

Linkitetyn Datan automaattinen analysointi ja dokumentointipalvelu

vocab.at

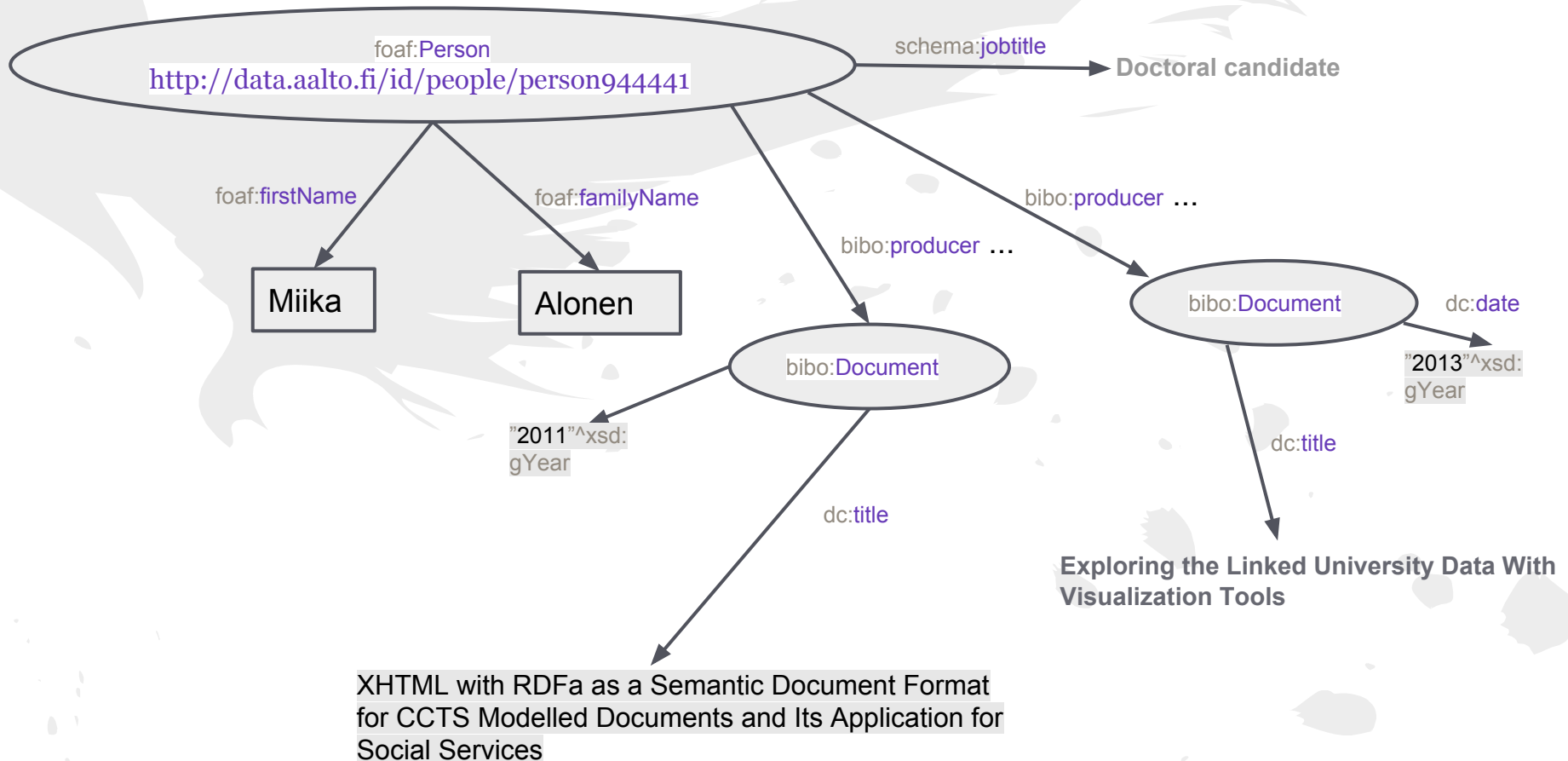


Miika Alonen

24.1.2014



About me



Motivaatio

Quality of Linked Data

2289 datasets

62,224,812,703 triples from **870 datasets**:

- 1,563,527,469 triples from **717 dumps**
- 60,661,285,234 from **156 datasets via SPARQL**

Problems with **1416 datasets = 61.9%**

(1061 dumps having errors, 355 SPARQL endpoints with errors)

Status from:
<http://datahub.io/>

at 27.9.2013

analyzed by <http://stats.lod2.eu/>

What's Wrong with Linked Data?

“The last issue of Semantic Web journal received 27 submissions of Linked Data publications and half of them had to be rejected due to issues with dataset quality or usefulness”

<http://blog.semantic-web.at/2012/08/09/whats-wrong-with-linked-data/>

- Osaamisen puute?
 - > 30 linkitettyyn dataan liittyvää suositusta
- Ei laadunhallintaan liittyviä W3C suosituksia
 - > Rakenteen validointi riittämätöntä
- Yksi suuri ongelma on laatupoikkeamien löytämisessä suuresta määrästä dataa

Linkitetyn datan sanastot

- Käytettävissä oleva sanasto määritellään RDF-muodossa ja dokumentoidaan ihmisille HTML-sivustoina
- Datassa käytetty sanasto ei kuitenkaan aina vastaa mitään yksittäistä sanastoa
- Objektiivisen dokumentaation puute johtaa helposti väärinymmärryksiin ja vääriin johtopäätöksiin tietosisällöstä

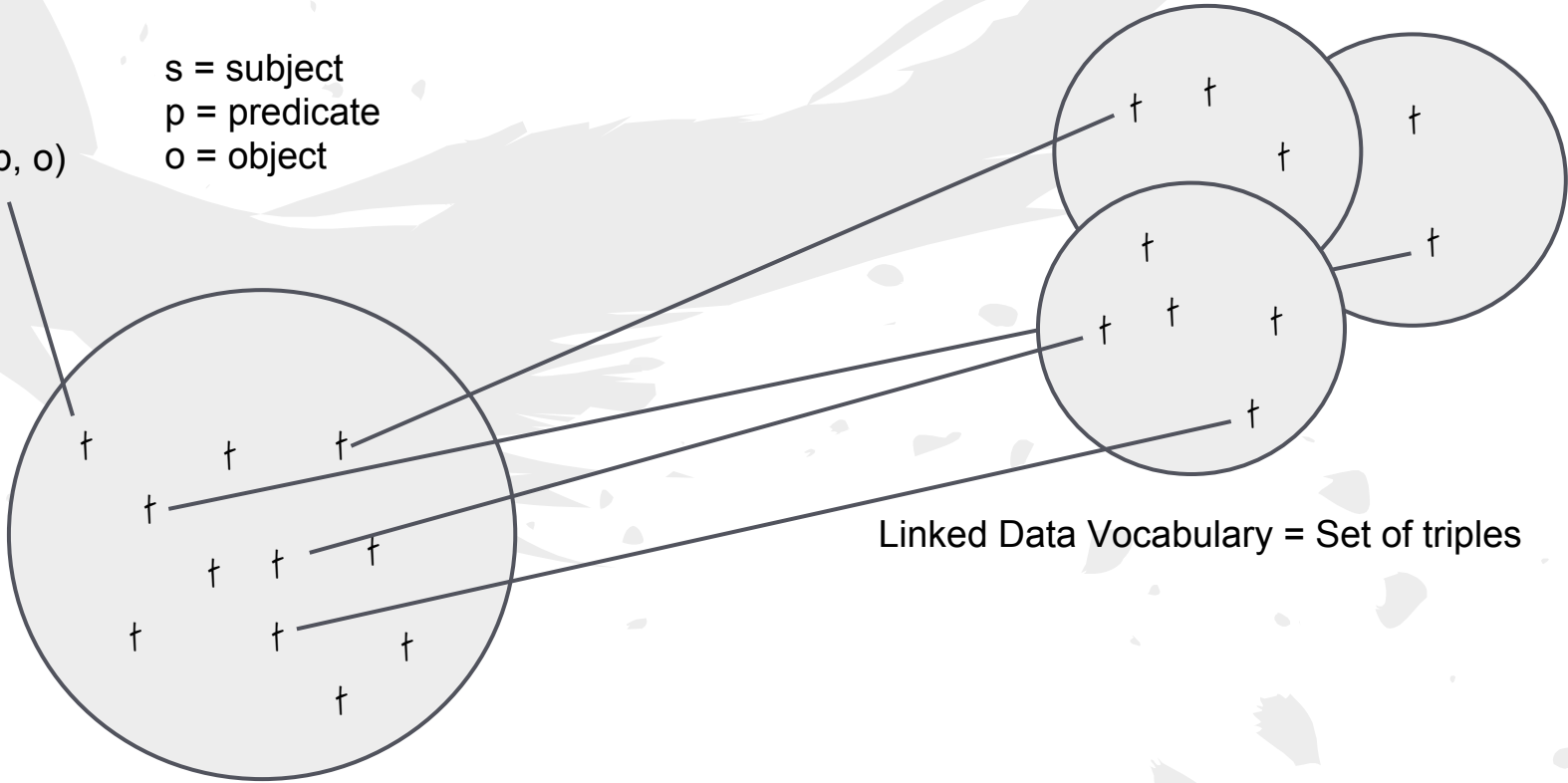
MIKSI?

Linkitetty data on koneluettavaa ja se voidaan analysoida ja dokumentoida automaattisesti!

Linkitetty data

s = subject
p = predicate
o = object

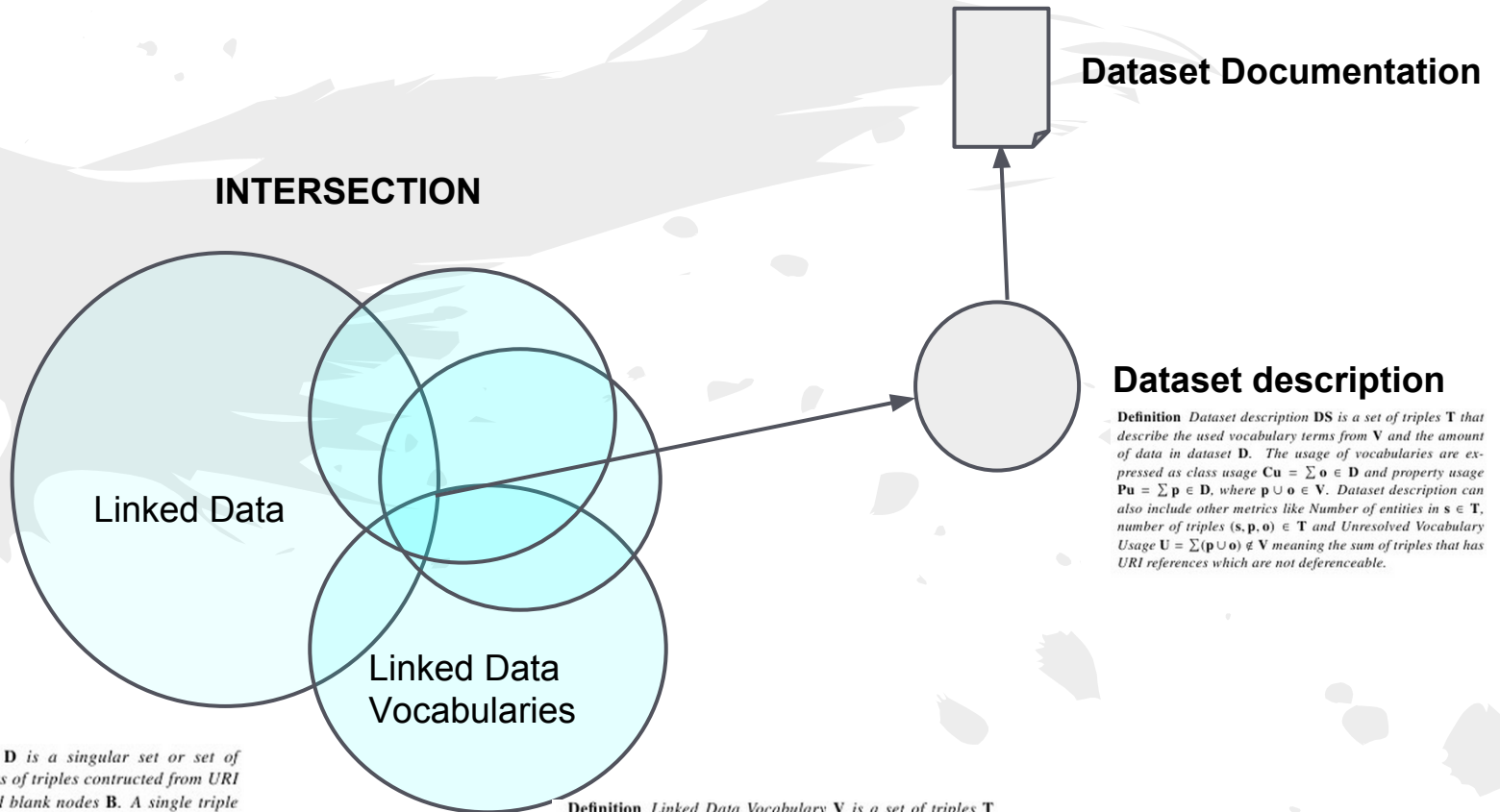
(s, p, o)



Linked Data Vocabulary = Set of triples

Linked Data = Set of triples

Linkitetyn datan analysointi ja dokumentointi



INTERSECTION

Linked Data

Linked Data Vocabularies

Dataset Documentation

Dataset description

Definition Dataset description DS is a set of triples T that describe the used vocabulary terms from V and the amount of data in dataset D . The usage of vocabularies are expressed as class usage $C_u = \sum o \in D$ and property usage $P_u = \sum p \in D$, where $p \cup o \in V$. Dataset description can also include other metrics like Number of entities in $s \in T$, number of triples $(s, p, o) \in T$ and Unresolved Vocabulary Usage $U = \sum (p \cup o) \notin V$ meaning the sum of triples that has URI references which are not dereferenceable.

Definition Linked Dataset D is a singular set or set of graphs G , which contain sets of triples constructed from URI references U , Literals L and blank nodes B . A single triple T is a construct that can be defined as $T = [s, p, o]$, where $s \in U \cup B$, $p \in U$ and $o \in U \cup B \cup L$.

Definition Linked Data Vocabulary V is a set of triples T that describe the classes C and properties P which are identified with URIs meaning $C \cup P \in U$. Any dataset D may use class $o \in C$ or property $p \in P$ from the vocabulary V to describe a resource.

http://vocab.at

Linkitetyn datan automaattinen dokumentointi ja analysointipalvelu

Example 1: Dataset description

```
1 @prefix dct: <http://purl.org/dc/terms/> .
2 @prefix prov: <http://www.w3.org/ns/prov#> .
3 @prefix void: <http://rdfs.org/ns/void#> .
4
5 <http://vocab.at/id/3Ivr> a <prov:Activity> ;
6 <dct:source> <http://vocab.at/info> ;
7 <prov:startedAtTime> "2013-09-24T08:00:46.122Z"^^
  xsd:dateTime ;
8 <prov:endedAtTime> "2013-09-24T08:00:48.249Z"^^xsd
  :dateTime ;
9 <prov:generated>
10 [ a <void:Dataset> ;
11   <dct:source> <http://vocab.at/info> ;
12   <void:classPartition>
13   [ <void:class> foaf:Document ;
14     <void:distinctSubjects> 1 ;
15     <void:propertyPartition>
16     [ <void:entities> 1 ;
17       <void:property> dct:subject ;
18       <void:triples> 4
19     ] ;
20   ];
21 <void:vocabulary>
22 <http://purl.org/dc/elements/1.1/> ,
23 <http://xmlns.com/foaf/0.1/> ;
24 # continues, see http://vocab.at/data/300v
25 ];
```

foaf:Person foaf:Person

Identifier: <http://xmlns.com/foaf0.1/Person>

Namespace: <http://xmlns.com/foaf0.1/>

This class is used by this dataset. Follow the identifier to the original specification. These are the used properties.

Properties at a glance

[foaf:topic_interest](#) [schema:jobtitle](#) [foaf:name](#) [org:memberOf](#) [foaf:workInfoHomepage](#) [foaf:familyName](#) [foaf:firstName](#) [foaf:plan](#) [foaf:skypeID](#) [foaf:jabberID](#)

Example query

Notice that this query is generated automatically from all of the properties used by foaf:Person-instances and some of the properties may be used optionally.

```
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX schema: <http://schema.org/>
PREFIX org: <http://www.w3.org/ns/org#>
SELECT ?topic_interest ?jobtitle ?name ?memberOf ?workInfoHomepage ?familyName ?firstName ?plan ?skypeID ?jabberID WHERE {
  ?s a foaf:Person .
  ?s foaf:topic_interest ?topic_interest .
  ?s schema:jobtitle ?jobtitle .
  ?s foaf:name ?name .
  ?s org:memberOf ?memberOf .
  ?s foaf:workInfoHomepage ?workInfoHomepage .
  ?s foaf:familyName ?familyName .
  ?s foaf:firstName ?firstName .
  ?s foaf:plan ?plan .
  ?s foaf:skypeID ?skypeID .
  ?s foaf:jabberID ?jabberID .
} LIMIT 1
```

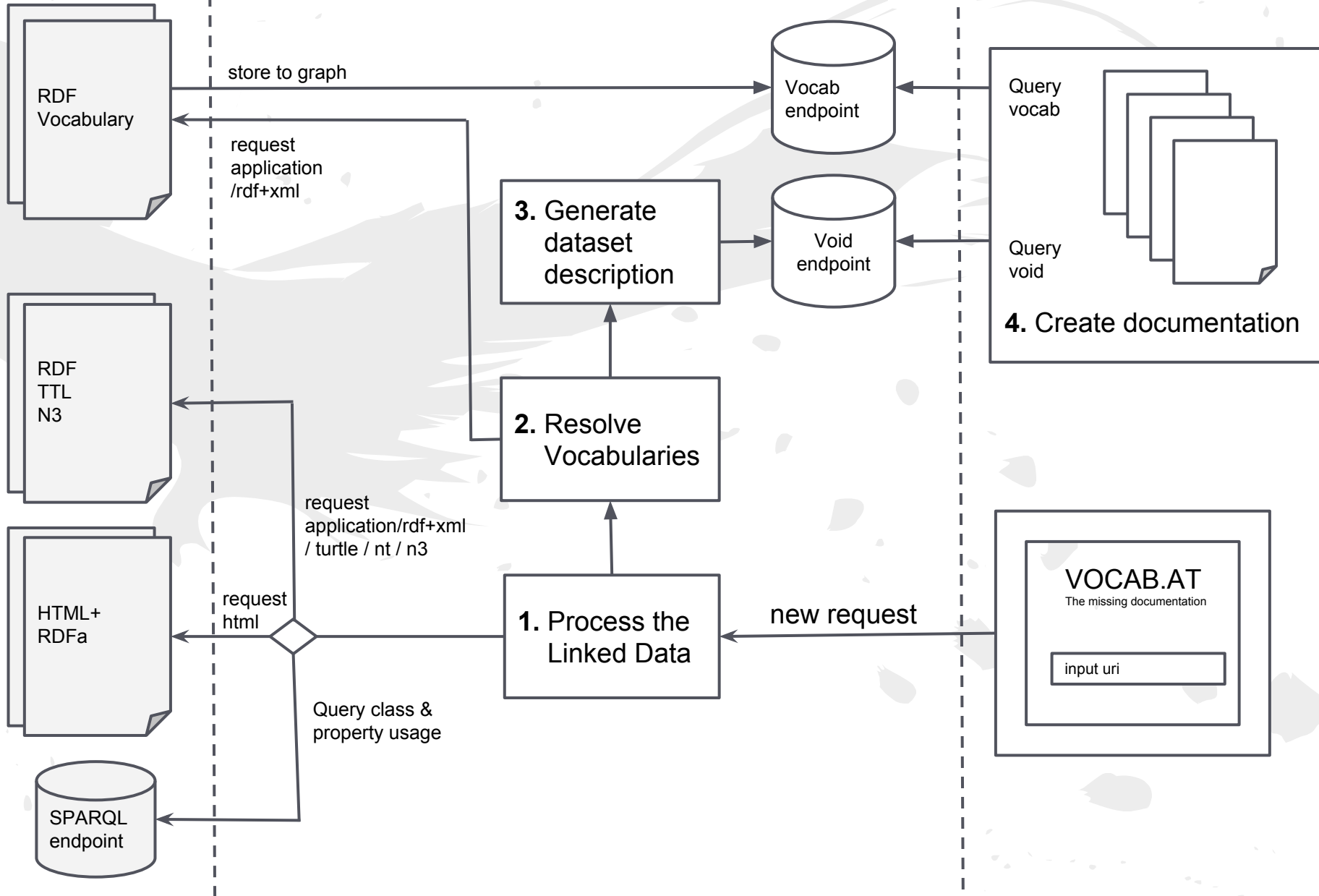
[\[back to class list\]](#)

Esim: <http://vocab.at?uri=http://vocab.at/info>

RDF out in the wild

VOCAB.AT services

VOCAB.AT interface



Preliminary Case Study

Table 1
Dataset Vocabulary Assesment

Dataset	Score	Issues
EARTH Dataset ¹⁰	1	0
Linked Amazon Dataset ¹¹	0.004 1	3
OGOLOD Dataset ¹²	0.269	40
LOD EUScreen Dataset ¹³	0.96	3
Kirjasampo Dataset ¹⁴	0.173	13
AEMET Dataset ¹⁵	0.556	1
TourMISLOD Dataset ¹⁶	0.17	2

Definition Vocabulary score $V = 1 - \frac{UT}{T}$ where UT is sum of properties and classes not defined in any vocabulary, $UT = \sum P \cup C \notin V$ and T is total number of triples in the dataset, $T = \sum (s, p, o) \in D$.



KIITOS!