# Re-using Cool URIs:
## Entity Reconciliation Against LOD Hubs
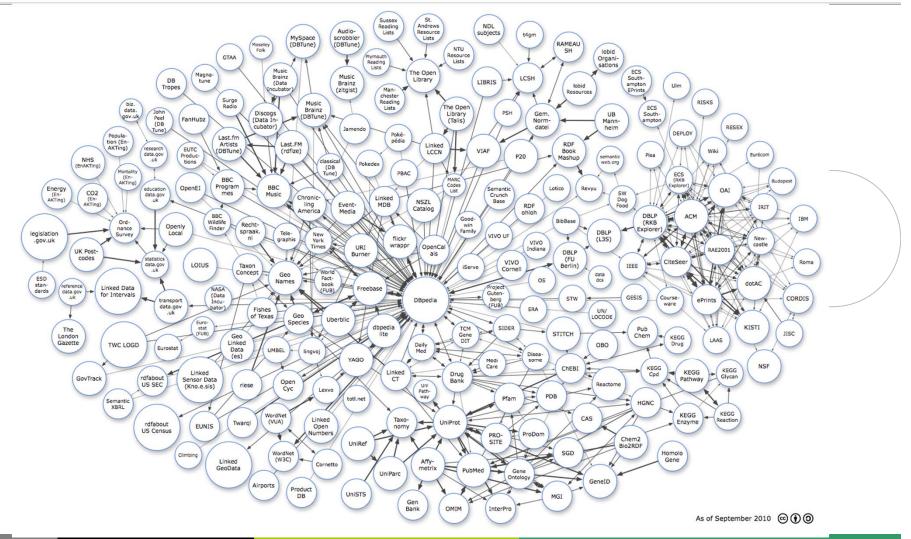
Fadi Maali, Richard Cyganiak, Vassilios Peristeras
LDOW 2011
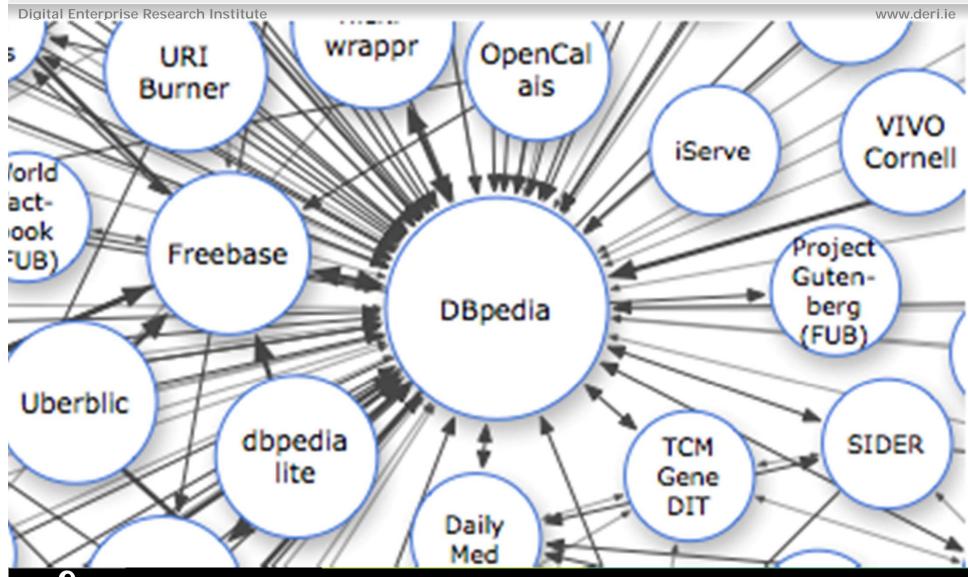
As of September 2010

# The Web of Data
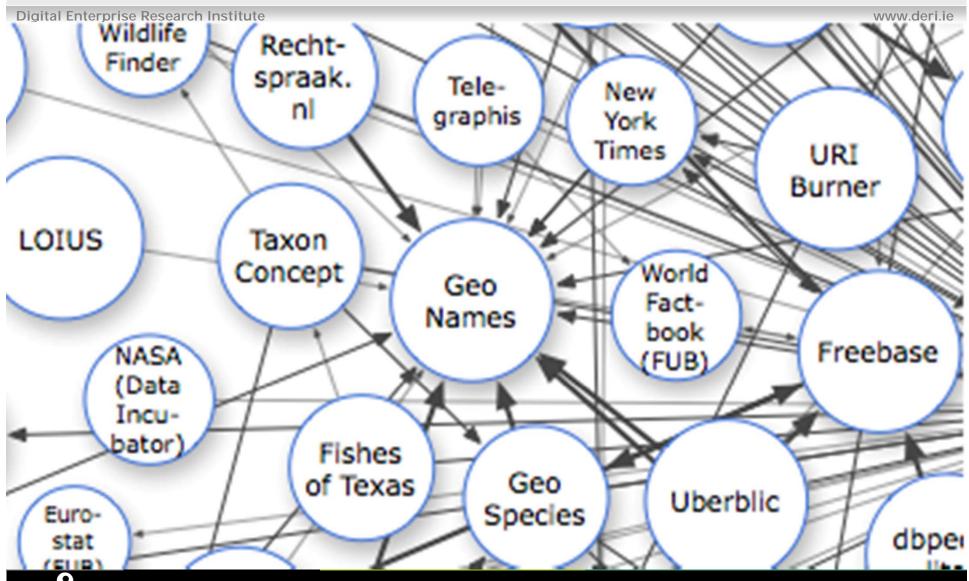
# The Web of Data

# The Web of Data

Enabling **networked** knowledge.

# The Web of Data

- LOD Hubs = datasets that attract many inlinks

- Music metadata community uses BBC Music identifiers

- UK government data community uses Ordnance Survey identifiers

- Library data community uses Library of Congress Subject Headings

Enabling **networked** knowledge.

# Standard identifiers

0  123456  789012

LOD Hubs are bar codes for a specific community.

# For example, government data

# For example, government data

# Reconciliation

| City | State | Country |
|------|-------|---------|
| Cambridge | Massachusetts | United States |

Enabling **networked** knowledge.

# Reconciliation

label=Cambridge

Cambridge Bay in
Canada

| City | State | Country |
|------|-------|---------|
| Cambridge | Massachusetts | United States |

# Reconciliation

label=Cambridge
type = City
*In the state of
Massachusetts*

Cambridge city in
Massachusetts

✔

| City | State | Country |
|------|-------|---------|
| Cambridge | Massachusetts | United States |

Enabling **networked** knowledge.

# Approaches

- SPARQL
- SPARQL + full-text search
- Silk Server
- Semantic Web search engines

- **Based on regular expressions**
- **Pros**
  - ☐ Standardised
  - ☐ Zero-effort approach
- **Cons**
  - ☐ Slow
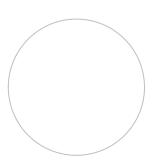  - ☐ Not good at text search
  - ☐ No ranked results

- Based on full-text extension for SPARQL

- Pros

  - More forgiving string matching

  - Ranking

  - Zero-effort (depending on your SPARQL store)

- Cons

  - Proprietary syntax

- **Pros**
  - Powerful declarative link specification
  - Variety of similarity functions
- **Cons**
  - Configuration needs to prepared
  - Silk Server needs to be deployed
  - Silk Server tightly couples its input and reference data

# Semantic Web Search Engine

- **Based on Sindice API**

- **Pros**
  - ☐ Zero-effort approach (if your dataset is indexed in Sindice)
  - ☐ Search distributed RDF datasets (e.g. FOAF profiles)

- **Cons**
  - ☐ Noisy

Enabling **networked** knowledge.

- **Data Interlinking benchmark (part of IM@OAEI2010)**
- **We reconciled DailyMed against:**
  - DBpedia SPARQL endpoint (http://dbpedia.org/sparql)
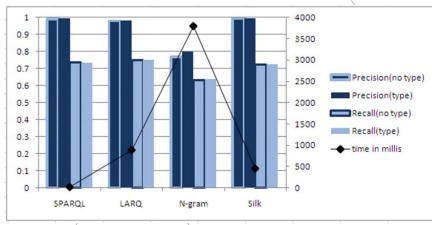  - Sider dump file (part of the benchmark)

- SPARQL with REGEX is unsuitable (performance)
- Except if labels are very consistent
- Type restrictions are very effective
- Silk has best recall (but requires custom link spec)

Services performance against DBpedia

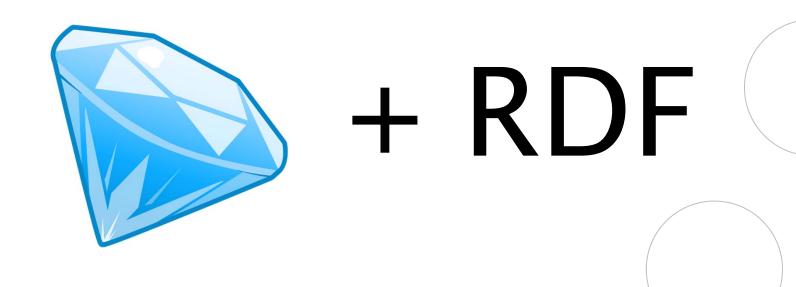Services performance against Sider RDF dump file

 + RDF

| |
| --- |
| Fadi Maali |
| Gofran Shukair |
| Souleiman Hasan |
| Richard Cyganiak |
| Michael Hausenblas |
| Manfred Hauswirth |
| Stefan Decker |
| Lukasz Porwol |
| Alexandre Passant |
| Owen |
| Maciej Dabrowski |

*List of people from DERI*

*List of people from DERI*

# Example

# Example

*Reconciliation result facets*

*Resource Preview*

- RDF Extension for Google Refine
  http://lab.linkeddata.deri.ie/2010/grefine-rdf-extension/

- Reconciliation will be in the upcoming next version

# Thanks!