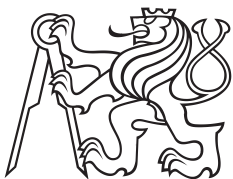


Diplomová práce



České
vysoké
učení technické
v Praze

F3

Fakulta elektrotechnická
Katedra kybernetiky

Minimální dokument

Jakub Podlaha

May 2014

/ Prohlášení

Prohlašuji, že jsem se neflákal.

Abstrakt / Abstract

Tento dokument je pouze pro potřeby testování.

This document is for testing purpose only.

Obsah /

1 Introduction	1	
1.1 Problem Statement and Motivation	1	
1.2 Current solution crOWler	2	
1.3 Proposed Solution and Methodology	2	
1.4 Specific goals of the project	2	
1.5 Work structure (XXX)	2	
2 Existing solutions	3	
2.1 Semantic and non semantic crawlers	3	
2.2 Advantages and pitfalls of Semantic crawler and linked data	3	
2.3 Research - existující řešení - platforma	4	
2.3.1 InfoCram 2000 - Jirka Mašek	4	
2.3.2 iMacros	4	
2.3.3 Selenium IDE	4	
2.4 crOWler	4	
2.4.1 zavislosti	5	
2.4.2 Classes of CrOWler	5	
2.4.3 Run configuration	5	
2.5 Strigil!	6	
2.5.1 What problems it solves? (Use cases)	6	
2.5.2 Architecture of Strigil platform	6	
2.5.3 What inspiration it brings for crawler	6	
3 Knowledge base, principles and technologies	7	
3.1 automatická extrakce dat	7	
3.2 RDF and RDFS	7	
3.3 OWL	7	
3.4 Linked Data	7	
3.5 Ontology repositories	7	
3.6 RDFa	7	
3.7 dalsi	7	
4 Program design	8	
4.1 Use Cases	8	
4.2 Workflow	8	
4.2.1 Main line	8	
4.2.2 Scenario creation	8	
4.2.3 Additional branches to Scenario Creation	8	
4.3 Model	9	
4.4 Implementation	9	
4.5 Issues - solved and unsolved	9	
5 Program Implementation	10	
6 Results and Tests	11	
6.1 Data	11	
6.1.1 Pamatky	11	
7 zaver	12	

Kapitola 1

Introduction

During past few years the Web went through bigger or smaller revolutions.

- WEB 2.0 and tag cloud
- HTML5 and semantic tags
- Smartphones, Tablets and mobile web everywhere,
- The run out of IPv4 addresses, nonexistent boom of IPv6,
- Cloud technologies and BigData,
- Bitcoin, Tor, anonymous internet,
- WikiLeaks, NSA, Heartbleed and security concerns
- Google Knowledge Graph, Facebook Open Graph, ...

That's only few examples of some of the biggest recent issues on the web in general. We live in an age, where so little can mean so much. The environment online is dramatically changing, mostly on a wave of some new, useful or frightening technology. The Semantic Web technologies have been described, standardized and implemented for several years now (XXX Example with linked data, rdf) and their tide seems to be near, though yet to come.

Semantic Web itself reates to several principles (along with their implementation) that allow users to add meaning to their data. This meaning brings not only a standardized structure, but also, as a consequence, the possibility to query and reason on data from multiple sources. Once given the structure, similar data can be joined in a form of a bigger cloud. This phenomena is called Linked Data.

In this work we'd like to bring the Semantic Web technologies closer to users. The approach is to propose a methodology for extracting structured data out of unstructured ones, designing and implementing an appropriate tool, to simplify the process of annotating (yet anonymous) data on a webpage, i.e. to bring structure and meaning into it.

1.1 Problem Statement and Motivation

Giving meaning, i.e. semantization of web pages gets more popular. Probably the most obvious example can be seen in the way the Google search engine serves it's results. Presenting not only the resulting pages but as well snippets of information scraped directly from the page content such as menu fields parsed directly from HTML5, contact information or opening hours, or even visualizing data from their own internal ontology, the Knowledge Graph.

XXX https://en.wikipedia.org/wiki/Google_Knowledge_Graph

XXX Strigil - <http://delivery.acm.org/10.1145/2540000/2539170/p453-starka.pdf>

What are the options for bringing sematic into a web?

One direction to go (XXX better) is to annotate data on **the server side**, i.e. at the time it is being created and/or published. The person or engine creating the

1. load the ontology
2. add selectors to specific resources from the ontology
3. run the crawling process according the above

■ 1.3 Proposed Solution and Methodology

To simplyfy the creation of guidelines, or scenarios for crOWLer, we propose a tool that allows user to select all the element directly on the web page being crawled, with all the necessary settings, pass the scenario created to the crOWLer and obtain the results in a form of a graphical feedback.

■ 1.4 Specific goals of the project

- design the semantic data creation use-cases
- implement extension for a browser
- load and visualise ontology
- create scenario for crOWLer
- serialize scenario and ontology
- parse it by crOWLer creating it's configuration
- run crOWLer
- visualize the extracted data (feedback)

■ 1.5 Work structure (XXX)

TBD

Kapitola 2

Existing solutions

2.1 Semantic and non semantic crawlers

By researching existing solutions, there is currently no open source or openly available solution to solve this task. Rumor goes there is proprietary tool in IBM.

Existing tools named as **Ontology-based Web Crawlers** refer mostly to crawlers that **rank** pages being crawled by guess-matching them against some ontology. In those programs user can't specify data that are being retrieved. Moreover, there is no way to get involved in the crawling process. It is solely used to automatically rank the relevance of documents. They are solving different task where input is several documents and possibly an ontology and output is the best matching document.

In case we are trying to solve the input is one or more documents and one or more ontologies and the result is data obtained from the documents and annotated with resources from the ontologies.

2.2 Advantages and pitfalls of Semantic crawler and linked data

The simplest approach is manual searching for keywords, or even simple browsing the web. That might be useful in some cases, but when there is a lot of data, it becomes exhausting.

Crawling data using simple tools like 'wget -mirror' allows us to load data and then write a program or script to retrieve a relevant information. This approach takes a lot of energy for one time only solution of a given problem.

By storing such crawled data into database we obtain persistent database, possibly automatically obtained by the script from pervious case. Such data is static, but can be queried over and over and possibly re-retrieved when becomes obsolete. It's structure is, however, based on programmers imagination and needs to be described in order to understand and handle the data properly.

When using Ontology-based solution, tailor made for crawling and annotating data from web, we obtain several benefits **for free**. The tool designed specially for this purpose makes it easy. Once the data is annotated, we can not only query on them, but also automatically reason on them and obtain more or more specific/narrow results than with general data. The attributes and relations within ontology, that allow reasoning, are usually part of the ontology definition and as such comes, again, **for free**.

Last for benefits: using ontology from public resource as a schema for our data can give us correct structure without need for XXX making it up or building it from scratch. Also by using some common ontology, we can join together any accessible data structured according to this ontology and simply query on resulting super set. With this approach we can utilize the power of linked data cloud (XXX reference).

Semantic crawling is not a silver bullet. The technology is only finding it's place and uses and it's being shaped by the needs of it's users. In current it's mostly used on accademic field XXX.

There is always a threat of inconsistency of an ontology when some data don't fit the rules or breaks structure of an ontology. (XXX more)

Just like with **hardcoded** crawling technique, the semantic crawling is tightly connected (XXX better) with the structure of the web being crawled and selectors (XXX explain term) used for matching data on the web. Any change on a webpage structure can lead to broken selectors or links during the crawling process (XXX and make the scenario useless, more on self-repairing of scenarios?).

A lot of web pages loads their data dynamically using AJAX queries. Some pages simply changes it's content frequently (XXX typically news pages, forums: rt.com, vimeo.com, ...) which would require almost constant crawling and growth into an massive ontology (XXX any suggestions on that? =).

Stating that, the semantic crawling is an usefull way to effectively obtain and query on (otherwise anonymous) data from the web, but it still have it's challenges to overtake.

2.3 Research - existující řešení - platforma

2.3.1 InfoCram 2000 - Jirka Mašek

- zalozeny na Aardwark ¹⁾

2.3.2 iMacros

- http://wiki.imacros.net/Command_Reference
- http://wiki.imacros.net/iMacros_for_Firefox
- http://wiki.imacros.net/iMacros_for_Chrome

2.3.3 Selenium IDE

- IDE - <http://www.seleniumhq.org/projects/ide/>
- plugins - <http://www.seleniumhq.org/projects/ide/plugins.jsp>
- current commands - <http://release.seleniumhq.org/selenium-core/1.0.1/reference.html>
- documentation - <http://docs.seleniumhq.org/docs/index.jsp>
- extending selenium API (blog, tutorial) - <http://adam.goucher.ca/?s=selenium&paged=2>
- randomString example - <http://adam.goucher.ca/?p=1348>

2.4 crOWLer

¹⁾ <https://addons.mozilla.org/en-US/firefox/addon/aardvark/>

- **2.5 Strigil!**
- **2.5.1 What problems it solves? (Use cases)**
- **2.5.2 Architecture of Strigil platform**
- **2.5.3 What inspiration it brings for crawler**

Kapitola 3

Knowledge base, principles and technologies

Linked Data, RDFa, ...

informativni cast, teorie

Seznamte se technologiemi pro automatickou extrakci dat z webových stránek a s jazyky sémantického webu RDF, RDFS a OWL.

3.1 automatická extrakce dat

3.2 RDF and RDFS

- https://en.wikipedia.org/wiki/Resource_Description_Framework

3.3 OWL

- <http://www.w3.org/TR/owl2-primer/>
- https://en.wikipedia.org/wiki/Web_Ontology_Language
- <http://www.w3.org/TR/2012/REC-owl2-quick-reference-20121211/>

3.4 Linked Data

- <http://linkeddata.org/guides-and-tutorials>
- <http://linkeddatabook.com/editions/1.0/>
- <http://lov.okfn.org/dataset/lov/>

3.5 Ontology repositories

- http://www.w3.org/wiki/Ontology_repositories

3.6 RDFa

- <https://www.sio2.cz/web/psiotwo/publications>
- <http://rdfa.info/play/>

3.7 dalsi

- <https://en.wikipedia.org/wiki/SPARQL>
- [https://en.wikipedia.org/wiki/Turtle_\(syntax\)](https://en.wikipedia.org/wiki/Turtle_(syntax))

Kapitola 4

Program design

4.1 Use Cases

- NPU
- RLP
- beerborec.cz
- citybee.cz

4.2 Workflow

4.2.1 Main line

- user loads/creates ontology using sowl
- user opens webpage with data
- user creates scenario using sowl
- sowl sends scenario to crawler
- crawler crawls the web according to scenario and stores results in repository
- + crawler sends data to sowl which embeds them in original web page (XXX)

4.2.2 Scenario creation

- user starts scenario creation in sowl
- loop until finished:
 - user selects an element on page
 - user select action on element (perform and record event, i.e. click on link, narrow HTML context, assign element - object or property according to situation, ...)
 - sowl records the action in scenario

4.2.3 Additional branches to Scenario Creation

- user can navigate through scenario by clicking scenario steps
- user can navigate through scenario by clicking ontological context
- user can navigate through scenario by clicking areas on webpage covered by scenario
- when user clicks on a hyperlink:
 - existing template can be assigned to the action (no need to actually follow the link)
 - new template can be created for resulting action (resulting page loaded, new template created, click through shown in breadcrumbs)

■ 4.3 Model

■ 4.4 Implementation

■ 4.5 Issues - solved and unsolved

- error handling (non existent selector, missing data, ...)



Kapitola **5**

Program Implementation

Kapitola 6

Results and Tests

6.1 Data

6.1.1 Památky

- <http://onto.mondis.cz/resource/page/npu/>
- <http://monumnet.npu.cz/pamfond/list.php?hledani=1&KrOk=&HiZe=&VybUzemi=1&sNazSidOb=&Adresa=&Cdom=&Pamatka=&CiRejst=&Uz=B&PrirUbytOd=3.5.1958&PrirUbytDo=10.12.2013>
- <http://dominanty.cz/pamatky-cihana.php>



Kapitola 7

zaver

TBD