



*Keith Alexander (Talis), Richard Cyganiak (DERI),  
Michael Hausenblas (DERI) and Jun Zhao (University of Oxford)*

# Describing Linked Datasets

On the Design and Usage of **void**,  
the 'Vocabulary Of Interlinked Datasets'

Linked Data Workshop at WWW09, 2009-04-20, Madrid, Spain

# Agenda

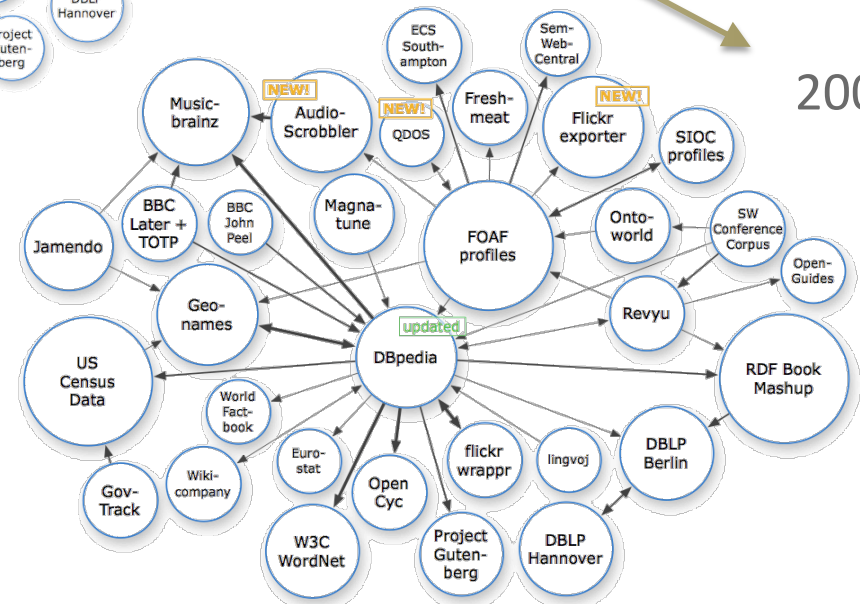
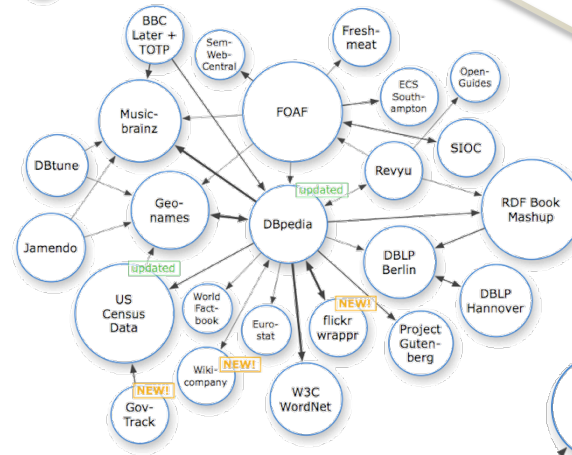
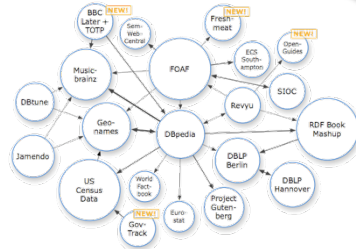
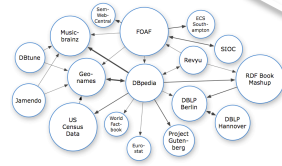


- The Problem
- Our Proposal – void
- Applications
- Next Steps



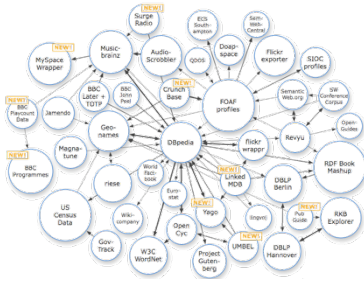
2007

# The Problem

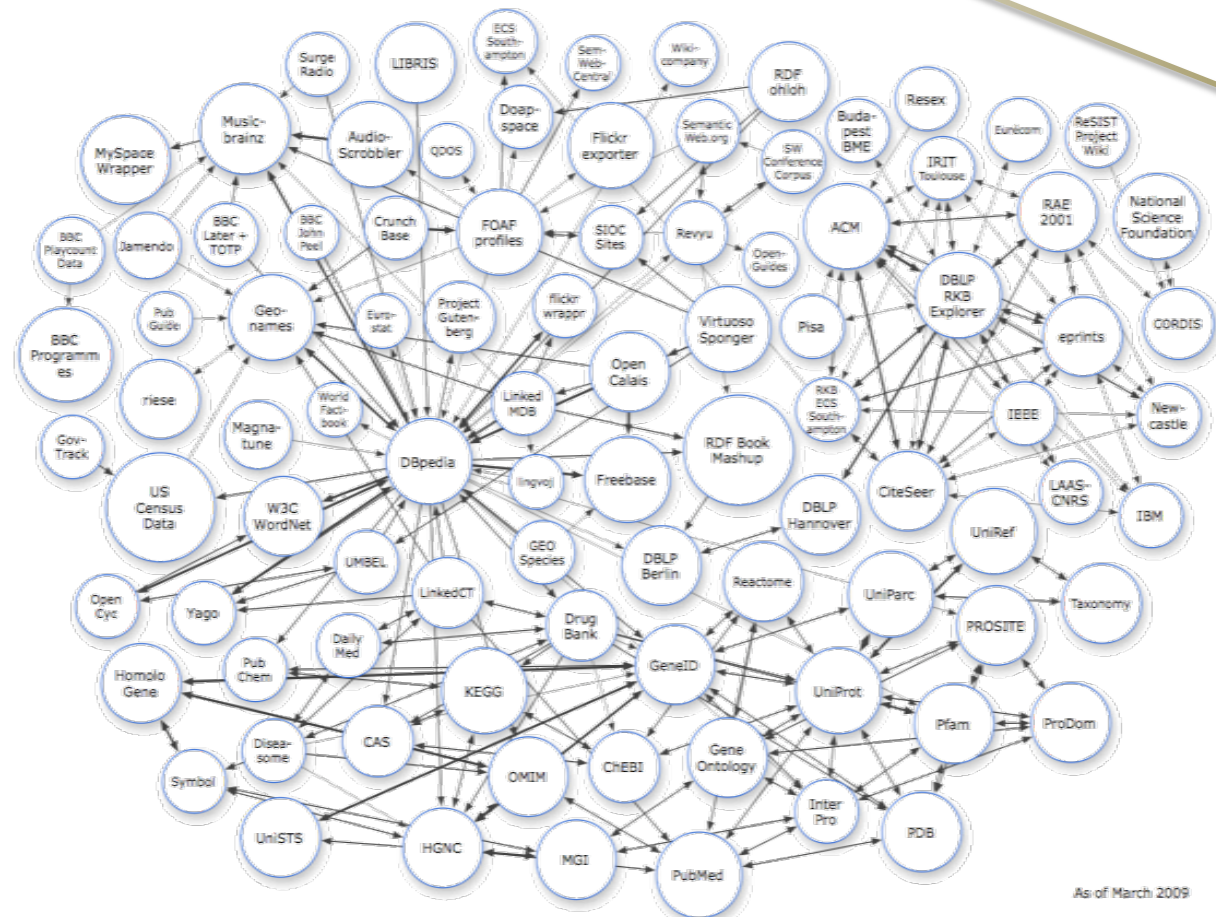


2008

**Describing Linked Datasets** – On the Design and Usage of voiD, the “Vocabulary Of Interlinked Datasets”,  
Linked Data Workshop at WWW09, 2009-04-20, Madrid, Spain



# The Problem



2009

**Describing Linked Datasets** – On the Design and Usage of *voID*, the “Vocabulary Of Interlinked Datasets”, Linked Data Workshop at WWW09, 2009-04-20, Madrid, Spain

# The Problem



- The Linking Open Data (LOD) cloud gathers currently roughly the same momentum as the **Web** in the early **1990s**
- How did people deal with the consequences of having a decentralized system, back then?

# The Problem



# The Problem



- From 2007 on, we have been doing it in the *Yahoo!-catalog-style*: **manually collecting** and **representing** data about the Linking Open Data cloud:
  - In the LOD cloud diagram, we give a qualitative view in form of a visual graph
  - In various ESW Wiki pages we create HTML tables:
    - <http://esw.w3.org/topic/TaskForces/CommunityProjects/LinkingOpenData/DataSets/Statistics>
    - <http://esw.w3.org/topic/TaskForces/CommunityProjects/LinkingOpenData/DataSets/LinkStatistics>



# The Problem



<http://esw.w3.org/topic/TaskForces/CommunityProjects/LinkingOpenData/DataSets/LinkStatistics>

## Statistics on links between Data sets

This page collects statistics on links between [Data sets that are available as Linked Data](#).

The link statistics are also used to update the LOD cloud diagram for the main page of the project. So if you publish a dataset yourself or if you know detailed statistics for datasets that you use, please add them to the table and we will include them in the next revision of the LOD cloud.

### Links between Data sets

- Outdated [CSV version of this table](#)
- Note that **Link count (range)** originally had 3 valid values – “> 100”, “> 1,000”, “> 100,000” – based on **Link count (actual)**.
  - Now has 5 levels – “> 100”, “> 1,000”, “> 10,000”, “> 100,000”, “> 1,000,000” – being updated (and corrected) as **Link count (actual)** is updated. (Linkages fewer than 100 are not counted here nor shown in the Cloud diagrams.)
- Also note – some links (e.g., that for RAE 2001) have been percent-escaped to get around spam-blocks in the ESW Wiki configuration. De-escaping these should allow loading, if they fail when simply clicked in this table.
- re “Type” column – moment, so counts Equivalent/Broader

## Statistics on Data sets

This page collects statistics on Data sets that are available as Linked Data.

The statistics are also used to update the LOD cloud diagram for the main page of the project. So if you publish a dataset yourself or if you know detailed statistics for datasets that you use, please add them to the table and we will include them in the next revision of the LOD cloud.

### Data set sizes

- “Wrapper” denotes any “data set” which is not available as an RDF Dump, for which size cannot be accurately assessed, or for which triples are dynamically produced.
- Note – some links (e.g., those for RAE 2001 and “R e s e x”) have been percent-escaped to get around spam-blocks in the ESW Wiki configuration. De-escaping these should allow loading, if they fail when simply clicked in this table.

Data set	Size of the data set (number of triples)	Wrapper?	endpoint?	RDF dump?
<a href="#">ACM (RKB)</a>	12,644,652	N	Y	Y
<a href="#">AudioScrobbler</a>	600,000,000	Y	Y	N
<a href="#">BBC John Peel</a>	277,000	N	Y	Y
<a href="#">BBC Later + TOTP (link not responding - 2009-04-01)</a>	10,000	N	Y	Y
<a href="#">BBC Playcount Data</a>	10,000	N	Y	Y
<a href="#">BBC Programmes</a>	10,000,000	N	Y	Y
<a href="#">Budapest BME (RKB)</a>	42,064	N	Y	Y
<a href="#">CAS</a>	100,000	N	Y	Y
<a href="#">ChEBI</a>	510,866	N	Y	Y
<a href="#">CiteSeer (RKB)</a>	8,294,523	N	Y	Y
<a href="#">CrunchBase</a>	955,676	Y	Y	N
<a href="#">Daily Med</a>	170,000	N	Y	Y

<http://esw.w3.org/topic/TaskForces/CommunityProjects/LinkingOpenData/DataSets/Statistics>

**Describing Linked Datasets** – On the Design and Usage of void, the “Vocabulary Of Interlinked Datasets”,  
Linked Data Workshop at WWW09, 2009-04-20, Madrid, Spain



# The Problem



- Currently, only **human comprehensible descriptions** (the LOD cloud, Wiki pages) available
- **We can't automate tasks**, such as
  - Efficient & effective search
  - Selection of dataset (for apps, interlinking targets)
  - Generation of maps, etc.

# The Problem



- We **can't apply our tools** and methods we have experiences with, such as editors, engines, stores, etc.
- Even worse, it **doesn't scale**
  - We'd need a Google-style approach that scales like hell and is powerful enough to enable the above mentioned
  - Providing **metadata** about the **LOD cloud** in a **machine-comprehensible** way

# Agenda



## ✓ The Problem

- Our Proposal – voidD
- Applications
- Next Steps

## Our Proposal - void



- Solution: providing a formal description of
  - What a dataset is about (topic, technical details)
  - How and under which conditions to access it
  - How the dataset is interlinked with other datasets
    - Qualitative level: type of interlinking
    - Quantitative level: number of links, resources, etc.
  - How to discover the metadata
- **void**, the “Vocabulary of Interlinked Datasets” provides precisely this

## Our Proposal - void

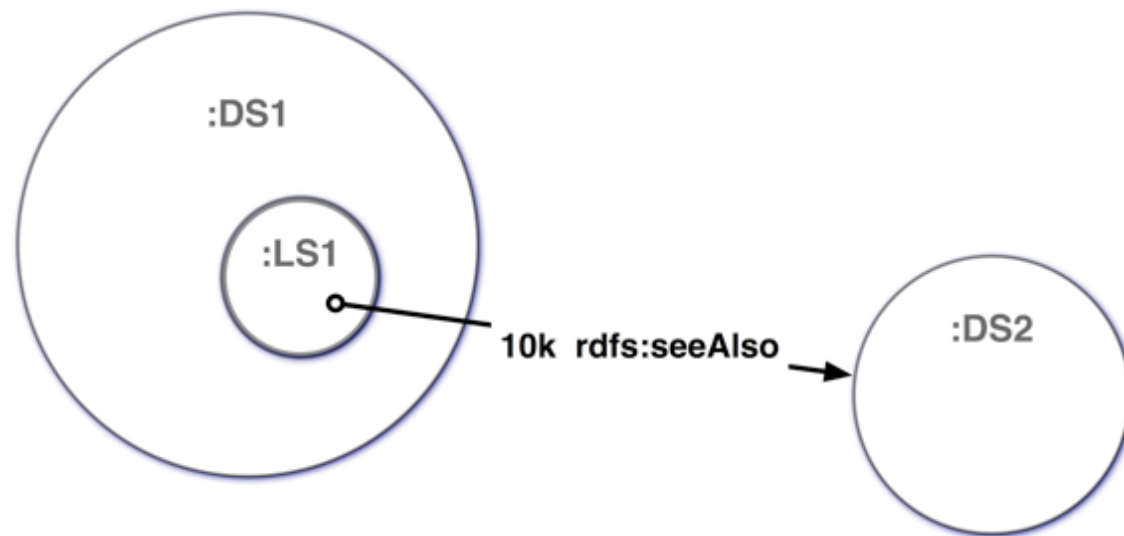


- A **dataset** is a set of RDF triples that are published, maintained or aggregated by a single provider.
- A **dataset** is **authoritative** with respect to a certain URI namespace if it contains information about resources named by URIs in this namespace, and is **published** by the **URI owner** (→ [URI ownership as of the AWWW1](#))

## Our Proposal - voidD

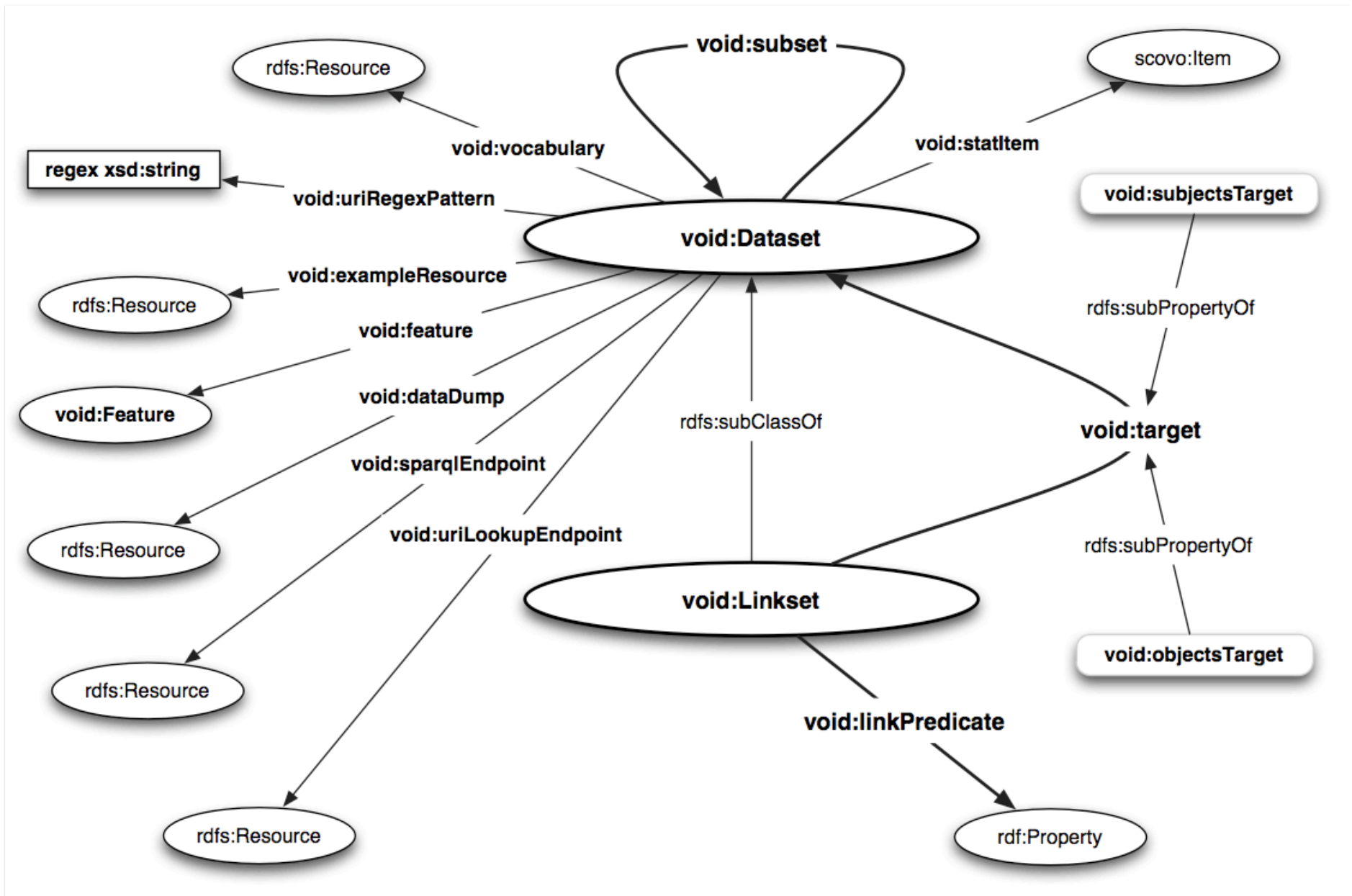


- A **linkset**  $LS$  is a set of RDF triples where for all triples  $t_i = \langle s_i, p_i, o_i \rangle \in LS$ , the subject is in one dataset, i.e. all  $s_i$  are described in  $DS_1$ , and the object is in another dataset, i.e. all  $o_i$  are described in  $DS_2$ .





# Our Proposal - voidD





## Our Proposal - void

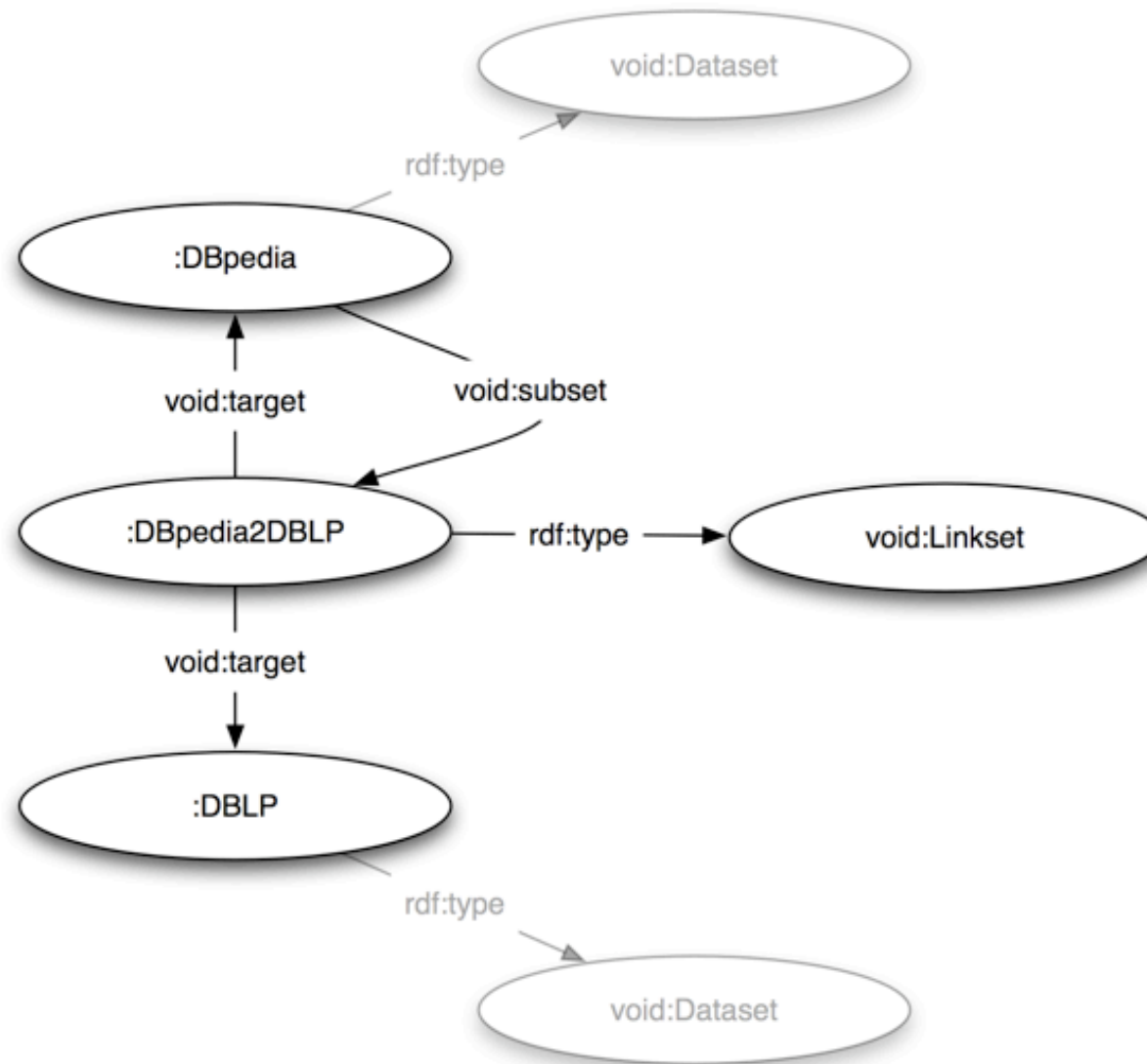


void offers two orthogonal interlinking types:

- **classic LOD** vs. **3rd-party**, differing in where the interlinking statements are kept. In the first case the interlinking triples, i.e. a linkset, are hosted in one of the two involved datasets, while in the latter case there is a third dataset involved that contains the interlinking triples, i.e. the linkset;
- **non-directed** vs. **directed**, which addresses the issue if someone is interested in stating the direction of the interlinking or not (for example with owl:sameAs)

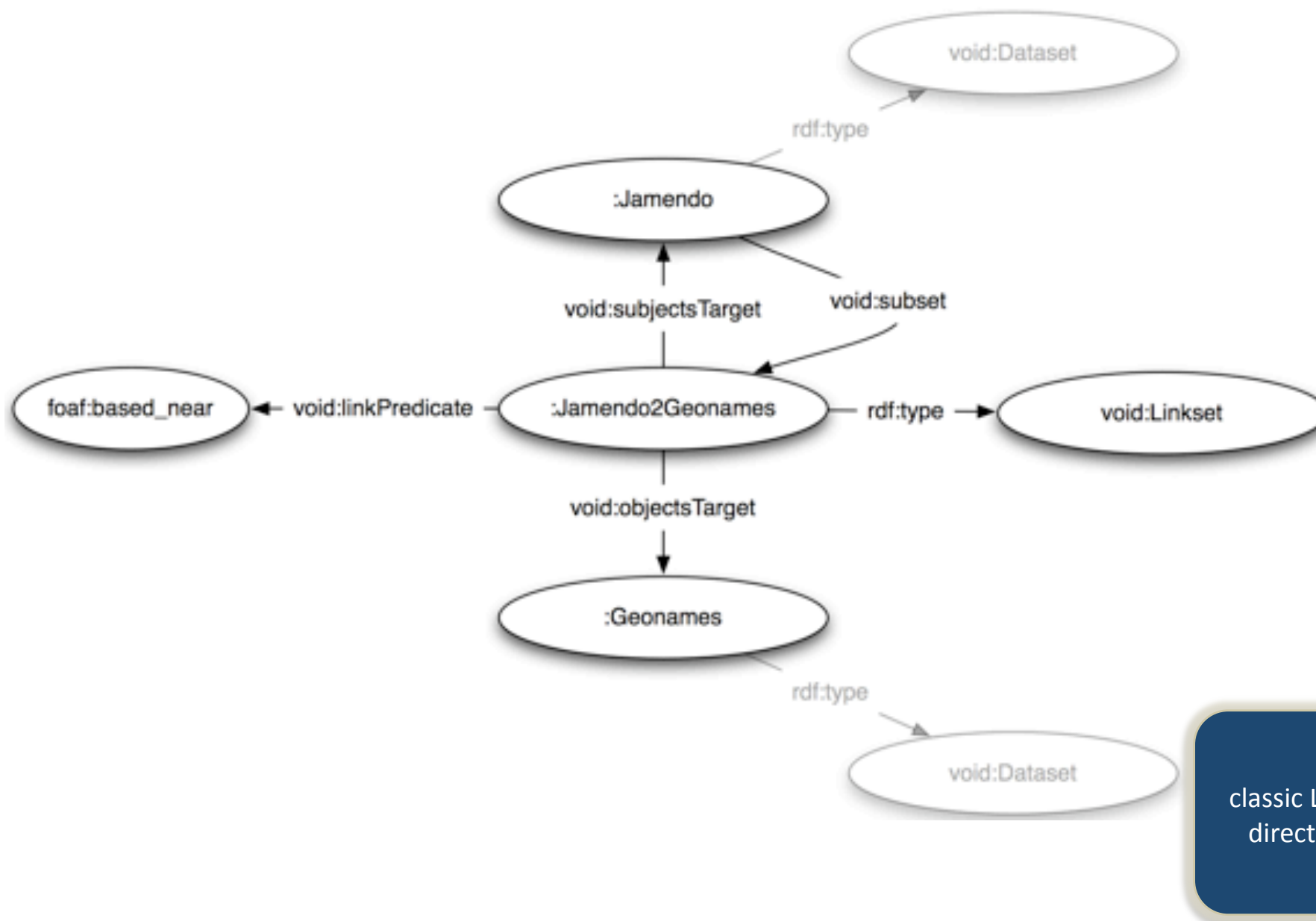


# Our Proposal - voidD

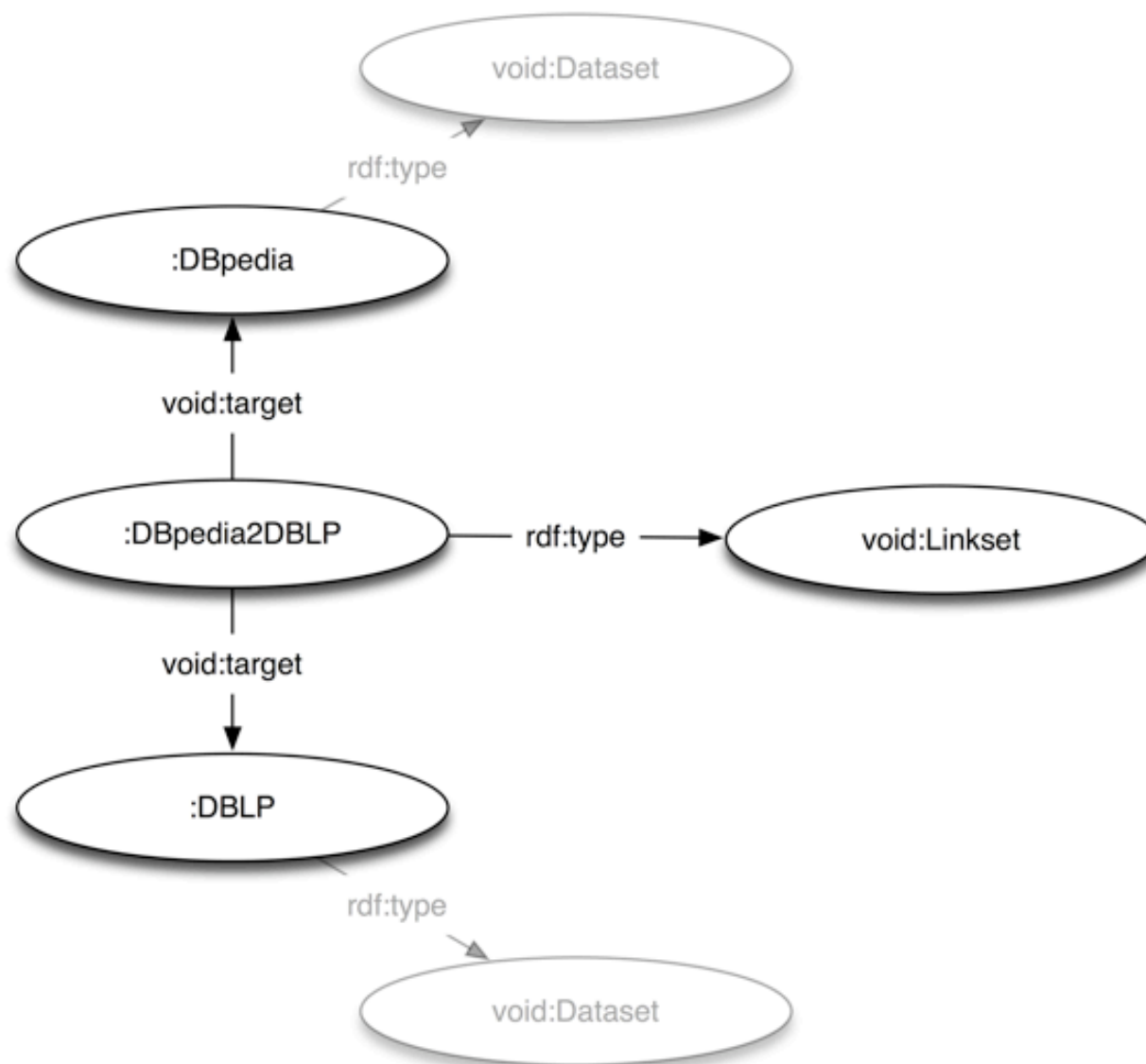


classic LOD,  
non-directed

# Our Proposal - voidD

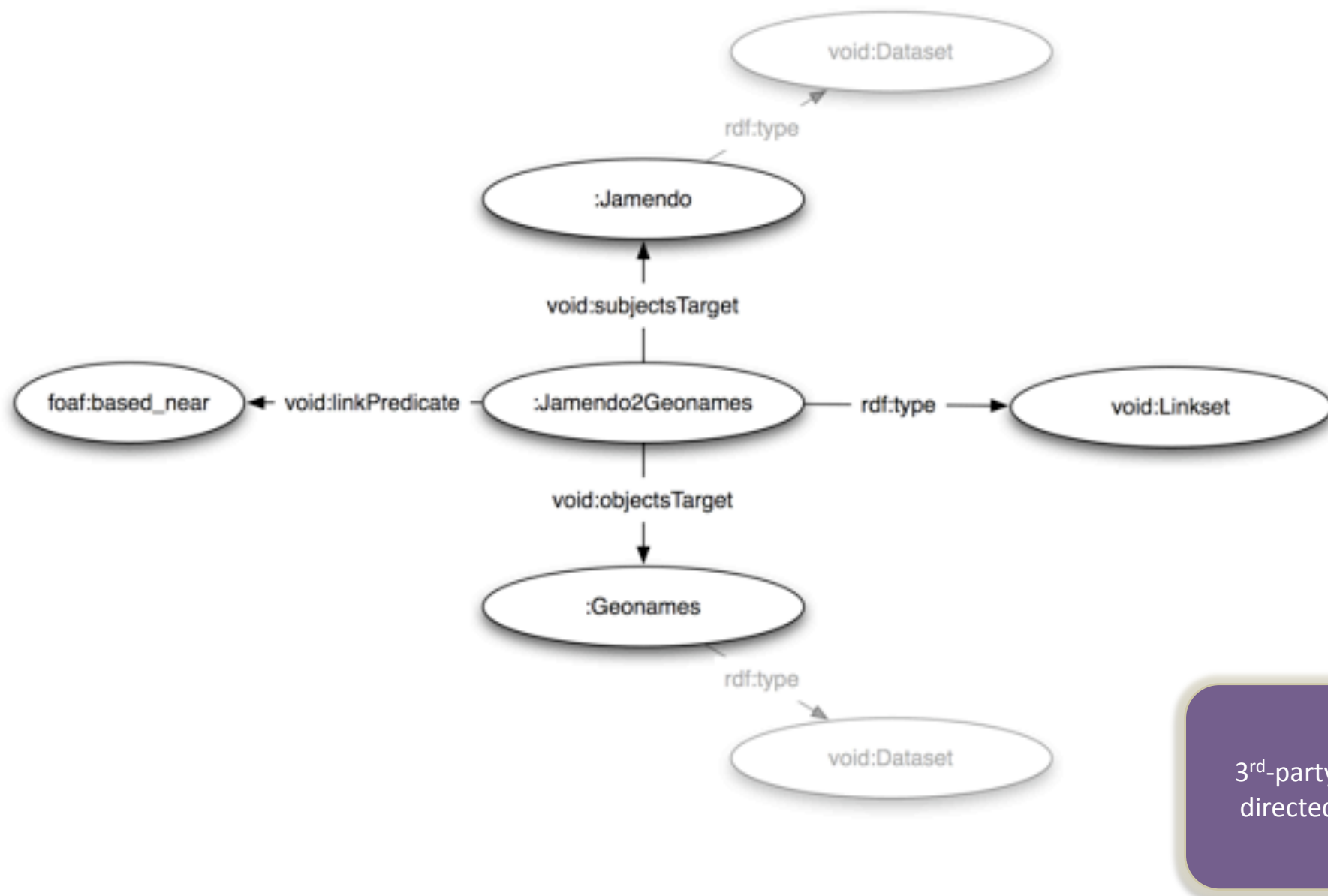


# Our Proposal - voidD



3<sup>rd</sup>-party,  
non-directed

# Our Proposal - voidD



3<sup>rd</sup>-party,  
directed



## Our Proposal - void

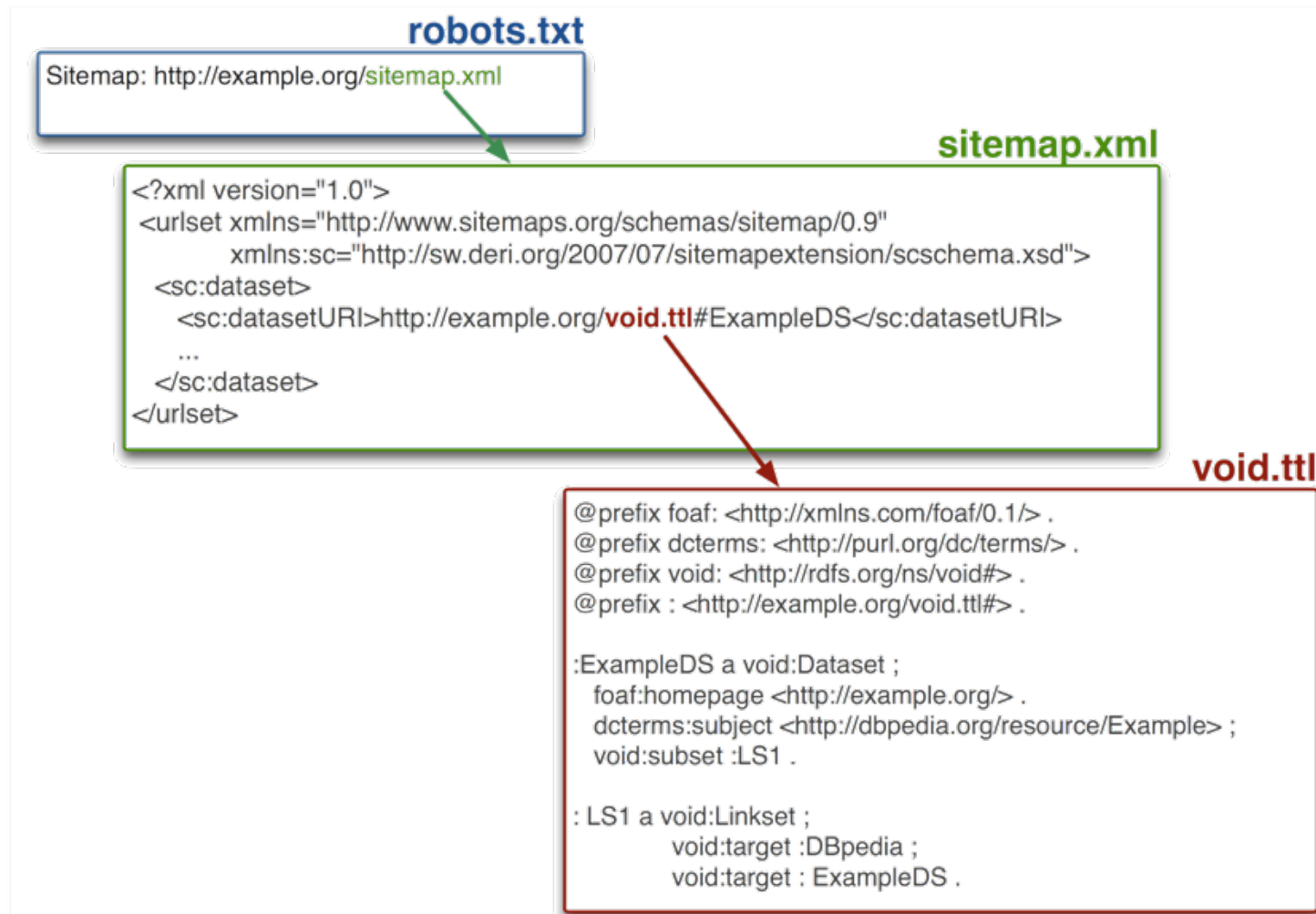


- Reusing terms from other vocabularies
  - foaf:homepage/IFP
  - dcterms:subject along with DBpedia URIs  
[http://dbpedia.org/resource/ XXX](http://dbpedia.org/resource/XXX)
  - [SCOVO](#) for statistics about triples, links, etc



## Our Proposal - voidD

- Publication & discovery via sitemaps and/or backlinks (dcterms:isPartOf)



## Our Proposal - void



- Once dataset providers have published their void description in RDF along with their dataset, one can address the following issues:
  - How to find **some** datasets?
  - How to **efficiently** find a specific dataset?
  - How to **effectively** find datasets?
  - How to **dynamically** select datasets?
  - How to select datasets based on certain **preferences**?

# Agenda



- ✓ The Problem
- ✓ Our Proposal – void
- Applications
- Next Steps

# Applications



- Generation (ve, liftSSM, NX parser)
- Vocabulary Management (Talis)
- Explorer (RKB, LDE)
- Query Federation (Clarck-Parsia, OpenLink)
- Dataset ranking (→ DING! talk)
- Potential Applications
  - Map of data (Sindice)
  - Dynamic Meshups for Application

# Applications



ve - the void editor

http://localhost:8888/ve/

ve - the void editor

## ve - the void editor

"... vi was yesterday"

[reset](#) [help](#)

Dataset Description

URI

Subject(s)

Found <<http://dbpedia.org/resource/Cars>> as a subject [[use this](#)].

There are aliases that might have more information:

- <http://dbpedia.org/resource/Automobile> [[use this](#)]

Interlinking

Following seed datasets are available:

- [DBpedia. A community effort to extract structured information from Wikipedia](#) [[use this](#)]

generate voidID

Output

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
@prefix dcterms: <http://purl.org/dc/terms/> .
@prefix void: <http://rdfs.org/ns/void#> .

:NEWDS rdf:type void:Dataset ;
        foaf:homepage <http://example.org/mycars> .
        dcterms:subject <http://dbpedia.org/resource/Cars> ;
        void:subset :LS1 .

:DBpedia rdf:type void:Dataset ;
        foaf:homepage <http://dbpedia.org/> .

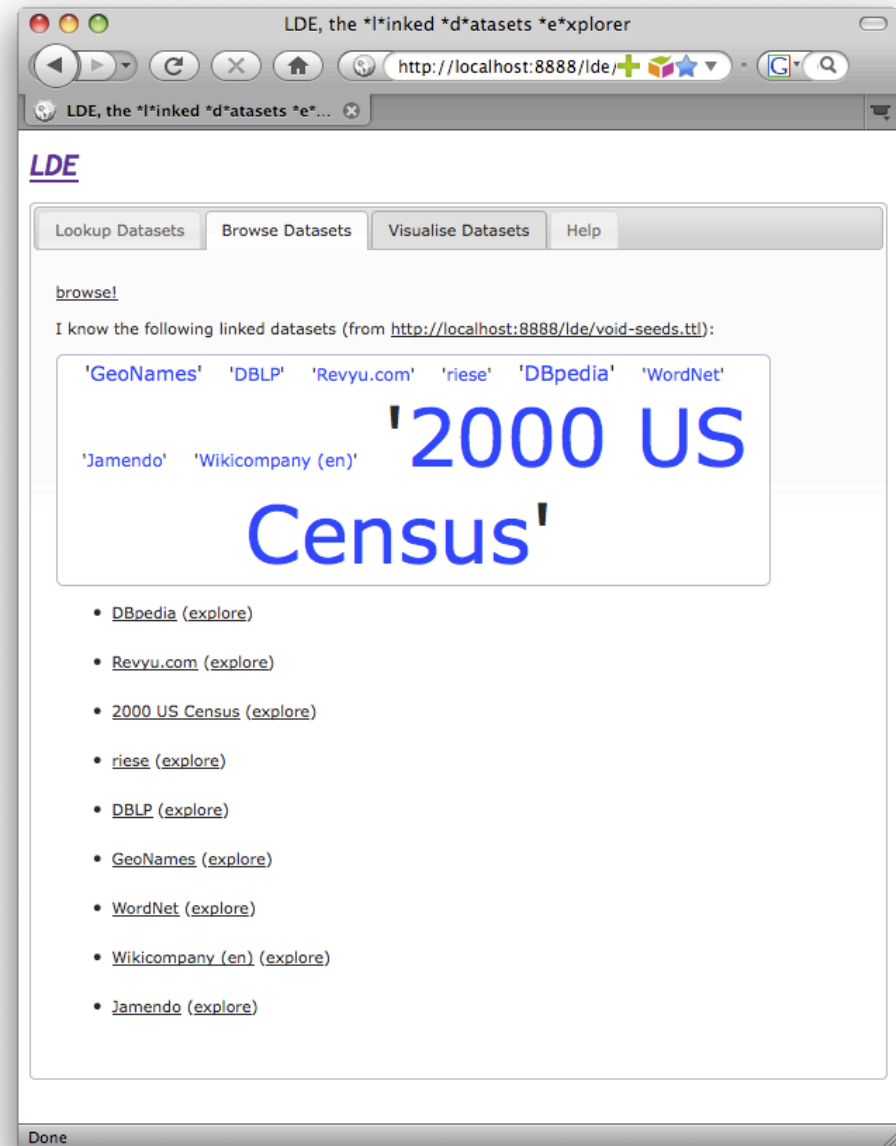
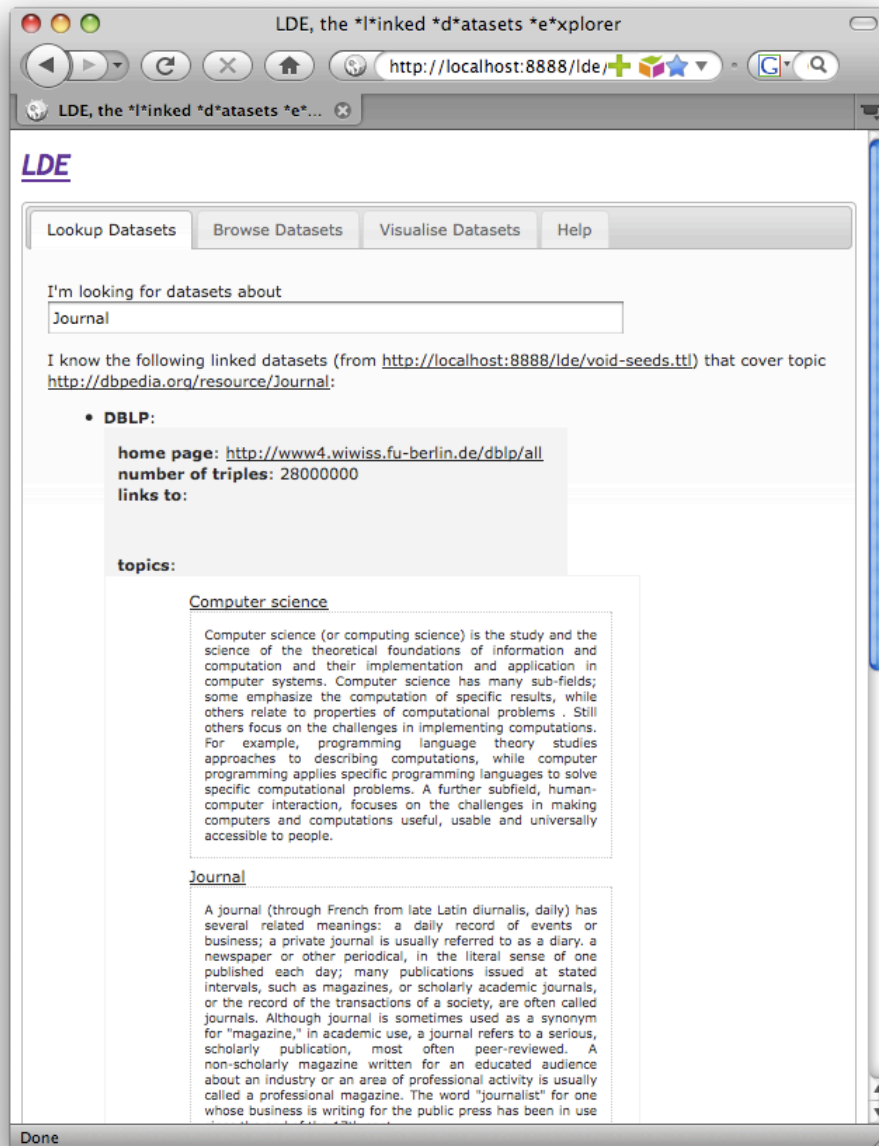
:LS1 rdf:type void:Linkset ;
      void:target :NEWDS ;
      void:target :DBpedia .
```

Done

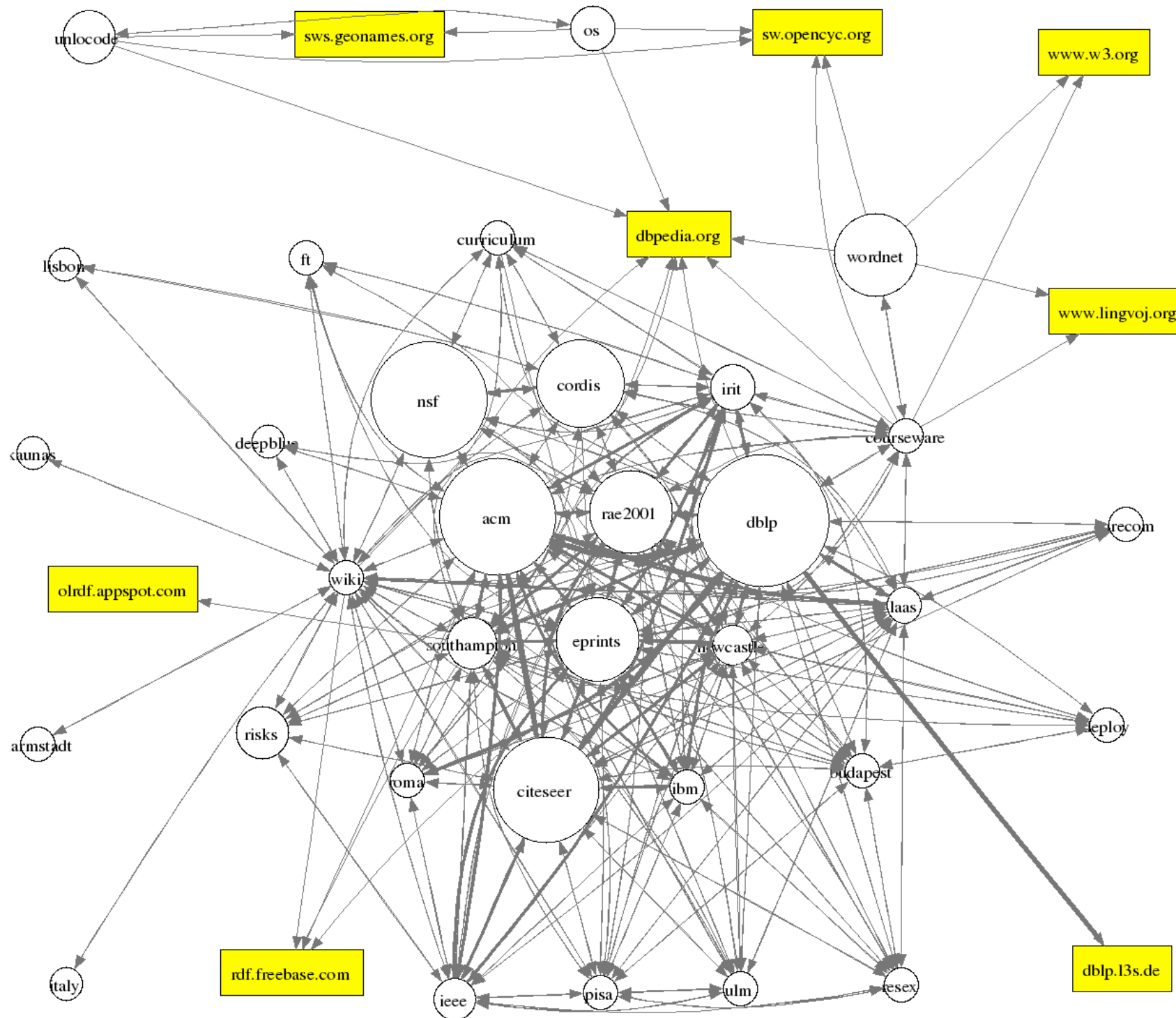
<http://ld2sd.deri.org/ve>



# Applications



<http://ld2sd.deri.org/ldc>



<http://dblp.rkbexplorer.com/models/void.ttl>

# Applications



## About: <http://twitter.com/mhausenblas#Dataset>



An Entity in Data Space: [linkeddata.uriburner.com](http://linkeddata.uriburner.com)

Property	Value
<a href="#">void:sparqlEndpoint</a>	▪ <a href="http://linkeddata.uriburner.com/sparql">http://linkeddata.uriburner.com/sparql</a>
<a href="#">void:statItem</a>	▪ <a href="http://twitter.com/mhausenblas#Stat">http://twitter.com/mhausenblas#Stat</a> ▪ <a href="http://twitter.com/mhausenblas#PersonStat">http://twitter.com/mhausenblas#PersonStat</a> ▪ <a href="http://twitter.com/mhausenblas#BoardPostStat">http://twitter.com/mhausenblas#BoardPostStat</a> ▪ <a href="http://twitter.com/mhausenblas#DataSourceStat">http://twitter.com/mhausenblas#DataSourceStat</a> ▪ <a href="http://twitter.com/mhausenblas#ContainerStat">http://twitter.com/mhausenblas#ContainerStat</a> ▪ »more«
<a href="#">rdf:type</a>	▪ <a href="#">void:Dataset</a>
<a href="#">rdfs:seeAlso</a>	▪ <a href="http://twitter.com/mhausenblas">http://twitter.com/mhausenblas</a>

Explore using: [OpenLink Data Explorer](#) | [Zitgist Data Viewer](#) | [Marbles](#) | [DISCO](#) | [Tabulator](#) Raw Data in: [N3](#) | [RDF/XML](#) [About](#)



This work is licensed under a [Creative Commons Attribution-Share Alike 3.0 Unported License](#).

<http://linkeddata.uriburner.com/>

# Agenda



- ✓ The Problem
- ✓ Our Proposal – void
- ✓ Applications
- **Next Steps**

## Next Steps



- void 2.0 see issues at <http://code.google.com/p/void-impl/issues/list>
- statistics module (fix/extend re SCOVO)
- SPARQL endpoints
- provenance, trust (?)
- Assist people in publishing void