

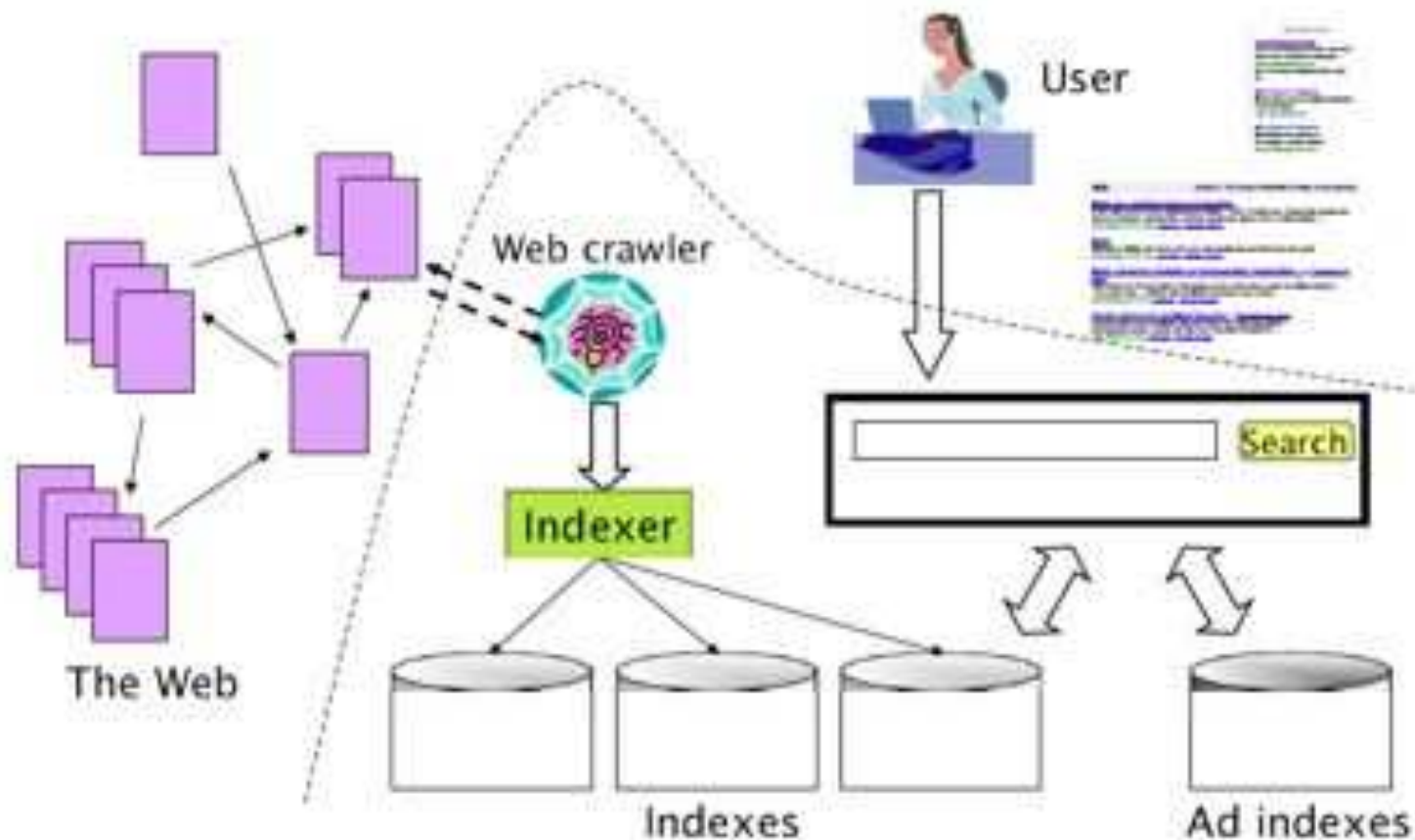
Sistemi informativi per la gestione d'azienda 2013

LABORATORIO
ESERCITAZIONE
WEB CRAWLER 2

Info

- Contacts
 - sandro.labruzzo@isti.cnr.it
 - claudio.atzori@isti.cnr.it
- WEB sites
 - <http://www.nmis.isti.cnr.it/casarosa/SIA/>

Architecture of a Search Engine



Open Source Crawlers

<http://java-source.net/open-source/crawlers>

- [WebSPHINX - http://www-2.cs.cmu.edu/~rcm/websphinx/](http://www-2.cs.cmu.edu/~rcm/websphinx/)
- [Jspider - http://j-spider.sourceforge.net/](http://j-spider.sourceforge.net/)
- [WebEater - http://sourceforge.net/projects/webeater](http://sourceforge.net/projects/webeater)
- [WebLech - http://weblech.sourceforge.net/](http://weblech.sourceforge.net/)
- [Arachnid - http://arachnid.sourceforge.net/](http://arachnid.sourceforge.net/)
- [JoBo - http://www.matuschek.net/software/jobbo/index.html](http://www.matuschek.net/software/jobbo/index.html)
- [Web-Harvest - http://web-harvest.sourceforge.net/](http://web-harvest.sourceforge.net/)
- [Crawler4j - http://code.google.com/p/crawler4j/](http://code.google.com/p/crawler4j/)
- [Ex-Crawler - http://ex-crawler.sourceforge.net](http://ex-crawler.sourceforge.net)
- [Bixo - http://openbixo.org](http://openbixo.org)
- [Nutch - http://nutch.apache.org/](http://nutch.apache.org/)

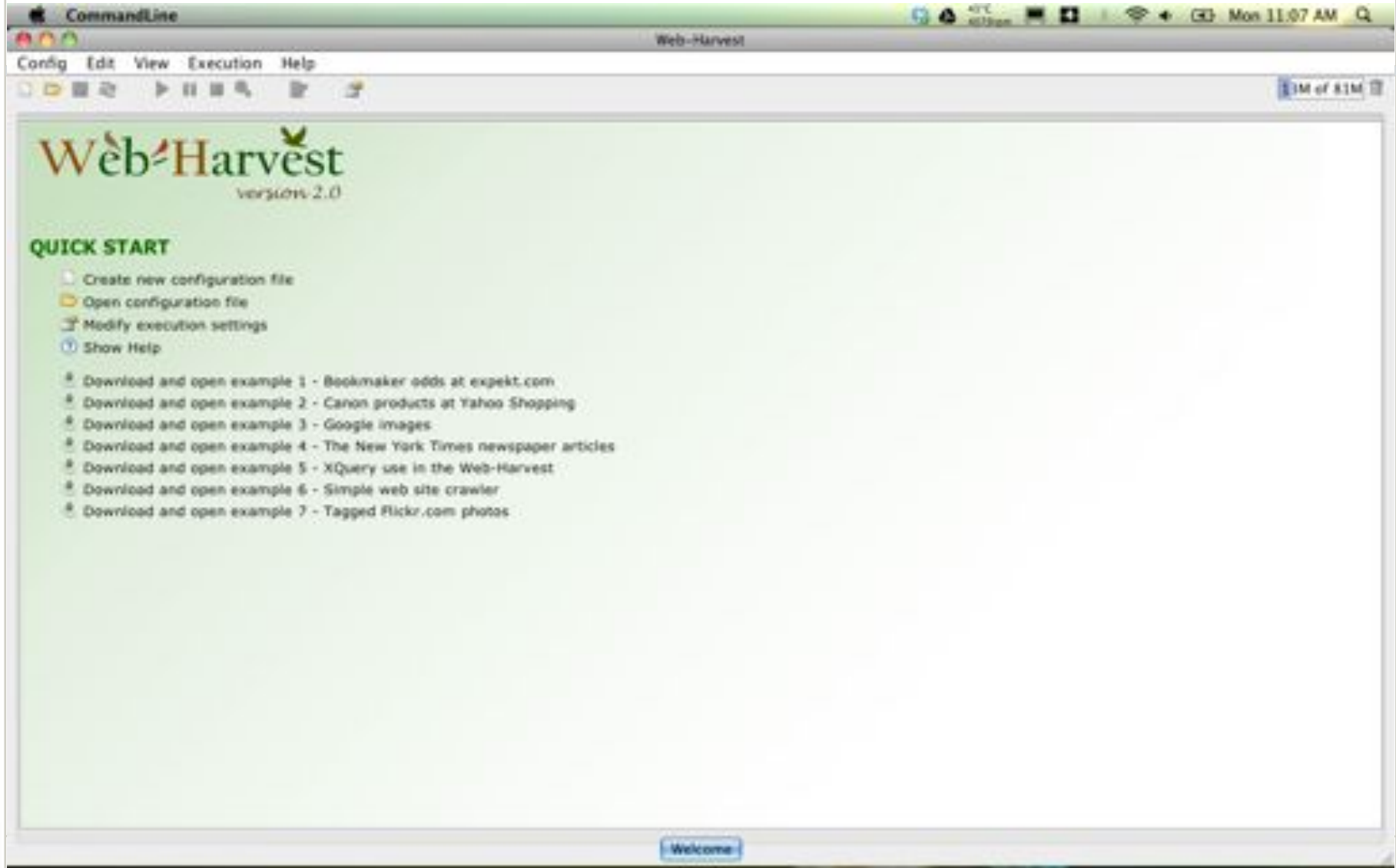
Web-Harvest

- <http://web-harvest.sourceforge.net/>
- Web-Harvest is an *Open Source Web Data Extraction tool* written in Java.
- Allows to collect desired Web pages and extract useful data from them.
- Leverages on well established techniques and technologies for text/xml manipulation such as
 - *XSLT*
 - *Xquery*
 - *Regular Expressions*
- Process of extracting data from Web pages is also referred as *Web Scraping* or *Web Data Mining*

Web-Harvest

- Work defined by “recipes”, no more need to code Java (yay!)
- Recipes are described in a specific configuration language (XML)
- Recipes reference manual
 - <http://web-harvest.sourceforge.net/manual.php>
- Download URL:
 - <http://web-harvest.sourceforge.net/download/webharvest2b1-exe.zip>
 - extract
 - Run (java -jar <webHarvest.jar> | double click... | whatever...)

Web-Harvest



Web-Harvest

The screenshot shows a web browser window titled "Web-Harvest" displaying the configuration for a web-harvesting tool. The browser's address bar shows "Welcome google_images.xml". The main content area displays XML code for configuring the tool. The code includes comments and XML tags for defining variables, including a search expression, a URL template, and a loop for downloading images. A sidebar on the left shows a tree view of the configuration elements, including "include", "var-def", "template", "call", "loop", "list", "body", and "file".

```
<?xml version='1.0' encoding='UTF-8'?>
<!--
  Supports following initial variables:
  search - search expression
  -->
<!-- Updated on February, 8th, 2011 -->
<config charset="UTF-8">
  <include path="functions.xml"/>
  <!-- defines search keyword and start URL -->
  <var-def name="search" overwrite="false">banana</var-def>
  <var-def name="url"><template>http://images.google.com/images?q=${search}&amp;hl=en</template></var-def>
  <!-- retrieves all image URLs -->
  <var-def name="imglinks">
    <call name="download-multipage-list">
      <call-param name="pageUrl"><var name="url"/></call-param>
      <call-param name="nextXPath">//a[@id="pbnext"]/@href</call-param>
      <call-param name="itemXPath">//img[contains(@src, 'images?img')]</call-param>
      <call-param name="maxloops">5</call-param>
    </call>
  </var-def>
  <!-- download images and save them to the files -->
  <loop item="link" index="1" filter="unique">
    <list>
      <var name="imglinks"/>
    </list>
    <body>
      <file action="write" type="binary" path="google_images/${search}_${i}.gif">
        <http uri="${sys.fullUri(url, link)"/>
      </file>
    </body>
  </loop>
</config>
```

Web-Harvest

```
<config charset="UTF-8">
  <var-def name="urlList">
    <xpath expression="//img/@src">
      <html-to-xml>
        <http url="http://news.bbc.co.uk"/>
      </html-to-xml>
    </xpath>
  </var-def>

  <loop item="link" index="i" filter="unique">
    <list>
      <var name="urlList"/>
    </list>
    <body>
      <file action="write" type="binary" path="images/${i}.gif">
        <http url="${sys.fullUrl('http://news.bbc.co.uk', link)}"/>
      </file>
    </body>
  </loop>
</config>
```

Web-Harvest

This configuration contains two pipelines.

```
<config charset="UTF-8">
  1 [ <var-def name="urlList">
      <xpath expression="//img/@src">
        <html-to-xml>
          <http url="http://news.bbc.co.uk"/>
        </html-to-xml>
      </xpath>
    </var-def>

  2 [ <loop item="link" index="i" filter="unique">
      <list>
        <var name="urlList"/>
      </list>
      <body>
        <file action="write" type="binary" path="images/${i}.gif">
          <http url="${sys.fullUrl('http://news.bbc.co.uk', link)}"/>
        </file>
      </body>
    </loop>
  </config>
```

Web-Harvest

The first pipeline performs the following steps:

- HTML content at *http://news.bbc.co.uk* is downloaded,
- HTML cleaning is performed on downloaded content producing XHTML,
- XPath expression is searched for, giving URL sequence of page images,
- New variable named "urlList" is defined containing sequence of image URLs.

```
1 [ <var-def name="urlList">
    <xpath expression="//img/@src">
        <html-to-xml>
            <http url="http://news.bbc.co.uk"/>
        </html-to-xml>
    </xpath>
</var-def>
```

Web-Harvest

The second pipeline uses result of the previous execution in order to collect all page images:

- *Loop* processor iterates over URL sequence and for every item:
- Downloads image at current URL,
- Stores the image on the file system.

```
2 [ <loop item="link" index="i" filter="unique">
    <list>
      <var name="urlList"/>
    </list>
    <body>
      <file action="write" type="binary" path="images/${i}.gif">
        <http url="${sys.fullUrl('http://news.bbc.co.uk', link)}"/>
      </file>
    </body>
  </loop>
```

Web-Harvest: task

- foreach example in examples
 - UNDERSTAND** example
 - RUN** example
 - WRITE** description(example)
- function description
 - foreach tag in example
 - UNDERSTAND** tag (<http://web-harvest.sourceforge.net/manual.php>)
 - return **detailed pseudocode** of the example

Email examples.zip to claudio.atzori@isti.cnr.it
[subject: SIA webharvest <name.surname>]

Hint: examples provide some comments ... consider integrating yours