# Heurist: a Web 2.0 Approach to Integrating Research, Teaching and Web Publishing

Ian JOHNSON

Senior Research Fellow
Archaeological Computing Laboratory/Digital Innovation Unit, University of Sydney
ian.johnson@sydney.edu.au

## Abstract

This paper argues for the integration of bibliographic data and other forms of research data in a single integrated web 2.0 social database. This approach contrasts the use of separate systems for different classes of data (references in a bibliographic system, excavation data in a database, photographs in a photographic catalogue, notes and interpretation in wordprocessing files, wikis or blogs, and so forth). The fragmentation of information in non-interoperable 'silos' impedes cross-referencing and often leads to poor workflow, redundancy and data currency problems. An integrated web 2.0 approach promotes accessibility, teamwork, data currency and cross-referencing of information. The paper describes the philosophy behind Heurist (HeuristScholar.org), a free academic social bookmarking, bibliographic and general-purpose database, which provides rich data handling in a single integrated web service. Sourcing of information, interpretation and discussion can be developed through record cross-references to which notes, annotation and discussion can be attached. Records store geographic and temporal information and can be shared and tagged to allow cooperative creation and peer-group rating of information. Filters and output transformations allow subsets of the database to be generated in a wide variety of formats including reference lists, maps, timelines and XML feeds, for incorporation into project web sites, teaching resources and mashups.

## Keywords

databases, data structures, web 2.0, social computing, linking, bibliographic systems

## 1. Introduction

Bibliographic databases traditionally record a narrow class of materials – formally published books and journal articles, theses, grey literature reports and so forth – which typically synthesize material around a specific issue and are instantiated as physical objects (printed matter).

The main sources of complexity in implementing such bibliographic databases are the variety of ways items can be structured in relationship to other items (eg. within books, conference proceedings, journals and various types of series), the lack of standardized versioning systems, and the disambiguation of entities (including the bibliographic entries themselves and their creators). While this complexity is addressed by the Functional Requirements for Bibliographic Records (FRBR) conceptual model (Tillet 2004), most software implementations view each bibliographic item as an independent database record complete in itself – in relational database terms they are very poorly normalized – and pay little or no attention to the relationships between items. Yet it is precisely the relationships between bibliographic items (eg.

chapters belong in books, articles in journal issues, both form part of a series/serial from a publisher), and between bibliographic items and other entities (sites, research projects, fields of research, theories and so forth) which are the focus of research and scholarship. This central importance of linking information items was one of the fundamental tenets of the initial proposal for the World Wide Web (Berners-Lee and Cailliau 1989).

New forms of publication have exploded over the past few years driven by the success of the web - multimedia resources (increasingly without physical manifestation), web sites of various genres, pre-print repositories, open access e-journals and the deep web of searchable online databases, not to mention the many tools which leverage and extend these resources (Friedberg 2009, 150). This increase in the complexity of what we might call bibliographic data – citeable resources, not necessarily textual or printed – owes much to the rise of the web, of open content, of digitization, of semi-automated and automated methods of data collection (field computers, digital cameras, video recorders, laser

scanners) and increasingly widespread access to database technology.

The increase in complexity and variety of published resources raises the question of where to draw the line between bibliographic resources, *sensu lato*, and research data; or indeed, whether such a line should, or can, be drawn. While few people would disagree that a printed descriptive catalogue of pottery specimens is a bibliographic item, is this also the case when the same information is delivered in an ever-evolving web-accessible database? The changing nature of digital information forces us to reconsider the convenient pigeonholing of certain types of data record as bibliographic data, to be handled separately from other types of data using specialized bibliographic software such as EndNote.

## 2. The need for integration

The arbitrary distinction between 'bibliographic' data and other types of data is a hangover from a simpler past, when there was a clear divide between published syntheses (in printed form) and the data from which they were derived (including plans and sections, handwritten notes and forms, tabulations, standalone databases and photographs). Today published material grades into other types of data resulting from analysis and interpretation (such as annotated site catalogues and identification keys), new forms of presentation (such as online journals and compendia), and raw data collected in field surveys, excavation and laboratory analysis at ever increasing levels of granularity.

Drawing artificial distinctions between published material and the data on which these syntheses depend, hinders cross-referencing of information which should be linked. For example, syntheses should be linked directly to the original sources and data on which they depend. Bibliographic referencing in published material rarely links directly to the item cited – in general its function is to claim authority for an assertion and provide a finding aid which would allow verification of the assertion. In the same way, image collections for teaching often include references to the source of images as a formatted bibliographic reference or note in a text field, but this is stored redundantly as text for every record and does not generally link to the actual source. The desktop silo nature of commonly used bibliographic software has hindered tighter integration of bibliographic data with research databases.

Similarly, project databases – whether survey data, excavation data, laboratory analysis, or syntheses of information around a specific issue – are often bespoke systems which admirably (or less than admirably) fulfill particular needs but rarely connect with the wider network world (due to the extra, unrewarded, effort required to expose databases as an effective web service and perceived specificity to a small group of users). They are therefore often in desktop or local area network based 'silos', inaccessible to anyone but the users of their host computer, other than through sharing of copies (with all the attendant problems of redundancy, maintenance of currency and merging of changes). When accessible they often lack multi-user capabilities and/or are locked down to modification by a small group of users on a local network because of the difficulties of monitoring and rolling back erroneous or hostile changes. Nichols (2009) further makes the case that many digital projects live in sub-disciplinary silos.

Even when available more widely, bespoke databases are generally accessed through a web interface (HTML forms) which allows human access but not machine access, and cannot therefore be linked programmatically with other data to create an integrated system for analyzing larger problems. Even where an API is provided, we are still largely at the stage of case-by-case connections – the semantic web information which would allow automatic discovery and use of a database through a web service is only really available within particular constrained domains, such as the geographic data domain (and even then, is only partially implemented).

## 3. A web 2.0 approach

This pessimistic view of the 'silo' problem in research information may start to find a solution in increasing adoption of web 2.0 (O'Reilly 2005; Anderson 2007) approaches to content creation through crowdsourcing and social software such as Wikipedia, Delicious or Open StreetMap, where content is built by the users. These systems demonstrate that it is possible to build substantial and highly usable databases – subject to intelligent assessment and some benevolent control – without the costs normally associated with building significant knowledge resources, by harnessing the collaborative enthusiasm of large numbers of people for data collection and data correction, and through data mining of collective behavior (Howe 2008).

Such social systems offer a great potential for collaborative development and sharing of information. However, there are a number of potential challenges for academic use. First, unlike conventional published material, content can be continuously added to or modified, so that there is no fixed versioning of content to be referenced; while versions can be snapshot, current publication practices do not cope well with such dynamic resources. Second, in most systems, contribution is anonymous or based on virtual web identities rather than real identities, and there is no editorial control other than filtering of inappropriate material. This model does not fit with existing models of peer review and weighting of content through knowledge of the author's/editor's/publisher's body of work and reputation there is therefore no opportunity for building webs of citation which highlight significant contributions (simple counts can represent notoriety as much as quality).

Third, and most serious, there is a significant danger (already apparent) that the free, or low-cost availability of a plethora of systems for different purposes – wikis, blogs, bookmarking, social networking, calendaring, to-do lists, website builders, image collections and mapping, to mention only the most obvious – leads to exactly the same siloing of information in separate, unconnected and incompatible systems that occurs on the desktop (with the added danger that they are outside the control of the user and can shift locations or disappear altogether according to the whim of commercial expediency).

Instead of this short-term convenience-driven approach to information, a coherent approach is needed. This would be based on the removal of the artificial distinction between bibliographic data and other forms of data, and between different types of non-bibliographic data. The approach needs to integrate data, discussion, interpretation and synthesis within a single integrated knowledge system which is agnostic towards different types of objects and draws little or no distinction between documents, multimedia, text, quantitative data and metadata[1]. This approach embodies the principles of

the web: "The dream behind the Web is of a common information space in which we communicate by sharing information. Its universality is essential: the fact that a hypertext link can point to anything, be it personal, local or global, be it draft or highly polished" (Berners-Lee 1998).

An integrated knowledge system which is agnostic about the types of information it can store is the aim of the Heurist eResearch database (*Fig. 1*; HeuristScholar.org), developed under my direction at the Archaeological Computing Laboratory, University of Sydney. Heurist combines many of the characteristics of different social web systems: bookmarking, tagging, discussion, wikis, blogs, multimedia, georeferencing, workgroups, messaging – with cross-linking of bibliographic references and other data categories, in a system with authorial identity.

By integrating all these different forms of information in a single system, we can start to leverage technical development – such as a programming API, web services, geographic contextualization and mapping, XML feeds and publishing transformations, rich text annotation and discussion tools, across all classes of information rather than developing incompatible functions in different systems. In this way we are able to build on the real strength of online databases, standards and web services to develop a
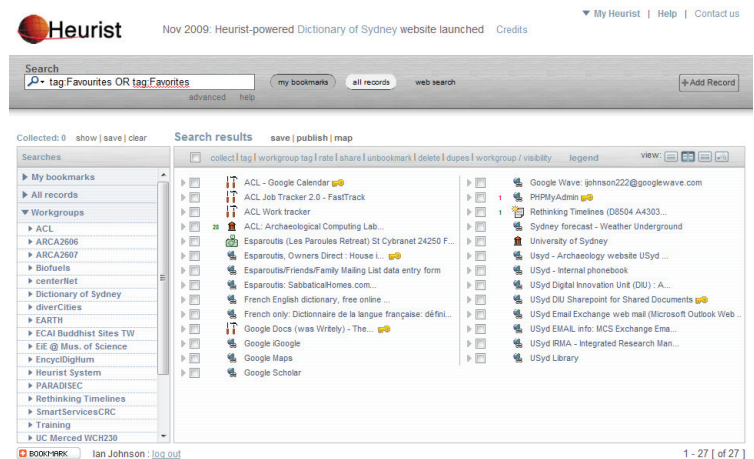


*Fig. 1. Heurist search page. The Publish function provides HTML code to embed the current search as live results in a web site. The user can choose to display a variety of content delivery formats, including HTML, XML data feeds, bibliographic formats, TimeMap and Google Map /Earth with timeline.*

---

[1] Similar approaches have been adopted in the design of archaeological information systems such as IDEA (Andresen and Madsen 1996) and INFRA (Schloen 2001), which use high-level abstractions of entities ('contexts', 'finds' and 'constructs' in IDEA, or simply 'items' in INFRA) to represent the full range of information in a few generic tables, rather than conventional object type categorizations and relationships manifest as separate tables.

Fig. 2. Results of a Heurist search for historical events rendered as an interactive, linked Google Map and Simile timeline.

system which flexibly supports research and teaching using distributed web-accessible resources which can be edited, extended and remixed on the fly. Nichols (2009) argues that competition for development resources is a zero sum game, and further suggests that 'tool-agnostic and use-agnostic approaches [have the] potential to encourage collaboration'.

Heurist handles an easily extensible set of more than 70 record types, ranging from bibliographic references and internet bookmarks, through encyclopaedia entries, seminars and grant programs, to C14 dates, archaeological sites and spatial databases, all uniquely referenced within a single database service and cooperatively editable. It allows users to attach geographic features (via a simple Google Maps digitizer), files, photographs, multimedia resources and rich text annotations to each entity in the database, with granular control of workgroup access to annotations at the paragraph level. Some entries, and parts of entries, can be locked off as authoritative content, while others can be left open to all comers.

## 4. Publishing data

Heurist is intended to remove the obstacles to web-based publishing of data by allowing users to publish subsets of the database to the web, in a variety of formats, without programming. Formats available



Fig. 3. University of Sydney Archaeology Department web site: research projects list with thumbnail images, published live from Heurist to the University CMS.

294

include interactive maps and timelines (*Fig. 2*) and XML feeds, as well as a variety of text layouts which can be extended by writing XSLT files to control formatting. As with many web applications today, it automatically generates the required HTML code to embed live Heurist output, allowing anyone with basic web page editing skills to create live published search results and interactive maps in their web pages.

Several web sites have been built using Heurist to populate the pages live with lists of courses, people, projects, seminars, publications and so forth (see *Fig. 3*). Over the last two years we have been assembling a wealth of background archaeological information in Heurist, ranging from useful web sites and sources of data to annotated bibliographic references. Dr Arianna Traviglia (University of Venice) has recently used Heurist to develop a comprehensive bibliography of remote sensing applications in Archaeology during a research fellowship at the University of Sydney. We intend to continue this work, using collaborative editing and update to maintain currency of this bibliography (and other archaeological resources in Heurist). Heurist is also being used to generate bibliographies live in project

and teaching web sites, and as a bookmarking and blogging tool in undergraduate classes; students create their assignments, bibliographies and other data within the system and publish them in their Heurist blog.
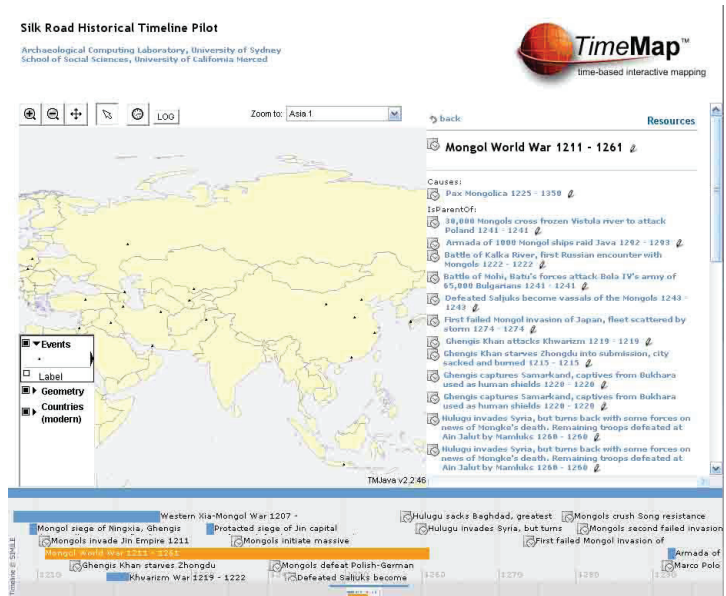


*Fig. 4. Rethinking Timelines project: pilot historical event visualization with linked interactive TimeMap (left), Simile timeline (bottom) and related events browser (right) generated from Heurist by following relationship records.*
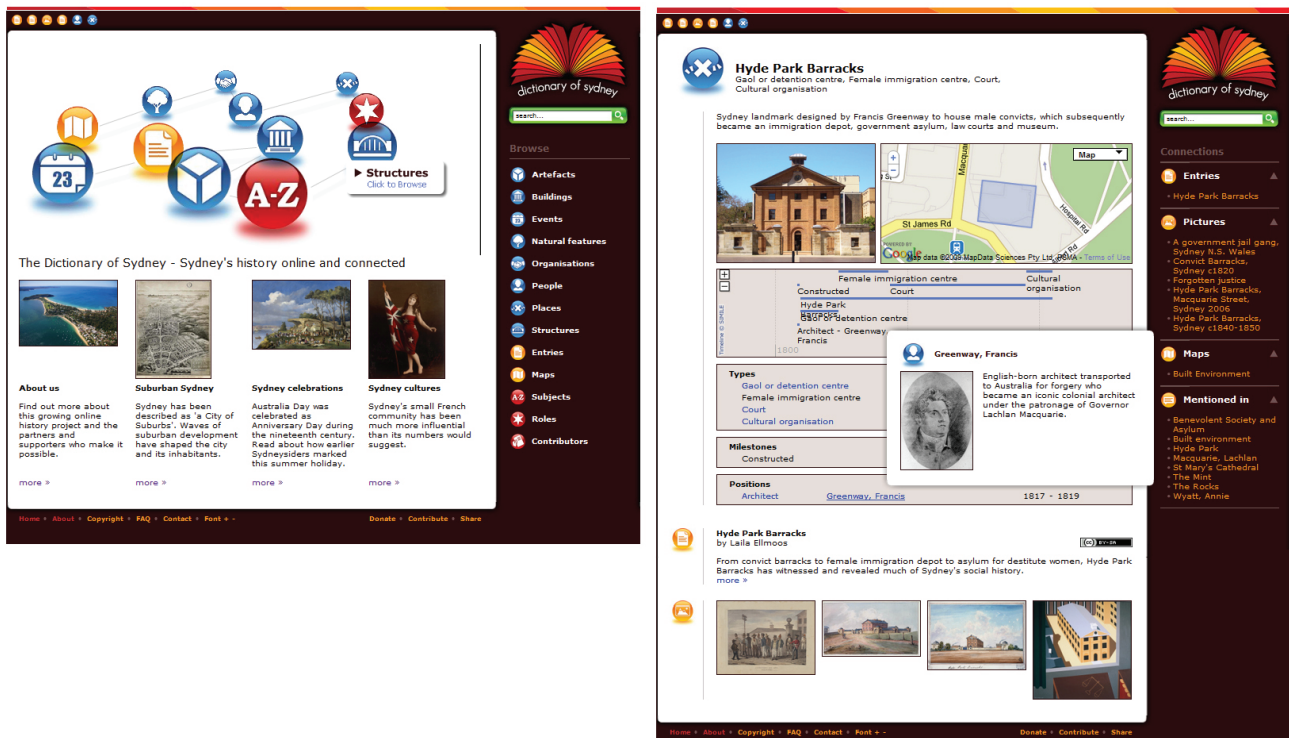


*Fig. 5. Dictionary of Sydney: encyclopaedia entries are annotated with links to entities (people, places, structures, events etc.) and to multimedia resources using Heurist relationship records. Maps, timelines and popups are automatically constructed from annotation data records.*

## 5. Relationships

A critical capability in Heurist is the ability to store annotated, date-stamped, typed relationships between any two entities in the database, allowing derived information to be linked back to bibliographic sources (such as sites to the documents in which they are referenced or field notes to the photographs, artefacts or trenches which they describe). The relationship capability was originally developed to allow chapters to be linked to the books containing them, journal articles to the journal volumes in which they occur, books to their publishers, and so forth. Since these linked entities are of different types, the relationship capability was necessarily type-agnostic.

The benefits of a type-agnostic relationship system (unlike the type-specific relationships typical of relational databases) have become increasingly apparent as we have built a variety of different applications based around Heurist. Type-agnostic relations are the key to linking bibliographic entries to other types of entity and building, browsing and visualizing networks of related entities such as historic events (*Fig. 4*) or contemporary stories. The Dictionary of Sydney (*Fig. 5*; dictionaryofsydney.org) illustrates the potential of relationships to act as the organizing principle and navigation metaphor for a rich collection of heterogeneous entities.

## 6. Next steps

The first-order typed relationships implemented in Heurist have proved pivotal for a number of projects. They are currently being further developed as part of the Rethinking Timelines project, funded by the Australian Research Council, which aims to develop new ways of modeling history through networks of interconnected historical events (Mostern and Johnson 2008). In particular, this will allow the allocation of periods to archaeological material using relationships with dated events which can nest within broader periods.

The Rethinking Timelines project is also extending the time stamping methodology to handle temporal objects which have uncertain and diffuse temporal limits, and investigating methods for delivering multiple interpretations and fuzzy dating in coordinated map-timeline visualisations. We are also investigating methods for generating self-documenting archival packages and RDF output

for semantic web applications. It is our belief that such richly interlinked, integrated databases of heterogeneous information, exposed as semantic web services, offer new opportunities for research across collections and the dissemination of research results.

## Appendix I: Technical background

Heurist is built in MySQL, PHP, Javascript and Cocoon, and easily customised through CSS and XSLT. Live search results and maps can be embedded in most CMS as well as static web pages. Multiple instances of Heurist can operate from the same set of record type definitions and user records.

Users can self-register, bookmark web pages, add, edit, tag, rate, share, relate, annotate, and provide commentary on records. Administrative functions allow the creation and management of user groups, controlling access to subsets of the database.

The service at HeuristScholar.org is available free of charge. Current development aims to make the software available as Open Source in 2010, to support installation of standalone instances, and to provide cross-instance software update, searching and data sharing.

## Bibliography

Anderson, Patricia (2007). What is Web 2.0? Ideas, technologies and implications for education. *JISC Technology and Standards Watch.*

Andresen, Jens and Torsten Madsen (1996). Dynamic classification and description in the IDEA. *Archeologia e Calcolatori* VII, 591–602.

Berners-Lee, Timothy (1998). The World Wide Web: A very short personal history (online: http://www.w3.org/People/Berners-Lee/ShortHistory.html, accessed 2 Feb 2010)

Berners-Lee, Timothy and Robert Cailliau (1989). WorldWideWeb: Proposal for a HyperText Project. Proposal to CERN. (archived at: http://info.cern.ch/, accessed 2 Feb 2010)

Friedberg, Anne (2009). On Digital Scholarship. *Cinema Journal* 48(2), 150–154

Howe, Jeff (2008). *Crowdsourcing: Why the Power of the Crowd Is Driving the Future of Business.* Random House.

Mostern, Ruth and Ian Johnson (2008). From Named Place to Naming Event: Creating Gazetteers for

History. *International Journal of Geographical Information Science* 22 (10), 1091–1108.

Nicholls, Stephen (2009). Time to Change Our Thinking: Dismantling the Silo Model of Digital Scholarship. *Ariadne* 58.

O'Reilly, Tim (2005). What is Web 2.0: Design Patterns and Business Models for the Next Generation of Software. (online: http://oreilly.com/web2/archive/what-is-web-20.html, accessed 2 Feb 2010)

Schloen, J. David (2001). Archaeological Data Models and Web Publication Using XML. *Computers and the Humanities* 35, 123–152.

Tillett, Barbara (2004). What is FRBR? A conceptual model for the bibliographic universe, Library of Congress Cataloging Distribution Service. (online: http://www.loc.gov/cds/downloads/FRBR.PDF, accessed 2 Feb 2010)