

Projektdokumentation im Modul Semantic Web

Anzahl der Publikationen pro Mitarbeiter der HTWK-Leipzig

Marcel Kirbst

30. Juni 2014

Recherchefragestellung: Eine Liste aller Angestellten der HTWK-Leipzig, geordnet nach der Anzahl der bisherigen Publikationen im deutschsprachigen Raum.

1 Inhaltliche Interpretation der Fragestellung

Neben der Lehre stellt die Forschung und die Publikation deren Resultate eine wichtige Kernaufgabe der Professoren und wissenschaftlichen Mitarbeiter an Hochschulen dar. Die Aufgabe dieser Projektarbeit besteht daher zum einen darin, alle Mitarbeiter der HTWK-Leipzig zu ermitteln und in einer semantischen Datenbank zu erfassen.

Weiterhin sind Daten aller im deutschsprachigen Raum erschienen Publikationen aufzubereiten und einzubinden. Diese Daten stellt die Deutsche Nationalbibliothek kostenfrei zur Verfügung.

Die Umsetzung des Projektes erfolgt auf einer virtuellen Maschine unter Verwendung der Virtualisierungslösung VirtualBox der Firma Oracle. Auf der virtuellen Maschine wurde als Gastbetriebssystem die Linux-Distribution Kubuntu 14.04 von Canonical eingesetzt.

2 Relevante Datenquellen

An dieser Stelle erfolgt eine Auflistung aller relevanten Datenquellen und deren Beschreibung.

2.1 Mitarbeiterkatalog der HTWK-Leipzig

Die aktuellste, öffentlich zugängliche Auflistung aller Mitarbeiter der HTWK-Leipzig besteht aus dem HTWK-Telefonverzeichnis. Diese HTML-Seite lässt sich unter der Webadresse [1] abrufen. Eine aufbereitete Version dieser Daten im JSON-Format wird von Herrn Henri Knochenhauer, B.Sc. und Herrn Roy Meisser, B.Sc unter [2] zur Verfügung gestellt und in diesem Projekt als Datenquelle genutzt.

Link	http://141.57.21.45:8080/info/staff
Datenformat	JSON
Schnittstelle	HTTP Rest-API
Lizenz	Daten: unbekannt. Datenquelle: GPL.

Die Daten dieser Datenquelle umfassen die eindeutige Mitarbeiteridentifikationsnummer, Nachname, Vornamen, akademischen Grad sowie die Fakultät jedes Mitarbeiters.

2.2 Bibliothekskatalog der Deutsch Nationalbibliothek

Die Deutsche Nationalbibliothek (DNB) sammelt alle deutschen Publikationen seit dem Jahr 1913. Darüber hinaus bietet sie umfangreiche Dienstleistungen für Bibliotheken und Wissenschaftler an. Mittels einer durch die DNB vergebenen Normdatensatz (GND)¹ lassen sich Werke, Personen oder Körperschaften identifizieren. Diese Identifizierung ermöglicht es Personen unterschiedlicher Quellen zu verlinken.

Link	http://datendienst.dnb.de/cgi-bin/mabit.pl?userID=opendata&pass=opendata&cmd=login
Datenformat	RDF
Schnittstelle	Linked Data, Dump, Rest-API
Lizenz	CC-SA

Die Daten lassen sich unter der angegebenen Webadresse direkt herunterladen. Für dieses Projekt wurde die Datei GND.rdf.gz herunter geladen, welche die gesammelten Normdaten zu den Autoren enthält und entpackt derzeit etwa 9 GByte groß ist. Weiterhin wurde

¹Gemeinsame Normdatei

die Datei DNBTitel.rdf.gz heruntergeladen. Diese Datei enthält die Namen und Autoren aller Werke der Deutschen Nationalbibliothek und ist entpackt derzeit etwa 12 GByte groß.

3 Extraktion relevanter Daten und import in einen Triplestore

Die Extraktion aller 3 Dateien erfolgt über eine für dieses Projekt erstellte JAVA-Anwendung, welche im Repository zu diesem Projekt hinterlegt ist. Das Repository kann unter folgender Webadresse abgerufen werden[3].

3.1 Extraktion der HTWK-Mitarbeiter

Die Daten der HTWK-Mitarbeiter liegen im JSON Format vor. Die exakte Struktur zeigt das folgende Listing auszugsweise.

```
1 [
2   - {
3     cuid: "597",
4     name: "Siebeck, Andrea",
5     degree: "Dipl.-Angl.",
6     faculty: "AAA"
7   },
8   - {
9     {
10      cuid: "8",
11      name: "Engel, Heike",
12      degree: "Dipl.-Wirtschaftsinf.",
13      faculty: "DF"
14    }
15 ]
```

Listing 1: Format der importierten HTWK-Mitarbeiterdaten

Die JAVA-Anwendung bezieht das JSON-Dokument direkt von der Webressource. Die JAVA-Anwendung generiert aus den importierten Daten eine RDF-Datei, welche nach folgendem RDF-Schema generiert wird.

```
1 @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
2 @prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
3 @prefix owl: <http://www.w3.org/2002/07/owl#> .
4 @prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
5 @prefix : <https://raw.githubusercontent.com/mkirstb/semanticweb/master/scheme#>↵
6 .
```

```

7 :Employee a rdfs:Class .
8 :Publications a rdfs:Class .
9
10 :Int a rdfs:Datatype; owl:onDatatype xsd:integer; xsd:minInclusive 1 .
11 :String a rdfs:Datatype; owl:onDatatype xsd:string .
12
13 :cuid a rdf:Property; rdfs:domain :Employee; rdfs:range :Int .
14 :lastname a rdf:Property; rdfs:domain :Employee; rdfs:range :String .
15 :firstname a rdf:Property; rdfs:domain :Employee; rdfs:range :String .
16 :degree a rdf:Property; rdfs:domain :Employee; rdfs:range :String .
17 :dnbautorid a rdf:Property; rdfs:domain :Employee; rdfs:range :Int .
18 :birth a rdf:Property; rdfs:domain :Employee; rdfs:range :Int .
19
20 :dnbpubid a rdf:Property; rdfs:domain :Publication; rdfs:range :Int .
21 :dnbautorid a rdf:Property; rdfs:domain :Publication; rdfs:range :Int .
22 :title a rdf:Property; rdfs:domain :Publication; rdfs:range :String .
23 :pubdate a rdf:Property; rdfs:domain :Publication; rdfs:range :Int .

```

Listing 2: RDF-Schema der exportierten HTWK-Mitarbeiterdaten

Da diese Datei im Format N3 vorliegt, wurde diese mit einem RDF-Translator [4] in das XML-Format überführt.

Die vom Java-Programm generierte RDF-Datei kann nun direkt über die Webschnittstelle in OntoWiki importiert werden und hat folgendes Format.

```

1 <rdf:RDF
2   xmlns:semweb="https://raw.githubusercontent.com/mkirbst/semanticweb/master/↵
   scheme#"
3   xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
4 >
5
6   <semweb:Employee rdf:about="https://raw.githubusercontent.com/mkirbst/↵
   semanticweb/master/htwkstaff.json#cuid597">
7     <semweb:cuid>597</semweb:cuid>
8     <semweb:lastname>Siebeck</semweb:lastname>
9     <semweb:firstname> Andrea</semweb:firstname>
10    <semweb:degree>Dipl.-Angl.</semweb:degree>
11    <semweb:faculty>AAA</semweb:faculty>
12    <semweb:dnbautorid>0</semweb:dnbautorid>
13    <semweb:birth>0</semweb:birth>
14  </semweb:Employee>
15  <semweb:Employee rdf:about="https://raw.githubusercontent.com/mkirbst/↵
   semanticweb/master/htwkstaff.json#cuid8">
16    <semweb:cuid>8</semweb:cuid>
17    <semweb:lastname>Engel</semweb:lastname>
18    <semweb:firstname> Heike</semweb:firstname>
19    <semweb:degree>Dipl.-Wirtschaftsinf.</semweb:degree>
20    <semweb:faculty>DF</semweb:faculty>
21    <semweb:dnbautorid>0</semweb:dnbautorid>
22    <semweb:birth>0</semweb:birth>
23  </semweb:Employee>
24 </rdf:RDF>

```

Listing 3: RDF-Format der exportierten HTWK-Mitarbeiterdaten

3.2 Extraktion der Daten der Deutschen Nationalbibliothek

Unter der Webadresse [5] lassen sich die benötigten Daten direkt im RDF-Format herunterladen. Die heruntergeladenen Dateien sind jedoch mit 9 und 12 GByte so groß, dass diese nicht direkt über die Webschnittstelle des OntoWiki importiert werden können.

Aus diesem Grund filtert die JAVA-Anwendung aus den Datensätzen nur diese heraus, deren Autor den Vor- und Nachname von HTWK-Mitarbeitern enthält. Um diese Operationen effizient durchführen zu können, werden in der JAVA-Anwendung Hashmaps erzeugt, die als Schlüssel die betreffenden RDF-Tags so wie Vor- und Nachname als Hashwert enthalten. Durch diese Implementierungsvariante werden die Vorteile der Hashmap-Datenstruktur genutzt. Außerdem werden so wenig wie möglich rechenintensive Stringmanipulationen durchgeführt.

```
199 Map<String, emp>hEmps = new HashMap<String, emp>();
```

Listing 4: Auszug JAVA-Anwendung: Hashmap zur schnelleren Filterung der Normdaten

```
233 //DNB-GND.rdf name format as hasmap key value
234 hEmps.put("<rdf:li>"+name+"</rdf:li>", tmpemp);
```

Listing 5: Auszug JAVA-Anwendung: Einfügen von Objekten in die Hashmap

```
135 int hits = 0;
136 String begintag = "<rdf:Description";
137 String endtag = "</rdf:Description";
138 String nametag = "<gndo:preferredNameForThePerson>";
139 StringBuilder t3 = new StringBuilder();
140 String ts3 = "";
141 boolean found = false;
142
143 BufferedReader br = new BufferedReader(new FileReader(gndfilepath));
144 String line;
145
146 while ((line = br.readLine()) != null) {
147     String linetrimmed = line.trim(); // performance: trim only once
148
149     //linetrimmed begins with begintag OR nametag OR endtag OR other
150     if(linetrimmed.startsWith(begintag))
151     {
152         t3 = new StringBuilder();
153         t3.append(line+"\n");
154         found = false;
155     } else if (linetrimmed.startsWith(nametag))
156     {
157         t3.append(line+"\n");
158         if(hEmps.containsKey(linetrimmed)) { //hEmps == hashmap containing htwk ←
            employee objects
159             found = true;
160             hits++;
161         }
162     }
163 }
```

```

162 } else if (linetruncated.startsWith(endtag)) {
163     t3.append(line+"\n\n");
164     if(found == true) {
165         System.out.println(hits + " " + t3.toString());
166         writer.write(t3.toString());
167     }
168 } else {
169     t3.append(line+"\n");
170 }
171 }

```

Listing 6: Auszug JAVA-Anwendung: RDF-Parser

Das resultierende JAVA-Programm verarbeitet die gesamten Datensätze beider Dateien innerhalb einer Laufzeit von circa 2 Minuten. Die zu Grunde liegende Hardware ist eine Laptop der Marke HP EliteBook 8470w (CPU:Intel(R) Core(TM) i7-3740QM CPU @ 2.70GHz, RAM 16 GB, Primärspeicher: Micron SSD 250 GB).

Die resultierenden beiden Dateien haben exakt das gleiche RDF-Format wie die Eingabedateien, jedoch bereinigt um alle Datensätze, deren Autoren nicht einen Namen eines HTWK-Mitarbeiters tragen. Diese Dateien lassen sich nun über die Webschnittstelle von OntoWiki importieren.

4 Verlinkung von Ressourcen

Mit den 3 verwendeten Datenquellen lassen sich die Daten auf 2 Arten verknüpfen um die gewünschten Ergebnisse zu erhalten. Die Daten der Deutschen Nationalbibliothek sehen für jeden Autor eine eindeutige Identifikationsnummer vor. Als problematisch haben sich im Verlauf des Projektes dabei 2 Umstände erwiesen.

Die Deutsche Nationalbibliothek weist darauf hin das sich die für das Projekt verwendeten Daten noch in einem Migrationsprozess befinden und noch nicht vollständig konsistent sind. Daraus resultiert, dass manche Autoren mehrfach mit unterschiedlichen Publikationen im Datenbestand vorhanden sind.

Das zweite Problem besteht darin, dass von der Datenquelle HTWK-Telefonbuch nur Nachname und Vornamen für die Verlinkung einbezogen werden können. Weitere hilfreiche Attribute wie beispielsweise das Geburtsjahr sind in dieser Datenquelle nicht hinterlegt. Für ein wirklich tragfähiges Ergebnis des Projekts sollten folglich konsistente Datenquellen heran gezogen werden.

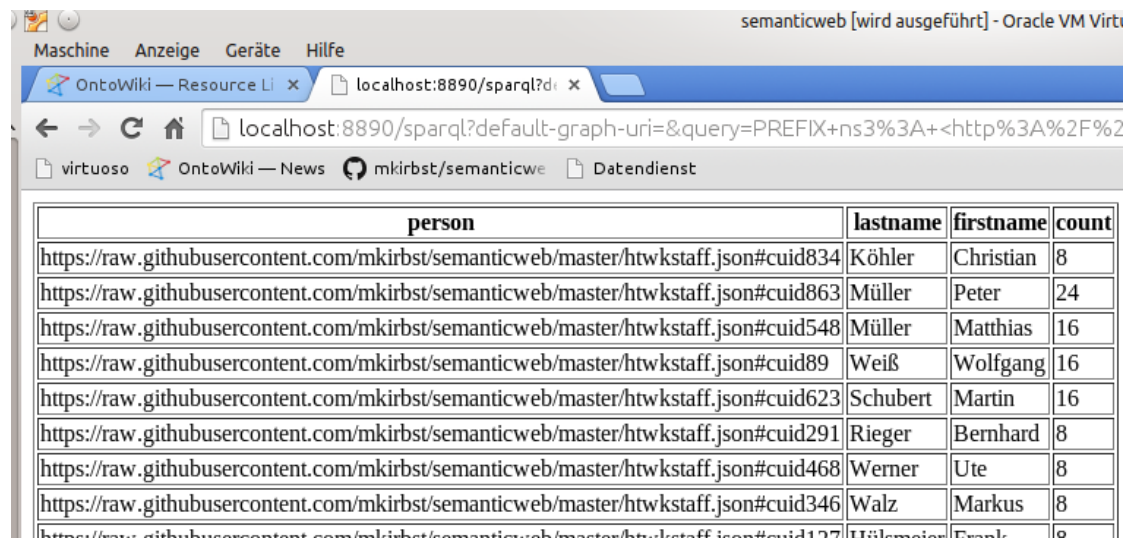
5 Anfrage an die Forschungswissensbasis

5.1 SPARQL-Anfrage

```
1 PREFIX bibo: <http://purl.org/ontology/bibo/>
2 PREFIX ns3: <http://purl.org/ontology/bibo/>
3 PREFIX owl: <http://www.w3.org/2002/07/owl#>
4 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
5 PREFIX ns2: <http://rdvocab.info/>
6 PREFIX dct: <http://purl.org/dc/terms/>
7 PREFIX dc: <http://purl.org/dc/elements/1.1/>
8 PREFIX ns1: <http://iflastandards.info/ns/isbd/elements/>
9 PREFIX ns0: <http://id.loc.gov/vocabulary/relators/>
10 PREFIX semweb: <https://raw.githubusercontent.com/mkirbst/semanticweb/master/↵
    scheme#>
11 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
12 PREFIX marcRole: <http://id.loc.gov/vocabulary/relators/>
13
14 select ?person ?lastname ?firstname (count(?document)+count(?article) as ?count)
15 where {
16     ?person rdf:type semweb:Employee;
17             semweb:lastname ?lastname;
18             semweb:firstname ?firstname .
19
20     ?document rdf:type ns3:Document;
21               marcRole:edt ?editor .
22
23     ?article rdf:type bibo:Article;
24              marcRole:edt ?aeditor;
25              dc:creator ?acreator
26
27     Filter (
28         strstarts(concat(?editor, " "), ?lastname)
29         AND strends(concat(?editor, " "), ?firstname)
30         AND strstarts(concat(?aeditor, " "), ?lastname)
31         AND strends(concat(?aeditor, " "), ?firstname)
32     )
33 }
34 }
35 LIMIT 200
```

Listing 7: SPARQL-Abfrage

5.2 Ergebnis der Anfrage



person	lastname	firstname	count
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid834	Köhler	Christian	8
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid863	Müller	Peter	24
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid548	Müller	Matthias	16
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid89	Weiß	Wolfgang	16
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid623	Schubert	Martin	16
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid291	Rieger	Bernhard	8
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid468	Werner	Ute	8
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid346	Walz	Markus	8
https://raw.githubusercontent.com/mkirbst/semanticweb/master/htwkstaff.json#cuid127	Hilmeier	Frank	8

Abbildung 1: Ergebnis der SPARQL-Abfrage

6 Interpretation und Zusammenfassung

Zusammenfassend lässt sich feststellen, dass die semantische Verknüpfung der Daten auf diese Weise möglich ist. Mit dem derzeitigen Datenbestand lassen sich jedoch automatisiert keine zufriedenstellenden Resultate erzielen. Sollten die zu Grunde liegenden Daten der Deutschen Nationalbibliothek in Zukunft in einem vollständig konsistenten Zustand befinden, ist mit aussagekräftigeren Ergebnissen zu rechnen. Ein manueller Abgleich der einzelnen Personendaten oder die Zuhilfenahme weiterer Autoreneigenschaften, wie beispielsweise des Geburtsdatums, lassen ebenfalls eine höhere Qualität der Ergebnisse erwarten.

Alle Ressourcen zu diesem Projekt sind im Repository unter der Webadresse [3] hinterlegt.

Projektdokumentation

- [1] <http://www.htwk-leipzig.de/de/hochschule/telefonverzeichnis/>
abrufbar am 11.06.2014
- [2] <http://141.57.21.45:8080/info/staff> abrufbar am 22.06.2014
- [3] <https://github.com/mkirbst/semanticweb>
abrufbar am 12.06.2014
- [4] <http://rdf-translator.appspot.com/>
abrufbar am 25.06.2014
- [5] <http://datendienst.dnb.de/cgi-bin/mabit.pl?userID=opendata&pass=opendata&cmd=login>
abrufbar am 25.06.2014