



Avec le soutien de l'Union européenne et de la Région wallonne



## Web Data Extraction with



**Fabrice Estiévenart ([fabrice.estievenart@cetic.be](mailto:fabrice.estievenart@cetic.be))**

*Rencontres Mondiales du Logiciel Libre - Bordeaux - 8<sup>th</sup> July 2010*



[www.cetic.be](http://www.cetic.be)

Your connection to ICT research

# CETIC, overview

- ICT research centre
- Created in 2001
- Initiated by 3 universities
- Connection between Industry & Research
- 3 departments, 40 researchers
- Contribution to Regional Economic Development
- International focus – European Research Area



# CETIC, the ICS team

- Knowledge Extraction from Unstructured Content

- Web wrapping
- Document clustering
- Text mining

- Search Engines

- Crawling
- Text extraction
- Analysis / Indexing
- Search

- Semantic Web

- Ontology and terminology engineering
- Micro-formats

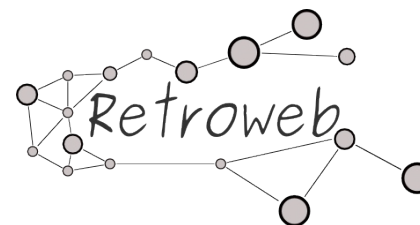






# Biggest challenges in web wrapping

- Collecting/organizing relevant documents
  - Intelligent crawling
  - Web services
- Locating data of interest within documents
  - RegExp
  - Element structure
  - External resources



## Retroweb, overview

Tool for web data extraction (web wrapping)

**Retroweb-GUI:** semi-automated definition of extraction rules, visual process

**Retroweb-Wrapper:** web data extraction from extraction rules

# Retroweb, terminology

## Web Page

<http://www.imdb.com/title/tt0821642>

<http://www.imdb.com/title/tt0858486>

<http://www.imdb.com/title/tt0458525>

## Page Type

imdb-movie

## Page Component

title

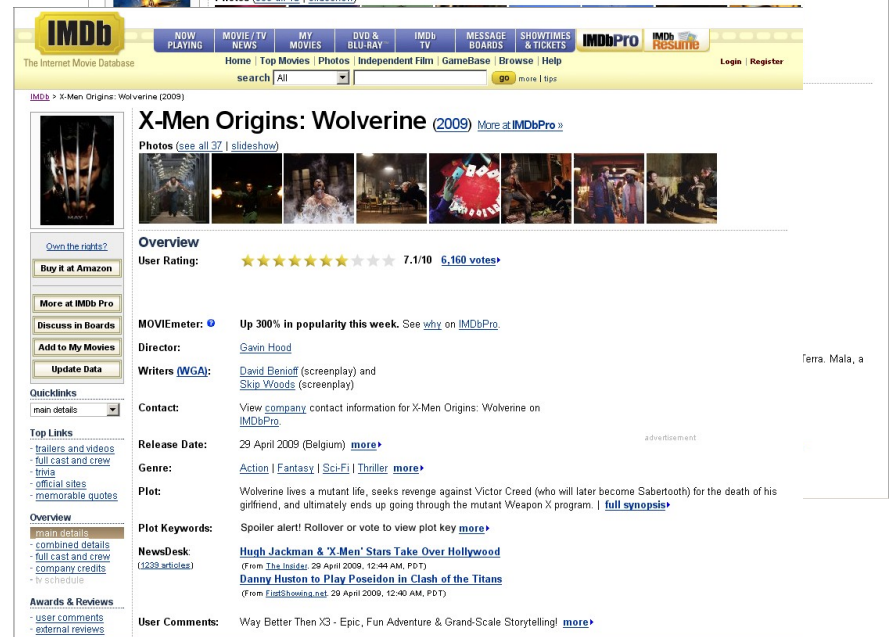
tagline

actors

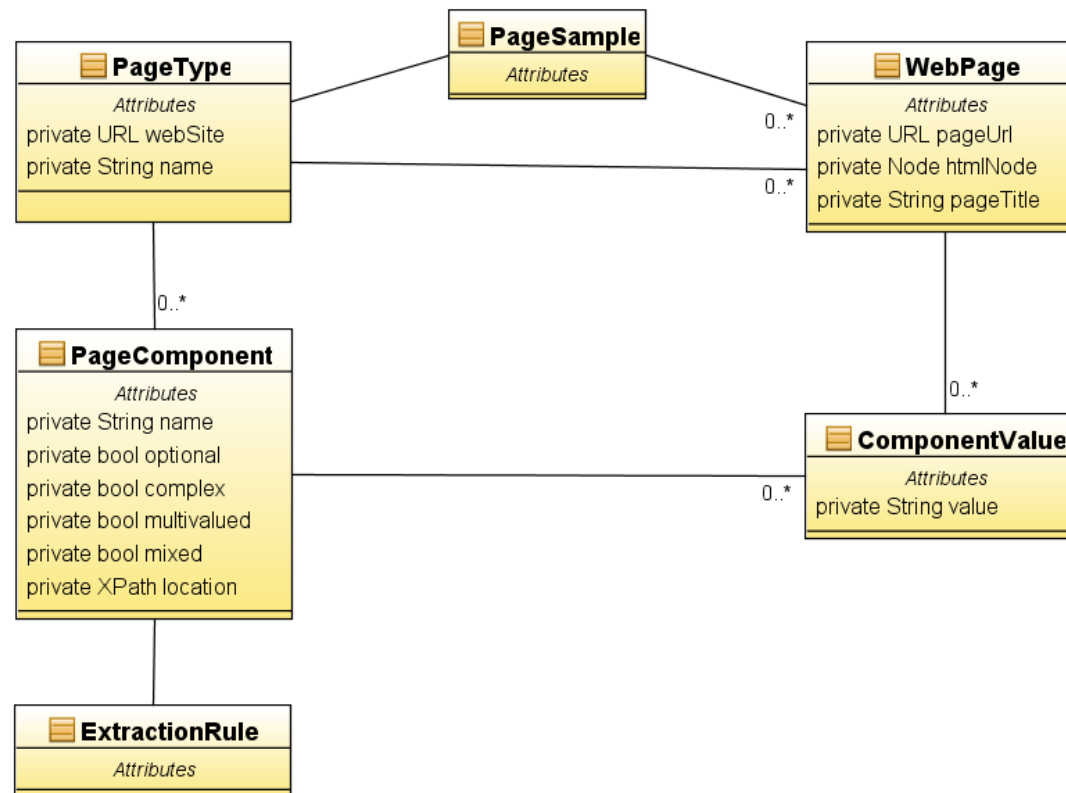
## Component Value

The soloist

[www.cetic.be](http://www.cetic.be)



# Retroweb, model





# The wrapping process

## Retroweb GUI

Select a page sample

Select values and name page components

Refine extraction rules

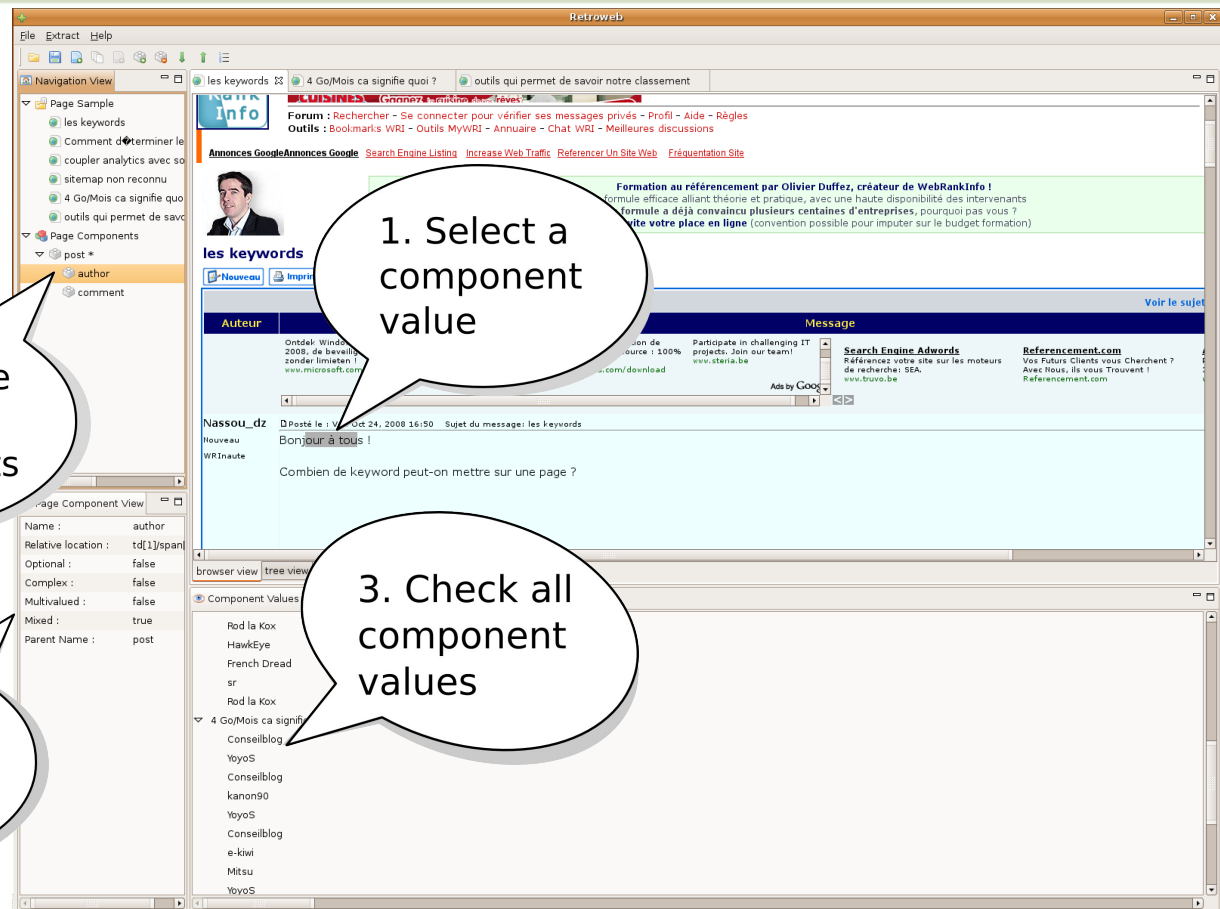
(Re)structure page components

Check component values

## Retroweb Wrapper

Extract to XML

# Screenshot



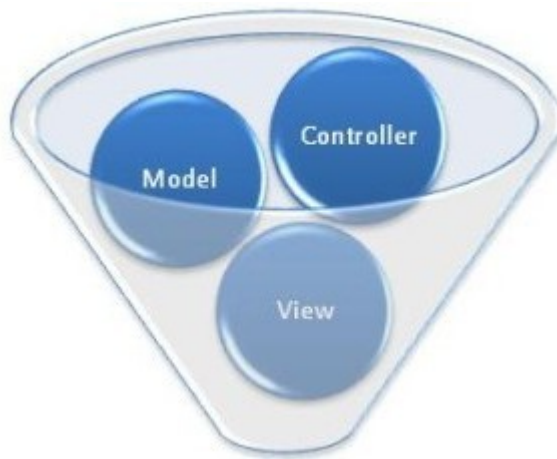
4. Structure page components

2. Visualize rules and properties

1. Select a component value

3. Check all component values

# Technologies



## Retroweb, benefits

Easy: no need to learn a specific language

Flexible: only data of interest are extracted

Robust: extraction rules are defined from a set of pages

Extensible: based on standards (XML, XPath) and open-source (Affero GPL v3)

Portable: GNU/Linux, MS-Windows



## Some applications

Web sites reverse engineering  
Search engines  
Competitive intelligence  
Semantic annotation of corpus

## Some similar applications

- Lixto Visual Developer
  - Spin-off from the Wien University
  - Services in web intelligence
  - Visual tool based on Eclipse-RCP
- Dapper
  - Free web application for data extraction
  - Not extensible

## Wanna be involved ?

- Centre d'Expertise en Logiciel Libre à Vocation Industrielle (CELLAVI)
- <https://forge.pallavi.be>



C E L L A V I

## Next steps

### Search engine integration

- Web pages collection and clustering

- Semantic indexation

- Advanced search

### Error detection

- Rules maintenance

### Semantic import/export

- Micro-formats (RSS, FOAF, ...)

- Ontology population





Avec le soutien de l'Union européenne et de la Région wallonne



**Thanks for your attention**  
**Questions?**



[www.cetic.be](http://www.cetic.be)

Your connection to ICT research