# Hahooa   Arabic   Web   Directory

-----------------------------------------------------------------

# Natural   Language   Processing for  Arabic  Information  Retrieval

## ACL  Workshop  on  Arabic  Language  Processing
## July 6 2001

Malek Boualem, France Telecom R&D

Régina Sneifer, Wanadoo

France Telecom R&D

(Arabic Search Engine)

# Contents

1. **Arabic Web**
2. **Arabic directories and portals**
3. **Arabic search engines**


4. **Hahooa Arabic Web directory**


5. **Some problems with Arabic information retrieval**
6. **Arabic linguistic features**


7. **Natural language processing**
8. **Arabic NLP activities at France Telecom R&D**

France Telecom R&D

# 1. Arabic Web

- **I**nternet users (2002) :

12 M / 240 M

- Arabic Web :

High speed development.

MS, Word, FrontPage, IE.

- Character encoding :
  - Microsoft CP1256
  - ISO-8859-6
  - UNICODE-UTF8

# 2. Arabic directories and portals

**Directories and Portals (Arabic & English) :**

- *http://www.albawaba.com*
- *http://www.arabia.com*
- *http://www.arabicseek.com*
- *http://www.ayna.com*
- *http://www.konouz.com*
- *http://www.maktoob.com*
- *http://www.naseej.com*

**Directories and Portals (Arabic, English & French) :**

- *http://www.hahooa.com*

France Telecom R&D

# 3. Arabic search engines

- **Search engines for directories and portals**

  - **Extense-Voila** (Echo-Wanadoo, France)
    **http://www.hahooa.com**

  - **Ayna** (Ayna, USA)
    **http://www.ayna.com**

  - **Al-Idrisi** (Sakhr, Egypt)
    **http://www.sakhr.com/**

  - **Konouz** (Alladin, Australia)
    **http://www.konouz.com/**

- **Search engines for the Web**

  - **Arabvista** (Emirates Internet & Multimedia and Compaq)
    **http://www.arabvista.com**

France Telecom R&D

# 4. Hahooa Arabic Web directory

## Hahooa : *Here it is*
### *http://www.hahooa.com*

➢ Echo-Wanadoo, France Telecom (Voila search engine *http://www.voila.fr*).

➢ Hahooa is trilingual : Arabic, English and French.

➢ More than 7600 URLs (Arabic around 3000).
➢ 418 themes (commerce and economy : more than 42%).

➢ Navigation languages : Arabic 49% ,  English,  French
➢ Users : France, Saudia, Lebanon, Morocco, Jordan, Canada, Belgium, etc.
➢ Frequent themes : Arts & Culture (music, …), Arab world, News & Media.

➢ Hahooa is regularly updated.

# Hahooa by countries

France Telecom R&D

# Hahooa by themes

**Répartition des catégories**



Education & Formation
6%

Loisirs & tourisme
8%

Organisations & Institutions
6%

Art et culture
4%

Regional
2%

Science & Technologie
1%

Annuaires & références
5%

informatique & Internet
9%

Société & Culture
8%

Sciences humaines
1%

Actualité & médias
7%

Economie& commerce
43%

# Hahooa by languages

**Répartition des sites pa langue**

| | |
|---|---|
| fr | 1019 |
| en | 4028 |
| en-fr | 192 |
| ar | 1054 |
| ar-fr | 32 |
| ar-en | 586 |
| ar-en-fr | 69 |
| autre | 154 |

Legend:
- fr
- en
- en-fr
- ar
- ar-fr
- ar-en
- ar-en
- autre

# 5. Some problems with Arabic information retrieval

| Keyword | Answers or documents |
|---|---|
| مجتمع (society) | 3012 |
| المجتمع (the society) | 20695 |
| مجتمعات (societies) | 599 |
| | |
| علم (science) | 12410 |
| علوم (sciences) | 3925 |
| | |
| إمرأة (woman) | 193 |
| نساء (women) | 11970 |
| | |
| كتب (book/write) | 44133 |
| كتابة (writing) | 2990 |
| | |
| تعليم (education) | 46728 |
| تدريس (education) | 770 |

Some tests using one of the search engines

France Telecom R&D

# 6. Arabic linguistic features

*Morphology*
*Syntax*
*Vowels*

*…*

سَتَكْتُبِينَهُمْ

| Enclitic | Suffix Morphem | Radical | Prefix Morphem | Proclitic |
|----------|----------------|---------|----------------|-----------|
| <hom'> | <yna> | < ktobi> | <ta> | <sa> |
| Object Masculine Plural | Subject Feminine singular | Root "ktb" | Subject-Pronoun "you" | Temporal affix |

## You (féminine) will write them (masculine)

# Examples of morphological features

**Inflected morphology by prefix (or by determination) :**

مجتمع (society) : المجتمع (the society)

**Inflected morphology by suffix :**

مجتمع (society) : مجتمعات (societies)

**Inflected morphology by infix :**

علم (science) : علوم (sciences)

**Irregular plural :**

إمرأة (femme) : نساء (femmes)

**Derivational morphology :**

كتب (book/write) : كتابة (writing)

# Arabic vowels

Le mot سلم

- signifie "paix" lorsqu'il est voyellé ainsi : سِلم

- signifie "échelle" lorsqu'il est voyellé ainsi : سُلّم

- signifie "a transmis" lorsqu'il est voyellé ainsi : سَلَّم

- signifie "est guéri" lorsqu'il est voyellé ainsi : سَلِم

- *Arabic texts that are published on the Web could be with/without/partially-with vowels.*
- *Queries on search engines should not contain vowels.*
- *Search engines should take this important morpho-syntactic feature into account.*

France Telecom R&D

# Cross-Language Information Retrieval

## *Most of the Arabic countries are bilingual*

**User**
Arabic language

Query in Arabic
Query in English
Query in French

Choice of target language
AR / EN / FR

**Query Translation**

**Search engine**
Query in Arabic
and / or
in English / French

**Monolingual Web sites**

English / French

**Monolingual Web sites**

Transliterated Arabic
Transcribed Arabic
Arabic characters

**Documents**
Arabic
and / or
English / French

**Choice of**
Document

**Choice of**
Document
+
Target language

**Document**

**Document Translation**

**Multilingual set of web sites**

# 7. Natural language processing

*The **eaats** mouse **tthe** cat*

Orthography :  **NO**  ➡ **Lexicons and morphology**


*The eats mouse the cat*

Orthography :  **YES**

Syntax :  **NO**  ➡ **Grammar**


*The **mouse** eats the cat*

Orthography :  **YES**

Syntax :  **YES**

Semantic :  **NO**  ➡ **Meaning**

    *(animal or computer ?)*


*The mouse eats the cat*

Orthography :  **YES**

Syntax :  **YES**

Semantic :  **YES**

Pragmatic :  **NO**  ➡ **Domain, real word**

    *(mousses do not eat cats !)*

France Telecom R&D

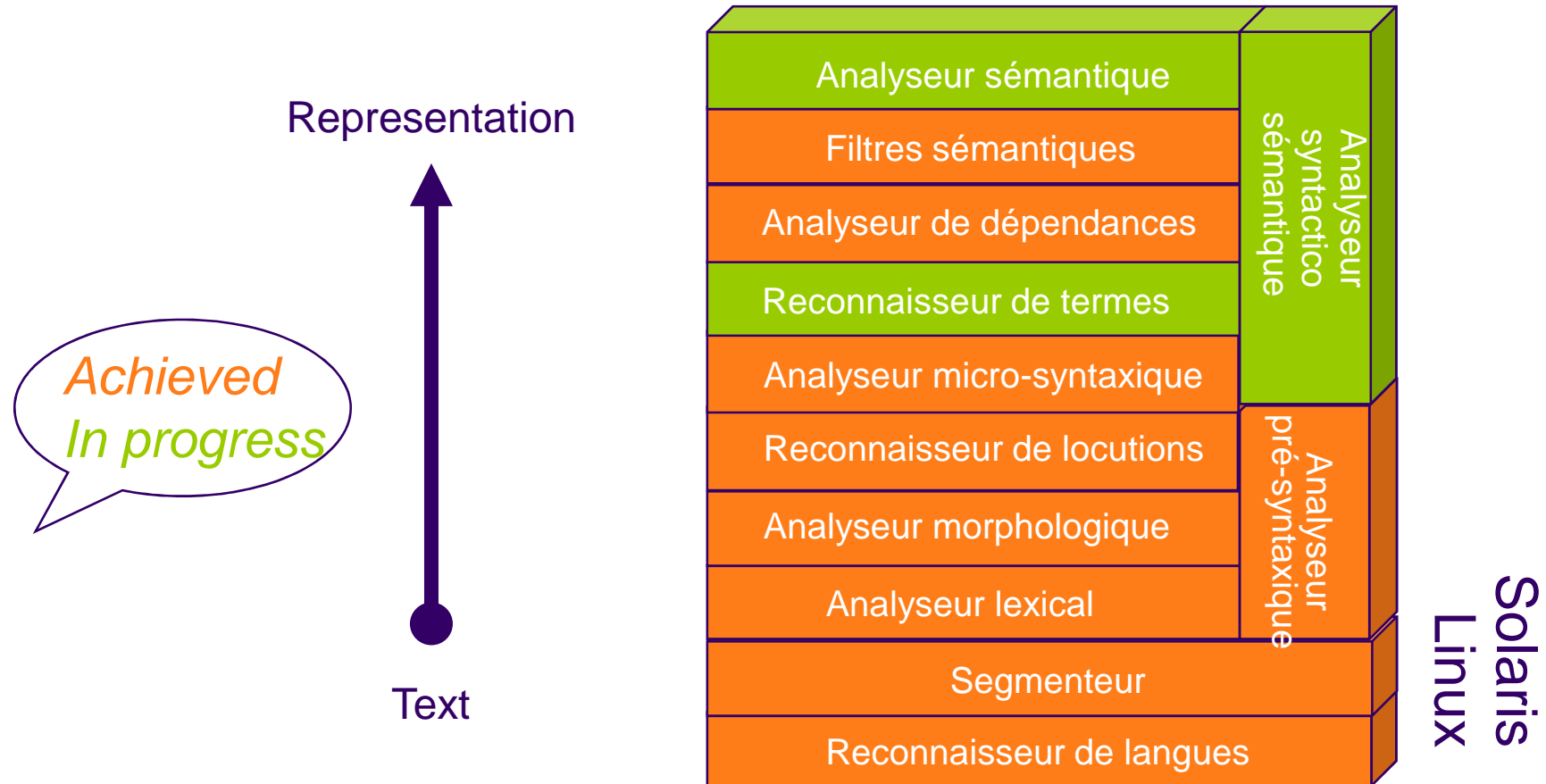# 8. Arabic NLP activities at France Telecom R&D

Natural Language Processing Group
14 members + 10 non members
Responsible : Jérôme Vinesse

France Telecom R&D
DMI Division
GRI laboratory (Data management and information retrieval)

**Address :**
France Telecom R&D
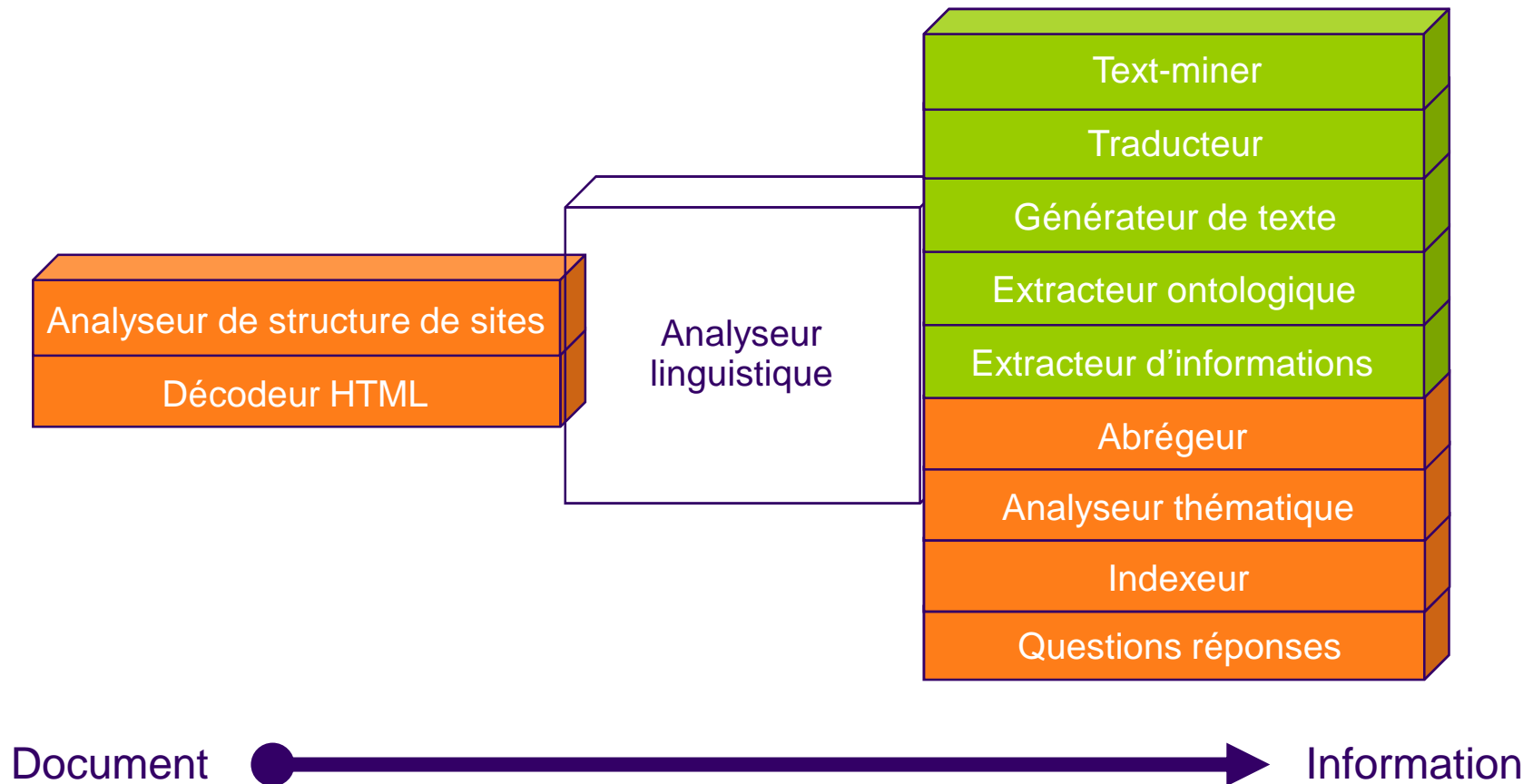2, avenue Pierre Marzin
22307 Lannion
France

France Telecom R&D

# Linguistic analysis tool : TiLT

Representation

↑

Text

Achieved
*In progress*

| Analyseur sémantique |
| Filtres sémantiques |
| Analyseur de dépendances |
| Reconnaisseur de termes |
| Analyseur micro-syntaxique |
| Reconnaisseur de locutions |
| Analyseur morphologique |
| Analyseur lexical |
| Segmenteur |
| Reconnaisseur de langues |

Analyseur syntactico sémantique

Analyseur pré-syntaxique

Solaris Linux

# Application modules

Text-miner

Traducteur

Générateur de texte

Extracteur ontologique

Extracteur d'informations

Abrégeur

Analyseur thématique

Indexeur

Questions réponses

Analyseur de structure de sites
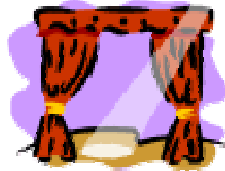
Décodeur HTML

Analyseur linguistique

Document ⟶ Information

# Catalogue

**Language and charset Identification**

**Thematic classification**

**Automatic generation of linguistic exercises**

**Linguistic filtering**

**Question / Answering**

**Linguistic localization**

**Summarization**

**HTML parsing**

# Languages

**French**

**English**

**Spanish**

**German**

**Polish**

**Arabic**

# Contact

Malek Boualem

Natural Language Processing Group
France Télécom R&D
DMI/GRI Laboratory
2, avenue Pierre Marzin
22307 Lannion
France

Phone:          (33)(0)2.96.05.29.83
Fax:            (33)(0)2.96.05.32.86
Email:          malek.boualem@rd.francetelecom.com

France Telecom R&D