

Using Linked Data as a basis for a Learning Resource Recommendation System

Nadeem Shabir, Chris Clarke

Talis Group Limited: nadeem.shabir@talis.com, chris.clarke@talis.com

Abstract. Resource List Management Systems (RLMS) allow the electronic publication of course reading lists. Aside from electronic access, existing systems in this area provide little utility for teachers and learners above and beyond the traditional paper based reading lists. Our vision is that resource lists could in actual fact become Open Educational Resources that can be shared, re-mixed and re-used across institutions and borders. This paper introduces how we used linked data to architect a RLMS to meet this vision. However, in implementing this system, questions arose around the provenance, sustainability, licensing and reliability of today's linked data cloud. This paper documents the steps we took to address these criticisms in our implementation. The paper goes on to discuss how the ecosystem of learning data managed by this application opens the way for future work, which involves leveraging typed relationships between learning goals, educational resources and system actors to provide recommendation-like services for academics creating new content.

1 What is Linked Data

When Sir Tim Berners-Lee originally expressed his vision for the Semantic Web, he was imagining a Web of Data[1]:

I have a dream for the Web [in which computers] become capable of analyzing all the data on the Web the content, links, and transactions between people and computers. A Semantic Web, which should make this possible, has yet to emerge, but when it does, the day-to-day mechanisms of trade, bureaucracy and our daily lives will be handled by machines talking to machines. The intelligent agents people have touted for ages will finally materialize.

Berners-Lee would later go on to define some of the properties of this Web of Data, and in doing so coined the term 'Linked Data', which simply refers to a set of best practices for publishing and connecting structured data on the Web. Berners-Lee expressed these a set of simple rules[2]:

1. Use URIs as names for things
2. Use HTTP URIs so that people can look up those names

In the past 10 years, several specialist online systems have emerged to assist in the management of such lists, such as Sentient Discover³, Talis List⁴, LORLS⁵ and Blackwell Reading Lists online⁶. Such systems offered an online representation of the paper-based lists, as well as providing tools for the library to assist in stock acquisition.

However, current Resource List Management Systems (RLMS) have provided limited extended functionality for the teacher over and above that of the paper-based solution - save online access for students - this may explain their limited adoption, and the continued proliferation of paper handouts. Some allow linking to journal articles via institutional link resolvers, and for items the library physically holds, most allow linking to the library catalogue. However, these systems are simply signposting solutions, providing none of the added services that users of Web 2.0-like systems might expect, such as recommendation services, rich user interface metaphors or the integration or in-lining of the resources themselves, including full text, into the list.

Our research showed that often teachers construct these lists in a style which reflects both the chronological order and/or the major topic areas the course unit covers. Thus the structure of these lists, and the relative position of resources on it tells us something about their intended usage and how they relate to each other.

We also know that in authoring these lists, teachers are either explicitly or implicitly influenced by similar works by their peers. An example of explicit influence is the teacher that seeks out similar syllabi when trying to author their own. An example of implicit influence is where peers discuss the availability and quality of educational resources, which may later lead to their use (or not) within the classroom. In essence, the latter could be described as a professional variant of the water cooler effect⁷.

This made us consider the impact of creating a system that enhances the authoring process of the lists by making it possible to formally harness the existing work of peers, thus supporting either the creation of derived works (with appropriate attribution), or using them as the basis for content suggestions for authors of new lists within comparable subject areas.

In the development of our new system, Talis Aspire⁸, our vision was to create a system which would allow resource lists themselves to be considered and operated on as Open Education Resources (OERs)⁹. This means that they can be re-used, remixed, shared and collaborated on easily, supporting the notion that open access to knowledge is in the interests of all.

³ <http://www.sentientdiscover.co.uk>

⁴ <http://www.talis.com/list>

⁵ <https://lorls.lboro.ac.uk/>

⁶ <http://www.readinglists.co.uk>

⁷ <http://www.wordspy.com/words/watercoolereffect.asp>

⁸ <http://www.talis.com/aspire>

⁹ http://en.wikipedia.org/wiki/Open_educational_resources

3 Benefits of a Linked Data approach

A key objective that makes our vision workable is that a user should be able to easily discover appropriate content (or have it recommended) to re-use and remix. It follows then that resource lists must be richly and homogeneously described in order that lists from different authors are comparable. In addition, combination with datasets outside the system boundary become important in the creation of a data ecosystem which supports discovery and recommendation - or in other words, ease of discovery of new relationships.

The nature of RDF-based systems, such as those that underpin the datasets on the linked data web, make it easy to re-combine graphs of data from multiple sources, allowing these new relationships to be discovered.

We concluded early that not only would the system have to merge resource metadata from multiple and incompatible sources, but that each individual customer implementation of the system should be able to publish the resulting resource lists in a way that could be re-combined at a later date to enable re-use, remixing and sharing of data within a multi-institution ecosystem. Without the resulting scale that combining data from multiple institutions provided, any discovery or recommendation features within a particular subject domain would be of limited use.

Our experience with RDF and specifically linked data indicates suitability for richness of description, standardised publication, interlinking and interoperability between disparate sets of data. By settling on linked data principals, as described in an earlier paper by Clarke[10], the team were able to unify not only the description of resources using shared ontologies such as Bibliographic Ontology¹⁰, Resource List Ontology¹¹, SIOC¹² and FOAF¹³, but also on how the resource lists were to be published, allowing them to be combined with other data sources at a later date.

This approach is supported by one of the challenges that a recent JISC-funded report[11] suggests semantic technologies can address:

Information in UK HE/FE institution seems to be fragmented and in formats that makes it often inaccessible. Discovery of relevant information over a large number of sources needs to be supported. Information that is publicly available on the institutions Web pages is not available in machine processable formats making it difficult to compare programmes of study, syllabuses or research angles.

¹⁰ <http://bibliontology.com/>

¹¹ <http://vocab.org/resourcelist/>

¹² <http://rdfs.org/sioc/spec/>

¹³ <http://xmlns.com/foaf/0.1/>

4 Leveraging the data ecosystem to support the discovery and recommendation of content in an Open Education context

Given that one could interconnect resource metadata used to construct resource lists between departments, schools and even institutions, one could discover which modules cite the textbook *Financial Accounting and Reporting* (Elliot & Elliot). Unifying the description of those modules, one could discover if the textbook was largely being cited on 1st year Business Studies courses, or if it was actually a core text on most MBA programmes. What resources usually appear alongside Elliot & Elliot on resource lists? Combine with this knowledge about how students actually use the text, (for example, do they purchase it or do they ignore it) and multiply this knowledge across all disciplines and resources used for learning and it is conceivable that one could create advanced, context-aware recommendation systems for a multitude of use cases.

For example, when building the list, teachers can use the system to help them predict the impact of including a text on a particular resource list. They can discover, and aggregate, which resources are routinely included by their peers, or how the majority of students choose to use resources. If they choose to include an item which is not held by the library, the system could suggest similar items that are held. Additionally, it is conceivable that they could browse a repository of lists in a related subject area, licensed as OERs, and use them as a basis for their own work.

Our work to date has focused on seeding an ecosystem of resource list data as a basis for future work in developing discovery and recommendation functionality we describe above. To date we have six UK HEIs using the system with plans to expand to a further twenty five during 2009/10. It is our assumption that this is around the lower limit required before the functionality described is viable for the end user.

What follows is the set of techniques we intend to blend together to realise the above scenarios, thus unlocking the potential of the ecosystem of data we have built.

4.1 Explicit hierarchical classification

Chandrasekaran stated that hierarchical classification was one of the generic tasks that must be addressed when designing expert systems[13]. To build effective recommendation systems, it is important to know the context of educator-selected resources - for example, a list for level one students for a module entitled *Introduction to Clinical Psychology*, is unlikely to have much in common with data from lists around the topic of *Organisational Behavior* at the same institution, so similarity between resources should be weighted much lower than those from equivalent level one *Clinical Psychology* courses delivered at other institutions.

The system uses the AIISO ontology¹⁴ to organise lists into a tree hierarchy at the institutional level, allowing one to trace the module, programme, department and institution that a list and its resources belong to. However, when taking a view across the whole ecosystem, one cannot use this mechanism to map equivalent lists at different institutions. By augmenting the AIISO descriptions with data from the Joint Academic Coding System (JACS)¹⁵ we now have a basis for comparison at the course level. Mixing in the level of each academic programme gives us further data to complete our hierarchical classification system.

In addition, it is our intention to ask teachers to optionally indicate their principal subject areas in their profile. This allows other users to locate their profile and subscribe to updates about resources they have recently pulled into the system (their bookmarks) and lists they have made available under OER-compatible licenses. The aim is to replicate the watercooler effect inside the application - teachers can follow the implicit recommendations made by their peers - an example being the inclusion of a resource on a resource list.

4.2 Deriving similarity without explicit classification

Where no explicit classification of a resource list is made, it is possible to derive similarity between a set of given lists by analyzing the pattern of resource usage (inference). The results can be used to populate a similarity index to aid the discovery of lists in a related topic area, or to recommend individual resources within a particular domain.

A very naive example is as follows: List 1 at location A contains a resource X. List 2 at location B contains a resource Y. List 3 at location C contains resource X and states Y as an alternative should X be unavailable. By merging data from all locations, we can discover new relationships between list 1 and list 2, even though they contain no shared resources and thus no direct links, even in the merged result set. We can suggest to the author of list 1 that Y could be a relevant resource, and given a high density of equivalent resources, we could make an assertion that lists 1 and 2 are potentially similar to a teacher pursuing a discovery use case.

The power of linked data here is that although the sophistication of the inference algorithm can be increased, the merging of datasets across institutions remains trivial.

4.3 Wisdom of crowds

By allowing teachers on-mass to re-mix, reuse and share lists to create derivative works, we introduce the wisdom of crowds into the system. The level of an individual teacher's impact in his subject area could be formulated as a product of the quantity of derivative works in the ecosystem. In essence, we enable the most attributed content to float to the top of the pile.

¹⁴ <http://vocab.org/aiiso/>

¹⁵ http://www.hesa.ac.uk/dox/jacs/JACS_complete.pdf

When combined with the ability to follow the actions of others, we provide the context for wider social network functionality within the system.

4.4 Other inputs

Taking into account the behavior of students can further feed into our recommendation algorithm. Both explicit and implicit inputs can give some indication as to similarity between resources. The application allows users to explicitly rate their intention to an item (a range from ‘Intend to purchase’ to ‘Won’t use’). In addition, a feed of loan history from the library can be used to discover potentially equivalent items, loaned at the same time as borrowing the resources on a given list.

As we have less context information for these inputs, they can only be used as supplementary indicators in the recommendation algorithm.

Other inputs can be taken from links to other points in the Linked Data cloud - section 6 discusses how resource lists link into the Library of Congress Subject Headings and the Linked Periodicals datasets - these links can provide further hints to our algorithms.

5 Critique of the state of the art of Linked Data

To date, the linked open data movement has conducted excellent work and demonstrated that publishing linked open data on the web is technically feasible. However, for those wanting to develop commercial enterprise systems, such as Talis Aspire, that operate over the data contained within the linked open data cloud, several several social and legal hurdles remain. These are especially relevant to our context, where we wish to leverage linked data to enhance learning through the machine-driven discovery and recommendation of OER content.

1. Sustainability - Many linked open data endpoints are maintained by hobbyists or by projects or programmes with time-limited funding. Developers of applications that either link to these data sources or extract data from them, and will rely on them, need to be confident that data source will continue to be available or that the contents of the service will be available for others to run, or mirror, if the original host disappears.
2. Provenance - How does a user trace the provenance of a piece of data represented on the linked data web? Currently when we view a graph of data it could be the aggregation of information from a number of different data sets. This ability to combine data is one of the unique and key properties of RDF. However, once aggregated together, it is not easy to identify how each data set contributed to the aggregated view. This becomes problematic where users want to be able to weight assertions according to the party that has made them, or indeed exclude them all together. For example, to a student trying to achieve a particular learning goal, learning resources recommended by a first year student would carry less weight than those recommended

by the professor of the course, assuming that the latter is considered more authoritative.

3. Licensing - Re-use of OERs relies on clear licensing to explain to consumers can and can't do with the original work. However, a large proportion of the data in the linked data cloud is not specifically licensed at all. This is exacerbated by the fact that there is currently no agreed protocol which allows a machine agent to interpret license terms attached to a given set of data.
4. Reliability - The nature of the document web is that sites frequently disappear, leading users to experience, expect and work around 404 errors. However, developers of linked data applications must be sure that the core data sets required for main flows will be available and able to meet the demands of their application. Developers of linked data applications need to ensure their systems will still operate if linked sources are temporarily or permanently unavailable.

6 Addressing sustainability

What happens to a data set if the company that published it suddenly becomes bankrupt? What can other businesses and their applications that are dependent on that data do?

Where a formal relationship exists between a service provider and a customer, the two parties can enter into an escrow agreement. In the software world these agreements are based around source code or data. It is unclear how this can apply when there is no formal relationship between the publisher of the data and its consumers, as is often the case with linked data.

The team was forced to find a different solution to this problem, one example of which we detail here. As mentioned earlier Talis Aspire contains vast amounts of bibliographic data, many of these entities describe periodicals, such as journals. Each Talis Aspire implementation might describe the same periodical, as such there is value in linking these descriptions together so that we know they are referring to the same entity.

We decided that whilst it was technically feasible to create linkages between each of the Talis Aspire datasets, it would make more sense for each of the Talis Aspire implementations to link to a single, authoritative, dataset (or linking hub) that only contained title level information about periodicals. These two approaches are contrasted in Fig.2. Crucially, whilst this proposed dataset would be of benefit to us it would remain application agnostic and would therefore have value to others.

At the time there was no dataset available to link to, but there were a number of available sources of data, mostly provided in the form of CSV files that could be converted and published as linked data, because the data was already in the public domain or with permission from the owners. We immediately recognised the value in engaging with the data owners and others as part of a community interested in curating and using this dataset.

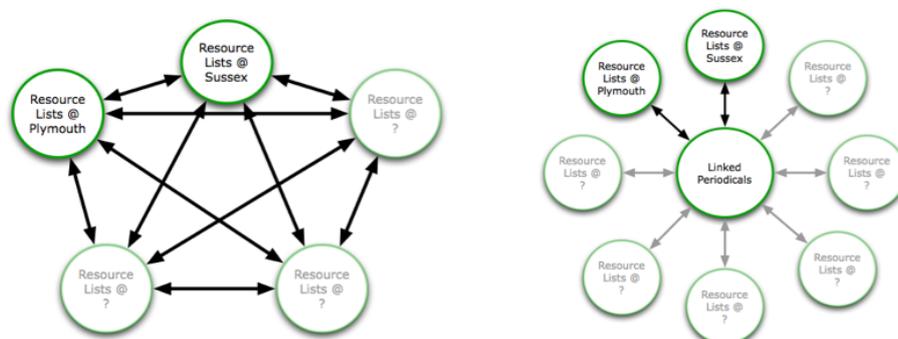


Fig. 2. Linking Talis Aspire tenancies directly to each other versus linking them to Linked Periodicals

We initiated an incubation project on Data Incubator¹⁶, called ‘Linked Periodicals’¹⁷. Data Incubator helps organise communities around particular datasets and leverages the skills and experience within the community to convert the original data and publish it as linked data. The code to perform the conversion is open sourced so that the community and indeed the data owner can repurpose it. In fact, one of the the major goals of each of these individual incubation projects is to provide tools that the original data owner can take, adapt with ease, and then use to emit the linked data themselves. If they are unwilling, or unable to do so themselves they can engage others in the community who are willing to do so. As with most community led initiatives an open license for use of the data would be recommended.

Data Incubator provides a temporary home for the dataset whilst the community works on the conversion and discusses¹⁸ use cases for the data as well as ways to link the dataset to other existing datasets in the linked data ecosystem thereby, both, demonstrating and increasing its utility. For example we recently discussed, and succeeded in linking the Linked Periodicals dataset to the Library of Congress Subject Headings¹⁹ data, as illustrated in Fig.3.

Interestingly, the LCSH subject headings were originally provided as linked data²⁰ by a single individual who believed there was utility in making the data available as linked data. The service became popular amongst the community, however the Library of Congress asserted its rights as the data owner and forced the service to be shut down²¹. The Library of Congress then went on to publish the data itself. However we believe that the history around how the Library of

¹⁶ <http://dataincubator.org/>
¹⁷ <http://periodicals.dataincubator.org/>
¹⁸ <http://groups.google.com/group/dataincubator>
¹⁹ <http://id.loc.gov/authorities/>
²⁰ <http://lcsch.info/>
²¹ <http://lcsch.info/comments1.html>

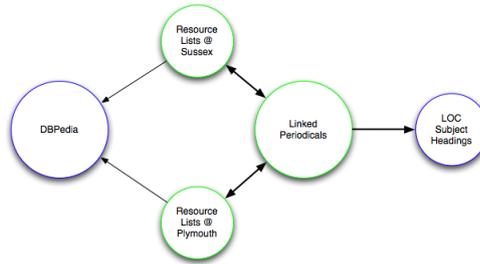


Fig. 3. Linking Talis Aspire, Linked Periodicals, DBPedia and LCSH

Congress Subject Headings were made available as linked data, serves to illustrate the need for initiatives/processes like Data Incubator, where data owners and communities interested in that data can work together to agree how and under what terms it should be made available.

Whilst this does not address all the issues of sustainability, our belief is that a motivated, open, community interested in that dataset is likely to be willing to curate, maintain and indeed, demonstrate the value of that dataset. Members of this community could enter into a ‘living will’ style arrangement for the data, so that if the original host, does disappear, or is no longer able to, then members of community could continue to provide the data, in perpetuity.

7 Addressing data provenance

With any published data there is a need to provide additional meta information describing the dataset itself. This meta information enables data consumers to make informed decisions about the quality of the data and to determine whether they want to trust and use that data[8]. Recently, The Vocabulary Of Interlinked Datasets (voID)[9] was published to specifically define the terms and best practices needed in order to categorise and provide statistical meta information about data sets as well as the ‘linksets’ connecting them. Fig. 4 illustrates an example of how you might use voID to describe a data set.

In addition to voID, one can infer provenance via the domain name used to form the resource URIs within the datasets. By using URIs stemmed from the registered domain name of the institution, which are in turn sub-domains of top level domains such as .edu and .ac.uk that can only be issued to academic institutions, the data generated by the application is essentially watermarked to an officially recognised educational body. Fig. 5 is an example of some basic data about a resource list using University of Plymouth’s domain. In this example, the data uses the plymouth.ac.uk domain, and according to linked data principles, an agent obtains this RDF/XML description by resolving the URI <http://lists.lib.plymouth.ac.uk/lists/abf203.rdf>, which can be safely assumed to only be under the ultimate control of The University of Plymouth.

```

:DBpedia a void:Dataset ;
foaf:homepage <http://dbpedia.org/> ;
void:subset :DBpedia2DBLP
dcterms:license <http://www.gnu.org/copyleft/fdl.html>
void:statItem [
  rdf:value 20000;
  scovo:dimension void: numberOfWorkResources ;
  scovo:dimension foaf:Person ;
  dcterms:source <http://wiki.dbpedia.org/> ;
] .

:DBLP a void:Dataset ;
foaf:homepage <http://dblp.13s.de/d2r/> ;
dc:subject dbp: Computer_science ;
dc:subject dbp:Journal ;
dc:subject dbp:Proceedings .

:DBpedia2DBLP a void:Linkset ;
void:subjectsTarget :DBpedia ;
void:objectsTarget :DBLP ;
void:linkPredicate owl:sameAs .

```

Fig. 4. An example void description of a data set including statistical information and licensing

```

<http://lists.lib.plymouth.ac.uk> a list:List ;
sioc:name "Financial Accounting and Reporting";
sioc:parent_of <http://lists.lib.plymouth.ac.uk/sections/abf203-1> ;
sioc:parent_of <http://lists.lib.plymouth.ac.uk/sections/abf203-2> .

```

Fig. 5. Watermarking data using the institution's domain name

The Talis Aspire application only allows users authenticated by the institution's Devolved Authentication (DA) infrastructure to write new descriptions that can be published back out as linked data. This therefore closes the provenance loop - that is, University of Plymouth would only allow educators (presumably in the Business School), authenticated using the institution's DA, to access the Talis Aspire system and create resource lists about Financial Accounting.

8 Addressing data licensing

As more and more data sets are published and made available online as linked data, there is an increasing need for the owners of those datasets to make clear what the terms of use are for that data, and make explicit which rights they are willing to exert or to waive[6].

As described in an earlier paper, the importance of attaching specific licensing terms to published data cannot be understated[10]. Specific licenses have been developed for data, The Open Data Commons Public Domain Dedication and Licence[7] (ODC PDDL) is available for publishers of data to explicitly gift data to the commons.

In the specific example of Linked Periodicals, we are working with the data owners to ensure that we can apply ODC PDDL license, opening up the way for Talis, and others, to innovate with this data. In choosing which other datasets

may inform our recommendation algorithms, we seek clarification first as to the appropriate use of data.

In our application, we provide a mechanism for explicitly licensing OERs with a set of pre-defined licenses. The system then can enforce some of the terms of these licenses because it has specific knowledge (in the form of application code) that governs both the operations a user can perform and the propagation of attribution.

One major barrier is that the user cannot arbitrarily assign a license that is unknown to the system - this is because there is no agreed protocol for machine interpretation of license terms, allowing software agents to determine at run time how users can operate over, and attribute linked data.

9 Addressing reliability

Individuals or organisations that are seriously committed to publishing linked data must realise that like any other web service their linked data endpoints must be reliable, by that we mean they must be available and able to meet increasing demand, particularly if the data set is popular. Not all organisations have the infrastructure internally to support services like these which may be expensive to host and maintain.

Software as a service and cloud computing models obviously have a role to play in addressing the reliability, availability and scaling issues, by enabling data publishers to rapidly build out a publishing infrastructure that will support these operations without significant and prohibitive capital investments. Whilst computing power²² and data storage services²³ are widely available, services that directly provide support for this particular niche, linked data publishing, have not yet become common place. However, the increasing maturity of semantic web standards and techniques has led to the emergence of managed platforms, such as the Talis Platform²⁴, which can enable mainstream data owners such as the BBC to publish linked data²⁵.

10 Summary

It is our firm belief that without embracing linked data, we cannot fully realise our vision of an RLMS that considers resource lists as OERs, providing seamless discovery of and access to education resources.

We have described Talis Aspire, a RLMS that is built on linked data principles. This application contributes to an ecosystem of open learning resource data, as well as interlinking to other datasets. This ecosystem can now be used as a basis for building advanced discovery and recommendation systems.

²² <http://aws.amazon.com/ec2/>

²³ <http://aws.amazon.com/s3/>

²⁴ <http://www.talis.com/platform>

²⁵ <http://blogs.talis.com/n2/archives/569>

We identified several major obstacles around the use of linked data in commercially deployed learning systems. In describing these problems, and the steps we have currently taken to address them, we recognise that they have not been completely resolved. There is still much work that needs to be done, for example the provisioning of machine readable licenses that can be embedded in data.

We hope that we can continue to engage with the community on forwarding the debate and building on our work in these areas, and support future initiatives that unlock and interlink relevant sources, in addition to resource list data, to deliver on the wider aspiration of a linked web of open education data.

References

- [1] Berners-Lee, T.: Weaving the Web: The past present and future of the World Wide Web by its inventor. Orion Business Books, London (1999).
- [2] Berners-Lee, T.: Linked Data - Design Issues. Retrieved June 27, 2009, <http://www.w3.org/DesignIssues/LinkedData.html> (2006).
- [3] Berners-Lee, T.: Linked Open Data. Retrieved June 27, 2009 [http://www.w3.org/2008/Talks/0617-lod-tbl/#\(3\)](http://www.w3.org/2008/Talks/0617-lod-tbl/#(3)) (2008)
- [4] Bizer, C., Cyganiak, R., Heath, T.: How to publish Linked Data on the Web. Retrieved June 27, 2009, <http://www4.wiwiw.fu-berlin.de/bizer/pub/LinkedDataTutorial/> (2007).
- [5] Bizer, C., Lehmann, J., Kobilarov, G., Auer, S., Becker, C., Cyganiak, S., Hellmann, S.: DBpedia - A Crystallization Point the Web of Data. Journal of Web Semantics (JWS), Special Issue on the Web of Data. Preprint: <http://www.wiwiw.fu-berlin.de/en/institute/pwo/bizer/research/publications/Bizer-et-al-DBpedia-CrystallizationPoint-JWS-Preprint.pdf> (2009).
- [6] Miller, P., Styles, R., Heath, T.: Open Data Commons, a License for Open Data. Proceedings of the 1st Workshop about Linked Data on the Web (LDOW) (2008)
- [7] The Open Data Commons Public Domain Dedication and Licence. Retrieved June 27, 2009, <http://www.opendatacommons.org/odc-public-domain-dedication-and-licence/>
- [8] Hartig, O.: Provenance Information in the Web of Data. Proceedings of the 2nd Workshop on Linked Data on the Web (LDOW2009). (2009)
- [9] Alexander, K., Cyganiak, R., Hausenblas, M., Zhao, J.: Describing Linked Datasets On the Design and Usage of void, the Vocabulary Of Interlinked Datasets. Proceedings of the 2nd Workshop on Linked Data on the Web (LDOW2009) (2009)
- [10] Clarke, C.: A Resource List Management Tool for Undergraduate Students based on Linked Open Data Principles. Proceedings of the 6th European Semantic Web Conference (ESWC2009) (2009)
- [11] Tiropanis, T., Davis, H., Millard, D., Weal, M., White, S., Wills, G.: JISC - SemTech Project Report. (2009)
- [12] Jackson, P.: Introduction to Expert Systems (3rd Edition). Addison Wesley. (1998)
- [13] Chandrasekaran, B.: Generic Tasks in Knowledge-Based Reasoning: High-Level Building Blocks for Expert System Design. IEEE Expert, 1(3):23-30. (1986)

14 Nadeem Shabir, Chris Clarke

[14] Dix, A., Beale, R., Wood, A.: Architectures to make Simple Visualisations using Simple Systems. Proceedings of AVI2000, ACM Press:51-61. (2000)