



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Going for GOLD - Adventures in Open Linked Metadata

Citation for published version:

Reid, J, Waites, W & Butchart, B 2011, 'Going for GOLD - Adventures in Open Linked Metadata' Paper presented at AGI GeoCommunity 2011, Nottingham, United Kingdom, 20/09/11 - 22/09/11, .

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Early version, also known as pre-print

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Going for GOLD – Adventures in Open Linked Metadata

James Reid, Business Development and Projects Manager , EDINA, University of Edinburgh; Will Waites, University of Edinburgh; Ben Butchart, Senior Software Engineer, EDINA, University of Edinburgh;

Introduction

EDINA¹ as a JISC² national data centre makes available, *inter alia* many geographical datasets and services for use by the UK academic community. To support discovery of these resources it operates metadata catalogues that one can search and browse to find specific items of interest. These metadata are typically made available both in HTML form for human consumption and in a machine readable form, often as ISO XML for compliance with international obligations.

Catalogue systems which publish geospatial metadata (that is, information held in structured records that describe the characteristics of data and services), and which provide interfaces for their query and maintenance have been promoted by the Open Geospatial Consortium³ (OGC) for some time. The Catalogue Services (CS) specification which provides just such a standard method of defining, querying and organising stores of geospatial metadata (Nebert *et al.* 2007) is the *de facto* standard for discovering information about geospatial data and services in a standards compliant fashion. The CS specification defines an HTTP protocol binding named Catalogue Services for the Web (CSW) which underpins resource discovery across and within Spatial Data Infrastructures (SDIs) such as the UK Location Information Infrastructure. The UK academic sector has its own thematic SDI and already has its own geospatial discovery service through GoGeo⁴.

Independent of this, there is a growing trend to make such metadata catalogues available in a uniform way using techniques from Linked Data⁵ and the semantic web. Reasons for this include the ability to correlate metadata across different catalogue implementations that may use very different internal representations, to facilitate linking and annotation of the metadata by third parties and to provide a basic level of referential infrastructure on top of which more difficult questions of provenance and quality can be addressed.

Context

An example of the broader relevance of CSW and geospatial metadata for discovery purposes, is the recommendation issued in the context of INSPIRE⁶ by the INSPIRE Network Services Drafting Team (2008) to SDIs in the European Union to derive the base functionality of discovery services from the ISO profile of CSW. However, CSW is, arguably, not ideally suited for the modern Web where search engines are the users' default gateway to information⁷. The reasons why CSW might

¹ edina.ac.uk

² www.jisc.ac.uk

³ www.opengeospatial.org

⁴ www.gogeo.ac.uk

⁵ linkeddata.org/

⁶ INSPIRE defines the legislative framework for the establishment of a pan-European geospatial resource infrastructure. One of its key aims is to improve the 'discoverability' of geospatial resources by publication of resource metadata.

⁷ Several studies of research behaviour in the UK HFE sector support the view that 'googling' is a default reflex when seeking out resources. As noted by in 'Information behaviour of the researcher of the future', UCL, 2008, "they [resource providers] need to make their sites more highly visible in cyberspace by opening them up to search engines"

be regarded as sub-optimal from a purely 'discoverability' perspective are (adapted after Lopez-Pellicer *et al.* 2010):

- Search engines are poorly optimised to index Deep Web databases. The term 'Deep Web' refers to the database content that is effectively hidden from search engines behind Web forms and back-office applications. Surfacing Deep Web content is a research problem that has concerned the search engine community since its description by Bergman (2001). From this perspective, the underlying content of SDI metadata repositories are opaque and hidden behind catalogue application interfaces; as such, SDI metadata forms part of the Deep Web. Hence, the 'findability' via search engines depends on the success of crawling processes that require the analysis of the Web interface, and then the automatic generation of appropriate queries.
- Applications are increasingly becoming Linked Data friendly. The Linked Data community, which has grown significantly over the last few years, promotes a Web of data based on the architectural principles of the Web (Bizer *et al.*, 2008). Linked Data is less a technology than a set of best practices for publishing, sharing and connecting data and information using Uniform Resource Identifiers (URIs) that are resolved to Resource Description Framework (RDF) documents⁸. RDF is a W3C recommendation for modeling and exchanging metadata (Miller *et al.*, 2004). As an aside, the original Geography Markup Language (GML) model was based upon RDF and vestiges of this ancestry is still evident today. In the UK, the Cabinet Office has published guidance on the design of URIs⁹ which forms the basis for the UK Location Programmes approach to publishing Linked Data¹⁰ (see below).
- The evolution of metadata vocabularies towards RDF based models. Well known metadata vocabularies have evolved to models based on RDF with an emphasis on the linking of metadata descriptions. The abstract data models of Dublin Core Metadata Initiative (DCMI)¹¹ and the Open Archive Initiative (OAI)¹² have evolved alongside development of the RDF data model. This has resulted in abstract models based on the RDF data model (Nilsson *et al.* 2008; Lagoze *et al.* 2008) which emphasis the use (and reuse) of entities rather than using plain literals as the value of properties. This evolution ultimately enables the effective hyperlinking of metadata and traversal queries using query languages and protocols, such as SPARQL (Seaborne *et al.*, 2008).

Whilst CSW are undoubtedly useful to enable the discovery and provide access to geographic information resources within the geographic community (Nogueras *et al.*, 2005) and indeed are essential to the development of regional, national and global SDIs, they are nevertheless disjoint with the operational model of Deep Web crawlers. Popular search engines have developed several techniques to extract information from Deep Web databases without previous knowledge of their interfaces - Lopez-Pellicer *et al.* (2010) note that:

"The operational model for Web crawlers described in Raghavan (2001), based on (1) form analysis, (2) query generation and (3) response analysis is widely accepted. It models queries as functions with n named inputs $X_1..X_n$. where the challenge is to discover the possible values of these named inputs that return most of the content of the database. This approach is suitable for CSW HTTP GET requests. However, the constraints are encoded in a single named input as a CQL string (see Nebert *et al.*

⁸ <http://www.w3.org/RDF/>

⁹ "Designing URI Sets for the UK Public Sector" downloadable at:
<http://www.cabinetoffice.gov.uk/sites/default/files/resources/designing-uri-sets-uk-public-sector.pdf>

¹⁰ See e.g. <http://location.defra.gov.uk/wp-content/uploads/2011/03/INSPIRE-UK-Location-and-Linked-Data.pdf>

¹¹ <http://dublincore.org/documents/dcmi-terms/>

¹² <http://www.openarchives.org/>

2007), or an XML Filter (Vretanos, 2004). This characteristic is incompatible with the query model of the Deep Web crawlers. Researchers working for search engines, such as Google (see Madhavan et al. 2008), discourage the alternative operational model that is based on the development of ad-hoc connectors as non-sustainable in production environments.”

The upshot of this, is that geospatial metadata is not as 'open' as it potentially could be because of the formalised constraints imposed by the (pre-)existing geospatial metadata query standards. For example, the GoGeo CSW metadata repository effectively resides behind an opaque interface¹³ from the point of view of other (non-geospatial) communities. The CSW interface does not define a simple Web API to query and retrieve metadata. Some communities that potentially could use CSW are accustomed to simple APIs and common formats for purposes of 'mash-up'. For example, many geo-mashups and related data services (see Turner, 2006) use Web APIs to access and share data built following the REST architectural style (Fielding, 2000) . These APIs are characterized by the identification of resources by opaque URIs, semantic descriptions of resources, stateless and cacheable communication, and a uniform interface based on the verbs of the HTTP protocol which sits in opposition to the style adopted by CSW.

Approach

The work described below addressed the perceived shortcomings in the by producing Linked Data from extant ISO19115/19139 records. To do so we evaluated a number of alternative strategies for producing the RDF:

1. Metadata crosswalking. There are several geographic metadata crosswalks to the Dublin Core vocabulary which may be viewed as the lowest common denominator metadata baseline. We adopted the use of a Dublin Core crosswalk to implement uniform mappings from geographic metadata schemas to the RDF data model. This approach consists of three steps:
 - Apply a metadata crosswalk from the original metadata schema (ISO) to the Dublin Core vocabulary.
 - Add additional metadata such as provenance of the record, original information model or crosswalk identification.
 - Apply the profile for expressing as RDF the metadata terms.
2. An alternative approach was to publish direct from the relational data store underpinning the metadata resource , direct to RDF. An extension to this approach was to produce the RDF direct from the Unified Modelling Language (UML) representations of the underlying schemas using a visual modeling approach.

Results

For the purposes of the project, we worked with two types of metadata catalogues:

1. Catalogue Services for the Web,(CSW) services, important particularly because they are required for the implementation of the EU INSIPRE directive. We harvested, amongst others, the Scottish Location Information Discovery Service CSW for this purpose.

¹³ GoGeo also supports a Z39.50 target interface but this is not heavily used.

2. Customised catalogues that themselves aggregate data from various sources as exemplified by GoGeo and are typically implemented on top of a relational database schema.

These two separate catalogue implementations, lend themselves conveniently to the two alternative strategies for RDF production. In the case of CSW the crosswalk approach is the obvious solution whilst for database based schemas, schema mapping and UML derivation approaches seemed more appropriate.

For the CSW "metadata crosswalk" approach, we applied XSLT transforms which are generally appropriate where data is available in a predictable, standard form and is ideally suited for a circumstance where administrative or other elevated network privileges are not available or practicable – for example the data is available via an HTTP API but connections to the underlying (SQL) database are not possible i.e. the CSW is the proxy interface to the Deep Web. For the purposes of establishing a Linked Data production flow-line, the metadata are harvested and stored in an intermediate database (triplestore) and it is against this intermediate database that queries are performed and representations (RDF and others) are published.

A Note on URI Governance

Implicit in any approach that mints URIs are certain assumptions. Firstly, that a URI can be clearly and unambiguously defined and assigned to identify particular resources. Secondly that URIs are stable and persist. Axiomatically, the essence of a *persistent* URI is its immutability. Permanence implies long term commitment and ownership which presupposes some established mechanism for governance i.e. some authority has to set the rules and own the URIs which are used to identify things. The aforementioned Cabinet Office paper, "Designing URI Sets for the UK Public Sector", endeavoured to address this issue but its adoption and adaptation to location (geospatial) information in "[Designing URI Sets for Location](#)"¹⁴ highlighted an issue of political sensitivity – specifically the fact that in its original presentation the guidance neglected to allow for nationally determined and nationally specific URI schemes. At time of writing, the guidance is being recast to allow for Scottish, Welsh and Northern Irish URI naming schemes to be adopted if required by the devolved administrations¹⁵.

A Note on Target Representation¹⁶

As far as possible we worked with well known and widely used vocabularies (a set of agreed terms). The core vocabulary is, Dublin Core Terms¹⁷, hereafter simply referred to as DC. Unfortunately, DC does not contain a notion of a collection except indirectly as an underspecified range for `dc:isPartOf` and does not contain any particular way to refer to a representation of an abstract dataset such as a downloadable resource like a CSV or Shape file. It also contains no particular mechanism either to talk about people and organisations who might be authors or maintainers of datasets or to represent geographical or temporal extent. For all else, bar the question of geographical extent there are, fortunately, solutions that are more or less well

¹⁴ http://location.defra.gov.uk/wp-content/uploads/2010/04/Designing_URI_Sets_for_Location-Ver0.5.pdf

¹⁵ Scottish Government are adopting a `http://{resource}.data.scotland.gov.uk` pattern where `{resource}` maps to the United Nations Classification of Functions of Government (COFOG) headings, using a `cofog01`, `cofog02` approach rather than the descriptive label e.g. "general public services", "defence", etc. To improve the overall quality of this approach work is being undertaken with the UN and Rensselaer (<http://rpi.edu/>) to produce a SKOS RDF version of the COFOG.

¹⁶ Throughout, in the examples given below, Turtle notation is used (<http://www.w3.org/TeamSubmission/turtle/>). Except where confusion and ambiguity might arise, prefix declarations are generally omitted and well known or standard prefixes are used.

¹⁷ We have studiously avoided the use of the legacy Dublin Core Elements namespace.

For this reason where the prefix `dc:` appears it is taken to refer to the Dublin Core Terms namespace.

established. For describing data catalogues there is a specific vocabulary which is specifically designed to augment DC called the Data Catalog Vocabulary or DCat. Of particular interest are the concepts of `dcat:Catalog`, `dcat:CatalogRecord` and `dcat:Distribution`. Also used is `dcat:Dataset` which is simply an alias for `dc:Dataset`.

Further predicates are used for expressing other metadata such as keywords and spatio-temporal granularity. For referring to people, either natural or corporate, common practice is to use the Friend-of-a-Friend or FOAF vocabulary¹⁸. Where the precise nature of the entity in question is unclear `foaf:Agent`¹⁹ was used. Where it is clear that a natural person is being referred to, the more specific `foaf:Person` was used and if it was an organisation we used `foaf:Organisation`. For more involved descriptions of people and organisations and their relationship, the Organisation Ontology[?] may be used. This vocabulary provides facilities for describing the role that a particular person may fill in an organisation, for example.

Before addressing questions of how to represent spatial data, (which are far from settled), consider the fictitious example catalogue in Figure 1.

```
@prefix ex: <http://example.org/>

ex:catalogue a dcat:Catalog;
  dc:description "An example catalogue";
  dcat:keyword "Examples";
  dc:record ex:rec1, ex:rec2, ex:rec3.

ex:org a foaf:Organisation;
  foaf:name "Fictitious Inc.";
  foaf:mbox <mailto:someone@example.org>.

ex:sa a foaf:Organisation;
  foaf:name "Space Agency of Somewheria".

ex:rec1 a dcat:CatalogRecord;
  dc:modified "2011-07-01"^^xsd:date;
  dc:maintainer ex:org;
  dcat:dataset [
    a dcat:Dataset;
    dc:identifier "ABCD-EF12-3456-7890-DCBA";
    dc:modified "1984-03-23"^^xsd:date;
    dc:contributor ex:sa;
    dc:title "Pictures of crop circles from space";
    dcat:distribution [
      a dcat:Distribution;
      dc:description "Shape file download of crop circles";
      dcat:accessURL <http://download.example.org/crops.shp>
    ], [
      a dcat:Distribution;
    ]
  ].
ex:rec2 a dcat:CatalogRecord;
...
ex:rec3 a dcat:CatalogRecord;
```

Figure 1. An example catalogue

This example is not intended to be a complete representation of all the metadata that may be contained in such a catalogue but is intended to give a clearer idea of the structure. The distinction between the description of a catalogue record and the description of the dataset itself may seem pedantic but is in fact quite important. Frequently, the metadata record may be changed or updated in isolation, even when no change to the dataset itself has been made. The dataset may be maintained by one person or organisation and this person may have no influence whatsoever over the metadata catalogue. Separating out the concepts is important in order to be able to express this difference. This separation is well known and is already expressed through

¹⁸ <http://xmlns.com/foaf/spec/>

¹⁹ a subclass of `dc:Agent`

existing geospatial metadata fields. It is worth noting the lack of an explicit URI to identify the dataset itself. This is for notational convenience in the above example. In practice, where a dataset has no natural dereferenceable URI, it is given a non-resolveable one in the urn:uuid namespace (where it has an identifier that is a conformant UUID). This non-resolveable URI is essentially a constant that is used to refer to the dataset in third-party documents. However in the current project we were concerned with publishing *catalogue records* and not *dataset descriptions* as such and for expediency (elaborated on in the section on the use of named graphs in the metadata crosswalk section below) it was more straightforward to adopt this approach - although it is not a necessary feature of the representation.

Geographical Metadata

When we talk about geographical metadata we typically mean an expression of the coverage (spatial) area of a particular dataset. Most commonly this will be a (minimal) bounding box but it may be, in principle, a shape of arbitrary complexity. Whilst there is well established practice for representing point data in RDF, there does not appear to be any such consensus when it comes to representing even a shape as simple as rectangle²⁰. What is reasonably certain is that the dc:spatial predicate should be used to link from the dataset to its spatial extent. However the range of this predicate is simply a dc:Location and is not further specified other than such an entity must describe "a spatial region or named place".

One approach, that taken by NeoGeo²¹ is to attempt to create a completely granular "native" RDF description of the geometry. This is not incorrect by any means but was regarded as inappropriate for several reasons. Firstly, there is no support in any current software for presenting such data to the user as a map rendering. Secondly, correspondence with the authors of NeoGeo suggests that they are primarily interested in publishing such geometries as first class entities in themselves where as we assume that the geometries in our application are unlikely to be of primary interest outside of the context of the datasets which they annotate. Lastly, with an eventual eye to publishing these metadata using an RDF database that is aware of geospatial datatypes, the complexity involved in creating indexes for data expressed in this way can be considerable.

The Ordnance Survey on the other hand has opted to encode any geographical component on data which they publish simply by encoding fragments of Geography Markup Language (GML) as literals (Goodwin *et al.* 2009). This works well in that it is easy to index and many more tools exist that understand GML natively - at least relative to the "native" RDF approach.

A third approach is that presented in the GeoSPARQL W3C member submission²² which, though the specification document conflates vocabulary with query language extensions which could potentially be better presented as two separate documents, allows for both approaches.

In our implementation we opted to construct literals using the Well Known Text (WKT)²³ and annotate them provisionally with terms from the GeoSPARQL namespace in the hope that the useful parts of this vocabulary will be incorporated at a later date into a more widely used standard i.e. as the standards are still undergoing comment and revision cycles, it is too early to

²⁰ Where by "rectangle" is meant the shape described by geodesics in a standard coordinate system like WGS84 implied by the coordinates of two diagonally opposite points
aka MBR or minimum bounding rectangle

²¹ <http://geovocab.org/doc/neogeo.html>

²² See <http://geosparql.org/>

²³ See <http://www.geoapi.org/3.0/javadoc/org/opengis/referencing/doc-files/WKT.html>

say which approach represents the 'best' for future proofing purposes and consequently we have applied a 'best guess' logic. Considering the example in Figure 2, though we have used the WKT as the lowest common denominator representation of geospatial data, there is room alongside for other equivalent representations in e.g. GML as the Ordnance Survey does - or indeed, in expressly materialised granular RDF as with NeoGeo. Yet it also retains the properties of being easy to display to users using common tools and easy to index when importing into an RDF database.

```
[
  a dcat:Dataset;
  dc:spatial [
    a geo:Geometry;
    geo:asWKT "<http://www.opengis.net/def/crs/ogc/1.3/CRS84>
      POLYGON((12.9709 52.3005, 12.9709 52.6954,
        13.8109 52.6954, 13.8109 52.3005))"^^geo:WKTLiteral
  ]
].
```

Figure 2. Geographical metadata example.

Metadata Crosswalk

In some sense any transformation technique could be called a "metadata crosswalk". The specific meaning here is a transformation that is done on a particular document (in this case a catalogue record) rather than at the level of a database or collection of such documents. The approach turns on a 'harvesting' mechanism similar to that used more generally across the European spatial data infrastructure and elsewhere. Algorithm 1 is executed periodically against a data source.

```
procedure Harvest(source, start)
  for xml in source modified since start do
    rdf ← Transform(xml)
    Store(rdf)
  end for
end procedure
```

Algorithm 1. Harvesting algorithm

Retrieved documents are transformed using an XSLT transform and are then stored in an RDF database²⁴. For this project a specialised CSW client was written and this implements the query and actual fetching of ISO19139 documents from a server such as Geonetwork²⁵ (which is used in national and thematic based catalogues in the UK).

The storage step is also important in that it must extract from the intermediary RDF data a suitable graph name, a URI that will be used to both identify a catalogue record and to group the statements that belong in this record. It is clear from Figure 1 that, because a catalogue record is "several levels deep" it is not sufficient to just, say, consider all statements with a given subject in order to obtain a complete catalogue record. In order to save on the expense of complex queries, the data relating to a particular record is therefore grouped together into a named graph during harvesting. This also simplifies updating records that have changed since it is merely necessary to replace the entire graph rather than do complex queries to determine exactly which statements need to be retracted before any new ones are added.

The bulk of the logic, however, is done by the transformation step. In this case it uses an XSLT stylesheet²⁶ to transform ISO19139 XML data into RDF/XML which may then be stored. The

²⁴ In our case using 4Store (<http://4store.org/>)

²⁵ <http://geonetwork-opensource.org/>

²⁶ The current version of this stylesheet is available at https://bitbucket.org/ww/gold/src/tip/static/csw/md_metadata.xsl

structure of this stylesheet is straightforward, however it is perhaps appropriate to give some account of some of the specific mappings that we made. Salient mappings that we adopted are reproduced in Table 1.

ISO19139	RDF
gmd:MD_Metadata	dcat:CatalogRecord
gmd:fileIdentifier	dc:identifier
	also used in construction of URI
gmd:contact/gmd:CI_ResponsibleParty	foaf:Agent
gmd:identificationInfo/gmd:MD_DataIdentification	dcat:Dataset
gmd:identificationInfo/srv:SV_ServiceIdentification	dcat:Distribution
gmd:distributionInfo/gmd:MD_Distribution	dcat:Distribution
gmd:CI_Citation/gmd:title	dc:title
gmd:CI_Citation/gmd:abstract	dc:abstract
gmd:MD_Keywords	dcat:keyword
gmd:MD_LegalConstraints	dc:accessRights
	dc:rights
gmd:MD_Constraints	dc:rights
gmd:MD_Resolution	dcat:spatialGranularity
	dcat:temporalGranularity
gmd:EX_GeographicBoundingBox	dc:spatial
	geo:Geometry
gmd:EX_TemporalExtent	dc:temporal
	time:Interval

Table 1. Crosswalk mapping of ISO1939 elements to common RDF vocabularies.

It should be noted, however, that this table is greatly simplified for in many cases it makes sense to produce more than one RDF statement, particularly when one considers type assertions and necessary bits of indirection, for what appears as perhaps a single element in the ISO19139 (a good example of this is the temporal extent where the expression of a time interval in RDF can be quite verbose), and vice-versa where it takes many elements (at least eight) in the ISO19139 to express a geographical bounding box but because we opted for the simpler representation, the RDF requires only three.

In our testing we experimented not only with CSW services run and maintained by ourselves but also those run by others, we found some anomalies in some elements. For example, to distinguish between various types of dates, e.g. publication or modification date, it appears that the controlled vocabulary and localisation support in many catalogues is incomplete for this purpose. This required a number of special cases in the transform, e.g. "publication" vs. "publicatie" as found in Dutch catalogues for example as appears in the gmd:dateType/gmd:CIDateTypeCode/codeListValue elements. This further complicated the transformation requiring manual intervention where automated harvesting and RDF production proved problematic. Nevertheless, this approach worked surprisingly well and provided us with a

generic pipeline for harvesting geospatial metadata from any CSW and producing Linked Data outputs, albeit by making some default choices on the vocabularies used and representation simplification.

Database and UML to RDF Approaches

Database Mapping Approach

For our second approach to producing RDF Linked Data, the database mapping approach, we used the popular D2R software²⁷. D2R is a toolset for publishing relational databases on the Semantic Web. It enables RDF and HTML browsers to navigate the content of the database, and allows applications to query the database using a formal query language - SPARQL (*SPARQL Protocol and RDF Query Language*)²⁸. In this approach, queries are expressed in the SPARQL query language (whether expressly composed or implied by the act of fetching or dereferencing a resource identifier) and are transformed to more traditional SQL queries that are then executed against the underlying database, the results being returned according to the standard practice for SPARQL query results. The major aspect for consideration here is the mapping between the database schema and the RDF constructs expected in the outputs. These mappings can become quite onerous to hand-craft and we developed an approach to making this less tedious by extending the capabilities of a visual UML tool²⁹.

Transforming the catalogue to Linked Data on the fly by translating requests into SQL queries against the relational database in which it is stored and then transforming the results has some advantages. Though it is slower than querying native RDF storage directly, chiefly because SPARQL queries may entail expensive self-joins when translated into SQL, there is no need to provision a large RDF database server to host what is essentially a cache of the original data. There is also no need to coordinate updates or periodically harvest the source. It may also be, with a well-normalised relational schema and judiciously chosen set of indexes, that the relational database has very much smaller resource requirements than an equivalent RDF store - this is particularly relevant when there is a significant volume of data involved. Realistically though, geospatial metadata repositories are seldom of sufficient size to make this a *prima facie* concern. Data currency is more of a consideration and in instances where the underlying geospatial metadata store is being regularly changed, it may be more appropriate to use a D2R type approach than try to re-synchronise and lock-step harvested resources. Data custodians need to balance the trade-off between periodic harvesting and change frequency in order to ensure overall data currency. This is of course a generic issue and not limited solely to the type of data publication work explored here.

The D2R server is a single Java program that takes an input file describing the coordinates of the database to be translated and the translation rules. The process for configuring it is relatively simple:

1. Generate a crude configuration file using the generate-mapping program
2. Refine and customise the mapping (the hard part!)
3. Run the d2r-server

The result is an HTTP service that supports simple HTML and RDF representations of the resources in the database and a SPARQL endpoint through which arbitrary queries may be made. For production use a reverse proxy or embedding into a Java servlet system may be considered desirable.

²⁷ <http://sourceforge.net/projects/d2rq-map/>

²⁸ <http://www.w3.org/TR/rdf-sparql-query/>

²⁹ We used Enterprise Architect.

There were however two main difficulties encountered with this approach. The first is that whereas working from a widely published international standard as the source data format (say ISO), we can be reasonably confident that it will not change in the short term and therefore any mapping from that to Linked Data also will not need to change, the same is not true for information in a relational database where internal business rules and *ad hoc* schemas often prevail. The main use of relational databases is as internal data storage and as such there are frequently no external constraints on the schema changing. As it stands, if the schema changes for whatever reason, the mapping for the D2R server must, of necessity, be revisited. Whilst this did not happen during the lifetime of this project, it can be reasonably expected that it will happen from time to time in the future as the services are extended or modified.

The other difficulty is related to the first. A further property of stable, agreed international standards as data sources, is that they tend to be relatively well documented. On the other hand, it is much less common to have the same level of documentation for a relational database schema that is intended primarily for internal use. Consequently, the process of refining the mapping configuration for the D2R server is not trivial if the schema is at all complicated or subject to frequent changes. The relative merits of which approach to use – crosswalking or derivation direct from a database, are ultimately bound to issues of available infrastructure, skillsets, data currency etc. Table 2 provides a summary of the two approaches.

Approach used	Crosswalk	D2R
Data currency	Cached	Live
Computing resources	RAM-bound for RDF DB	CPU and network bound
Moving parts / Administrative complexity	<ul style="list-style-type: none"> • Harvesting / update machinery • RDF storage • HTTP/REST service 	Just D2R Software
Configuration complexity	Simple	Complex
Skillsets for customisation	<ol style="list-style-type: none"> 1. RDF (Turtle, RDF/XML) 2. XSLT (for customising the transformation) 3. HTML+JS+CSS (for customising human-readable output) 4. Go language (for customising the behaviour of this particular implementation) 	<ol style="list-style-type: none"> 1. RDF (Turtle, RDF/XML) 2. Java (customisation of behaviour, deployment, appearance) 3. HTML+JS+CSS (customisation of behaviour)
Look and Feel	Easy customisation of HTML files	Somewhat more complex
SPARQL queries	Faster	Slower

Table 2: Side-by-side comparison of Crosswalk and D2R based approaches

Production from UML – a UML Profile for D2R scripting

As already noted, the task of republishing data deposited in relational database as Linked Data principally involves two main technical challenges.

1. Converting the tables, columns and relationships in the relational model to appropriate RDF constructs (that is, RDF classes and properties)
2. Mapping columns and relationships to appropriate existing Linked Data vocabularies.

This is the approach described above and is actually a fairly common problem which has lead to the evolution of tools such as D2R. These tools have emerged to assist dataset providers in mapping relational models to RDF and tend to treat the relational database as a virtual RDF graph enabling access through APIs such as Sesame³⁰ and Jena³¹. While the D2R platform greatly facilitated semantic access to data in our geospatial relational database, the creation of the necessary D2R mapping scripts is tedious and error prone, as no visual authoring tools are currently available. As many UML authoring tools already have good support for relational database modelling , we extended our UML tool to support simple D2R constructs, thereby automatically generating D2r scripts from within a visual editing environment. Our approach was to define an XML profile with D2R constructs that could augment existing UML stereotypes used in standard database modelling (<table>, <column> etc). The data modeller first loads the relational model into the UML tool. Most UML tools have a simple import mechanism to load tables into the tool direct from the relational store. Then, the data modeller imports the 'D2R profile' (our GOLD extensions), to expand their toolbox. Theses imported D2R tools provide access to constructs such as "ClassMap" and "PropertyBridge" that can then be dragged onto the tables and columns in the UML model to define the mappings. Once the modeller has completed all the mappings they can export the UML model to XMI³². The output XMI contains all the information that the data modeller specified in the UML modelling tool, including the stereotypes and TaggedValues from the D2R profile. An XSLT stylesheet can the be applied to the XMI file to generate the D2R mapping in RDF format. In the case of Enterprise Architect it was possible to specify the XSLT stylesheet as part of the XMI export so the export and xsl transform can be combined into one step. Our experiences with this approach suggests that for Linked Data production direct from a relational data store, particularly ones where table structure is more than very simple, this visual editing support for mapping production is both more intuitive from the data modellers perspective and significantly, less error prone than conventional 'hand-crafting' approaches. It is also more adaptable and robust to underlying schema changes and more flexible where multiple target output schemas are required. Our D2R profile extensions are available on request.

Conclusion

In endeavouring to expose the Deep web, we have taken a Linked Data approach and published geospatial metadata as RDF. We have explored alternate options for RDF generation – from cross-walking through well known vocabularies such as Dublin Core, to RDF generation direct from a relational database. In either case, there is an expectation that the outputs will be consistent yet the vagaries of geospatial data (quality aspects) mean that establishing which approach is more flexible or robust in any particular application instance, is necessarily subject to trial and error. We have established a basic workflow and infrastructure which support CSW harvesting (from any CSW) and automated Linked Data publication of that geospatial metadata by adopting well known

³⁰ <http://www.w3.org/2001/sw/wiki/Sesame>

³¹ <http://jena.sourceforge.net/>

³² The XML Metadata Interchange (XMI) is an [Object Management Group \(OMG\)](#) standard for exchanging metadata information

and frequently used vocabularies e.g. FOAF, DCat. An open question remains as to whether or not Linked Data is the panacea to resource discovery and reuse that its proponents assert. A significant issue to overcome is to establish core, common vocabularies – particularly in respect to alternate and competing approaches to the representation of geometry information.

References

- Bergman, Michael K. (2001) White Paper: The Deep Web: Surfacing Hidden Value, The Journal of Electronic Publishing (JEP), Volume 7, Issue 1, August, 2001. Available from: <http://hdl.handle.net/2027/spo.3336451.0007.104>
- Bizer, C.; Heath, T.; Idehen, K. and Berners-Lee, T. (2008) Linked data on the web (LDOW2008) WWW '08: Proceeding of the 17th international conference on World Wide Web, ACM, 1265-1266
- Fielding, R. T. (2000) REST: Architectural Styles and the Design of Network based Software Architectures University of California, Irvine
- Goodwin, J.; Dolbear, C. and Hart, G. (2009) Geographical Linked Data: The Administrative Geography of Great Britain on the Semantic Web Transactions in GIS, doi: 10.1111/j.1467-9671.2008.01133.x
- Lagoze, C. and de Sompel, H. V. (2008) ORE User Guide – Resource Map Implementation in RDF/XML Open Archives Initiative. Available from: <http://www.openarchives.org/ore/1.0/rdfxml>
- Lopez-Pellicer F. J., Florczyk A. J., Nogueras-Iso J, Muro-Medrano P. R. and Zarazaga-Soria F. J. (2010) Exposing CSW catalogues as Linked Data, Lecture Notes in Geoinformation and Cartography (LNG&C). Geospatial Thinking. 2010, vol. , p. 183-200. ISSN 1863-2246
- Madhavan, J., Ko, D., Kot, L., Ganapathy, V., Rasmussen, A. and Halevy, A. (2008) Google's Deep Web crawl, Proceedings of the VLDB Endowment, VLDB Endowment, 1, 1241-1252
- Miller, E. and Manola, F. (2004) RDF Primer. W3C, Available from <http://www.w3.org/TR/2004/REC-rdf-primer-20040210/>
- Nebert, D., Whiteside, A. and Vretanos, P. A. (2007). Open GIS Catalogue Services Specification. OpenGIS Publicly Available Standard, Open GIS Consortium Inc.
- Nilsson, M.; Powell, A.; Johnston, P. and Naeve, A. (2008) Expressing Dublin Core metadata using the Resource Description Framework (RDF) [online] Dublin Core Metadata Initiative, DCM
- Nogueras-Iso, J., Zarazaga-Soria, F.J., Muro-Medrano, P.R. (2005) Geographic Information Metadata for Spatial Data Infrastructures: Resources, Interoperability and Information Retrieval. Springer-Verlag New York, Inc., Secaucus, NJ, USA
- Raghavan, S. and Garcia-Molina, H. (2001) Crawling the Hidden Web VLDB '01:Proceedings of the 27th International Conference on Very Large Data Bases, Morgan Kaufmann Publishers Inc., 129-138
- Seaborne, A. & Prud'hommeaux, E. (2008) SPARQL Query Language for RDF W3C. W3C Recommendation. Available from: <http://www.w3.org/TR/2008/REC-rdf-sparql-query-20080115/>
- Turner, A. (2006) Introduction to neogeography O'Reilly Media, Inc.
- Vretanos, P. A. (2004) OpenGIS Filter Encoding Implementation Specification Open Geospatial Consortium Inc.