

Master's thesis



Czech
Technical
University
in Prague

F3

Faculty of Electrical Engineering
Department of Computer Science and Engineering

Platform for semantic extraction of the web

Jakub Podlaha

January 2015

/ Declaration

Prohlašuji, že jsem se neflákal.

Abstrakt / Abstract

Tento dokument je pouze pro potřeby testování.

Překlad titulu: Platforma pro sémantickou extrakci webu

This document is for testing purpose only.

Contents /

1 Introduction	1
1.1 Problem Statement and Motivation	1
1.2 Current solution crOWLer	3
1.3 Proposed Solution and Methodology	3
1.4 Specific goals of the project	3
1.5 Work structure (XXX)	3
2 Knowledge base, principles and technologies	4
2.1 Technology of Semantic Web	4
2.2 RDF and RDFS	4
2.2.1 URI	4
2.2.2 RDF and RDFS vocabulary	5
2.3 OWL	5
2.4 Linked Data	5
2.5 Ontology repositories	5
2.6 RDFa	6
2.7 dalsi	6
2.8 automatická extrakce dat	6
3 Existing solutions	7
3.1 Semantic and non semantic crawlers	7
3.2 Advantages and pitfalls of Semantic crawler and linked data	7
3.3 Research - existující řešení - platforma	8
3.3.1 InfoCram 2000 - Jirka Mašek	8
3.3.2 iMacros	8
3.3.3 Selenium IDE	8
3.4 crOWLer	8
3.4.1 zavislosti	9
3.4.2 Classes of CrOWLer	9
3.4.3 Run configuration	9
3.5 Strigil!	10
3.5.1 What problems it solves? (Use cases)	10
3.5.2 Architecture of Strigil platform	10
3.5.3 What inspiration it brings for crawler	10
4 Program design	11
4.1 Use Cases	11
4.2 Workflow	11
4.2.1 Main line	11
4.2.2 Scenario creation	11
4.2.3 Additional branches to Scenario Creation	11
4.3 Model	12
4.4 Implementation	12
4.5 Issues - solved and unsolved	12
5 Program Implementation	13
6 Results and Tests	14
6.1 Data	14
6.1.1 Pamatky	14
7 zaver	15

Chapter 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 relates 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 semantic 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 data

have to use the right tool and spend time giving the data the appropriate annotation. There are standards covering this use case: HTML5 brings in tags for clearer specification of the page structure (such as `nav`, `article`, `section`, `aside`, ...). Microformats <http://microformats.org/> define specialized values for HTML class attribute to bring standardized patterns for several basic use cases with fixed structure, such as vCard or Event. The microformat approach is easy to implement as it doesn't impose any extra technology into the stack (XXX better). Last but not least, we can use RDFa to annotate data on a webpage with an actual ontology. This technology is part of the Semantic Web stack and we'll describe it closer in further chapter (see in separate section XXX).

Annotating data on the server side enables users to use tools highlight data they are specifically interested in, extract them and reason on them. Services can use annotated data, combine them and offer results from multiple sources.

Some examples of utilities for extracting and testing structured data:

<http://www.google.com/webmasters/tools/richsnippets>

`http://rdfa.info/play/`

To bypass the gap between anonymous data present on the web on one side and rich, linked, meaningful ontologies (XXX example) on the other, we can go the opposite direction as well. We can take the unannotated data already present on the web and retrieve them in a form, that is defined by some ontology structure. This process can be performed automatically or manually.

There've been several (XXX) attempts in the automatic data annotation. In principle, the goal of this approach is to automatically analyze a web page, find the most appropriate ontology for it and use this ontology to automatically annotate data on the page.

As an example, on a webpage of a music band this would find (XXX <http://www.musicontology.com/>) as the best matching schema, and it would annotate band name, genre, albums using classes and properties from this ontology.

While on well structured, simple and/or partially annotated pages this process can be very successful and produce useful results, on pages with unorganized data the confidence on results produced by this approach might drop to lottery (XXX). Unfortunately many online webpages and services are poorly structured. Pages containing many unrelated data, in form of advertisements or other external content might confuse such an engine. Old servers present their content in poorly structured or even invalid HTML in a form of multiple nested tables that serve for structuring only and makes the whole structure of the web unreadable (XXX). Social aspect of web brings in almost complete randomness making it even harder to automatically reason on page's content if it can't be distinguished and potentially left aside. To sum it up, there are many potential and real threads that prevent us from automatically annotate all data on web with confidence (XXX).

In some use cases the ontology of the desired data is yet to be created and the user is aware of the data structure and capable of manually spot and select the data on a web page. Currently there isn't many tools allowing this kind of operation. The ideal implementation and the vision here will allow user to partially identify the structure of a webpage while leaving the repetitive tedious work on crawler following the same structure on all data on webpage.

To allow such a process we need to create tools that allow users to annotate the, previously meaning-free, data with elements of existing ontology and/or create resources

on-the-go. By using existing ontologies we not only give the meaning to our data, but also valuable connection to any other dataset annotated using the same ontology.

1.2 Current solution crOWLer

The suggested base-technology is being developed on our faculty XXX. Crawler called crOWLer serves the needs of extracting data from web. In current implementation, both, the scenario, followed by the crawler, and the ontology structure/schema are hard-coded into the crOWLer code. This requires unnecessary load of work for each separate use case, whilst in practice all the use cases share the same workflow.

1. load the ontology
2. add selectors to specific resources from the ontology
3. implement the rules to follow another page
4. run the crawling process according the above

1.3 Proposed Solution and Methodology

To simplify 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

Chapter 2

Knowledge base, principles and technologies

2.1 Technology of Semantic Web

Wikipedia defines Semantic Web as a collaborative movement led by international standards body the World Wide Web Consortium (W3C). (XXX https://en.wikipedia.org/wiki/Semantic_Web) W3C itself defines Semantic Web as a technology stack to support a Web of data, as opposed to Web of documents, the web we commonly know and use (XXX <http://www.w3.org/standards/semanticweb/>). Just like with Cloud or Big Data the proper definition tends to vary, but the notion remains the same. It is collaborative movement led by W3C and it does define a technology stack. It also includes users and companies using this technology and the data itself. Technologies and languages of Semantic Web such as RDF, RDFa, OWL, SPARQL (XXX) are well standardized and will be described in following chapters.

As a general logical concept of the technology... (XXX) Technology of Semantic Web is used to take data and metadata, give them unique identifiers and form them into oriented graphs. The metadata part define a schema of types and properties that can be assigned to data and also relations between this types and properties possibly in a form of ontology. When some data are annotated by resources from such an ontology we gain power to reason on this data, i.e. resolve new relations based on known ones, and also to query on our data along with any data annotated using the same ontology.

In RDF(S) the is defined in a form of triples. Triple consists of subject, predicate and object, wich all are simply **resources** listed by their identifiers. In this very general form we can express basically any relationship between two resources. On a level of classes and properties, we can for example assign a type to an individual, or set a class as a domain of some property. On a level of ontologies we can specify author and date it was released. (XXX DELME)

2.2 RDF and RDFS

RDF is a family of specifications for syntax notations and data serialization formats, meta data modeling, and vocabulary used for it.

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

We will look closely on URI, the resource identifier, vocabularies and semantics defined by RDF, RDFS, and OWL, and serialization into Turtle and RDF/XML formats.

2.2.1 URI

In order to give each resource an unique identifier a Uniform Resource Identifier is used. This is mostly in a form of URL as we commonly know it as **web address** (e.g. <http://www.example.org/some/place#something>). In some cases URI can be a URN as well. URN is a complementary syntax for URL that allow us to identify resources without specifying their location. This way we can for example use ISBN codes when

working with books and records, or UUID identifier a Universally Unique Identifier widely used to identify technically any data instance.

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

■ 2.2.2 RDF and RDFS vocabulary

In order to work with data properly (XXX) RDF(S) vocabulary defines several basic URIs along with their semantics.

`rdf:type` is a property used to state that a resource is an instance of a class. A commonly accepted qname for this property is `a`.^[4]

`rdfs:Resource` - is the class of everything. All things described by RDF are resources.
`rdfs:Class` - declares a resource as a class for other resources.

`rdfs:Literal` – literal values such as strings and integers. Property values such as textual strings are examples of RDF literals. Literals may be plain or typed. `rdfs:Datatype` – the class of datatypes. `rdfs:Datatype` is both an instance of and a subclass of `rdfs:Class`. Each instance of `rdfs:Datatype` is a subclass of `rdfs:Literal`. `rdf:XMLLiteral` – the class of XML literal values. `rdf:XMLLiteral` is an instance of `rdfs:Datatype` (and thus a subclass of `rdfs:Literal`).

`rdf:Property` – the class of properties. `rdfs:domain` of an `rdf:property` declares the class of the subject in a triple whose second component is the predicate. `rdfs:range` of an `rdf:property` declares the class or datatype of the object in a triple whose second component is the predicate.

`rdfs:subClassOf` allows to declare hierarchies of classes. `rdfs:subPropertyOf` is an instance of `rdf:Property` that is used to state that all resources related by one property are also related by another.

`rdfs:label` is an instance of `rdf:Property` that may be used to provide a human-readable version of a resource's name. `rdfs:comment` is an instance of `rdf:Property` that may be used to provide a human-readable description of a resource.

These are the basic building blocks of our future RDF graphs. The semantics defined in the specification and slightly described here allow us to specify class hierarchy, properties with domain and range as well as use this structure on individuals and literals.

■ 2.3 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/>

■ 2.4 Linked Data

Wikipedia defines Linked Data as a term used to describe a recommended best practice for exposing

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

■ 2.5 Ontology repositories

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

■ 2.6 RDFa

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

■ 2.7 dalsi

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

■ 2.8 automatická extrakce dat

TODO in next section

Chapter 3

Existing solutions

3.1 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.

3.2 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 useles, 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.

■ 3.3 Research - existující řešení - platforma

■ 3.3.1 InfoCram 2000 - Jirka Mašek

- zalozeny na Aardwark ¹⁾

■ 3.3.2 iMacros

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

■ 3.3.3 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>

■ 3.4 crOWLer

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

■ 3.4.1 zavislosti

■ maven - apache project managing tool

- <https://maven.apache.org>
- <https://maven.apache.org/run-maven/index.html>
- <https://maven.apache.org/guides/mini/guide-ide-eclipse.html>

■ sesame

- <http://www.openrdf.org/download.jsp> ??

■ jena

- <https://github.com/ansell/JenaSesame> !!
- or <https://github.com/afs/JenaSesame> ??
- or <http://jena.apache.org/> ???
- or <http://sjadapter.sourceforge.net/> ????
- or <http://sourceforge.net/projects/jenasesamemodel/>
- might help <http://www.iandickinson.me.uk/articles/jena-eclipse-helloworld/>
- little hint http://spqr.cerch.kcl.ac.uk/?page_id=130
- another hit <http://answers.semanticweb.com/questions/20865/how-to-get-the-jena-sesame-adapter>
- wiki [https://en.wikipedia.org/wiki/Jena_\(framework\)](https://en.wikipedia.org/wiki/Jena_(framework))
- jena vs. sesame flame <http://answers.semanticweb.com/questions/1638/jena-vs-sesame-is-there-a-serious-complete-up-to-date-unbiased-well-informed-side-by-side-comparison-between-the-two>

■ 3.4.2 Classes of CrOWLer

■ ImmovableHeritageConfiguration extends MonumnetConfiguration implements ConfigurationFactory

- implements Configuration, which is parameter for FullCrawler.run() method

■ FullCrawler

- implements the whole crawling algorithm

■

■ 3.4.3 Run configuration

```
crowler cz.sio2.crowler.configurations.npu.ImmovableHeritageConfiguration
file results
```

```
crowler cz.sio2.crowler.configurations.kub1x.KbxConfiguration file re-
sults
```

```
crowler cz.sio2.crowler.configurations.parser.SeleniumConfiguration\
file results generated.html
```

- Class ImmovableHeritageConfiguration implements Configuration class.
- Folder jena_con will be created and all the rdf's will be stored in int with names derived from ontology uri

■ 3.5 Strigil!

- 3.5.1 What problems it solves? (Use cases)
- 3.5.2 Architecture of Strigil platform
- 3.5.3 What inspiration it brings for crawler

Chapter 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, ...)



Chapter 5

Program Implementation


Chapter 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>



Chapter 7

zaver

TBD