

## Technologies

The Discovery programme has been funded by JISC to improve access to collections that support research and education. This document is part of a series that describes the lessons from 8 JISC projects funded under the Discovery programme in 2011 to explore open metadata for libraries, museums and archives. More information about the projects can be found at: <http://www.jisc.ac.uk/whatwedo/programmes/inf11/infrastructureforresourcediscovery.aspx>. The other documents in the series can be found at: <http://discovery.ac.uk>

### • Scope •

There are a large, and ever increasing, number of technologies relevant to the manipulation and publication of data. This guide summarises some of the key technologies used by the eight Discovery projects, and other related work.

### • The Problem Space •

This guide does not tackle a specific problem space, but rather attempts to present a snapshot of technologies and tools of relevance to the manipulation, publication and aggregation of metadata.

### • Outcomes •

#### • Data storage and related technologies

	Example Use	Related Tools
<b>MongoDB</b> A schema-less 'document-oriented' database, which stores JSON-style documents. MongoDB is designed to scale easily and offer high availability.  <a href="http://www.mongodb.org">http://www.mongodb.org</a>	The Jerome project <sup>1</sup>	ElasticSearch, CouchDB <sup>2</sup>
<b>ElasticSearch</b> A schema-less data store, index and search engine, which uses the JSON data format. ElasticSearch is built on Lucene (see below), and designed to offer a fast, easy to configure and highly scalable search engine.  <a href="http://www.elasticsearch.org">http://www.elasticsearch.org</a>	BibServer, part of the Bibliographic Knowledge Network toolset <sup>3</sup>	MongoDB, CouchDB, Lucene

<sup>1</sup> <http://jerome.blogs.lincoln.ac.uk>

<sup>2</sup> <http://couchdb.apache.org>

<sup>3</sup> <http://bibserver.okfn.org>

	Example Use	Related Tools
<p><b>Redis</b></p> <p>A key-value store, supporting keys containing strings, hashes, lists, sets and sorted sets. Redis keeps datasets in memory, with periodic writes to disk, making it extremely fast.</p> <p><a href="http://redis.io">http://redis.io</a></p>	<p>JISC Open Bibliography project for data analysis<sup>4</sup>, and to enhance SPARQL results with geolocation data in a data visualisation<sup>5</sup>.</p>	
<p><b>ARC2 + MySQL</b></p> <p>A PHP library for RDF. It supports the use of MySQL as an RDF triplestore.</p> <p><b>ARC2:</b> <a href="https://github.com/semsol/arc2/wiki">https://github.com/semsol/arc2/wiki</a></p> <p><b>MySQL:</b> <a href="http://www.mysql.com">http://www.mysql.com</a></p>	<p>COMET<sup>6</sup> used ARC2 + MySQL to host their linked data expression of bibliographic records from the University of Cambridge. While it proved easy to setup, speed and scalability were both issues.</p>	<p>OWLIM<sup>7</sup>, 4Store<sup>8</sup>, Virtuoso<sup>9</sup>, the Talis Platform, Sesame<sup>10</sup></p>
<p><b>Talis Platform</b></p> <p>A Software as a Service (SaaS) RDF triplestore</p> <p><a href="http://www.talis.com/platform">http://www.talis.com/platform</a></p>	<p>SALDA<sup>11</sup>, LOCAH<sup>12</sup></p>	<p>ARC2, OWLIM, 4Store, Virtuoso, Sesame</p>

## Flexible Infrastructure

	Example Use	Related Tools
<p><b>Amazon Web Services</b></p> <p>A set of services offering flexible, scalable, infrastructure. Specifically of note in this context, the 'Elastic Cloud Compute' (EC2) service offering deployment of virtual servers, and 'Simple Storage Service' (S3) offering disk space.</p> <p><a href="http://aws.amazon.com">http://aws.amazon.com</a></p>	<p>OpenART used AWS EC2 for easy deployment of virtual machines</p>	<p>Rackspace Cloud Hosting<sup>13</sup>, Heroku<sup>14</sup>, Google App Engine<sup>15</sup></p>

<sup>4</sup> <http://openbiblio.net/2010/11/18/characterising-the-british-library-bibliographic-dataset>

<sup>5</sup> <http://benosteen.com/globe>

<sup>6</sup> <http://cul-comet.blogspot.com>

<sup>7</sup> <http://www.ontotext.com/owlim>

<sup>8</sup> <http://4store.org>

<sup>9</sup> <http://www.openlinksw.com/dataspace/dav/wiki/Main/VOSIndex>

<sup>10</sup> <http://www.openrdf.org>

<sup>11</sup> <http://blogs.sussex.ac.uk/salda>

<sup>12</sup> <http://blogs.ukoln.ac.uk/locah>

<sup>13</sup> <http://www.rackspace.co.uk/cloud-hosting>

<sup>14</sup> <http://www.heroku.com>

<sup>15</sup> <http://code.google.com/appengine>

## Information Extraction

	Example Use	Related Tools
<p><b>OpenCalais</b></p> <p>An information extraction web service, including named entity recognition<sup>16</sup>. The Calais webservice accepts unstructured data as an input and returns entities, facts and events contained in the the data as RDF.</p> <p><a href="http://www.opencalais.com">http://www.opencalais.com</a></p>	Open Metadata Pathfinder (OMP) <sup>17</sup>	DBpedia Spotlight, GATE <sup>18</sup>
<p><b>DBpedia Spotlight</b></p> <p>An named entity recognition service based on data in DBpedia (linked data extracted from Wikipedia).</p> <p><a href="http://dbpedia.org/spotlight">http://dbpedia.org/spotlight</a></p>	DiscoBro <sup>19</sup> (an entry to the Discovery developer competition <sup>20</sup> )	OpenCalais, GATE

## Indexing and Search tools

	Example Use	Related Tools
<p><b>Sphinx</b></p> <p>An open source search server. It aims to offer very high speeds of both indexing and retrieval, and to be scalable up to indexing billions of documents. Sphinx can index from both SQL and NoSQL data stores as well as index files directly.</p> <p><a href="http://sphinxsearch.com">http://sphinxsearch.com</a></p>	The Jerome Project	Solr/Lucene, Zebra <sup>21</sup>
<p><b>Solr/Lucene</b></p> <p>Open source indexing and search software.</p> <p><a href="http://lucene.apache.org">http://lucene.apache.org</a></p>	<p>As ElasticSearch is built on Lucene, by extension this is used by BibServer</p> <p>Solr/Lucene is used in a wide range of 'resource discovery' software and services including VuFind<sup>22</sup>, Blacklight<sup>23</sup> and CultureGrid<sup>24</sup>.</p>	SolrMARC <sup>25</sup> , VuFind, Blacklight, ElasticSearch, Sphinx, Zebra

<sup>16</sup> [http://en.wikipedia.org/wiki/Named\\_entity\\_recognition](http://en.wikipedia.org/wiki/Named_entity_recognition)

<sup>17</sup> <http://openmetadatapathway.blogspot.com>

<sup>18</sup> <http://gate.ac.uk>

<sup>19</sup> <http://people.kmi.open.ac.uk/mathieu/about/discobro-discovering-linked-data-resources-while-browsing>

<sup>20</sup> <http://discovery.ac.uk/developers/competition>

<sup>21</sup> <http://www.indexdata.com/zebra>

<sup>22</sup> <http://vufind.org>

<sup>23</sup> <http://projectblacklight.org>

<sup>24</sup> <http://www.culturegrid.org.uk>

<sup>25</sup> <http://code.google.com/p/solmarc>

## Data manipulation and transformation tools

	Example Use	Related Tools
<p><b>MARC code libraries</b></p> <p>There are code libraries specifically for the reading, writing and manipulation of MARC records for a wide variety of programming languages.</p> <p>More information: <a href="http://wiki.code4lib.org/index.php/Working_with_MaRC">http://wiki.code4lib.org/index.php/Working_with_MaRC</a></p>	<p>Comet used MARC/Perl<sup>26</sup> to create the 'Marc<sup>27</sup> to RDF triples conversion utility' and the 'Marc record sorter' tools .</p> <p>The Lucero project<sup>28</sup> used MARC4J<sup>29</sup> (a Java library) to convert MARC records to RDF</p>	<p>pymarc<sup>30</sup> (Python library for MARC), ruby-marc<sup>31</sup> (Ruby library for MARC), File_MARC<sup>32</sup> (PHP library for MARC)</p>
<p><b>XSLT</b></p> <p>XSLT is a language or mechanism to transform XML documents in various ways, including into HTML, plain text, RDF and new XML documents. An XSLT provides the rules for the transformation, which can then be applied to XML using an XSLT processor such as Saxon<sup>33</sup>, Xalan<sup>34</sup> and xsltproc<sup>35</sup>.</p> <p><a href="http://www.w3schools.com/xsl">http://www.w3schools.com/xsl</a></p>	<p>Discovering Babel used XSLT<sup>36</sup> to transform information from the TEI XML format into Dublin Core</p> <p>LOCAH and SALDA both used XSLT to transform information from the EAD XML format commonly used for Archival description, to RDF</p>	<p>Xpath<sup>37</sup></p>
<p><b>Google Refine</b></p> <p>Google Refine is a tool for "working with messy data". It can help clean up, transform and extend data. It also supports 'reconciling' (i.e. matching) entities from one dataset to other, remote, datasets.</p> <p><a href="http://code.google.com/p/google-refine">http://code.google.com/p/google-refine</a></p>	<p>OpenART<sup>38</sup> used Google Refine to manipulate information, and to connect parts of the data using the reconciliation feature.</p>	<p>Stanford Data Wrangler<sup>39</sup></p>

<sup>26</sup> <http://marcpm.sf.net/>

<sup>27</sup> <http://data.lib.cam.ac.uk/code.php>

<sup>28</sup> <http://lucero-project.info/lb>

<sup>29</sup> <http://marc4j.tigris.org>

<sup>30</sup> <http://pypi.python.org/pypi/pymarc>

<sup>31</sup> <https://github.com/ruby-marc/ruby-marc>

<sup>32</sup> [http://pear.php.net/package/File\\_MARC](http://pear.php.net/package/File_MARC)

<sup>33</sup> <http://saxon.sourceforge.net>

<sup>34</sup> <http://xalan.apache.org>

<sup>35</sup> <http://xmlsoft.org/XSLT/xsltproc2.html>

<sup>36</sup> <http://blogs.oucs.ox.ac.uk/martinw/2011/08/03/discovering-babel-technical-issues>

<sup>37</sup> <http://www.w3schools.com/xpath>

<sup>38</sup> <http://yorkdl.wordpress.com/2011/11/02/openart-final-report>

<sup>39</sup> <http://vis.stanford.edu/wrangler>

## Big data

	Example Use	Related Tools
<p><b>Mahout</b></p> <p>A machine learning and data mining library designed to work with very large datasets.</p> <p><a href="http://mahout.apache.org">http://mahout.apache.org</a></p>	<p>The AEIOU Project used of Mahout for 'recommender systems'<sup>40</sup></p> <p>The Library Innovation Lab at Harvard<sup>41</sup> has been experimenting with analysing bibliographic records using Mahout and Gephi<sup>42</sup></p>	<p>Hadoop<sup>43</sup>, Gephi<sup>44</sup></p>

## Recommendations

- Schema-less approaches to data storage are a good fit to some of the data issues in the library, archive and museum sectors
- Information Extraction software and services offer routes to link traditional metadata into other data sources on the web
- Issues of storing, manipulating and analysing data are common across many domains and sectors, so existing tools and technologies may well be useful to library, archive and museum data
- 'Big data' and 'Visualisation' technologies will become more relevant to libraries, archives and museums as they wish to build and exploit large aggregations of activity data and metadata

## Key Discovery Projects

- **Comet** - <http://cul-comet.blogspot.com>
- **Discovering Babel** - <http://blogs.oucs.ox.ac.uk/martinw/2011/03/23/discovering-babel>
- **Jerome** - <http://jerome.blogs.lincoln.ac.uk>
- **OMP** - <http://openmetadatapathway.blogspot.com>
- **SALDA** - <http://blogs.sussex.ac.uk/salda>

## Other Related Work

- **Linked and Open Copac and Archives Hub (LOCAH)** - <http://blogs.ukoln.ac.uk/locah>  
A project which addresses how Linked Data can support cross-domain discovery of L&A material and beyond
- **Linking University Content for Education and Research Online (LUCERO)** - <http://lucero-project.info/lb>  
A project to publish Linked Data at the Open University (including bibliographic data) and to provide guidance to those wishing to do the same
- **AEIOU** - <http://www.wm.aber.ac.uk/aeiou>  
A project to increase the visibility and usage of Welsh academic research
- **JISC Open Bibliography** - <http://openbiblio.net/tag/jiscopenbib>  
Project to publish bibliographic metadata as Linked Open Data
- **Bibliographic Knowledge Network tools** - <http://bibserver.okfn.org>

<sup>40</sup> <http://www.wm.aber.ac.uk/aeiou>

<sup>43</sup> <http://hadoop.apache.org>

<sup>41</sup> <http://librarylab.law.harvard.edu>

<sup>44</sup> <http://gephi.org>

<sup>42</sup> <http://danbri.org/words/2011/10/11/720>