

Revyu.com: a Reviewing and Rating Site for the Web of Data

Tom Heath and Enrico Motta

Knowledge Media Institute, The Open University,
Walton Hall, Milton Keynes, MK7 6AA, United Kingdom
{t.heath, e.motta}@open.ac.uk

Abstract. Revyu.com is a live, publicly accessible reviewing and rating Web site, designed to be usable by humans whilst transparently generating machine-readable RDF metadata for the Semantic Web, based on their input. The site uses Semantic Web specifications such as RDF and SPARQL, and the latest Linked Data best practices to create a major node in a potentially Web-wide ecosystem of reviews and related data. Throughout the implementation of Revyu design decisions have been made that aim to minimize the burden on users, by maximizing the reuse of external data sources, and allowing less structured human input (in the form of Web2.0-style tagging) from which stronger semantics can later be derived. Links to external sources such as DBpedia are exploited to create human-oriented mashups at the HTML level, whilst links are also made in RDF to ensure Revyu plays a first class role in the blossoming Web of Data. The site is available at <<http://revyu.com>>.

1 Introduction

Revyu.com is a live, publicly usable (and used!) reviewing and rating Web site developed using Semantic Web technologies and standards, and according to the principles of Linked Data [1]. Reviews and ratings are widely available on the Web and are one major form of Web2.0-inspired *user-generated content*. However, despite the availability of review data through Web2.0 APIs such as Amazon Web Services, these reviews largely remain in isolated silos and in formats that prevent their easy integration and interlinking with data from other sources. This presents considerable barriers to the aggregation of all reviews of a particular item from across the Web. As has been recognised by several previous authors [2, 3], the Semantic Web, or Web of Data, provides a technological platform with which to overcome this problem. Revyu takes a significant and concrete step towards solving this problem by exposing reviews as Linked Data using standards such as RDF and SPARQL. In doing so it helps to create an ecosystem of interlinked reviews and ratings on the Web, and to bootstrap the Semantic Web as a whole.

2 Revyu Overview

Revyu allows people to write reviews and give ratings of things simply by filling in a Web form. The site takes an open world view of the reviewing process by not constraining users to reviewing items from a fixed database; anything they can name can be reviewed, whilst links supplied with the review can disambiguate items thanks to inverse functional properties such as *foaf:homepage*. As of July 2007 Revyu has been live for 9 months, attracting 381 reviews from exactly 100 reviewers.

The mode of interacting with the site will be familiar to those who have read or written reviews at sites such as Epinions¹ or Amazon²: users create reviews by filling in Web forms; these are then available immediately on the site in HTML. Whilst this functionality is not especially novel, as a reviewing application Revyu improves significantly over other work in the area in the following ways: it goes well beyond the closed world 'silos' of sites such as Epinions and TripAdvisor, by exposing reviews in a reusable, machine-readable format; it improves upon the APIs of sites such as Amazon by using a more flexible data format (RDF), allowing more versatile queries via SPARQL, and linking to external data sources; lastly it is not restricted to reviews and ratings in one domain, as is the case with Golbeck's FilmTrust [3].

Revyu is built from the ground upwards on Semantic Web technologies. By also following Linked Data principles the site ensures that reviews it hosts can be fully connected into a Web of Data. This approach manifests itself in a number of ways. All site content, in addition to being available in HTML, is also published in RDF/XML that is interlinked with the corresponding HTML pages but available as separate crawlable documents. This creation and publication of RDF is invisible to the reviewer, enabling novice users to contribute data to the Semantic Web through a familiar, Web2.0-style mode of interaction. To date this approach has yielded over 12,000 RDF triples publicly available on the Semantic Web. Whilst not a large figure by many standards, it is significant that these triples have been generated primarily from direct user input, rather than by data mining or extraction from natural language.

In addition to the review data, RDF is published on the site describing the reviewers, tags that they associate with reviewed items, and the reviewed items themselves. These descriptions use the FOAF [4] and Tag [5] ontologies, in addition to properties and classes from RDFS and OWL. This data is also available via the Revyu SPARQL endpoint³, which enables programmatic query access to the underlying triplestore. Providing such a query interface allows third parties to retrieve reviews and related data in a flexible fashion, for reuse in their own applications. Whilst in some ways analogous to Web2.0 APIs that provide remote query capabilities, SPARQL endpoints afford many advantages to the developer: for example, common libraries can be used to query multiple RDF graphs yet return the results as one resultset, effectively allowing joins over multiple data sources.

Underlying these features of Revyu is a technical infrastructure that enables us to follow current best practices in serving Linked Data. In the following section we will detail the Revyu architecture and discuss decisions made in implementing the system.

¹ <http://www.epinions.com/>

² <http://www.amazon.com/>

³ <http://revyu.com/sparql/welcome>

3 Revyu Architecture and Implementation

Revyu is implemented in PHP, and runs on a regular Apache web server. RDF processing capabilities are provided by RAP, the RDF API for PHP [6], with RDF data persisted to a de-normalised MySQL database following the RAP database schema. The Revyu SPARQL endpoint relies on the RAP SPARQL engine, which operates against the same MySQL-based triplestore.

From the outset Revyu was designed to adhere to the four 'commandments' of Linked Data outlined by Tim Berners-Lee [1]: using URIs as names for things, using HTTP URIs so people can look up those names, providing useful information when someone looks up a URI, and linking to other URIs so more things can be discovered.

All things represented on Revyu are assigned URIs: reviews, people, reviewed things, tags assigned to things, and even the bundles that represent tags assigned by one person at one point in time. Providing URIs for all these things gives many items a presence on the Semantic Web which they would not otherwise have, and enables any third party to make reference to these items in other RDF statements. This opens the way for links between Revyu and other data sets, thereby helping to lay the foundations for a Web of Data.

All URIs in the Revyu URI-space can be dereferenced. Attempts to dereference the URIs of non-information resources receive an HTTP303 "See Other" response, along with a URI of a document containing a description of the resource. This configuration is inline with the W3C Technical Architecture Group's finding on the httpRange-14 issue [7], and serves to reinforce the distinction between a resource and a description of that resource. Content negotiation (implemented as Apache *rewrite rules* according to the recipes in [8]) is carried out on Revyu URIs, whereby the user agent is redirected to either an HTML or RDF document that describes the resource, depending on the value of the Accept header sent in the initial request.

4 Deriving Semantics from Tagging Data

When creating Revyu, a significant decision was taken to not require users to classify the items they were reviewing, but instead to associate keyword tags with the item. This decision was taken for several reasons: firstly there was seen to be a lack of sufficiently comprehensive classifications of items that users may want to review; secondly, we did not want to require all users to subscribe to a single classification scheme for reviewed items as this seemed unnecessarily constraining and against the spirit of the Semantic Web; thirdly, providing an interface for classification using arbitrary types discovered in ontologies on the Semantic Web was seen as being hard to implement in a way that would be usable for non-specialists; and lastly, the availability of ontologies on the Web was thought insufficient to provide adequate coverage and therefore was likely to result in a more closed world of reviewed items.

The availability of Yago [9] class URIs via DBpedia [10] in recent months has gone some way to addressing a number of these issues, and we will be investigating this in future work. However we believe that the tagging route retains the appropriate balance of usability whilst also providing sufficient data from which stronger

semantics can be derived. At present we use tagging data in two ways: to identify basic semantic relationships between tags and to derive type information about a reviewed item.

Tags that reviewers frequently associate with the same item are assumed to be related in some way. In the HTML pages about each tag, tags that co-occur above certain threshold are displayed to the user. This threshold is set low for HTML output, as human readers of the page are unlikely to infer erroneous information based on these relationships. In contrast however, relationships exposed in RDF descriptions of tags (using the *skos:related* property) are based on a more conservative threshold, in order to avoid erroneous inferences based on these assertions. Unfortunately at present we are unable to determine more precise relationships (such as superclass/subclass) between tags, as this would require substantial further processing. However, we hope to investigate this issue in more detail in the future.

We currently derive type information from tagging data in two domains, books and films, also relying on external data sources to help ensure accurate results. Firstly, where items are tagged 'book' we parse Web links provided by the reviewer that relate to the item, and attempt to extract ISBN numbers embedded in these links. Where we are able to extract an ISBN number in this fashion we conclude that the reviewed item is in fact a book, and assert a corresponding *rdf:type* statement into the triplestore.

If an item has been tagged 'film' or 'movie', we execute a query against the DBpedia SPARQL endpoint⁴ in order to find any entries of type *yago:Film* that have the same name as the reviewed item. If a match is found then we conclude this item is in fact a film, and add an *rdf:type* statement to this effect to the triplestore. These type statements for both books and films are exposed in the RDF descriptions of items on Revyu, and also used as the basis for showing additional relevant data in the HTML pages about an item, as detailed in the following section.

Validating Revyu data against external sources in this way not only allows the derivation of more robust type information than would be possible using tags alone, it also allows us to link items on Revyu with items elsewhere on the Web of Data. Where matches are found, we use the *owl:sameAs* property to assert that two URIs identify the same resource. Publishing these links in RDF connects Revyu in to a growing Web of Linked Data, signified in particular by initiatives such as the Linking Open Data community project [11].

5 Production and Consumption of Linked Data

As described above Revyu adheres to the principles of Linked Data. This is demonstrated by the linking of Revyu data with heterogeneous external data sources (such as DBpedia⁵, Open Guides⁶, and FOAF data) wherever possible. Doing so helps create a Web of Data rather than simply isolated islands of RDF; Revyu data is *in the Web*, not just *on the Web*. This provides a distributed, interlinked data set managed by

⁴ <http://dbpedia.org/sparql>

⁵ <http://dbpedia.org/>

⁶ <http://openguides.org/>

a wide range of parties, on which others can begin to build applications, and into which other data sets can also be linked.

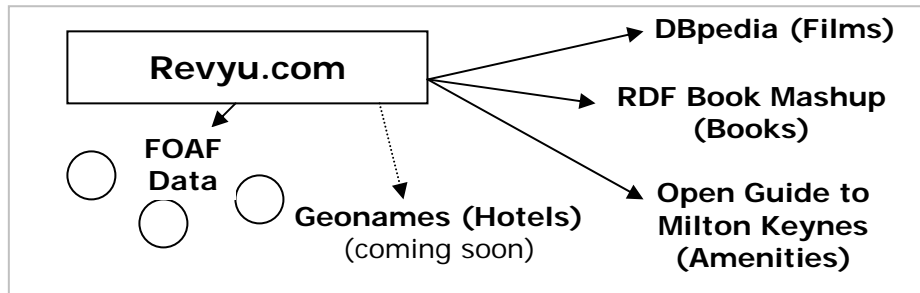


Fig. 1. Links from Revyu.com to external data sets

We also actively exploit the links we set between Revyu and external data sources to enhance the experience of our users, without placing an additional burden on reviewers. For example, users registering with the site are not asked to provide copious information to populate their user profile. Instead, where they have an existing FOAF description in an external location they may provide its URI, in which case Revyu dereferences this URI and queries the resulting graph for relevant information (such as a photo, location, home page address, and interests), which is then displayed on the their profile page, as illustrated in Fig. 2. This approach reduces the burden on the user by not requiring them to manage multiple redundant sets of personal information stored in different locations. Furthermore, where the user has assigned themselves a URI in their FOAF description, Revyu sets *owl:sameAs* links asserting that this URI identifies the same resource as the user's Revyu URI. Recently we have also implemented a feature allowing users to state that they know other Revyu reviewers, and add these people to their social network. This relationship is then recorded in the triplestore using the *foaf:knows* property, and exposed in the user's RDF description on the site.

The process of matching films and books on Revyu with those in external data sources, as described above, also allows *owl:sameAs* links to be set between these items. Based on *owl:sameAs* links from Revyu to DBpedia films we retrieve additional information about the film, such as the URI of the films promotional poster, and the name of the director. This information is displayed on the Revyu HTML page about the film (as shown in Fig. 3), thereby enhancing the value of the site for users without requiring this information to be manually entered into Revyu.

Similarly we exploit links between Revyu and the RDF Book Mashup [12] as the basis for retrieving book cover and author information which is also then displayed on the Revyu HTML page about the book. In the RDF descriptions of items we take a slightly different approach, choosing to simply expose the links between items but without republishing RDF data from external sources. This approach could be described as using Semantic Web data to produce Web2.0-style mashups at the human-readable, HTML level, whilst also creating linked data mashups at the RDF level. Not only does this linked data approach to mashups reduce issues with licensing of data for republication, it is also a more Web-like approach; duplicating data is of

much lesser value than linking to it, and the user agent of the future should be able to 'look ahead' to linked items and merge data accordingly.

About tom (Tom Heath)



[tom's Home Page](#)

tom's location:
[Borough of Milton Keynes](#)

tom's Interests

[Semantic Web](#)
[Word of Mouth](#)
[Trust](#)
[Beer](#)

Fig. 2. Excerpts from the first author's Revyu profile page, showing data sourced automatically from his external FOAF file⁷

Broken Flowers



directed by [Jim Jarmusch](#)

Broken Flowers is a 2005 comedy-drama film directed and written by

Fig. 3. Excerpts from the Revyu HTML page about the film *Broken Flowers*, showing the film poster, director information, and summary drawn from DBpedia⁸

It should be noted, we do not claim that the Revyu Web2.0-style mashups achieve something that could not have been done using conventional Web2.0 approaches. What distinguishes our approach are the following: the simultaneous publishing of data and human-oriented mashups, so that the data integration effort we have invested is not lost but can be reused by other parties; the ability to easily integrate

⁷ <http://revyu.com/people/tom/about/html>

⁸ <http://revyu.com/things/broken-flowers-film-movie-bill-murray-jim-jarmusch-sharon/about/html>

heterogeneous sources using RDF; and the substantially reduced development costs in producing human-oriented mashups through use of Semantic Web technologies.

Whilst we currently wait for new film reviews on Revyu and then attempt to automatically match them with entries in DBpedia, we will shortly be creating skeleton records in Revyu covering 12,000 films described in DBpedia. The skeleton records will simply include the title of the film, a statement indicating that this item is of type 'Film', a number of keyword tags, and links to the corresponding item on DBpedia. Not only will this provide a foundation on which new reviews can be created, it will also ensure that all films being reviewed in the future will already be interlinked with the corresponding DBpedia entry, and thus the Web of Data.

This approach was followed when linking Revyu to data from the Open Guide to Milton Keynes⁹, a member of the Open Guides family of wiki-based city guides that provide data in RDF. Milton Keynes is a city in south east England, and home of The Open University. Whilst some amenities in the city, such as pubs and restaurants, were already reviewed on Revyu, many more were listed in the Open Guide due to its longer history. Therefore, after identifying items existing in both locations and making the appropriate mappings to avoid duplication, we created skeleton records in Revyu for the remaining items, setting links back to their Open Guide URIs. This now enables latitude and longitude data for many items to be retrieved from RDF exposed by the Open Guide, and used to show a Google Map of the items location. The same approach can also be used to expose address, telephone, and opening time information held in the Open Guide. We will shortly be extending this generic method to linking Revyu with data from other Open Guides, such as London and Boston.

Whilst frequently suggested, at present there are no plans to import external review data into Revyu, for a number of reasons. Firstly little review data is available under a suitable license; secondly our ongoing research is predicated on the ability to combine review data with social networks, requiring some global identifier (such as *foaf:mbox_sha1sum*) to be available for each reviewer, which is rarely the case with traditional reviewing sites. Unfortunately, to the best of our knowledge Revyu is the only site serving reviews as Linked Data according to current best practices, which also limits our abilities to interlink Revyu with external review data sets.

6 Future Work and Conclusions

In addition to encouraging further user participation in order to increase the value delivered by the site, we plan to integrate Revyu with a number of additional data sets. Most notably we are preparing to create skeleton records in Revyu of 70,000 hotels worldwide, linked to their corresponding entry in the *Geonames* dataset. Additional data will be integrated as further relevant sources become available. It should be noted that our aim in linking to external datasets is not to constrain, but merely to seed, users conceptions of what can be reviewed. As we integrate further data sets we hope to achieve a more automated linking process by investigating generic similarity matching techniques for operation on the wider Semantic Web.

⁹ <http://miltonkeynes.openguides.org/>

Furthermore, by providing reviews in a reusable format that is easily integrated and interlinked with other data, Revyu provides core data for our ongoing work into information seeking, recommendation, and trust in social networks on the Web.

In conclusion, in this paper we have described Revyu.com, a human usable reviewing and rating Web site built on Semantic Web technologies, and fundamentally designed to contribute to the realization of a Web of Data. Whilst superficially not unique in functionality, the site is rare in its status as a publicly available service in daily use, that is oriented towards human users but also embodies current best practices in developing for the Semantic Web.

Acknowledgements

This research was partially supported by the Advanced Knowledge Technologies (AKT) project. AKT is an Interdisciplinary Research Collaboration (IRC), which is sponsored by the UK Engineering and Physical Sciences Research Council under grant number GR/N15764/01. Peter Coetzee did a superb job of turning data into skeleton records for import into Revyu. Lastly, the Open Guides and DBpedia communities and the RDF Book Mashup team deserve our special thanks.

References

1. Berners-Lee, T.: (2006) Linked Data. <http://www.w3.org/DesignIssues/LinkedData.html>
2. Guha, R.: Open Rating Systems. In: Proc. 1st Workshop on Friend of a Friend (2004)
3. Golbeck, J., Hendler, J.: FilmTrust: Movie Recommendations using Trust in Web-based Social Networks. In: Proc. IEEE Consumer Communications and Networking Conference (2006)
4. Brickley, D., Miller, L.: (2007) FOAF Vocabulary Specification 0.9. <http://xmlns.com/foaf/0.1/>
5. Newman, R., Russell, S., Ayers, D.: (2005) Tag Ontology. <http://www.holygoat.co.uk/owl/redwood/0.1/tags/>
6. Oldakowski, R., Bizer, C., Westphal, D.: RAP: RDF API for PHP. In: Proc. 1st Workshop on Scripting for the Semantic Web, 2nd European Semantic Web Conference (2005)
7. W3C Technical Architecture Group: (2007) httpRange-14: What is the range of the HTTP dereference function? <http://www.w3.org/2001/tag/issues.html#httpRange-14> (2005)
8. Miles, A., Baker, T., Swick, R.: (2006) Best Practice Recipes for Publishing RDF Vocabularies. <http://www.w3.org/TR/swbp-vocab-pub/>
9. Suchanek, F. M., Kasneci, G., Weikum, G.: Yago: A Core of Semantic Knowledge - Unifying WordNet and Wikipedia. In: Proc. 16th Intl. World Wide Web Conference (2007)
10. Auer, S., Lehmann, J.: What have Innsbruck and Leipzig in common? Extracting Semantics from Wiki Content. In: Proc. 4th European Semantic Web Conference (2007)
11. Bizer, C., Heath, T., Ayers, D., Raimond, Y.: Interlinking Open Data on the Web. In: Proc. Demonstrations Track, 4th European Semantic Web Conference (2007)
12. Bizer, C., Cyganiak, R., Gauss, T.: The RDF Book Mashup: From Web APIs to a Web of Data. In: Proc. 3rd Workshop on Scripting for the Semantic Web, 4th European Semantic Web Conference (2007)