

Development of Linked Data Services to Support Widespread Exposure of Data

CHRIS WOOD
BRITISH OCEANOGRAPHIC DATA CENTRE



**National
Oceanography Centre**
NATURAL ENVIRONMENT RESEARCH COUNCIL

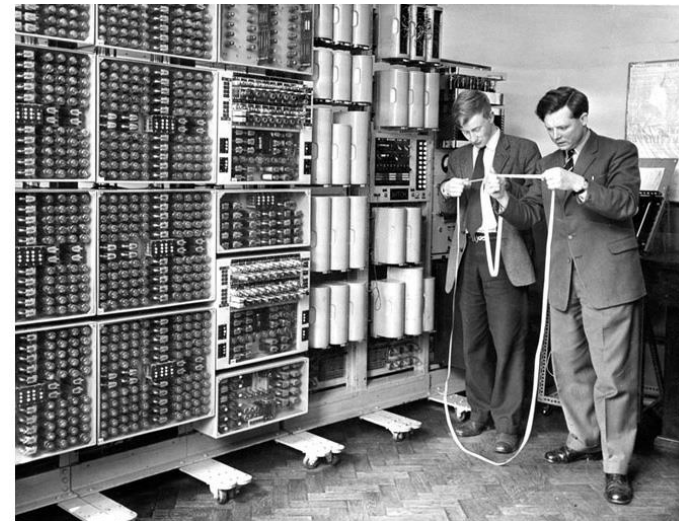
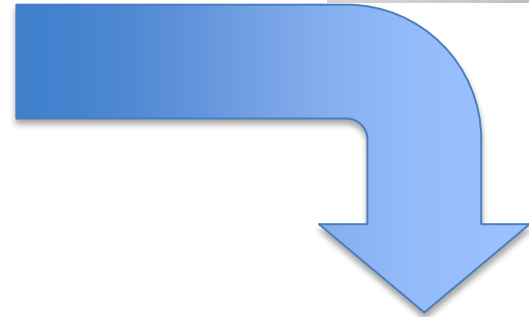
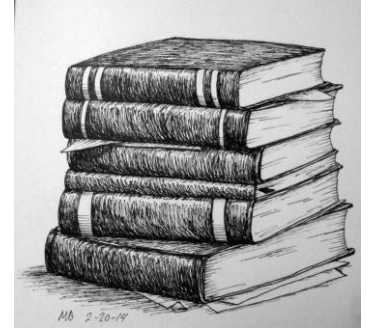
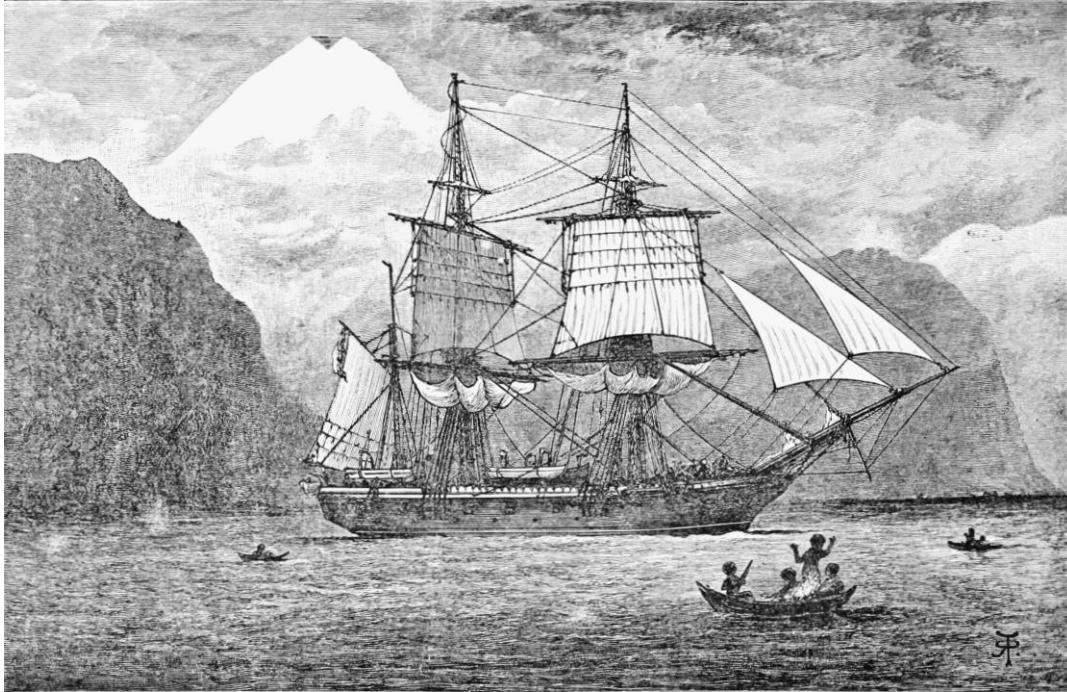


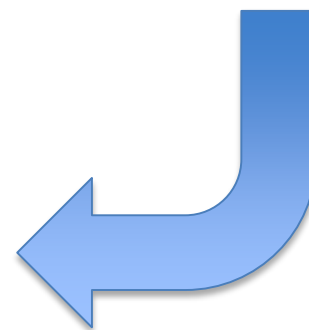
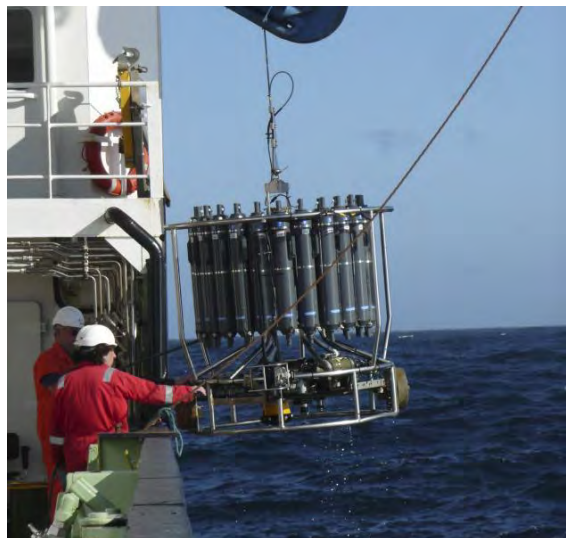
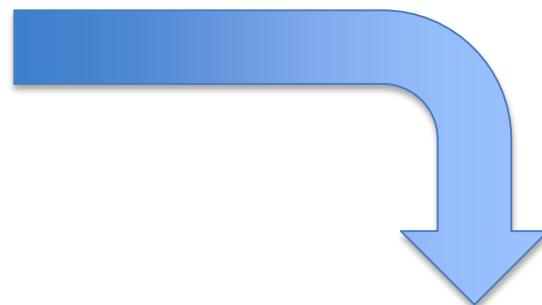
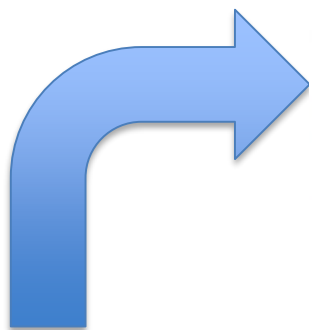
noc.ac.uk

OVERVIEW

- Evolution of data centres
- What is linked data (and why should we care?)
- What are the technical challenges?
- What are the cultural / societal challenges?
- Two case studies:
 - The BODC triplestore project
 - The Celtic Seas Partnership portal

Evolution of data centres





Linked data

- RDF: W3C standard

“*RDF is a **standard** model for data interchange on the Web. RDF has features that facilitate data merging even if the **underlying schemas differ**, and it specifically supports the evolution of schemas over time without requiring all the data consumers to be changed.*”

- RDF is a model, **not** a file format
- RDF models can be represented in various file formats (e.g. but not limited to, XML, JSON, turtle, n3)
- RDF models can be stored in various types of database

Technical details of RDF

- Individual expressions are called *triples*: anything can be described by a statement with three parts

Subject	Predicate	Object
Person	hasName	Chris
Chris	worksAt	BODC
BODC	locatedIn	Liverpool
Liverpool	hasCathedral	Anglican
Liverpool	hasCathedral	Catholic

- A database of triples is called a *triplestore*, and can be queried using SPARQL:

select ?s where {“Liverpool” hasCathedral ?s . }

The power of predicates

- Discovery of *objects* and filtering of *subjects* is much easier if known *predicates* are used
- Predicates come from external ontologies
- e.g: the *foaf* ('friend of a friend') ontology contains a property *name*

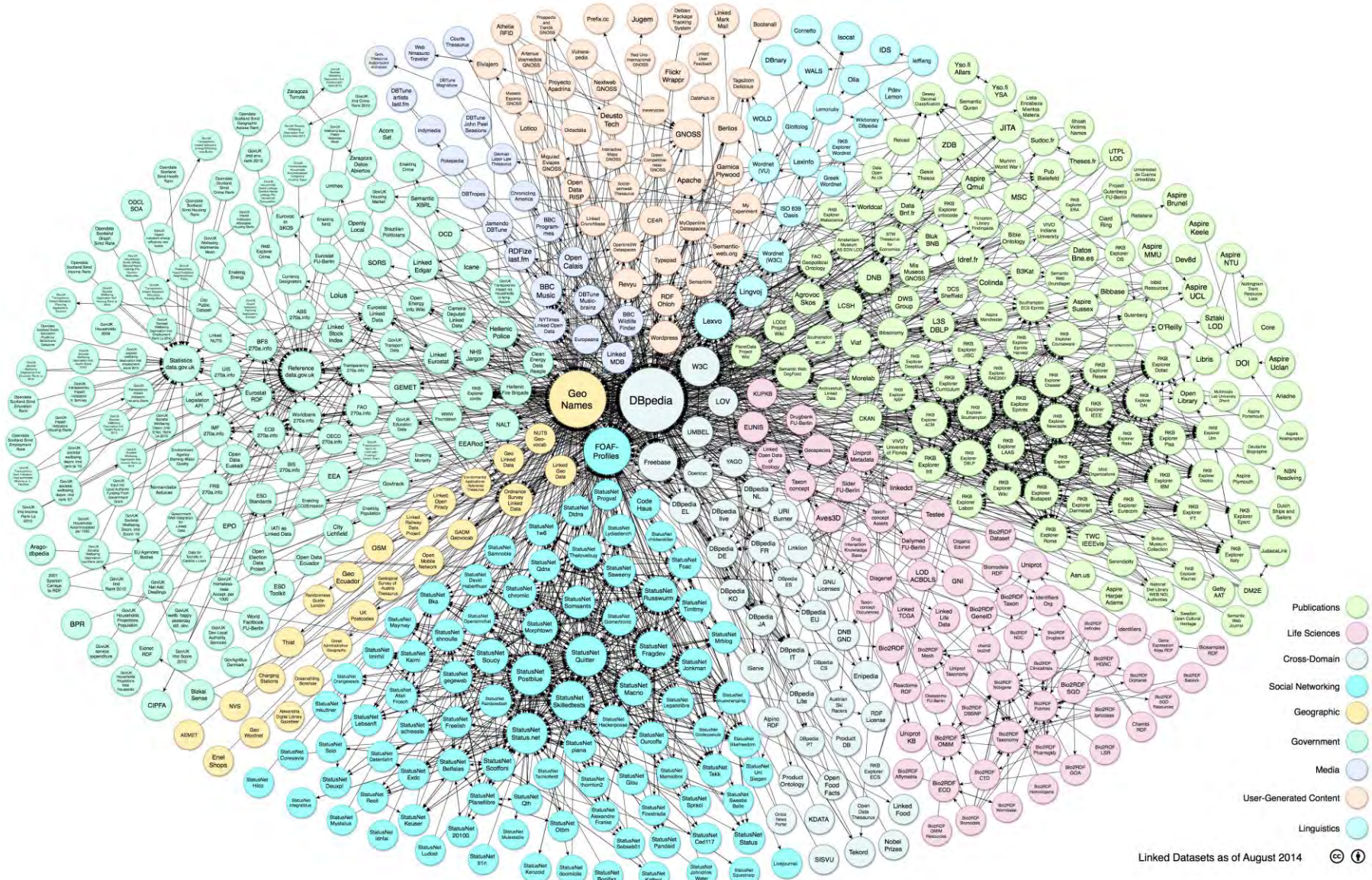
person foaf:name Chris

- Exploration and understanding of the dataset becomes easy!

Pros of ontologies: anyone can create them!

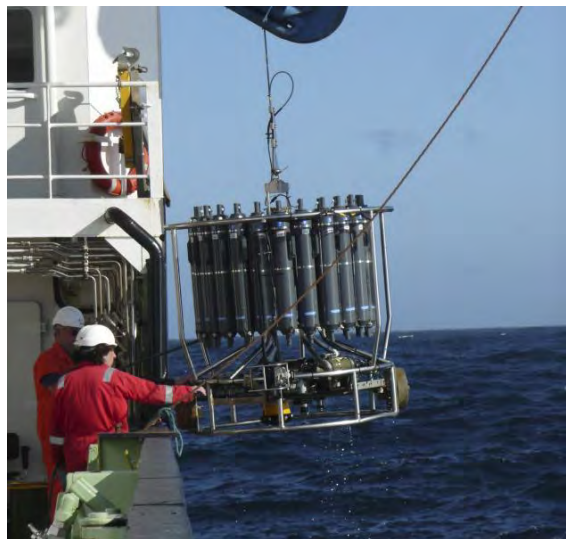
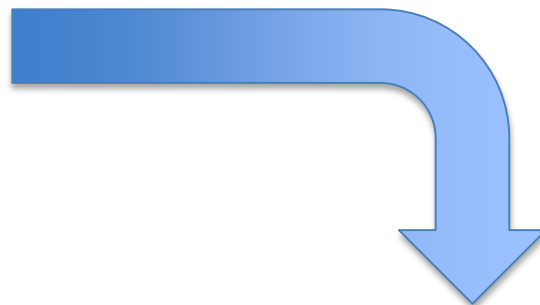
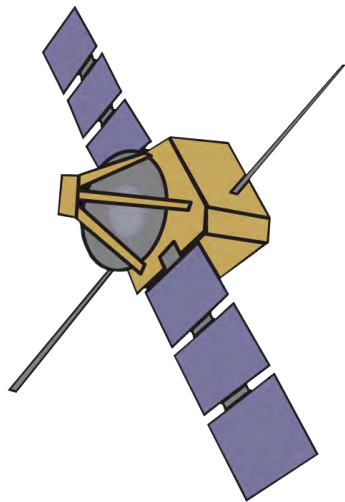
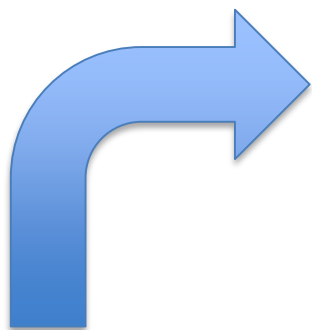
Cons of ontologies: anyone can create them!

Communication in the community is vital!



Reflections so far...

- Triplestores are designed to be exposed to the real world (via a SPARQL endpoint)
 - i.e. all data is exposed by default
- Traditional relational databases can normally only be exposed via software (e.g. a web interface)
 - For data to be exposed via traditional databases a custom API needs to be built
- **Standards** need to be followed for triplestores to be most effective
- Triplestores are just a different flavour of database: most **standard transactions** are available



Case study 1: The BODC architecture

- The main BODC schema holds metadata about data series (physical / biological / chemical data from CTD casts / transects / float deployments)
- Currently: ~108000 data series
- Discovery of data to date has been relatively labour intensive
- Triples needed to be created for all relevant metadata
- Triplestore is updated nightly
- Relevant coding for triplestore creation (e.g. data integrity & transactions)
- SPARQL endpoint software installation & configuration: stack of Jena (Fuseki & TDB), and elda

Real life triples...

```
<http://linked.bodc.ac.uk/series/26229><http://www.w3.org/2004/02/skos/core#notation>26229  
<http://linked.bodc.ac.uk/series/26229><http://mmisw.org/ont/ioos/biological#minimumDepthInMeters>"60.0"^^<http://www.w3.org/2001/XMLSchema#float/>0789/>
```

```
select ?depth where {  
  <http://linked.bodc.ac.uk/series/26229>  
  <http://mmisw.org/ont/ioos/biological#minimumDepthInMeters>  
  ?depth .  
}
```

```
-----  
|                depth                |  
=====
```

"60.0"^^<http://www.w3.org/2001/XMLSchema#float/>0789/>

```
-----
```

More advanced querying

← → ↻ linked.bodc.ac.uk

BODC Linked Open Data

linked.bodc.ac.uk/sparql/ ×

← → ↻ linked.bodc.ac.uk/sparql/?query=select+%3Fseries+%3Fdepth+where+%7B%3Fseries+<http%3A%2F%2Fmmisw.org%2Font%2F

series	depth
<http://linked.bodc.ac.uk/series/581431>	"5812.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/1137202>	"5572.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/1137214>	"5571.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/1137195>	"5541.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/1137183>	"5540.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/972689>	"5520.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/972690>	"5520.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/1092795>	"5517.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/945968>	"5512.0"^^<http://www.w3.org/2001/XMLSchema#float>
<http://linked.bodc.ac.uk/series/1137447>	"5502.0"^^<http://www.w3.org/2001/XMLSchema#float>

Force the accept header to text/plain regardless.

Get Results

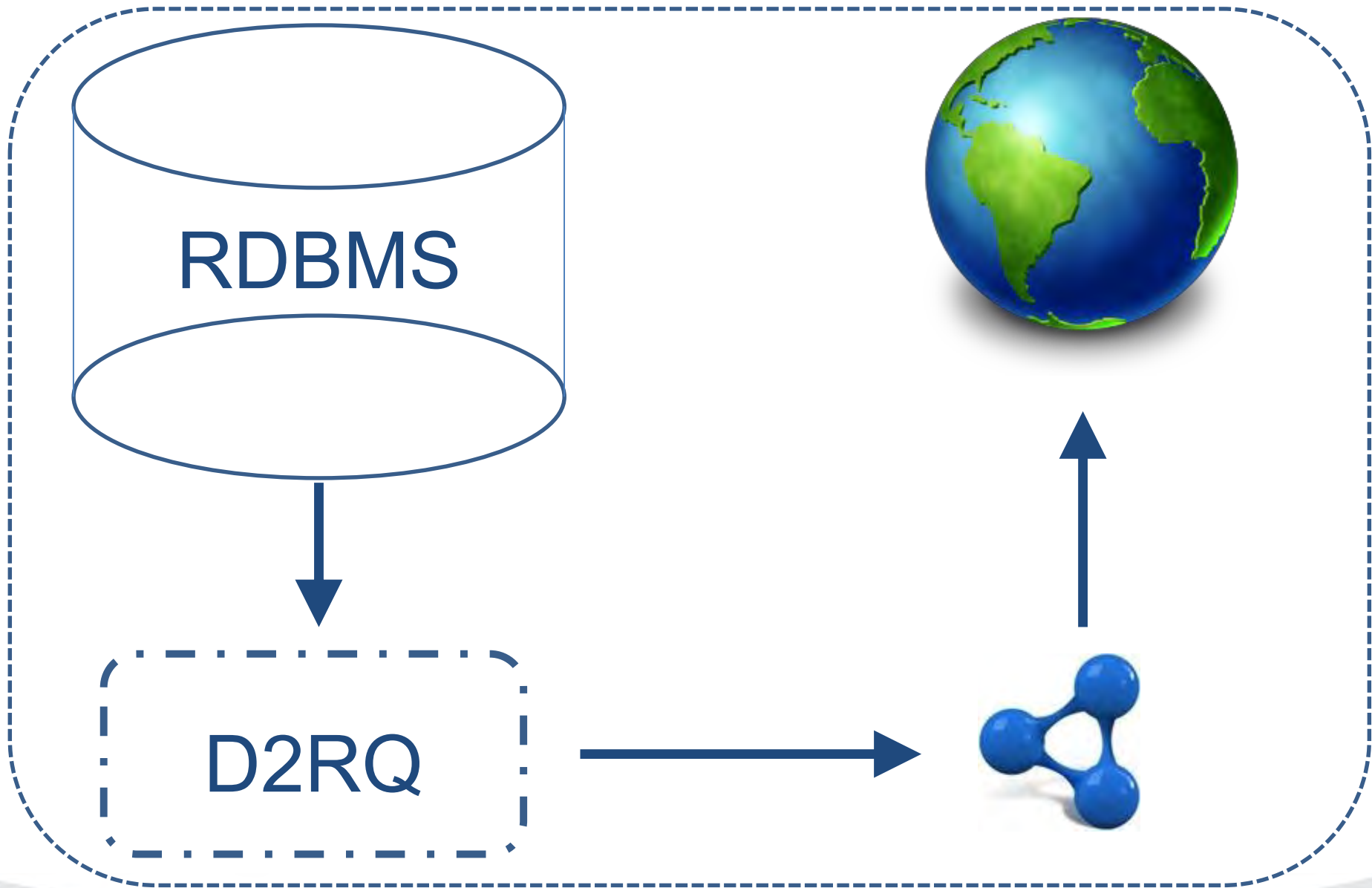
Case study 2: The Celtic Seas Partnership

- Marine Strategy Framework Directive (MSFD) requires EU member states to reach or maintain Good Environmental Status by 2020
- There are various frameworks to help implementation
- The Celtic Seas Partnership is a transboundary management framework to help implementation in the Celtic Seas



Celtic Seas Partnership Portal

- Key deliverable: a web-based portal for the marine community and interested stakeholders to expose datasets and documents that are of relevance
- Metadata stored in Oracle and exposed via a SPARQL endpoint using the D2RQ architecture
- D2RQ: dynamic mapping of an RDBMS to a pseudo-triplestore.
 - Mappings for ontologies can be applied in a configuration file
 - This allows the data to be queried as if it were stored in a native triplestore → the data is publicly exposed
 - SPARQL queries are transformed to SQL



FINDING CELTIC SEAS DATA

Use the filter menu to refine your search. Navigate back to choose a different descriptor or return to the homepage via the link in the top-left corner of the page.

CURRENTLY DISPLAYING
4 DATA SETS

CLEAR FILTERS

MSFD criteria

Concentration of contaminants

MSFD indicator

Concentration of the contaminants

Country

England

Ireland

Scotland

Part of a national marine monitoring programme

NO

YES

DESCRIPTOR 8 - Contaminants Do Not Produce Pollution Effects

Concentrations of contaminants are at levels not giving rise to pollution effects

Biota: Contaminants in the Marine Environment

Marine Institute (Custodian)

Concentration of contaminants > Concentration of the contaminants mentioned above, measured in the relevant matrix (such as biota, sediment and water) in a way that ensures comparability with the assessments under Directive 2000/60/EC



Contaminants in biota in the marine environment

Marine Institute (Pointofcontact)

Concentration of contaminants > Concentration of the contaminants mentioned above, measured in the relevant matrix (such as biota, sediment and water) in a way that ensures comparability with the assessments under Directive 2000/60/EC



Contaminants in sediments in the marine environment

Marine Institute (Pointofcontact)

Concentration of contaminants > Concentration of the contaminants mentioned above, measured in the relevant matrix (such as biota, sediment and water) in a way that ensures comparability with the assessments under Directive 2000/60/EC



Contaminants in water in the marine environment

Marine Institute (Pointofcontact)



Portal developed by the British Oceanographic Data Centre — Graphic design by POLAR 10, Cardiff
Centre for Environment, Fisheries and Aquaculture Science, Lowestoft Laboratory (Originator)
Portal developed by the British Oceanographic Data Centre — Graphic design by POLAR 10, Cardiff.

Conclusions & take home messages

- The RDF model (and associated technologies) are **important** and **powerful tools** for **discovery** and **delivery** of data and metadata
- Implementation can have a steep learning curve, & can be time consuming
- Until access to data via triplestores is widespread, end-users need to learn a new technology
 - Other types of API (e.g. pure JavaScript) will currently be more familiar to 3rd party developers but are often much less powerful
- Need to consider three user-groups: internal users, 3rd party developers, and external users
- The investment required to make and populate a triplestore of your data will be worth it – and your users will appreciate it in the long term!

Acknowledgments

- Justin Buck, Rob Thomas, Alexandra Kokkinaki
- Developers of D2RQ, Jena, Fuseki, and elda (and all the people who've answered various questions!)
- Funding: EU LIFE+ (LIFE011 ENV/UK/000392) via the Celtic Seas Partnership

Questions?

chwood@bodc.ac.uk
@c_wood

<http://linked.bodc.ac.uk>

<http://linked.bodc.ac.uk/documentation>

To go live in 2 weeks:

<http://resources.celticseaspartnership.eu>

Documentation

The screenshot shows a web browser window with the address bar containing the URL `linkeddev.bodc.ac.uk/documentation/`. The page title is "Basic documentation about BODC's NODB SPARQL API". The BODC logo is visible in the top left corner. The main content area includes a "Contents:" section with a list of links: "About the project", "Documentation", "Some basic queries", "Limitations of the service", and "Contact details for the project team". On the left side, there is a "Navigation" section with links to "About the project", "Documentation", "Some basic queries", "Limitations of the service", and "Contact details for the project team". Below this is a "Related Topics" section with a link to "Documentation overview" and a sub-link "Next: About the project". At the bottom of the page, there is a "Quick search" section with a search input field and a "Go" button. The Windows taskbar at the bottom shows the time as 12:20 on 07/07/2016.