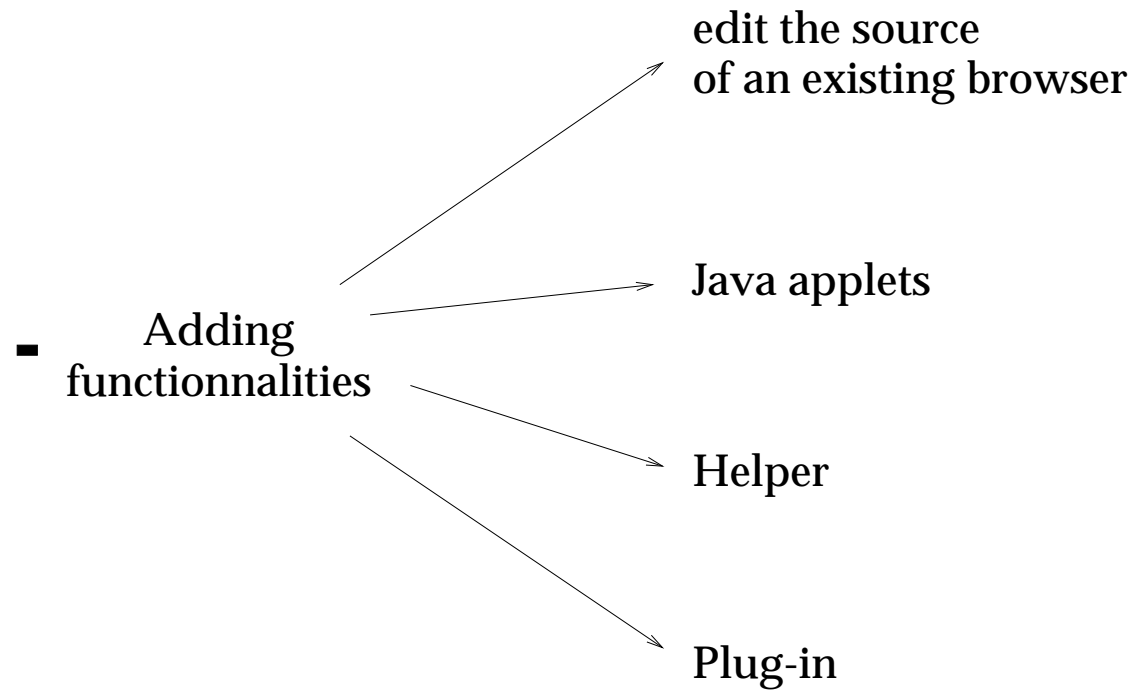
User agents

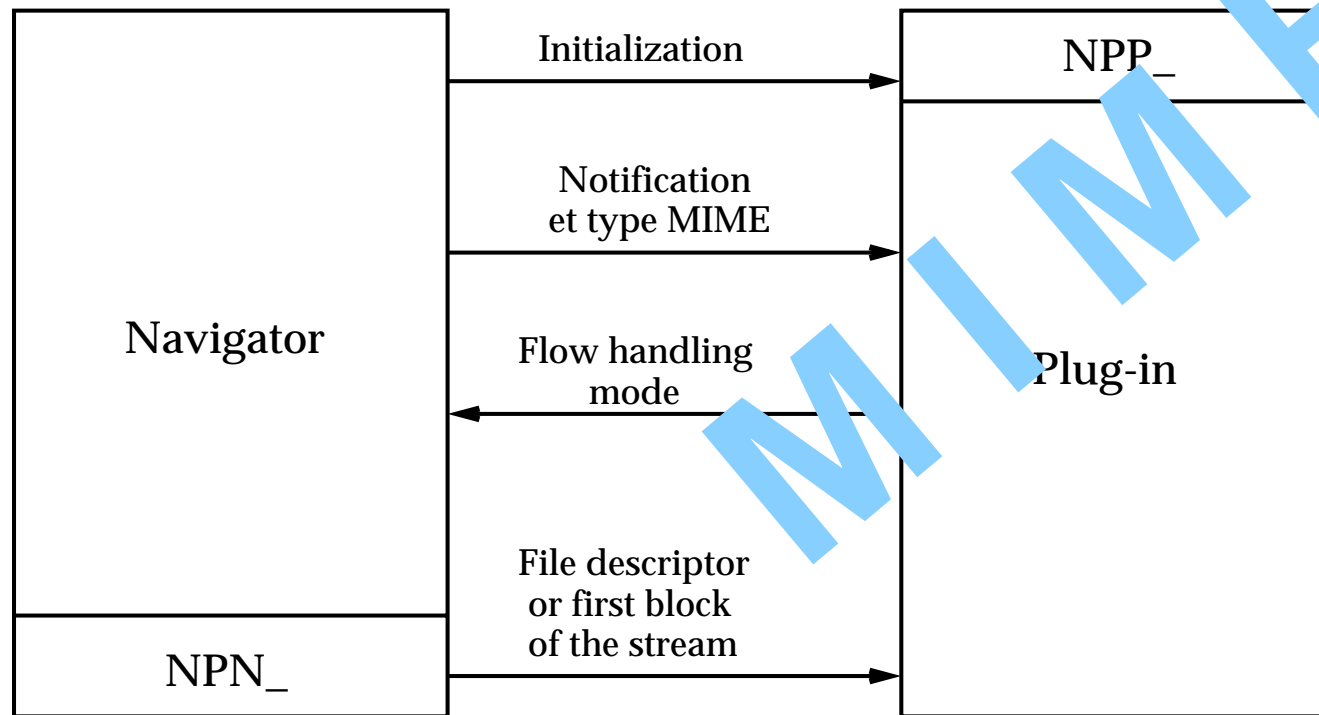
	CP 1256	ISIRI 3342	ISO 8895-6	UTF	Unix	Mac	Windows
Sakhr Sindbad			X	X			X
Alis Tango	X	X	X	X			X
Accent Multilingual Mosaic	X		X				X
Netscape Communicator				X	X	X	X
LangBox AraMosaic	X	X	X		X		
PMosaic		X			X		
Internet Explorer				X	X	X	X

Solving the problem

7

- "Rely" on Unicode





Registered MIME type

text/html; charset=iso-8859-6

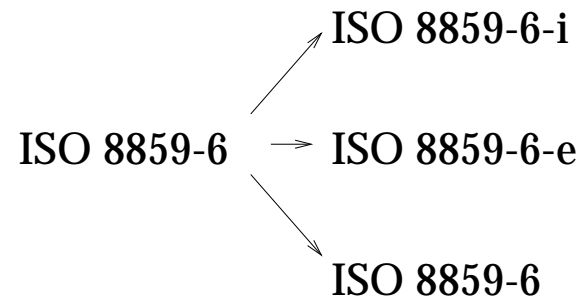
- a too large plug-in,
- not the objective,
- the solution does exist

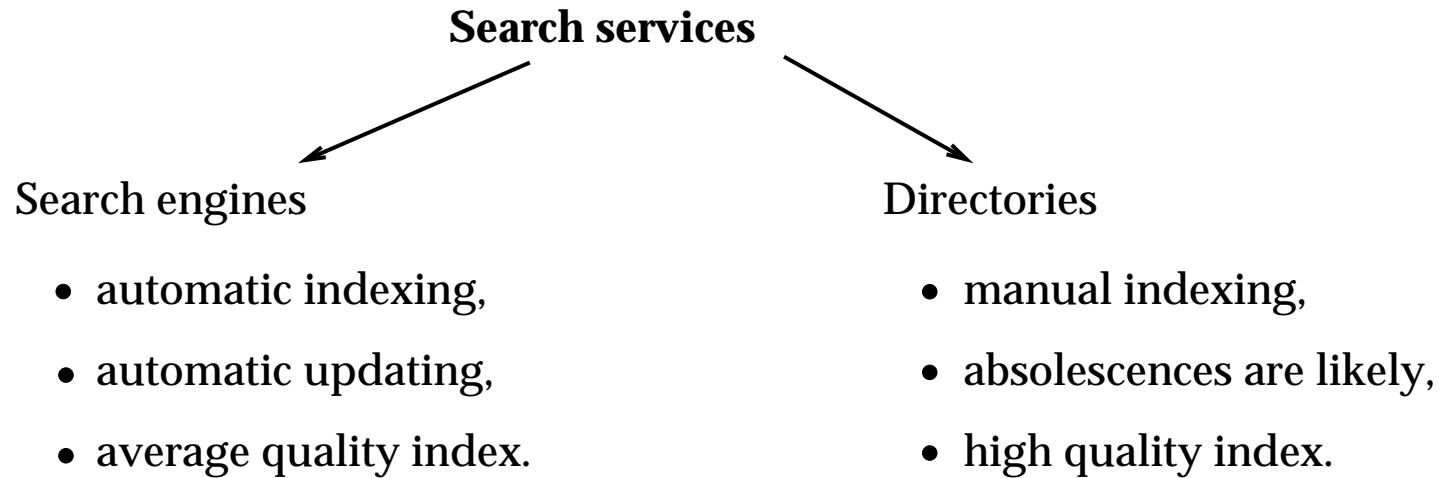
→ text/plain: unformatted text

The charset:

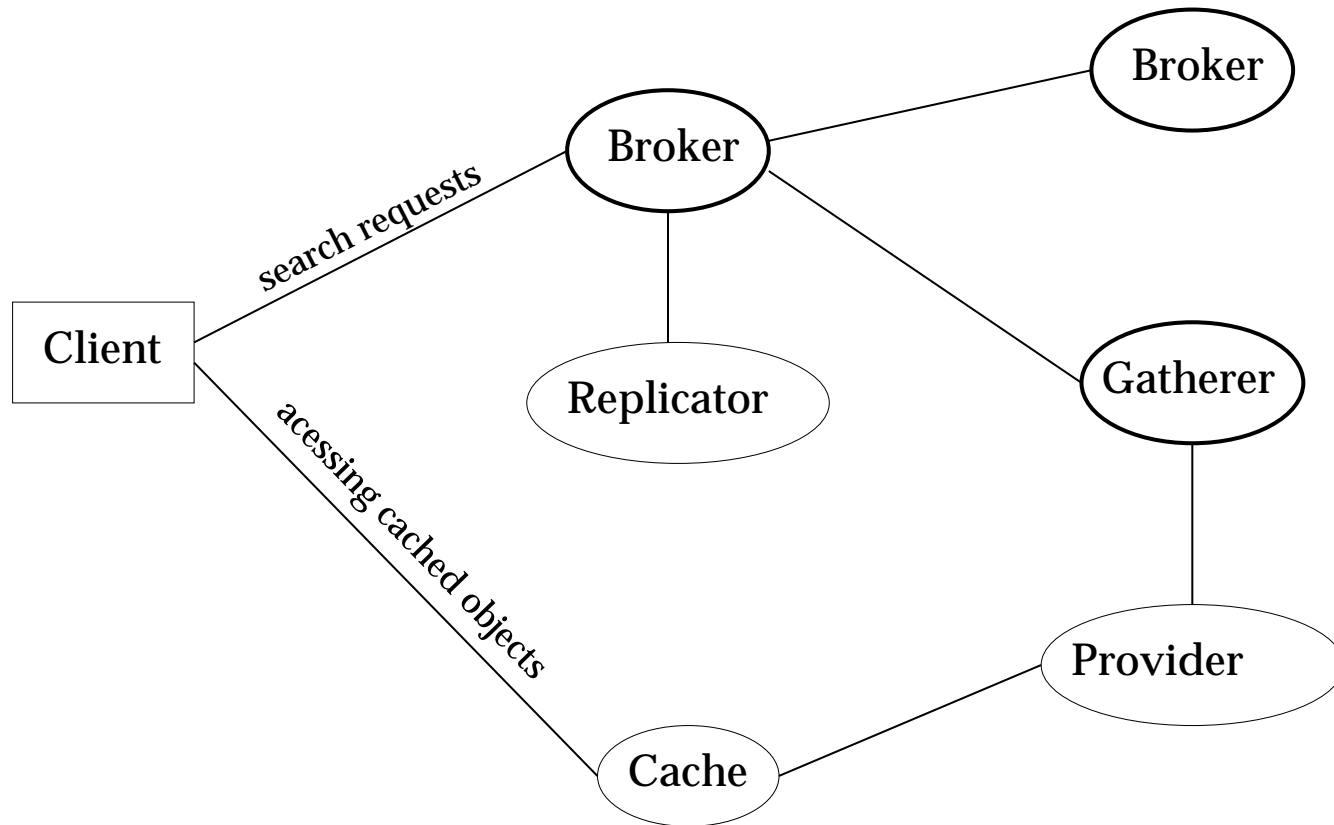
CP 1256 proprietary

ASMO_449
et ISIRI 3342 no Latin glyphs



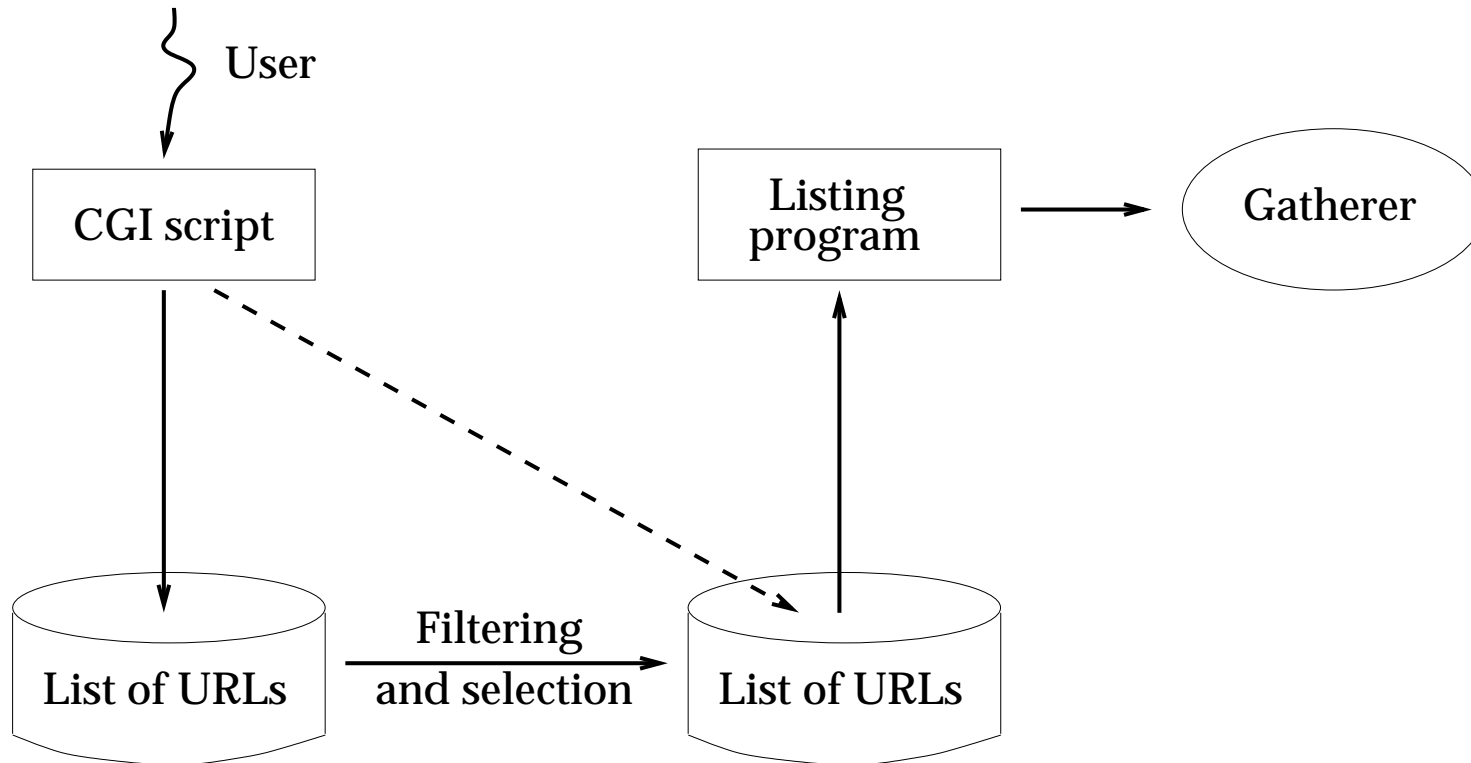


The Harvest system



The *Add URL* service

12



Case-sensitivity toggelling,

Search on word boundaries,

Searches with spelling errors,

Display of matched lines.

- What I did:
 - An overview of the existing standards and technologies,
 - A tool for Arabic documents visualization,
 - A functional search service.

- The ultimate solution for Arabic:

Unicode

- Plug-in
 - ports to other platforms,
 - adding other encodings support (especially UTF-8),
 - promote author systems in Arabic.

- The search service
 - distribute the gathering and brokering processus,
 - index documents in Arabic.