



# LumberJacz

Web data acquisition and  
on-top-of-Web applications

Viktor Zigo

Viktor Zigo, May 17th 2006

# About



Viktor Zigo

**lixto** Lixto Software GmbH.

Senior software architect



DBAI, TU Wien

PHD student

# Semantic Web



- **Web of data and relations** vs. web of proprietary documents
- Annotated content in machine readable representation and common ontology
- Data interchange, integration and inference

Uhm, coool but...

# Semantic Web ???



- Inherently agnostic Web segments  
prices, services, news, competition
- Google ignores it  
not visible = 80% fake
- No “semantic browser”
- Behavioral aspects  
navigation, interaction,
- Complexity  
WS are cool, REST and XmlRpc are way cooler
- Legacy Web, Web 2.0

# To Do



- Web-2-Sweb conversion
- Data integration and aggregation
- Adding new functionality upon Web  
largest information data-source

Web monitoring, data searching or meta-search, making Web more structured and semantic, integration of Web enabled apps, content aggregation, Web tests automation, screen scraping, web extraction, Web remixing or mash-up

# An Example – Simple Meta Search



## Flight Search Application:

“Search for the best flight offers according to user preference in real-time out of several airlines”

User : *From, To, Date*

**easyJet.com**

**SKY**  
EUROPE

**RYANAIR.COM**

197.00	EUR	Wed. 17.05.	07:15	Airberlin AB 8893
197.00	EUR	Wed. 17.05.	15:55	Airberlin AB 8855
197.00	EUR	Thu. 18.05.	07:15	Airberlin AB 8893
221.75	EUR (GBP)	Tue. 16 May 06	11:10.	Air Berlin
231.09	EUR (GBP)	Tue. 16 May 06	11:30.	Air Berlin
280.73	EUR (GBP)	Tue. 16 May 06	15:55.	Air Berlin
401.46	EUR (GBP)	Tue. 16 May 06	14:40.	Czech Airlines CSA
411.27	EUR (GBP)	Tue. 16 May 06	19:30.	Austrian

# Process Stages



1. User query mapping to portal search forms
2. Web navigation and query setting
3. Parallel execution
4. Results synchronization
5. Data extraction (navigation)
6. Data unification, merging, cleaning
7. Presentation and Syndication

# Query Mapping



## Map user query to Portal Specific Queries

- Hard-coded
- Automatic on-fly mapping

The screenshot shows a flight search interface with two main sections. The left section, titled 'Flugsuche', contains fields for 'von / nach' (departure/arrival), 'Hinflug' (outbound flight), 'Rückflug' (return flight), and 'Fluggäste' (passengers). The right section, titled 'bestprice.click', contains fields for 'Von' (origin), 'Nach' (destination), 'Erwachsene' (adults), and 'Outsicht' (view). Below these are fields for 'Departure' and 'Arrival' dates and times, and a 'Search' button. Two blue arrows point from the 'Flugsuche' section to the 'bestprice.click' section, indicating a mapping between the two.

## Problems :

- M :N mapping  
2006-01-03  
3<sup>rd</sup>, January, 2006
- syntactic, semantic
- Mapping of intervals (prices, age)
- Units
- Multilingual mapping, synonyms, free text



# Deep Web Navigation



## Types :

- Crawling
- Given navigation
- Auto navigation

## Obstacles :

- Dynamic content
- Sessions, state-full
- Authorizations,HTTPS
- Proxies
- AJAX
- Changing structure
- Page interaction
  - Forms
  - In-page navigation
  - Flash

# Common Wrapping Techniques



## ● Techniques

- Structural (tree)
- Syntactic (regexp)
- Tokenization
- Text analysis
- Visually based

## Wrapper Generation

- Hardcode
- Visually Design
- Supervised learning
- Automatic  
(pattern recognition,  
templates)

## Normalization, transformation and data-cleaning



- Inverse process of query mapping
- Common output schema (ontology)
- Value mapping

L.S.O, Symphonic Orch London, Orchester, London

- Duplicate entries, data clustering

# Syndication



- Format
  - XML, RDB mapping, RDF, PDF, text....
- Publishing
  - Portals, legacy app, CMS, RSS, DB ...



# LumberJacz

Axe the Web!

<http://lumberjacz.org>

# LumberJacz Project



***“An open-source technology that enables better ways to get and use data and information on the Web.”***

Pragmatic approach :

light-weight, flexible, end-to-end, get things done first

- Aug 2005 : Started
- May 2006 : usable, rolling out

# Tech background



- Mature and rich Mozilla code-base, open source
- Many similar applications :
  - MIT PiggyBank, Chickenfoot, GreaseMonkey
  - Test4Web, Selenium, Solvent
  - AJAX Toolkit Framework (IBM)
  - FF extensions
- Real-world experience with data extraction

# Appetizer – Personal Flight Search



- Running standalone **Flight Search** application
- Download & run
- 2 days (incl. GUI)

The screenshot displays two windows of the 'Personal Flight Search' application. The top window shows search results for a query from London to Vienna. The bottom window shows the search form with the same query parameters.

Price*	Currency	Departure	Carrier	Merchant
190.00	EUR	Mon. 15.05. 11:10	Airberlin HG 8553	http://www.airberlin.com
197.00	EUR	Mon. 15.05. 07:15	Airberlin AB 8893	http://www.airberlin.com
197.00	EUR	Mon. 15.05. 15:55	Airberlin AB 8855	http://www.airberlin.com
197.00	EUR	Tue. 16.05. 07:15	Airberlin AB 8893	http://www.airberlin.com
197.00	EUR	Tue. 16.05. 15:55	Airberlin AB 8855	http://www.airberlin.com
162.00	EUR	Tue. 16.05. 11:10	Airberlin HG 8553	http://www.airberlin.com
141.00	EUR	Wed. 17.05. 11:10	Airberlin HG 8553	http://www.airberlin.com
197.00	EUR	Wed. 17.05. 07:15	Airberlin AB 8893	http://www.airberlin.com

Search time 2006-4-15 17:58

From: london  
To: vienna  
Dep. Date: 2006-4-16  
Dep. Time: 00:24  
Adults: 1  
Children: 0

From: london  
To: vienna  
Departure: 15 April 2006 Anytime  
Adults: 1 Children: 0 Infants: 0



# What a heck is this Lumber ...?



A **generic workflow engine** and framework...

- ...primarily **client-side**

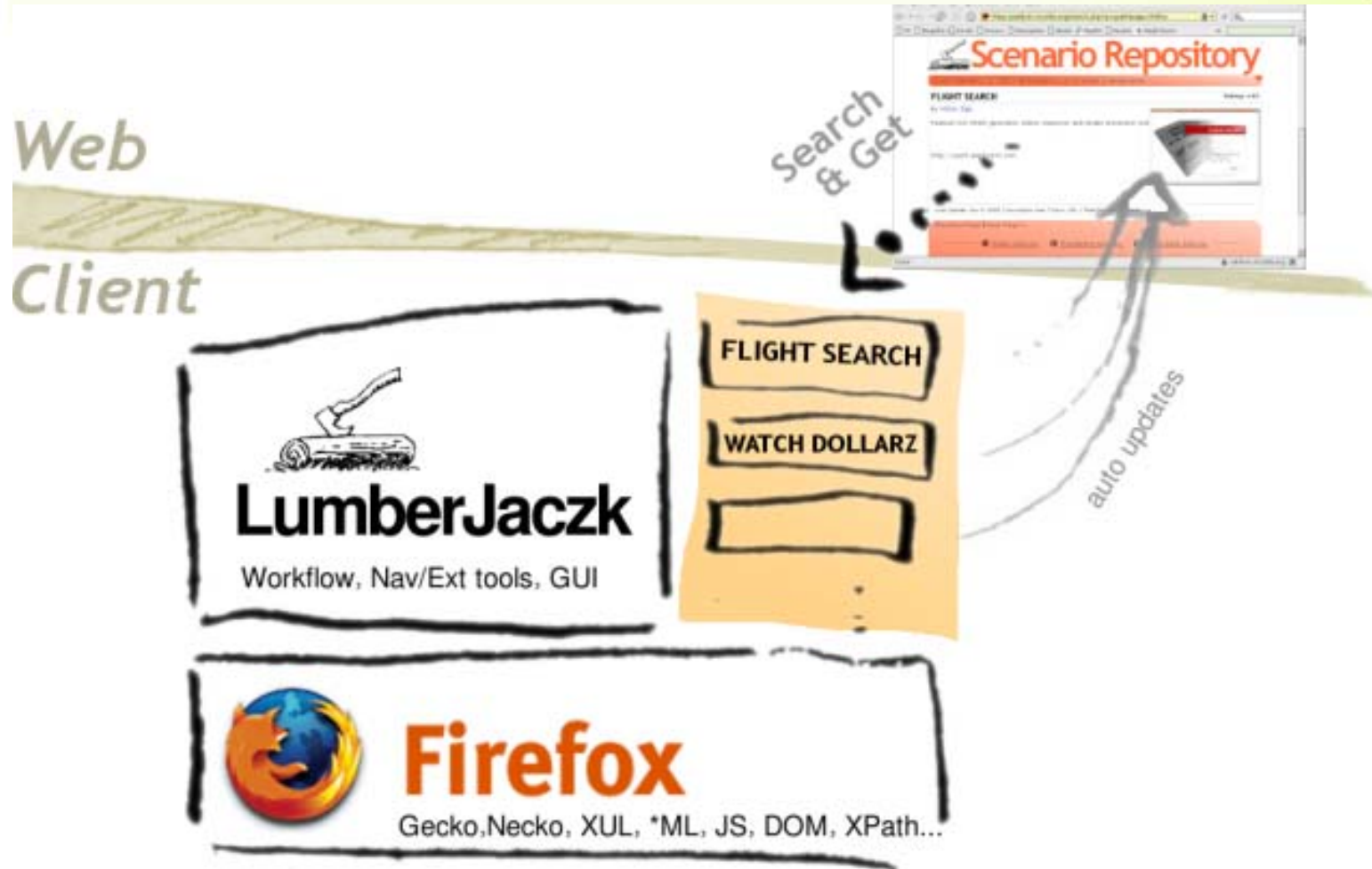
- ...for **execution of apps** (*scenarios*)

- ...operating **on top of the Web**

- ...provided in very **portable packages**

- ...with **rich GUI** and **user Interactions**

# Architecture and Principles



# Generic Player



Get the generic player once, either

- as a **standalone app**
- or in your Firefox **browser**

You can do this then :

1. Browse web repositories
2. Install them from web by single click
3. Customize, manage and run them from within your client (or browser)
4. Get automatically the updates

# Standalone Generic Player



- Launch and manage all the installed scenarios from your desktop
- Manage your recent runs and results



# Execution Modes



- standalone scenario (desktop)
- standalone generic player (desktop)
  - runs any installed scenarios
- single scenario browser extension
- generic player as a browser extension

# Scenario properties



It is a single application LumberJaczk runs.

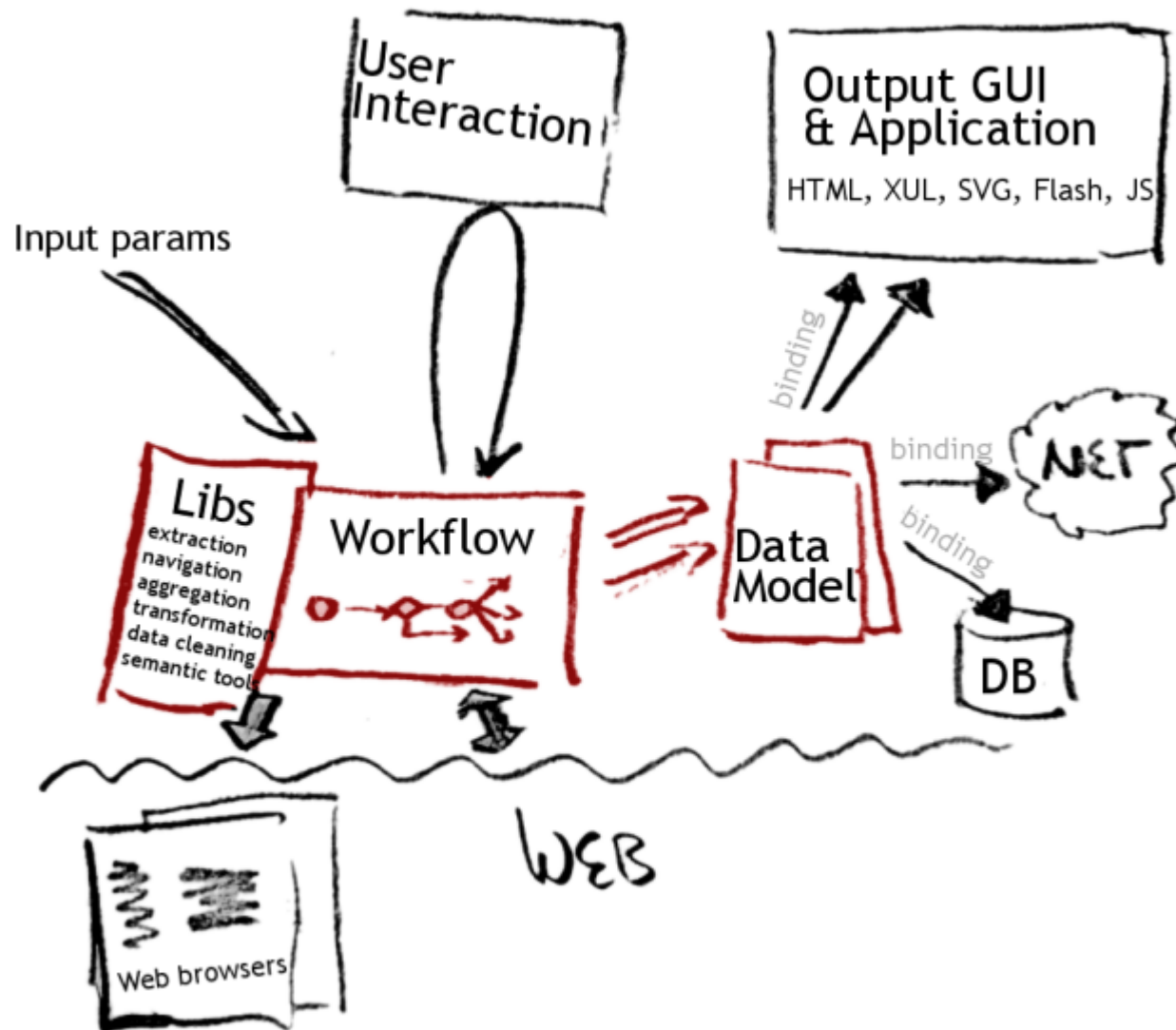
- small (kB), compact, and self-contained
- installs by finger-click
- installs directly from internet web page
- automatically updated
- dynamic scenario federations
  - e.g. simply add new sources to your app
- might be signed
- platform independent

# Scenario features



- application business logic
- configuration GUI
- output data schemas
- output application (GUI + logic)
  - presentation in GUI
  - filtering, graph generation, system integration, etc.
- löalizätion

# How does it work ?





# Scenario GUI ...



....is rich... more than rich !!!

It can be a blend of

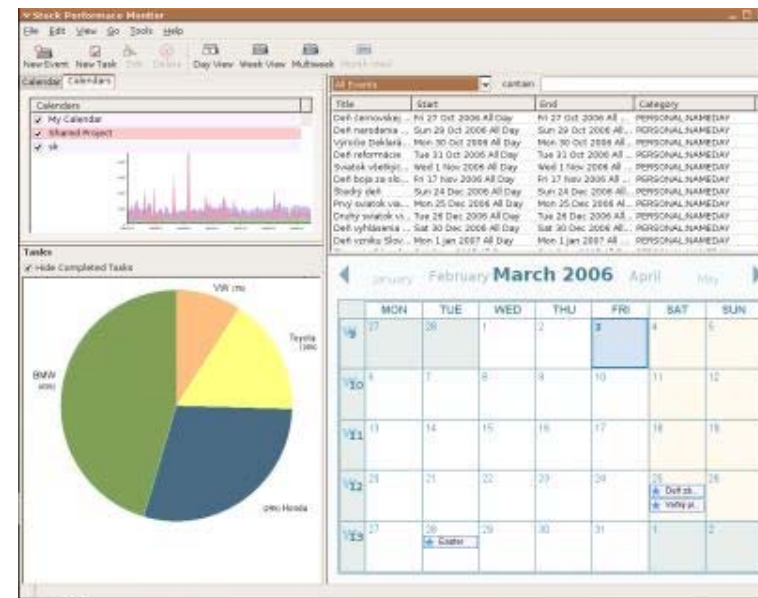
- **HTML, SVG**
- **Flash**
- **native GUI**
- **integration to existing Web page**

# Native GUI Scenario



- Native widgets
- Application logic
- Graphs and charts

## Business events monitoring example



# In-Browser Scenario



- up and running in a browser by 2 clicks
- launched from generic player
- HTML, native, Graphs possible
- Integration to 3<sup>rd</sup> party Web

## Car Search scenario example

The screenshot shows a Mozilla Firefox browser window with the following content:

- Browser Tabs:** chrome://cars...nt/output.xml
- Address Bar:** chrome://cars...nt/output.xml
- Page Title:** LumberJaczK
- Scenario Configuration:**
  - Scenario: Car Search
  - Run button
  - Car Search section with Make (HONDA) and Model (CIVIC) dropdowns.
  - Fulltext search:
  - Go button
- Logs:**
  - Choose a Scenario and run Engine...
  - CarSearch Scenario state START
  - CarSearch Scenario state FINISH
- Taskbar:** Task:autoviaExtractrems DONE
- Main Content:**
  - Pie chart: 33% AAAA, 2% B, 64% Autovia
  - Section: CAR SEARCH RESULTS
  - Table with columns: # Title, # Year, # Price, # Date

# Title	# Year	# Price	# Date
HONDA Civic 1.3 16V, Coupé.	1992	79.000	15. 5. 2008
HONDA Civic Sedan,	2005	484.000	30. 4. 2008
HONDA Civic 1.4, hatchback 5dv.,	1995	109.000	15. 5. 2008
HONDA Civic 1.5, hatchback 5dv.,	1996	169.000	15. 5. 2008
HONDA Civic 1.5, sedan,	1996	129.000	15. 5. 2008
HONDA Civic 1.4 16V 66kW 3dv.,	1999	169.000	15. 5. 2008
HONDA Civic 1.4 Limited Edition hatchback, P, M5,	1996	155.000	13. 5. 2008
HONDA Civic 1.4 LS hatchback, P, M5,	2003	339.000	15. 5. 2008
HONDA Civic 1.4 S hatchback, 5d.,	1995	122.000	11. 5. 2008
HONDA Civic 1.5i, Benzin, 4 dvore	1993	59.817	-
HONDA Civic 1.4 S hatchback, P, A4,	1995	139.000	12. 5. 2008
HONDA Civic 1.5i, Benzin, 2 dvore	1995	118.304	-
HONDA Civic 1.4 S hatchback, P, M5,	1998	195.000	15. 5. 2008
HONDA Civic 1.4i, Benzin, 5 dvore	1997	171.474	-
HONDA Civic 1.4 S hatchback, P, M5,	2001	289.000	11. 5. 2008
HONDA Civic 1.4i, Benzin, 5 dvore	1997	166.157	-
HONDA Civic 1.4 S hatchback, P, M5,	2002	325.000	15. 5. 2008
HONDA Civic 1.4i, Benzin, 3 dvore	1997	159.000	-

# Live and Dynamic



Execution experience can be very live and dynamic

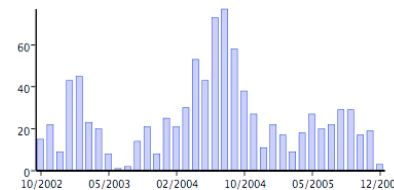
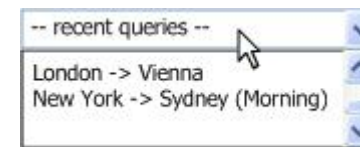
- watch the progress
  - search, workflow, communication, extraction
- live results modification
  - adding new, making more precise
- on fly generation of graphs
  - cool eh ?
- stop, pause, resume, inspect, and interact !

# It lives with you



The scenario player can :

- remember your inputs and configurations
  - you have run recently
  - you have predefined
- keep the state and results from the previous executions
  - and combine them for example (graphs)





# Getting Technical

Viktor Zigo, May 17th 2006

# Engine



- Generic workflow engine
- Scriptable actions
- Execution Context
- Output Data Models – XOutput
- Navigation, Extraction, transformation libraries

# Workflow engine



- dynamic flow decision
- synchronization on any custom event
- parallel execution (forking)
- native iterations (for-each)
- sub-processes (nesting)
- arbitrary scriptable action in any state (JS)
- interpreted
- arguments, shared variables
- static function, custom action libraries



# Actions



- JavaScript (Java)
- complete Mozilla control
- networking, windows, DOM, XPath, RegExp, XML, XSL, SOAP, WS, XmlRpc
- Libraries for navigation, extraction, and data integration
- XOutput
- interaction with Scenario GUI
- error handling



# Workflow Script

```
// output.set('in_carrier', in_carrier);
output.set('link', transformedLink);
output.set('merchant', 'http://opodo.co.uk');
output.link('flight', transformedLink, true);
}

//*****
//***** AIRBERLIN *****
function airberlin(){
  lib.addTab('http://www.airberlin.com/site/index.php?LANG=eng'); //ONLINE
  return syncOn( window , 'airberlinFillForm', 500)
}

function airberlinFillForm(){
  var root = lib.assertFirst( lib.xpathFrames( "/html/body/table[1]/tbody/tr[2]/td[2]/t

  var input = lib.assertFirst( lib.xpath(root, ".tbody/tr[1]/td/table/tbody/tr[2]/td[1]/s
  var option=lib.select(input, parameters.from, Lib.MATCH_TOP | Lib.MATCH_VISUAL

  input = lib.assertFirst( lib.xpath(root, ".tbody/tr[2]/td/table/tbody/tr[1]/td/table/tb
  lib.click(input);

  input = lib.assertFirst( lib.xpath(root, ".tbody/tr[5]/td/table/tbody/tr[2]/td[1]/selec
  option=lib.select(input, parameters.to, Lib.MATCH_TOP | Lib.MATCH_VISUAL);
```

# Data Output - XOutput



- native XML generation
- simplified model schemas
- linking with XLink
- unique XOutput technique
- parallel data generation
- per-partes flushing
- direct binding to output device (GUI)
  - HTML, XUL applications, HTTP, SOAP, XMLRPC, WS

# Further usage



## Bottom-up approach

- Generic client application workflow engine
- Rule-based extraction engine
- Hierarchical extraction engine (XML)

# Status

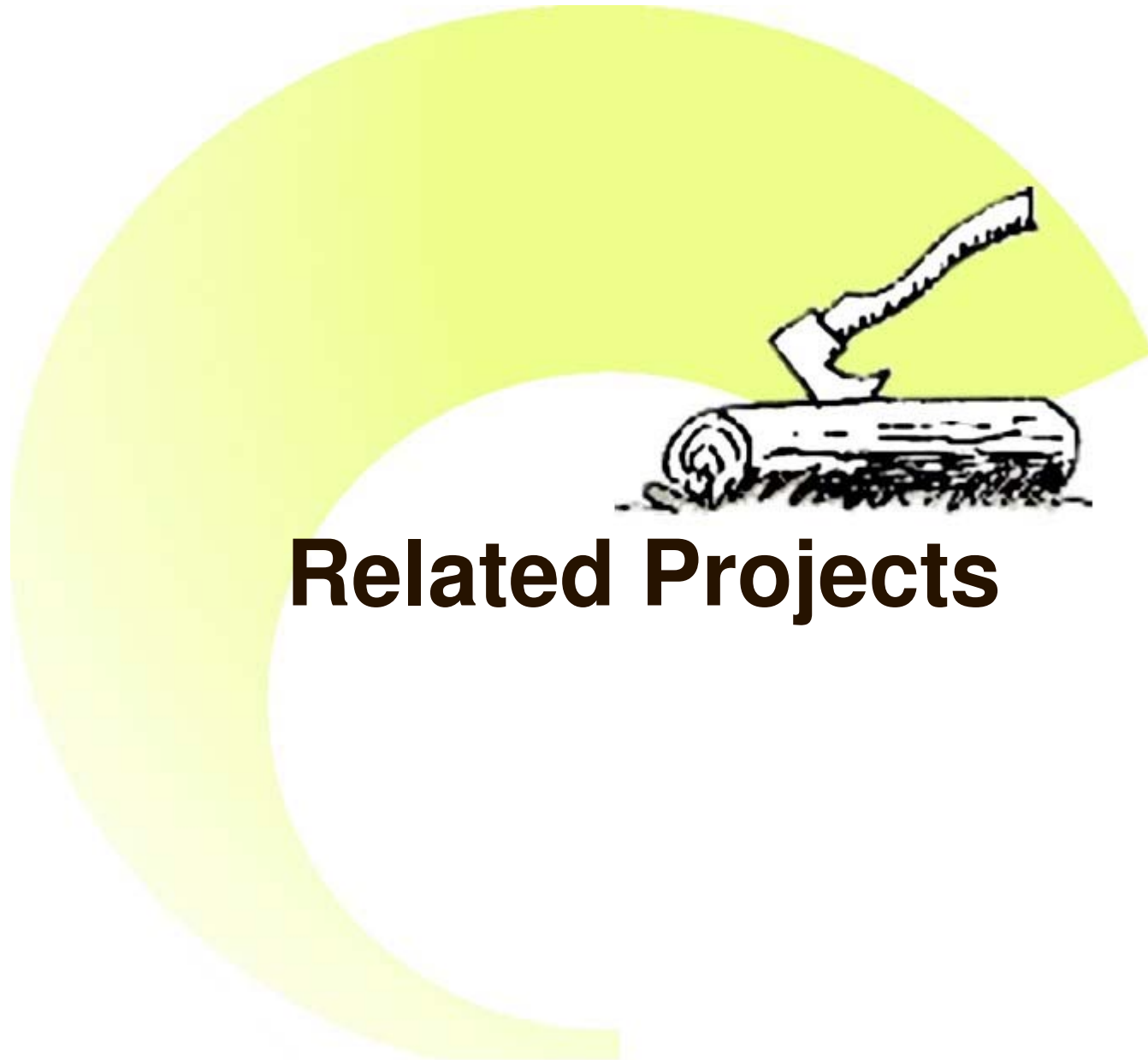


- Usable system
- Preparation for OS release
- Building web site and infrastructure

Web : <http://lumberjacz.org>

Interested ?

[interested@lumberjacz.org](mailto:interested@lumberjacz.org)



# Related Projects

Viktor Zigo, May 17th 2006

## Related and New Projects



- Public web repository of Scenarios
- Generic Wrapping
- Collaborative semantic database
- Server side solutions
- Data merging and cleaning toolset
- Visual Scenario Development

# Scenario Repository Portal



- Public repository for LJ scenarios
  - Publishing /review / share process
  - Multi-dimensional categorization, tagging
  - Scenario Marketplace

Example :

<https://addons.mozilla.org>





# Generic Wrapping I

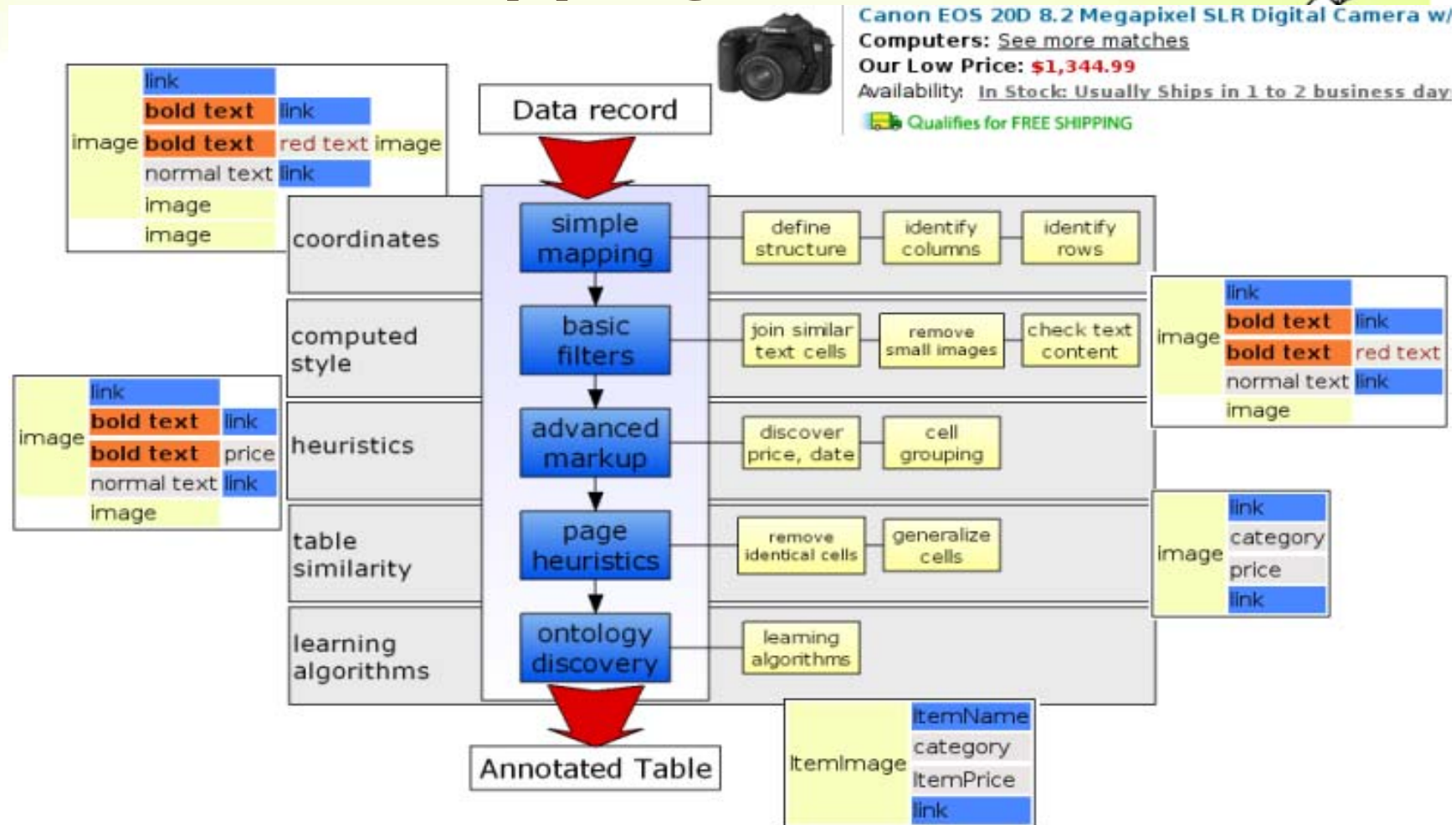


## *Visually based semantic annotation*

(Peter Szniek)

- Problems of tree-based /syntactic wrapping
  - Strong noise, dynamic, complex
- Visual Elements (WYSIWYW)
  - HTML, PDF → **meta-table**
  - Robust and simple structure for data extraction

# Generic Wrapping I - schema



# Generic Wrapping II



## *Inductive Wrapper Generation of Wrappers*

(Michal Ceresna)

- Applies machine learning
  - Tree-aligning
  - Attribute classification
  - Set covering machines

# Collaborative Semantic Database



- Scenario users supply and update the shared data storages
- In turn, the scenario users query also the data storages
- Data storage search engine
- Data schemas and ontologies mapping
- Data integration and inference

# Server-side solutions



- Farming
  - Scalability, load-balancing, failover
- Portal integration
  - Web access to server-side scenarios

# Misc



- Data cleaning toolset
- Advanced browser controls
- Internet Explorer hacking ☺
- Scenario development tools
- Core
  - Semantic Data Models – RDF, Sesame
  - Profiles, history support
  - XUL Runner
  - Scenarios Lego
  - Crawling

# Contact



<http://lumberjacz.org>

Viktor Zigo

Interested ?

[interested@lumberjacz.org](mailto:interested@lumberjacz.org)