

# Projektdokumentation im Modul Semantic Web

## Anzahl der Publikationen pro Mitarbeiter der HTWK-Leipzig

Marcel Kirbst

1. Juli 2014

# Inhaltsverzeichnis

- 1 Recherchefragestellung
- 2 relevante Datenquellen
- 3 Extraktion relevanter Daten
- 4 Verlinkung der Ressourcen
- 5 Anfrage an die Forschungswissensbasis
- 6 Interpretation und Zusammenfassung

# Recherchefragestellung

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

# Inhaltliche Interpretation

- erfassen alle Mitarbeiter der HTWK in semantischer Datenbank
- erfassen aller Publikationen im deutschsprachigen Raum
- semantische Verknüpfung dieser Daten

# Umsetzung

- Virtuelle Maschine Oracle VirtualBox
- Betriebssystem Linux-Distribution Canonical Kubuntu 14.04
- Datenbank Virtuoso + OntoWiki

# HTWK-Telefonverzeichnis

- aktuellste, öffentlich zugängliche Ressource die gefunden wurde
- abrufbar unter `http://www.htwk-leipzig.de/de/hochschule/telefonverzeichnis/`
- als JSON aufbereitete Daten von Herrn Henri Knochenhauer, B.Sc und Herrn Roy Meissner, B.Sc. verfügbar
- abrufbar unter `http://141.57.21.45:8080/info/staff`

# Normdaten der Deutschen Nationalbibliothek

- sind öffentlich zugänglich, halbjährlich aktualisiert, liegen direkt im RDF-Format vor
- jedoch sehr groß, (GND.rdf ca.9GB, DNB-Titel.rdf ca.12GB), teilweise inkonsistent
- abrufbar unter `http://datendienst.dnb.de/cgi-bin/mabit.pl?userID=opendata&pass=opendata&cmd=login`

## Importierbare Aufbereitung der Daten erforderlich

- HTWK-Mitarbeiterdaten: JSON  $\Rightarrow$  RDF
- DNB-Daten passen nicht in die VM, Extraktion der Daten die HTWK-Mitarbeiter betreffen
- Implementierung eines JAVA-Programms



## Verlinkung der Ressourcen direkt ueber Nachname, Vorname

- Vorteil: Mitarbeiterdaten lassen sich direkt mit den DNB-Titeldaten verknüpfen
- Nachteil: sehr viele falsche Ergebnisse aufgrund doppelter Namen

## Verlinkung der Ressourcen ueber eindeutige AutorenID

- (theoretischer) Vorteil: AutorenID ist theoretisch eindeutig, es existiert theoretisch nur eine AutorenID pro HTWK-Mitarbeiter
- Praxis: Daten (im Moment noch) inkosistent, teilweise mehrere AutorenIDs, jeweils unterschiedliche Publikationen zugeordnet

# SPARQL-Anfrage

siehe code

## Ergebnis SPARQL-Anfrage

siehe code

# Interpretation und Zusammenfassung

- Aufbereitung der Daten viel zeitintensiver als erwartet
- aufgrund teilweise inkonsistenter Daten noch keine tragfähigen Ergebnisse
- ggf manuelles Aufbereiten der HTWK-Daten (Geburtsjahr, DNB-ID hinzufügen) sollte auch zu besseren Ergebnissen führen