

# Semantic Web & Linked Open Data use cases

Dr David Croft

Coventry University

david.croft@coventry.ac.uk

2017

# Overview

## 1 Introduction

## 2 Fast recap

- Why LOD?

## 3 Use cases

- Elsevier
- Audi
- Swiss Life
- Personal experience

## 4 Benefits

## 5 Conclusion

## Fast recap

*“The Semantic Web is an extension of the current web in which information is given well-defined meaning, better enabling computers and people to work in cooperation.”*

–Tim Berners-Lee

## Here's one I made earlier

Latest attempt to represent human knowledge in a compatible, linkable, intuitively searchable format.

- 1895 Repertoire Bibliographique Universel (Universal bibliographic directory).
- 1927 Statistical Machine.
- 1939 Memex.
- 1987 Hypercard.
- 1989 Information management: A proposal (became the Web).



Repertoire Bibliographique Universel <sup>1</sup>

<sup>1</sup>[http://www.udcc.org/index.php/site/page?view=about\\_history](http://www.udcc.org/index.php/site/page?view=about_history)

Doesn't the World Wide Web (Web) do this already?

- Web is clearly the greatest collection of knowledge ever produced.
- Early search engines just did text matching.
  - Improvements made but are still basically relying on advanced text matching.
- Too much data.
  - Because you can find everything, is difficult to find anything (specific).

In order to be really effective at searching need to understand meaning of data being searched.

- Solutions?
  - 1 Solve general purpose Artificial Intelligence (AI).
  - 2 Cheat.

# Linked Open Data (LOD)

The current attempt to achieve the ideals of the semantic web.

- Massive collection of standards formatted data.
- Importantly the meaning of the data is also included.
  - Don't try and deduce the meaning of the data.
  - Just tell our search tools the meaning of everything.
  - Horray for cheating.

# Artificial Intelligence (AI)

If any of this sounds familiar to the expert systems of the 1970s.

- Far more limited claims regarding ability.
- No claims of general purpose AI.
- Old approaches failed due to difficulty of producing sufficient knowledge base<sup>2</sup>.
  - LOD offers distributed generation.
  - Each additional resource increases the value of everything else.

# eXtensible Markup Language (XML)

User definable markup language.

- Strong unicode support.
  - Non-english character sets.
- Often confused with HTML
  - Data presentation vs. data storage.
  - Both use angle bracket contained tags to denote structure.

```
<quiz>
  <question num="1">
    What is anatidaephobia?
    <answer>Fear of a duck watching you.</answer>
  </question>
</quiz>
```



# Resource Description Framework (RDF)

## Key features of Resource Description Framework (RDF):

- Representation of information as triples.
  - Object - Attribute - Value.
- Chained triples form a graph.
- Use of Uniform Resource Identifiers (URIs).
- Predicates/properties are also URIs.

## Alternative notations

RDF generally made available as XML. Fine for automated tools but unpleasant to read.

Alternative notations do exist.

- Turtle.
- N3.
- N-Triples.
- RDF/JSON.

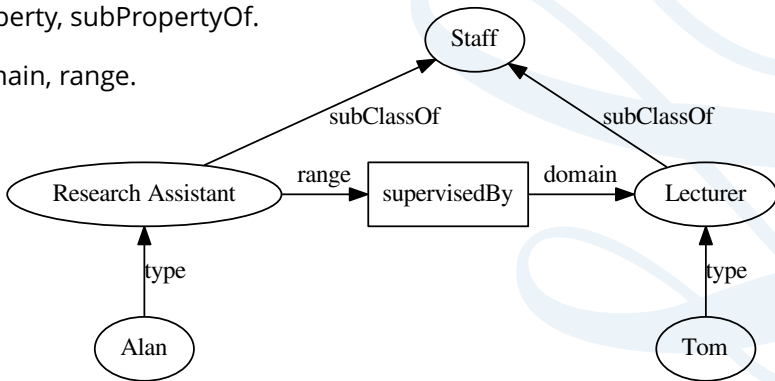
```
<#green-goblin>  
  rel:enemyOf <#spiderman>;  
  a foaf:Person;  
  foaf:name "Green Goblin" .
```

```
<#spiderman>  
  rel:enemyOf <#green-goblin>;  
  a foaf:Person;  
  foaf:name "Spiderman".
```

## RDF-Schema (RDFS)

Defined vocabulary for RDF.

- Hierarchical organisation.
  - Class, subClassOf, type.
  - Property, subPropertyOf.
  - domain, range.



# Web Ontology Language (OWL)

Don't get me started on the acronym.

- Describe term similarity and/or difference.
- Construct new classes.
  - Membership requirements.
  - As intersection, union, complement of existing ones.
- Reason about terms
  - if  $\langle \text{Person} \rangle \langle A \rangle$  and  $\langle \text{Person} \rangle \langle B \rangle$  have same  $\langle \text{foaf:email} \rangle$  then  $\langle A \rangle$  and  $\langle B \rangle$  are same individual.

# SPARQL Protocol and RDF Query Language (SPARQL)

Structured Query Language (SQL) like syntax for querying RDF sources.

- Querying and transformation only.
  - Not updating, deleting, database admin etc.
- Can get data from RDF graphs.
- Can create new graphs based on existing ones.

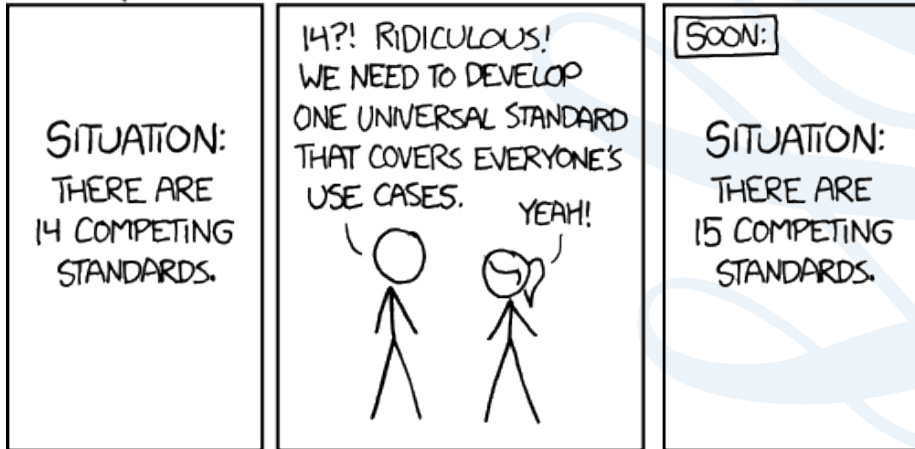
# Another set of standards, really?

The LOD tech stack succeeds where others have failed because:

- You don't have to migrate your data to a new format.
  - Keeps using whatever weird in-house developed format you want.
- Doesn't try to develop an ontology for everything.
  - Jack of all trades, master of none problem.
- Specify how individual elements of it map to other schemas.
  - Or let others map their schema onto yours.
- Can do conversion on the fly.

## Standards

HOW STANDARDS PROLIFERATE:  
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)



<https://xkcd.com/927/>

# Why LOD?

## Why use LOD?

- Do not need to use full LOD tech stack.
- Open data format.
  - Non-proprietary.
  - Easy of future migration.
- Includes the schema.
- Includes the ontology.
- Self documenting.
  - Just follow the URIs.



# Real world use cases

## Elsevier

- Major academic journal publisher.

## Problem

- Wide range of different domains.
  - Subscriptions to individual journals/domains.
- Increase in interdisciplinary research.
  - Information may be spread over multiple subject areas.
- Want everything related to topic e.g. Alzheimers.
  - Information spread - biology, chemistry, medicine, sociology.
  - Differing terminologies (keyword search problematic).



ELSEVIER

## Solution

- Mapping article topics to EMTREE thesaurus.
  - Unique Identifier (UID) for individual concepts allows cross domain searching despite differing domain terminology.
  - Mapped to existing domain specific ontologies e.g. Medical Subject Headings (MeSH).
  - So keep the existing annotations.
- EMTREE has RDF representation.
  - RDF acts and interoperability format between data sets.

## More use cases

### Audi

- $\approx 88\,000$  employees (2016).
- $\approx \text{€}60$  billion annual revenue.
- $\approx 2$  million cars sold annually.

### Problem

- 1000+ databases.
  - Isolated data.
  - Duplicated data.
- Can't be easily queried.
  - Reliance on manual code generation.
  - Extract, Transform & Load (ETL) processes.



## Solution

- You could create a gigantic data warehouse which everything should be stored in.
  - Data warehouse projects fail a lot i.e. the failed care.data NHS project.
- Or ontology based approach.
  - Annotate existing data fields with meaning.
  - Existing applications are unaffected (and unaware).

Imagine two data sources.

## Source A

```
<SLR rdf:ID="Olympus-OM-10">  
  <viewFinder>twin mirror</viewFinder>  
  <optics>  
    <Lens>  
      <focal-length>75-300mm zoom</focal-length>  
      <f-stop>4.0-4.5</f-stop>  
    </Lens>  
  </optics>  
</SLR>
```

## Source B

```
<Camera rdf:ID="Olympus-OM-10">  
  <viewFinder>twin mirror</viewFinder>  
  <optics>  
    <Lens>  
      <size>30cm zoom</size>  
      <aperture>4.5</aperture>  
    </Lens>  
  </optics>  
</SLR>
```

- No way for source A to know that source B's aperture is equivalent to its f-stop.
- No way for source B to know that an SLR is type of camera.

Application B starts processing some data from source A.

- What is an SLR?
- What is focal-length?
- What is f-stop?
- Could program these rules into each program.
  - 2 sources means explaining it 4 times.
  - 1000 sources means explaining it 1 000 000 times.
- Or use an ontology approach.
  - Encounters unknown term, look it up.
  - Oh that's what it means.

```
<owl:Class rdf:ID="SLR">  
  <rdfs:subClassOf rdf:resource="#Camera" />  
</owl:Class>
```

```
<owl:DatatypeProperty rdf:ID="f-stop">  
  <rdfs:domain rdf:resources="#Lens" />  
</owl:DatatypeProperty>
```

```
<owl:DatatypeProperty rdf:ID="aperture">  
  <rdfs:equivalentProperty rdf:resources="#f-stop" />  
</owl:DatatypeProperty>
```

```
<owl:DatatypeProperty rdf:ID="focal-length">  
  <rdfs:domain rdf:resources="#Lens" />  
</owl:DatatypeProperty>
```

```
<owl:DatatypeProperty rdf:ID="size">  
  <rdfs:domain rdf:resources="#focal-length" />  
</owl:DatatypeProperty>
```

## Unit conversion

How much you can achieve depends on level of detail in ontology.

Start including unit type information in the ontology and we can start doing unit conversion on the fly too.

- Quantities, Units, Dimensions & Types (QUDT), SPARQL Inferencing Notation (SPIN)
- Now not just understanding the equivlence of the fields.
  - Understanding the equivlence of the values.

```
<Lens>
```

```
  <size>30</size>
```

```
</Lens>
```

```
<owl:DatatypeProperty rdf:ID="size">
```

```
  <rdfs:domain rdf:resources="#focal-length" />
```

```
  <qudt:unit rdf:resource="http://qudt.org/vocab/unit#Centimeter"/>
```

```
</owl:DatatypeProperty>
```



## Even more use cases

Switzerland's largest life insurance group.

- $\approx 11\,000$  employees.
- $\approx \$14$  billion in premiums.
- Branches in 50 countries.

### Problem

- Distributed widely culturally and geographically.
- Creating company wide skills repository is tricky.
  - Who has what skills, qualifications.
  - How are different skills described in different countries?
  - Equivalent qualifications, UK GCSE  $\leftrightarrow$  US High School Diploma?



# SwissLife

## Solution

- Hand built ontology.
- Cover skills from 3 organisational units.
  - IT
  - Private insurance
  - HR
- 700 concepts + 180 educational concepts.
- 130 job function concepts overall.
  - I.e. Administration / Office.
  - Associated job titles; "executive assistant", "office manager"...
  - Responsibilities; basic clerical, purchasing, facilities management...
  - Skills; word processing, filing, customer relations...

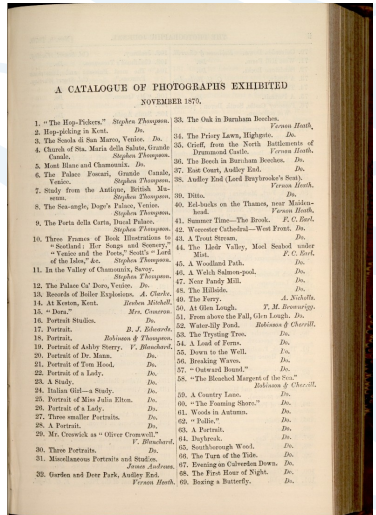
# Personal experience

Did my PhD and subsequent research in an area adjacent to Semantic web research.

- FuzzyPhoto project  
<http://fuzzyphoto.edublogs.org/>

Trying to do co-reference identification of photographic artefacts across disparate Gallery, Library, Archive & Museum (GLAM) collections.

- No digitised images so entirely metadata based.
- Success .... kind of.



# GLAMs problem

Galleries, Libraries, Archives & Museums (GLAMs) store cultural artefacts.  
Millions of items.

## Problem

- Large digitisation efforts for decades.
  - For conservation.
  - For organisation.
  - For access.
- But no standardisation.
- Every collection has its own schemas and ontologies.
  - If they have any at all.
- Metadata is itself a historical artefact.

# GLAMs improvement

Situation is better now.

- Awareness of need for standardisation and cross-collection compatibility.
  - Lightweight Information Describing Objects (LIDO) schema, CIDOC-CRM ontology.
  - ★★★ source but can be mapped over to RDF/linked data etc.
- Awareness of desirability of URIs for referencing/research etc.
- Desire to cross-reference collections.
  - For research, authentication.
- New collections may be created in Linked Data formats.
  - Or large institutions
    - British Library - <http://bnb.data.bl.uk/flint-sparql>
    - British Museum - <http://collection.britishmuseum.org/>
  - Unfortunately old data is not easily converted.

FuzzyPhoto, unable to use Semantic web technology.

- GLAM metadata is often too poor quality to be automatically converted to RDF.
- Direct transcription of human readable text with no standardised syntax.
  - I.e. date columns with 21 different date formats in them.
  - No standardised schema.
- Inability to trust the data.
  - Imprecise and/or incorrect.

Semantic web technology would have been brilliant.

- But data couldn't support it.
- Used Natural Language Processing (NLP), Fuzzy logic and clustering instead.

## Imprecise/Incorrect

Examples of issues. All examples are real field values.

- Data field says "03/05/10"
  - Most places - That's the 3<sup>rd</sup> of May 2010.
  - USA - 5<sup>th</sup> of March 2010.
  - China, Canada, Korea - 10<sup>th</sup> of May 2003.
  - But dealing with GLAM data, could be anywhen, 1903, 1803 etc.
- "Tuesday July 7<sup>th</sup>".
  - Narrows it to 1 day every 7 years across all of recorded history.
- "18<sup>th</sup>, 19<sup>th</sup> or 20<sup>th</sup> century".
  - 100 years less precise than if the field was blank.
- "Before the second world war".

Just the date information, comparable issues exist for all field types used.

Does demonstrate an issue with the Linked Data approach to achieving the Semantic Web.

- Requires high quality data.
  - Structured data
    - If pre-existing data is lacking, can be very expensive to add it in.
- Requires annotated data.
  - If pre-existing annotations are lacking, can be very expensive to add in.



# Unstructured problems

Limited benefits when dealing with unstructured data.

- I.e. Web 2.0 material.
  - Facebook posts may be annotated with date, location author.
  - No-one is going to annotate the meaning of the text.
- Data trust.
  - Authorative sources.
  - People lie (spam).

## LOD principles

Even if not using the entire LOD tech stack, consider LOD principles and ★ rating.

### Principle

- 1 Use URIs as names for things.
- 2 Use HTTP URIs so that people can look up those names.
- 3 When someone looks up a URI, provide useful information.
- 4 Include links to other URIs, so that they can discover more things.

### Advantages

- 1 URIs make it trivial to link objects.
  - DOI, ISBN, UPC etc.
- 2 Means that URIs become Globally Unique IDentifiers (GUIDs), don't have to worry about domain disambiguation.
- 3 When someone looks up a URI, provide useful information.
  - Why make it hard to learn out about an item/object/resource?
- 4 Include links to other URIs, so that they can discover more things.
  - Place things in context.

## Rating

- ★ Data is online.
- ★★ Data is online in a structured format.
- ★★★ Data is online in a non-proprietary structured format.
- ★★★★ Uses URIs so data can be referenced from elsewhere.
- ★★★★★ Contains links to data sources elsewhere.

## Advantages

- ★ Humans can read it.
- ★★ Computers can read it, if they have the right software.
- ★★★ Computers can read it, if you tell them what the data is.
- ★★★★ Data sources elsewhere can refer to your data.
- ★★★★★ Your data can be fully understood by computers.

# The End

## Any questions?

# Glossary

<b>AI</b>	Artificial Intelligence
<b>CIDOC</b>	ICOM International Committee for Documentation
<b>CIDOC-CRM</b>	CIDOC-Conceptual Reference Model
<b>DOI</b>	Document On Internet
<b>ETL</b>	Extract, Transform & Load
<b>GLAM</b>	Gallery, Library, Archive & Museum
<b>GUID</b>	Globally Unique Identifier
<b>HTTP</b>	HyperText Transfer Protocol
<b>ICOM</b>	International Council of Museums
<b>ISBN</b>	International Standard Book Number
<b>JSON</b>	JavaScript Object Notation
<b>LIDO</b>	Lightweight Information Describing Objects
<b>LOD</b>	Linked Open Data

<b>MeSH</b>	Medical Subject Headings
<b>NHS</b>	National Health Service
<b>NLP</b>	Natural Language Processing
<b>OWL</b>	Web Ontology Language
<b>QUDT</b>	Quantities, Units, Dimensions & Types
<b>RDF</b>	Resource Description Framework
<b>RDF/JSON</b>	RDF/JavaScript Object Notation
<b>RDFS</b>	RDF-Schema
<b>SPARQL</b>	SPARQL Protocol and RDF Query Language
<b>SPIN</b>	SPARQL Inferencing Notation
<b>SQL</b>	Structured Query Language
<b>UID</b>	Unique Identifier
<b>UPC</b>	Universal Product Code
<b>URI</b>	Uniform Resource Identifier
<b>Web</b>	World Wide Web
<b>XML</b>	eXtensible Markup Language