



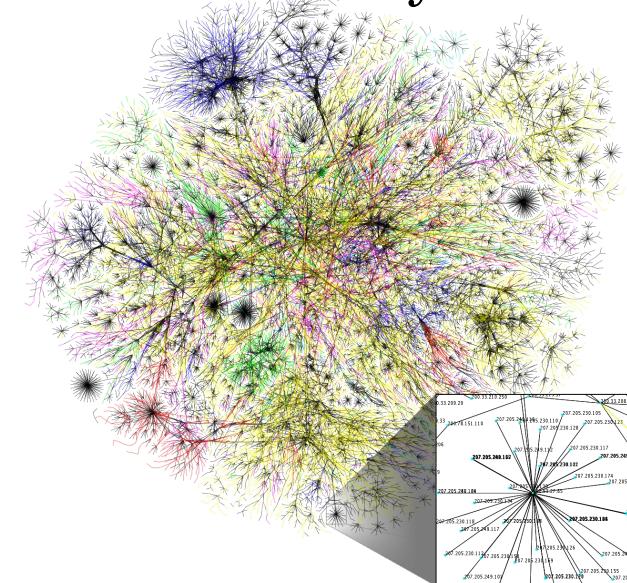
Linked Data

Craig Knoblock

DSCI 558: Building Knowledge Graphs

The Internet

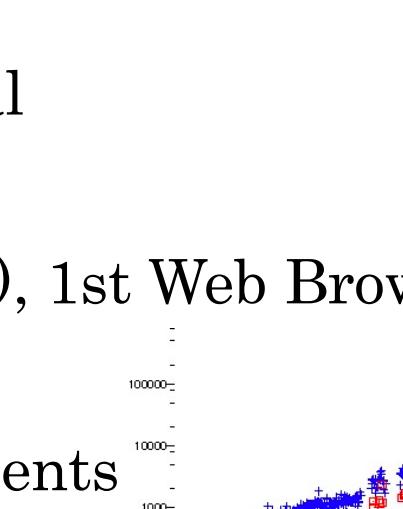
- 1962, Rand Corporation, communications systems for military
- 1965, Packet switching at the National Physical Laboratory
- 1969, ARPANET
- 1971, First Email
- 1981, CSNET & BITNET are born
- 1983, DNS is born
- 1983, TCP/IP is born
- 1985, FTP is standardized

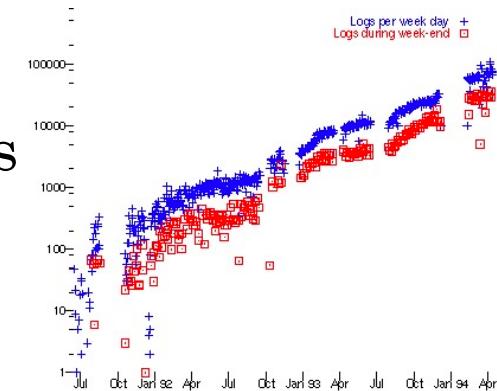


slide by Basel Shbita, picture by the Opte Project https://en.wikipedia.org/wiki/History_of_the_Internet#/media/File:Internet_map_1024_-_transparent,_inverted.png



The World Wide Web

- 1980s, Tim Berners-Lee @ CERN
 - 1989, HTTP
 - 1990, WWW Proposal
 - 1991, HTML
 - 1993, Mosaic (NCSA), 1st Web Browser
 - *Since*, Web of Documents
 - *Now*, Web of Data

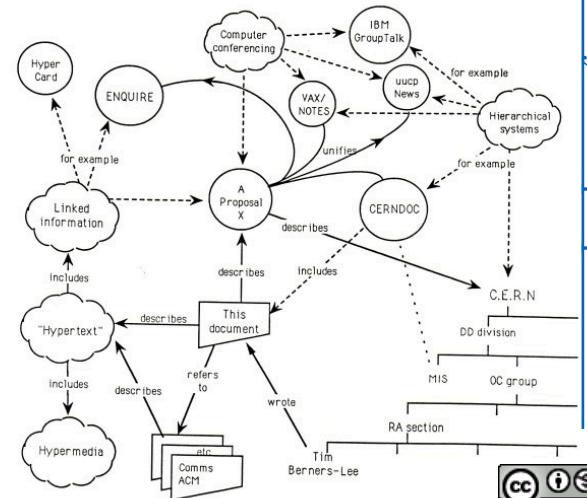


Information Management: A Proposal

Abstract

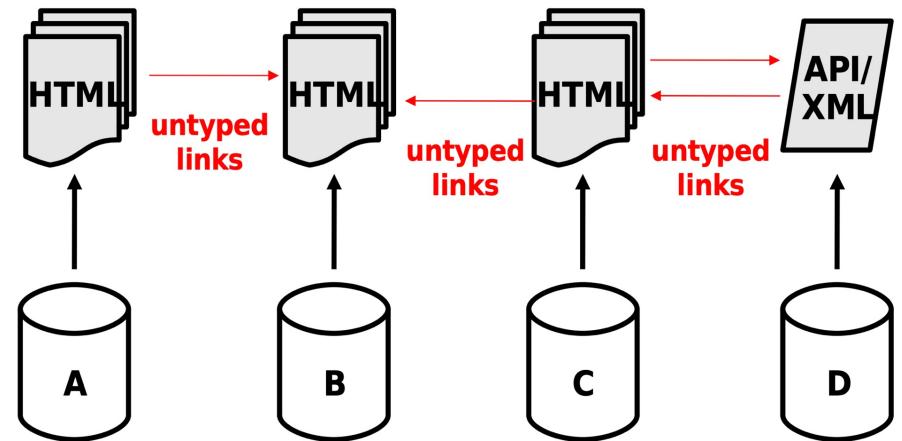
This proposal concerns the management of general information about accelerators and experiments at CERN. It discusses the problems of loss of information about complex evolving systems and derives a solution based on a distributed hypertext system.

Keywords: Hypertext, Computer conferencing, Document retrieval, Information management, Project control



The Web of Documents

- Analogy
 - a global filesystem
- Primary objects
 - documents
- Links between
 - documents (or sub-parts of)
- Degree of structure in objects
 - fairly low
- Semantics of content and links
 - implicit
- Designed for
 - human consumption

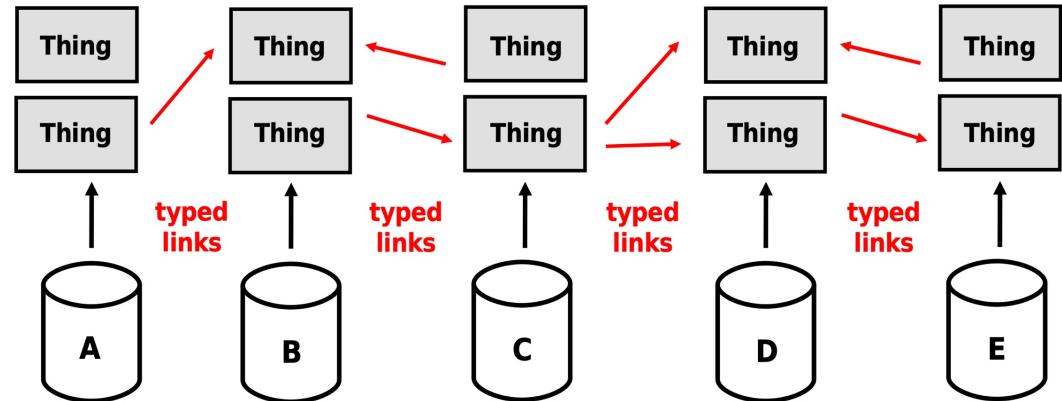


Disconnected Data ☹



The Web of Data

- Analogy
 - a global database
- Primary objects
 - things (or descriptions of things)
- Links between
 - things (including documents)
- Degree of structure in (descriptions of) things
 - high
- Semantics of content and links
 - explicit
- Designed for
 - machines first, humans later



Linked Data

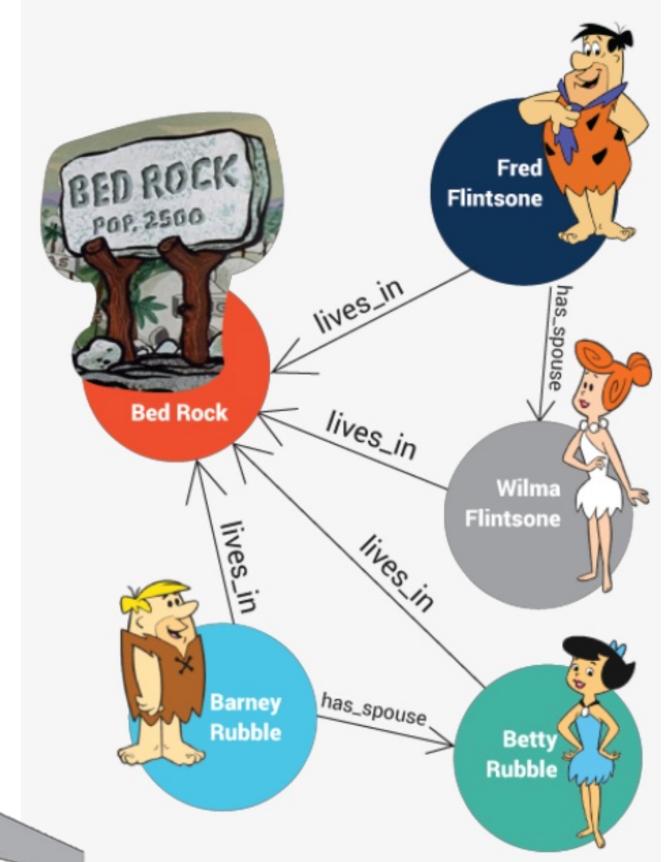
- A method of publishing **structured data** so that it can be **interlinked** and become **more useful**
- Web technologies:
 - HTTP
 - URIs
 - RDFto share information → processed automatically by computers



slide by Basel Shbita

Semantic Web

- W3C, extension of the WWW through new standards
- Encoding of semantics with the data
- Emergence of Ontologies
 - FOAF
 - OWL

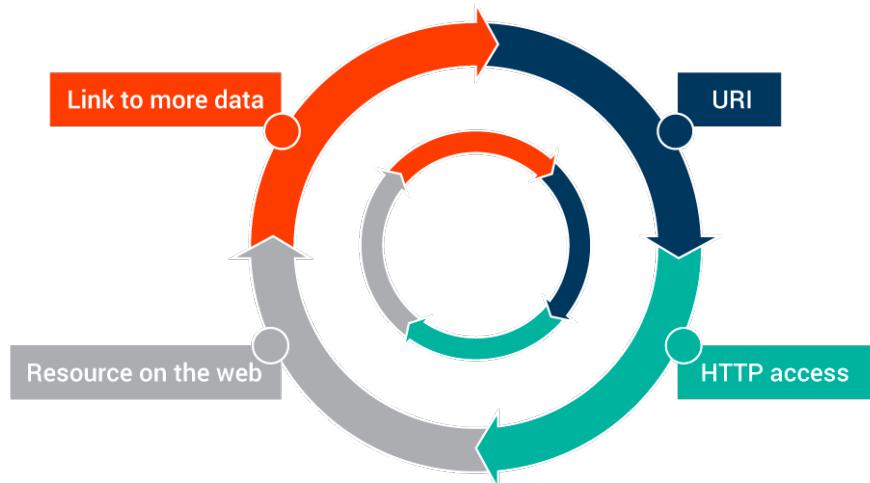


slide by Basel Shbita, images from <https://www.ontotext.com/knowledgehub/fundamentals/what-is-the-semantic-web/>



Linked Open Data

- Linked Data + Open Data
 - DBpedia
 - GeoNames
 - Wikidata
- Present day “Knowledge Graphs”
 - **across** vast amount of general importance
 - **alive** LOD
 - **graph-computing** techniques and algorithms



Linked Data Principles

- Use URIs as names for things
- Use HTTP URIs so that people can look up those names
- When someone looks up a URI, provide useful information, using the standards (RDF*, SPARQL)
- Include links to other URIs so that they can discover more things

http://youtu.be/OM6XIIcM_qo



Tim Berners-Lee
<http://www.w3.org/DesignIssues/LinkedData.html>

slide by Pedro Szekely

Principle 1:
Use URIs as names for things

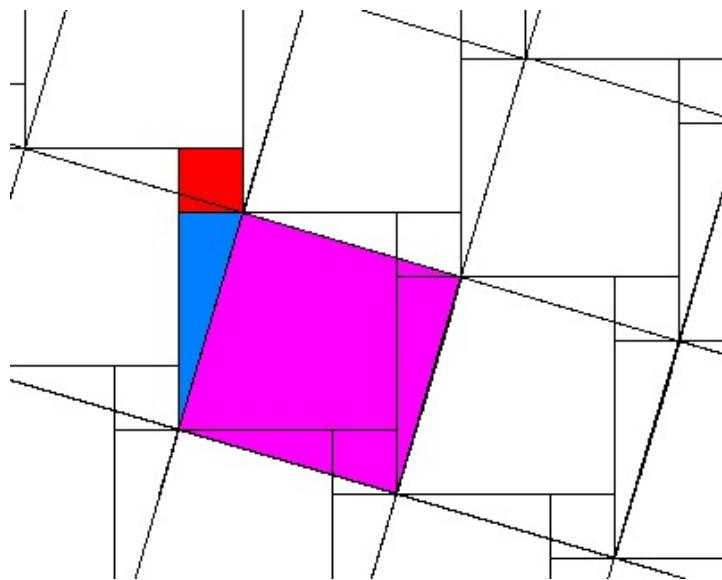
Principle 2:
Use HTTP URIs so that people can look
up those names

Can USC Have a URI?



slide by Pedro Szekely

Can the Pythagoras Theorem Have a URI?



slide by Pedro Szekely

My Dog: Can He Have a URI?



<http://szekelys.com/diego>

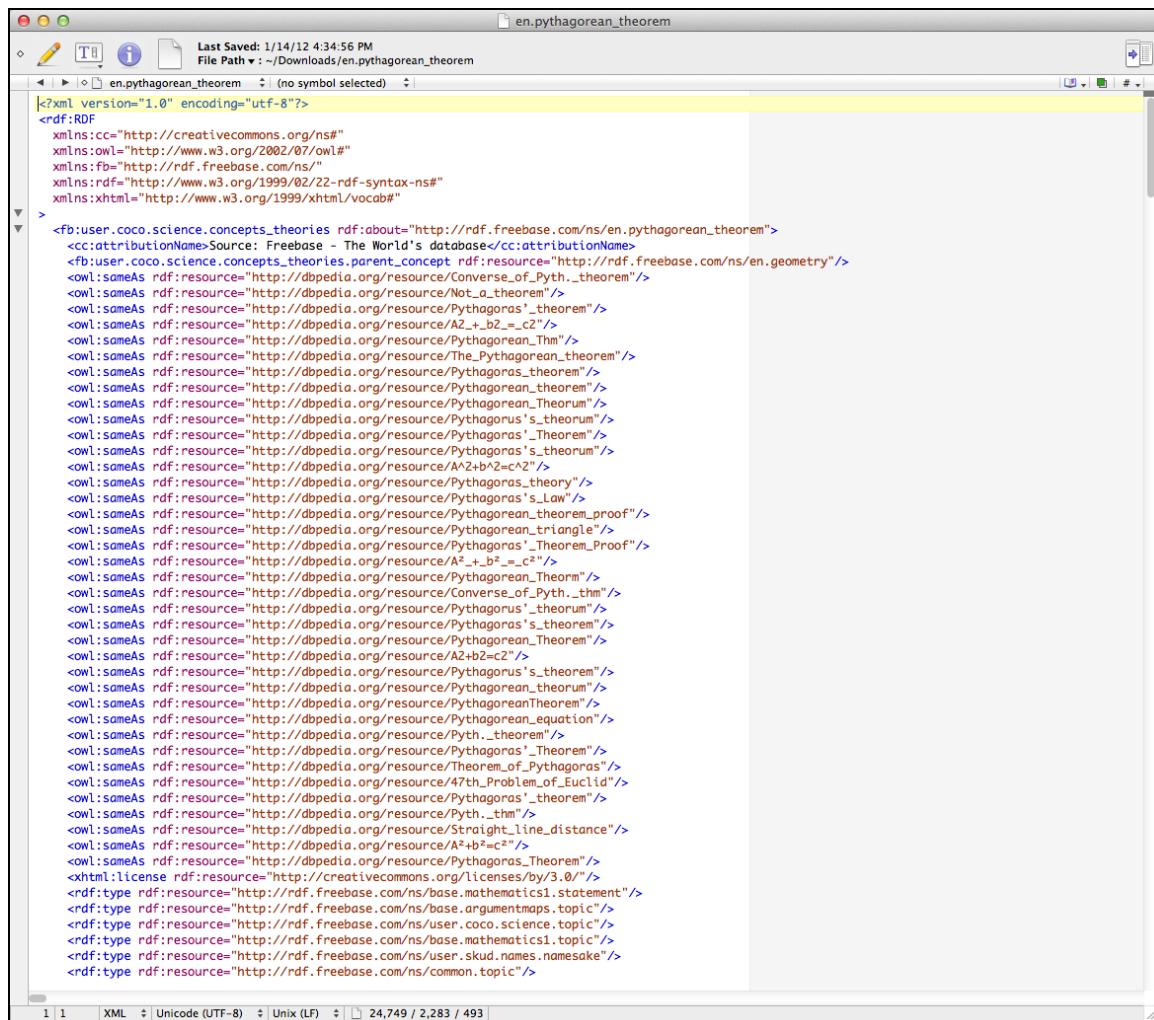
slide by Pedro Szekely

Principle 3:

When someone looks up a URI, provide
useful information, using the standards
(RDF*, SPARQL)

slide by Pedro Szekely

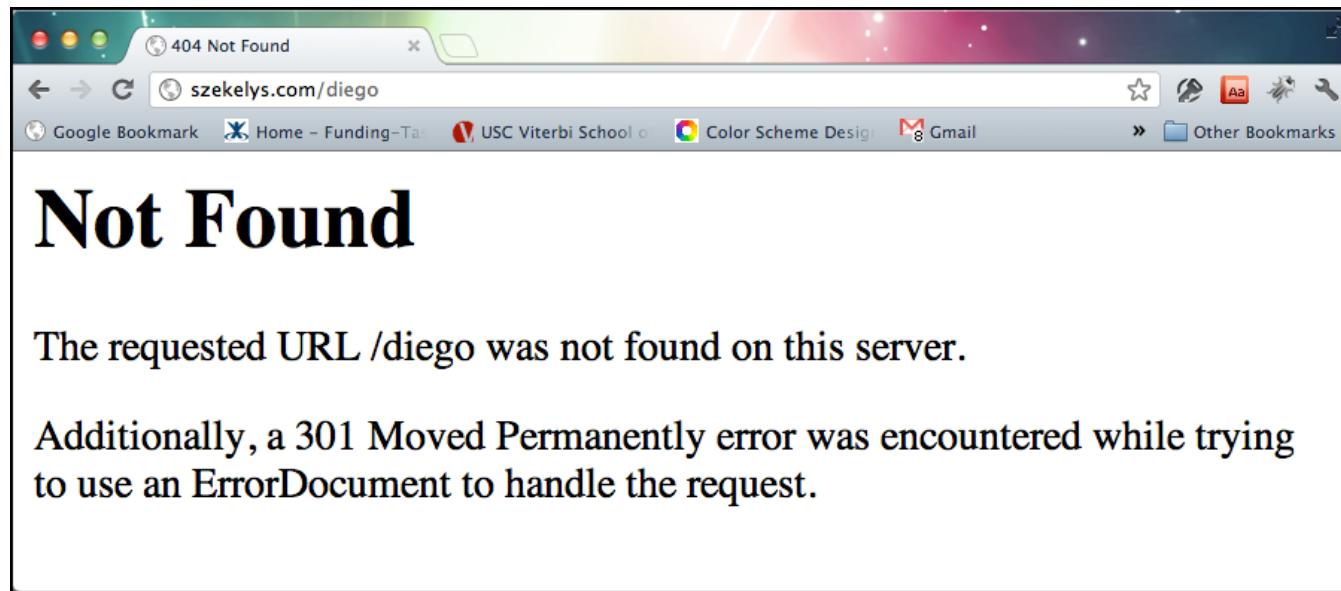
http://www.freebase.com/view/en/pythagorean_theorem



The screenshot shows a text editor window titled "en.pythagorean_theorem". The file path is listed as "~Downloads/en.pythagorean_theorem". The content of the file is an RDF XML document. The XML starts with the declaration "<?xml version='1.0' encoding='utf-8'?>" followed by the root element "<rdf:RDF>". The document contains numerous triples, primarily using namespaces such as "http://rdf.freebase.com/ns", "http://www.w3.org/1999/xhtml/vocab#", and "http://www.w3.org/2002/07/owl#". The triples describe various aspects of the Pythagorean theorem, including its converse, related concepts like the Pythagorean theorem, and its proof. The code is heavily nested, reflecting the complex relationships in the Freebase database.

slide by Pedro Szekely

<http://szekelys.com/diego>



~~Principle 3:~~
~~When someone looks up a URI, provide useful information, using the standards (RDF*, SPARQL)~~

Principle 4:
Include links to other URIs so that they
can discover more things

<http://szekelys.com/diego>

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .  
@prefix dbpprop:<http://dbpedia.org/property/> .  
@prefix dbpedia: <http://dbpedia.org/resource/> .  
@prefix dbpedia-owl: <http://dbpedia.org/ontology/> .  
@prefix fb: <http://rdf.freebase.com/ns/> .
```

<http://szekelys.com/diego>

rdf:type	“Dog” ;
http://szekelys.com/name	“Diego” ;
dbpedia-owl:species	“Labrador Retriever” ;
dbprop:country	“Canada” ;
dbprop:color	“Yellow” ;
fb:base.petbreeds.dog_size.gender	“Male” .

Linked Data?

slide by Pedro Szekely

<http://szekelys.com/diego>

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .  
@prefix dbpprop:<http://dbpedia.org/property/> .  
@prefix dbpedia: <http://dbpedia.org/resource/> .  
@prefix dbpedia-owl: <http://dbpedia.org/ontology/> .  
@prefix fb: <http://rdf.freebase.com/ns/> .
```

<http://szekelys.com/diego>

 rdf:type

<http://szekelys.com/name>

 dbpedia-owl:species

 dbprop:country

 dbprop:color

 fb:base.petbreeds.dog_size.gender "Male".

 "Dog" ;
 "Diego" ;
 "Labrador Retriever" ;
 "Canada" ,
 "Yellow" ;

Not Linked Data

slide by Pedro Szekely

<http://szekelys.com/diego>

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .  
@prefix dbpprop:<http://dbpedia.org/property/> .  
@prefix dbpedia: <http://dbpedia.org/resource/> .  
@prefix dbpedia-owl: <http://dbpedia.org/ontology/> .  
@prefix fb: <http://rdf.freebase.com/ns/> .
```

```
http://szekelys.com/diego  
    rdf:type dbpedia:Dog ;  
    http://szekelys.com/name "Diego" ;  
    dbpedia-owl:species dbpedia:Labrador Retriever ;  
    dbprop:country dbpedia:Canada ;  
    dbprop:color dbpedia:Yellow ;  
    fb:base.petbreeds.dog_size.gender fb:en.male"/> .
```

Principle 4:
Include links to other URIs so that they can
discover more things

<http://szekelys.com/diego>

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .  
@prefix dbpprop:<http://dbpedia.org/property/> .  
@prefix dbpedia: <http://dbpedia.org/resource/> .  
@prefix dbpedia-owl: <http://dbpedia.org/ontology/> .  
@prefix fb: <http://rdf.freebase.com/ns/> .  
  
http://szekelys.com/diego  
    rdf:type dbpedia:Dog ;  
    http://szekelys.com/name "Diego" ;  
    dbpedia-owl:species dbpedia:Labrador Retriever ;  
    dbprop:country dbpedia:Canada ;  
    dbprop:color dbpedia:Yellow ;  
    fb:base.petbreeds.dog_size.gender fb:en.male"/> .
```

Principle 3:

When someone looks up a URI, provide **useful information**, using the standards (RDF*, SPARQL)

slide by Pedro Szekely

<http://szekelys.com/diego>

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .  
@prefix dbpprop:<http://dbpedia.org/property/> .  
@prefix dbpedia: <http://dbpedia.org/resource/> .  
@prefix dbpedia-owl: <http://dbpedia.org/ontology/> .  
@prefix fb: <http://rdf.freebase.com/ns/> .  
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
```

<http://szekelys.com/diego>

rdf:type	dbpedia:Dog ;
foaf:name	"Diego" ;
dbpedia-owl:species	dbpedia:Labrador Retriever ;
dbprop:country	dbpedia:Canada ;
dbprop:color	dbpedia:Yellow ;
fb:base.petbreeds.dog_size.gender	fb:en.male"/> .

foaf is a widely used vocabulary

friend of a friend

slide by Pedro Szekely

Now we know what linked data is

What can go wrong?



slide by Pedro Szekely

Different URIs For the Same Thing

<https://www.wikidata.org/wiki/Q36107>

http://dbpedia.org/resource/Muhammad_Ali

https://yago-knowledge.org/resource/Muhammad_Ali

<http://data.nytimes.com/N13611972026987463463>



slide by Basel Shbita

Linked Data Challenges

- Different URIs for the same thing
 - ... makes it harder to link the data
- Timeliness
 - ... not up to date
- Provenance
 - ... not only a linked data problem
- Tools
 - ... slow performance compared to traditional data
 - ... search engines not yet mature
 - ... many RDF formats, not supported by all tools



slide by Pedro Szekely

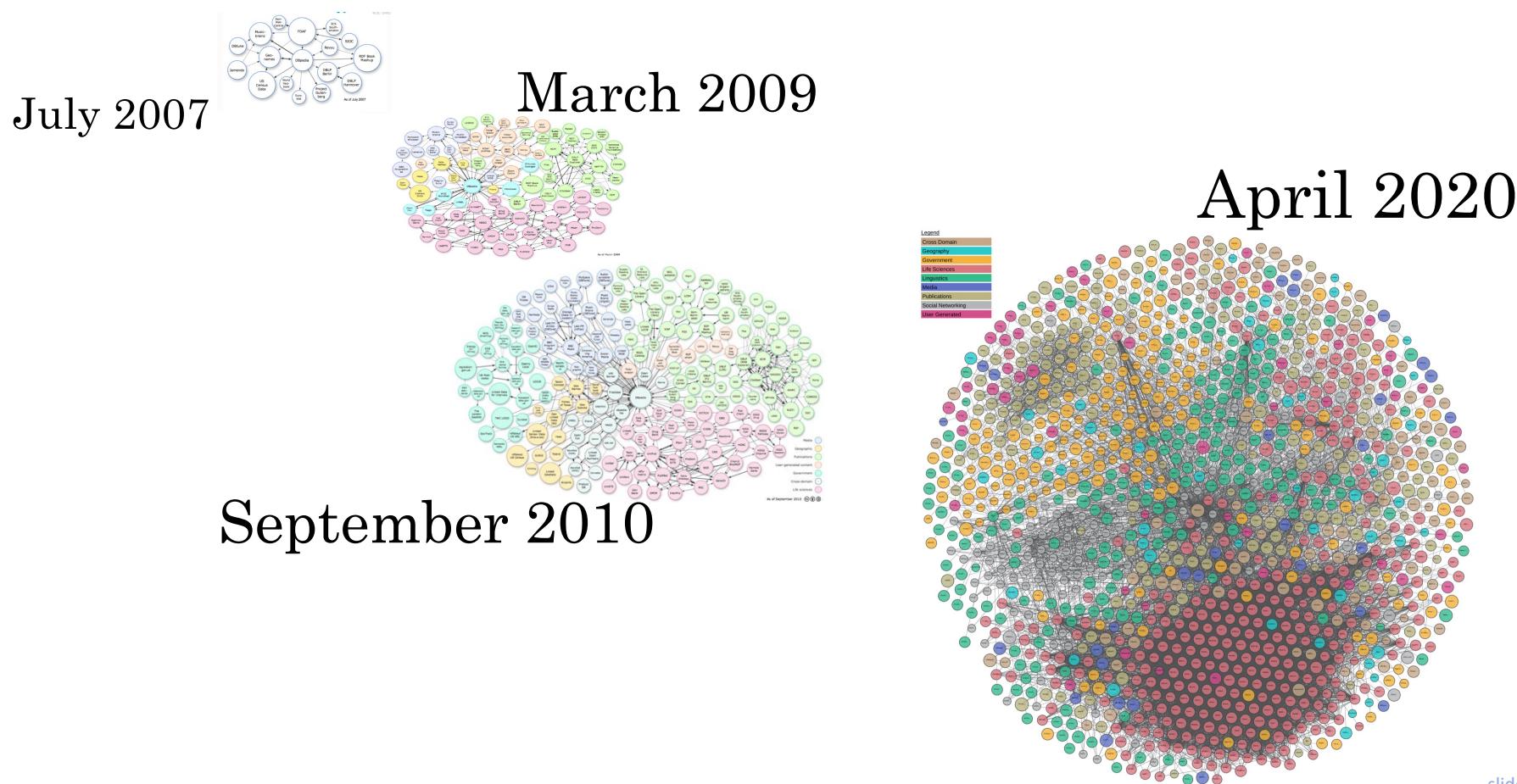
Linked Data Challenges

- Different URIs for the same thing
 - ... makes it harder to link the data
- Timeliness
 - ... not up to date
- Provenance
 - ... not only a linked data problem
- Tools
 - ... slow performance compared to traditional data
 - ... search engines not yet mature
 - ... many RDF formats, not supported by all tools



slide by Pedro Szekely

Getting Bigger Every Day



slide by Pedro Szekely

Working with Linked Data



RDF Dump



SPARQL Endpoint



URI Dereferencing



REST Service



slide by Pedro Szekely

YASGUI

yasgui.org

Query X Query X + http://yasgui.triply.cc/

http://dbpedia.org/sparql

```
1 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
2 PREFIX dbpedia-owl: <http://dbpedia.org/ontology/>
3 PREFIX dbpprop: <http://dbpedia.org/property/>
4 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
5 PREFIX dbprop: <http://dbpedia.org/property/>
6
7 select distinct ?c where {
8   ?s dbprop:hairColor ?c .
9 } limit 10 offset 0
```

Table Raw Response Pivot Table Google Chart

Showing 1 to 10 of 10 entries

Search: Show 50 entries

c

1	http://dbpedia.org/resource/Human_hair_color
2	http://dbpedia.org/resource/Red_hair
3	http://dbpedia.org/resource/Black
4	http://dbpedia.org/resource/Chestnut_hair
5	http://dbpedia.org/resource/Brown

slide by Pedro Szekely

Is your Linked Open Data 5 Star?



Available on the web (whatever format) *but with an open licence, to be Open Data*



Available as machine-readable structured data (e.g. excel instead of image scan of a table)



as (2) plus non-proprietary format (e.g. CSV instead of excel)



All the above plus, Use open standards from W3C (RDF and SPARQL) to identify things, so that people can point at your stuff



All the above, plus: Link your data to other people's data to provide context



slide by Basel Shbita

Best Practices for Data on the Web

<https://www.w3.org/TR/dwbp>

- [Best Practice 1: Provide metadata](#)
- [Best Practice 2: Provide descriptive metadata](#)
- [Best Practice 3: Provide structural metadata](#)
- [Best Practice 4: Provide data license information](#)
- [Best Practice 5: Provide data provenance information](#)
- [Best Practice 6: Provide data quality information](#)
- [Best Practice 7: Provide a version indicator](#)
- [Best Practice 8: Provide version history](#)
- [Best Practice 9: Use persistent URIs as identifiers of datasets](#)
- [Best Practice 10: Use persistent URIs as identifiers within datasets](#)
- [Best Practice 11: Assign URIs to dataset versions and series](#)
- [Best Practice 12: Use machine-readable standardized data formats](#)
- [Best Practice 13: Use locale-neutral data representations](#)
- [Best Practice 14: Provide data in multiple formats](#)
- [Best Practice 15: Reuse vocabularies, preferably standardized ones](#)
- [Best Practice 16: Choose the right formalization level](#)
- [Best Practice 17: Provide bulk download](#)
- [Best Practice 18: Provide Subsets for Large Datasets](#)
- [Best Practice 19: Use content negotiation for serving data available in multiple formats](#)
- [Best Practice 20: Provide real-time access](#)
- [Best Practice 21: Provide data up to date](#)
- [Best Practice 22: Provide an explanation for data that is not available](#)
- [Best Practice 23: Make data available through an API](#)
- [Best Practice 24: Use Web Standards as the foundation of APIs](#)
- [Best Practice 25: Provide complete documentation for your API](#)
- [Best Practice 26: Avoid Breaking Changes to Your API](#)
- [Best Practice 27: Preserve identifiers](#)
- [Best Practice 28: Assess dataset coverage](#)
- [Best Practice 29: Gather feedback from data consumers](#)
- [Best Practice 30: Make feedback available](#)
- [Best Practice 31: Enrich data by generating new data](#)
- [Best Practice 32: Provide Complementary Presentations](#)
- [Best Practice 33: Provide Feedback to the Original Publisher](#)
- [Best Practice 34: Follow Licensing Terms](#)
- [Best Practice 35: Cite the Original Publication](#)



slide by Basel Shbita

What's next?

