# Converting WHO's Global Health Observatory Data to RDF

Amrapali Zaveri

AKSW, Institut für Informatik

1

# Outline

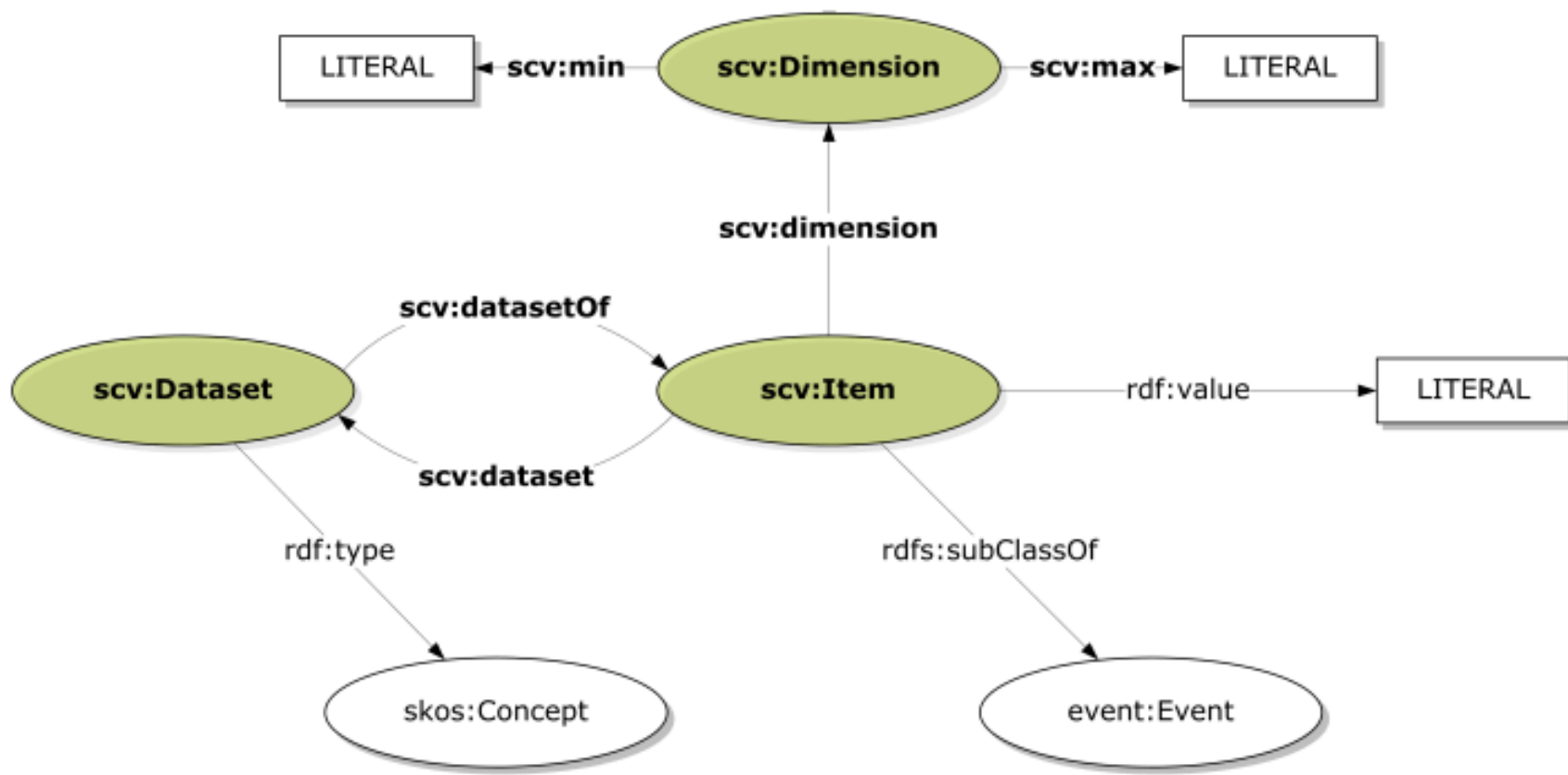- Background
- What is SCOVO?
- Semi-automated approach
- OntoWiki's CSV Import plug-in
- SCOVOfied WHO data
- Challenges and Future Work
- References

# Background

- Biomedical statistical data
  - Published as Excel sheets
- Advantage
  - Readable by humans
- Disadvantages
  - Cannot be queried efficiently
  - Difficult to integrate with other data (in different formats)
- Our approach
  - Converting data into a single data model - RDF
  - Using SCOVO (Statistical Core Vocabulary)*
    - designed particularly to represent multidimensional statistical data using RDF.

*Michael Hausenblas, Wolfgang Halb, Yves Raimond, Lee Feigenbaum, and Danny Ayers. Scovo: Using statistics on the web of data. In ESWC, pages 708–722, 2009.

Thursday, August 26, 2010

# What is SCOVO?



scv: <http://purl.org/NET/scovo#>
event: <http://purl.org/NET/c4dm/event.owl#>
skos: <http://www.w3.org/2004/02/skos/core#>
rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
rdfs: <http://www.w3.org/2000/01/rdf-schema#>

The Statistical Core Vocabulary (scovo)
http://purl.org/NET/scovo
v0.3@2008-05-15

**Fig.: The Statistical Core Vocabulary (scovo)**

Thursday, August 26, 2010

# Semi-automated approach

- Transforming CSV to RDF in a fully automated way is not feasible.
    - Dimensions may often be encoded in heading or label of a sheet
- Our semi-automatic approach:
    - As a plug-in in OntoWiki#
        - a semantic collaboration platform developed by the AKSW research group.
    - A CSV file is converted into RDF using SCOVO



*# Sören Auer, Sebastian Tramp (geb. Dietzold), Jens Lehmann, and Thomas Riechert: OntoWiki: A Tool for Social Semantic Collaboration In: Proceedings of the Workshop on Social and Collaborative Construction of Structured Knowledge CKC 2007 at the 16th*

Thursday, August 26, 2010

# 1. Create Knowledge Base

# 2. Import a CSV file

# 3. Define dimensions

Thursday, August 26, 2010

# 4. Define data range

Thursday, August 26, 2010

# 5. Save template, extract triples

10

# 6. Re-use template for similar files

# 7. View resources

Thursday, August 26, 2010

# SCOVOfied WHO's Global Health Observatory Data

```
prefix ex:<http://example.org/who-data>
prefix scv:<http://purl.org/NET/scovo>

ex:Country            rdfs:subClassOf      scv:Dimension;
                      rdf:type             rdfs:Class;
                      dc:title             "Country".
ex:Disease            rdfs:subClassOf      scv:Dimension;
                      rdf:type             rdfs:Class;
                      dc:title             "Disease".
ex:CountryCode        rdfs:subClassOf      scv:Dimension;
                      rdf:type             rdfs:Class;
                      dc:title             "CountryCode".


ex: Afghanistan       rdf:type             ex:Country;
                      dc:title             "Afghanistan" .
ex:Tuberculosis       rdf:type             ex:Disease;
                      dc:title             "Tuberculosis" .
ex:3010               rdf:type             ex:CountryCode;
                      dc:title             "3010".


ex:c1-r6              rdf:type             scv:Item;
                      rdf:value            127;
                      scv:dimension        ex:Afghanistan;
                      scv:dimension        ex:Tuberculosis .
                      scv:dimension        ex:3010
```

After converting a file containing 5 dimensions and 22384 statistical data items, an RDF model containing [13]

# Challenges and Future Work

- There may be some Excel sheets that contain taxonomies only readable by humans.

| | | | | |
|---|---|---|---|---|
| | All Causes | | | |
| I. | Communicable, maternal, perinatal and nutritional conditions | | | |
| | A. | Infectious and parasitic diseases | | |
| | | 1 | Tuberculosis | |
| | | 2 | STDs excluding HIV | |
| | | | a. | Syphilis |
| | | | b. | Chlamydia |
| | | | c. | Gonorrhoea |
| | | 3 | HIV/AIDS | |
| | | 4 | Diarrhoeal diseases | |
| | | 5 | Childhood-cluster diseases | |
| | | | a. | Pertussis |
| | | | b. | Poliomyelitis |
| | | | c. | Diphtheria |
| | | | d. | Measles |
| | | | e. | Tetanus |
| | | 6 | Meningitis | |
| | | 7 | Hepatitis B (g) | |
| | | | Hepatitis C (g) | |
| | | 8 | Malaria | |

14

# Future Work

- Converting other WHO datasets
    - WHO Global InfoBase Online
    - Global Health Atlas
    - Regional Statistics
- Evolution patterns$
    - Facilitate seamless evolution of knowledge-base

*$ Christoph Rieß, Norman Heino, Sebastian Tramp (geb. Dietzold), and Sören Auer: EvoPat -- Pattern-Based Evolution and Refactoring of RDF Knowledge Bases. In: Proceedings of the 9th International Semantic Web Conference ISWC2010*

Thursday, August 26, 2010