# Beyond Digital Incunabula: Modeling the Next Generation of Digital Libraries⋆

Gregory Crane, David Bamman, Lisa Cerrato, Alison Jones, David Mimno, Adrian Packel, David Sculley, and Gabriel Weaver

Perseus Project, Tufts University

**Abstract.** This paper describes several incunabular assumptions that impose upon early digital libraries the limitations drawn from print, and argues for a design strategy aimed at providing customization and personalization services that go beyond the limiting models of print distribution, based on services and experiments developed for the Greco-Roman collections in the Perseus Digital Library. Three features fundamentally characterize a successful digital library design: finer granularity of collection objects, automated processes, and decentralized community contributions.

## 1 Introduction

Potentially massive digital libraries such as the Google Library Project [17], the Open Content Alliance [20] and the EU i2010 Initiative [21] emphasize high volume digitization based primarily on automatically generated output of page images. While digitizing page images should be the first step in any digitization project and must comprise the core of any million book library, we need as well a sound infrastructure to augment our ability to search, browse and analyze the collection. Too great an emphasis on quantity can reinforce usage models that perpetuate limits from print distribution. This paper argues for a more aggressive, but still conservative, design strategy aimed at providing customization and personalization services that go beyond limiting models based on print distribution. While customization involves the user making explicit choices about the interface or system they are using, personalization involves the system adapting itself automatically to the user [35]. We base our argument on existing services and experiments developed for the Greco-Roman collections in the Perseus Digital Library [8, 6, 7]. The underlying methods are broader in application, have contributed to work completed for the National Science Digital Library [19], and lay the foundation for a range of services in Fedora and other digital repositories [27].

This paper has the following components. First, it describes several incunabular assumptions that impose upon early digital libraries the limitations drawn from print [34]. Second, it describes three features that fundamentally characterize emergent digital libraries. Third, it provides examples of customization and personalization built upon these three features.

## 2 Incunabular Models of Early Digital Libraries

New media begin by solving well known problems but also by imitating the forms that precede them. Consider three habits of thought drawn from paper libraries that constrain the design of

---

digital libraries. First, academic publications are based on coarse chunks of information, usually PDF or HTML files with heavily structured data (e.g., section headers, notes, bibliography). Pre-structured documents perpetuate the primacy of the author, leaving readers to do what they can with the structure and content that the author has chosen to include. Second, the emphasis on metadata carries forward the card catalogue of the print library. While metadata is important, metadata repositories that do not also include content are of limited use and have imposed an elegant modularity that constrains, as much as it has enhanced, intellectual life [26]. Third, print libraries are static. We may replace the books on the shelves with new editions, but those old books do not generate new versions of themselves. Print libraries cannot learn about their holdings and generate new content on their own. Information retrieval systems, which reindex libraries as they grow, constitute only a partial step in this direction, for they do not cycle over and over generating new knowledge by learning from their collections and from their users.

Hand-crafted collections, their contents marked up according to XML DTD/schemas with RDF metadata, incorporate major advantages over, but still narrowly replicate, their print antecedents. Some projects have, however, begun to move beyond these limitations. Figure 1 (left) is drawn from the most current and best documented survey of Athenian democracy now available: an electronic publication called Demos, published on-line as part of the Stoa publishing consortium [18].

First, Demos combines traditional and emerging structural designs: it can be downloaded as PDF chapters resembling conventional publications, but it is also available as logical chunks, each containing a semantically significant unit of text, providing greater precision than chapter heading and greater accuracy than page numbers in a book.

Second, Blackwell based a densely hypertextual work on an academic digital library that contained most of the primary evidence about Athenian democracy. The digital library shapes the form of Demos. Every major statement contains links to the primary evidence: where print reference works avoid visual clutter and save space, style sheets can turn these links on and off in a digital publication. More substantively, documented writing diminishes authorial claims of authority, offering readers an opportunity to compare conclusions with their evidence and enabling discussion. Unlike their print antecedents, citations in a digital library can point not only to pre-existing documents but also to services such as morphological analysis of Greek words or the plotting of geographic locations on high-resolution maps. The author combines links to the digital library and contextualizing materials within Demos (one of which is pictured in the figure). These internal links include discussions of the primary sources and the issues that they raise. Demos does not, however, directly address the model of the static library: each chunk of Demos lists its last modification dates, and each date testifies to the fact that Demos, for all its strengths, is not changing.

Figure 1 (right) shows the Wikipedia article on the Athenian council (boule) as it appears on March 2, 2006 [14]. Taken together, these demonstrate the possibilities of both existing and emergent digital libraries. While the Demos discussion of the boule is rigorous, it is also dated (January 23, 2003) and grows slowly out of date with each passing day. Second the Demos article reflects the synthesis of a single author interacting with a small editorial community, and such a publication method could omit important perspectives from an authoritative discussion. The Wikipedia article, by contrast, is subject to constant modification. Wikipedia can thus capture broader perspectives and remain current if opinion shifts [5]. The Wikipedia article, however, contains no source citations: modeled on contemporary encyclopedias and reference works and their relatively superficial bibliographic apparatus, Wikipedia articles contain high statement/evidence ratios. They thus claim credibility, and the lack of systematic pointers to evidence is problematic. If every statement were linked to its source, gross misrepresentations would be much more readily identified and corrected.

**Fig. 1.** (Left) A page from Demos: Classical Athenian Democracy, ed. Christopher Blackwell. (Right) Wikipedia article on the Athenian Boule.

## 3 Three Features that Characterize Post-Incunabular Digital Libraries

If we combine the scholarly rigor of Demos with the self-organizing qualities of Wikipedia, we can begin to see emerging a new model not only for digital libraries but also for the disciplined intellectual discourse which digital libraries should support. At least, three strategic functions distinguish emerging digital libraries from their predecessors.

1. Finer granularity: while many documents (like this paper) cite other publications as chunks, users often do not want to plough through an entire information object but would rather use an overview or particular sub-objects (proposition, bibliographic references, etc.) [37]. As digital libraries evolve, these structures go beyond those implicit in traditional publication models and begin to include explicit propositional statements.

2. Autonomous learning: Improved granularity implies this second fundamental characteristic of emergent digital libraries. Digital libraries should be constantly learning and becoming more intelligent as they scan their contents – a phenomenon already visible in rudimentary form with existing search engines. In a model digital library, the articles should be constantly scanning for new secondary sources, new editions of the primary materials and, to the extent possible, shifts in language that suggest perspectives that differ from the content of the current document. Thus documents and their sub-components should be in constant, autonomous communication with the libraries of which they are a part, scrutinizing new materials submitted and rereading the rest of the collection as automated processes evolve [39].

3. Decentralized, real-time community contributions: Wikipedia presents one of the most important practical experiments in the history of publication. For all the criticism that it may deserve, the English language Wikipedia has generated more than 1,000,000 entries in five years and supports several million queries a day. If we go beyond the higher level prose and examine individual verifiable propositions, the accuracy is remarkable. A recent study of relational statements in Wikipedia demonstrated that 97.22% of basic propositional statements and 99.47% of disambiguating links prove to be correct [38]. Thus, even if we reject the expository prose of Wikipedia for bias, the propositions within Wikipedia demonstrate that decentralized communities will accumulate immense amounts of highly accurate propositional data [28].

The following sections describe concrete, if rudimentary, steps Perseus has taken toward all three of the above principles, and addresses the issues that digital libraries in general must confront in order to transcend the limitations of print distribution.

## 3.1 Granularity

Cultural heritage documents often have complex contents that do not lend themselves to simple hierarchical representation [3]. While modern publications pre-define form and structure for the sake of simplifying system design, cultural heritage documents often have multiple, overlapping hierarchies (e.g., Thucydides' *History of the Peloponnesian War*, for instance, can be organized by book/chapter or by alternations of third person narrative and speeches). Digital libraries need to be able to address these various parts of a document that users want to work with. Figure 2 illustrates a dynamically generated report on what the Perseus DL knows about a particular chapter in Thucydides (Thuc. 1.86).



**Fig. 2.** Information about Thucydides, *History of the Peloponnesian War*, book 1, chapter 86, with user focus upon one of several translations. The right hand of the screen illustrates other dynamically collected resources extracted from digital objects and organized into an ontology of document types.

First, information display depends upon an authority list of meaningful citations: "Thuc. 1.86" provides a common designation with which references to Thucydides are aligned. Using this citation, we can identify which digital objects mention this particular chunk of text. Second, the individual passages that cite this passage have also been classified. Using this classification, we can distinguish Greek source texts from English translations and annotations specifically about Thuc. 1.86, and both from texts that only mention this passage in passing. Third, all documents in the collection have been broken down into structural units. Thus, we can extract multiple Greek versions or foreign language translations of the precise chunk designated by Thuc. 1.86. This allows us to identify not only that Thuc. 1.86 shows up in a particular Greek grammar, for instance, but that it also appears specifically in the discussion on "dative of interest."

Such fine grained chunking can transform the value of information: the main Greek-English lexicon entry on the Greek preposition "pros" mentions Thuc. 1.86, but few readers would scan
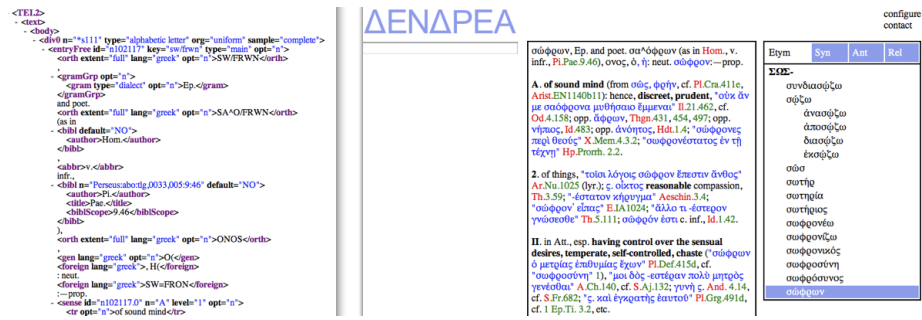
all fifty-two senses for the particular sense that cites Thuc. 1.86. Because we have captured the tree-structure of the dictionary, we know that this citation occurs at the third level down (as sense C.I.6), and we can extract this precise chunk from the much larger article. The right hand column aggregates and organizes these citations into a single report that could, in turn, be filtered and personalized for particular users.

Automatically organized reports on chunks of text (or museum objects, archaeological sites or other entities) build upon well structured documents which were designed for reference and subsequent citation. In order to address the complexity of cultural heritage documents that do not lend themselves to such representation, digital libraries need to confront the following issues, none of which are glamorous but each of which demands resolution:

1. Consistent markup for complex documents: For all of the progress that has been made, we do not have large, interoperable, richly structured documents in TEI or any other markup. Capturing chapters, sections, notes and similar well-defined elements is not the problem. Rather, we need consistent ways of managing documents within documents. Some carefully marked up collections (such as DocSouth [13] and American Memory [15]) may contain indices with accurate transcriptions of citations but have not included markup that captures the structure of the index or expands often idiosyncratic abbreviations into machine actionable data.

2. Mechanisms, automated and semi-automated, to identify meaningful chunks of individual documents: This desideratum addresses the need to generate large quantities of markup scalably as digital libraries with millions of books emerge: tables, notes, address/sender/ receiver/dateline tags for embedded letters and documents, indices, etc. are difficult to extract even from well transcribed documents. Such tools need to support both markup (e.g., where they add valid elements to existing structures) and extraction (where TEI texts might provide the source for generating external ontologies). Pilot projects integrating information extraction into digital libraries have begun to emerge (GATE/Greenstone), but developed solutions will be complex for large, heterogeneous collections and much needs to be done [40].

3. Digital library systems that can represent complex documents: We need as a starting point systems that support multiple, overlapping hierarchies within the same document [9]. We also need to be able to manage partial structures (e.g., browse all quoted speech, excerpts of poetry, personal names). Digital libraries need to model complex documents early in their development rather than concentrating on structurally simpler data types.

4. Conventions for exchanging complex content: Even if we address the first three issues, we still need conventions whereby we can exchange and recombine chunks of data from multiple collections [23]. These conventions include citation schemes, standard credit lines for author, institutions and funders, and an infrastructure for redundant, robust access in the present and for preservation into the future. Figure 3 illustrates an initial version of such a service, with a Perseus dictionary entry delivered as XML fragment (left) and a third party (Dendrea) representation of that data (right), with services added that are not currently present in the home digital library of the dictionary.

## 3.2  Automated Processes

While granularity can give us the ability to bring together related sections from different pre-existing documents, automatic analysis gives digital libraries the opportunity to create entirely new documents rather than quietly wait for new acquisitions.

**Fig. 3.** (Left) A well-formed fragment of XML representing an entry from the Liddell Scott Jones lexicon served over the Web. This service supports third party added-value services that exploit the rich underlying structure. Note that this document has a carefully captured structure, with each sense definition possessing a unique identification number. Thus, individual senses can be precisely extracted and reused in third-party hybrid documents. (Right) A third-party representation of the same dictionary article. Dendrea.org not only provides a different front end but also information extraction services that scan for etymological data, related words, and semantic relationships such as synonymy and antonymy.

Digital libraries need to include a range of automated processes that add value to their collections, including both classification (matching pre-determined patterns) and data mining (discovering new patterns) [1, 25]. While fields such as machine learning and data mining are major topics in their own right, integrating such technologies into digital libraries raises a range of problems [11, 10]. We need systems that can draw upon the contents of the digital library, continually applying new knowledge sources (e.g., gazetteers, machine readable dictionaries) as these become available, recalculating its results and assessing its performance [31, 41].



**Fig. 4.** Named entity analysis for classical texts. Notice the lack of culturally appropriate markup, with the TEI SURNAME tag used problematically to capture the primary name, as we begin adapting instruments for modern sources to Greco-Roman documents.

Not only do we need scalable methods to identify the semantically significant document chunks such as tables, embedded letters, and notes, we need more ways to analyze raw text and classify it as propositional knowledge, bibliographic citations, quotations, and named entities (e.g., is "London" a person or a place, and, which person or place?). People, places, dates and organizations are fundamental data which any mature digital library must track. Figure 4 displays the results of a named entity recognition system in place at Perseus since 2000. In an excerpt from Thucydides 8.108 ("About the same time Alcibiades returned with his thirteen ships from Caunus and Phaselis

to Samos, bringing word ...”), “Alcibiades” has been automatically annotated as a person, and “Caunus,” “Phaselis” and “Samos” have been annotated as places, two of them matched with the Getty Thesaurus of Geographic Names via TGN identification numbers.

Automated analysis also includes multi-lingual services such as cross-language information retrieval and machine translation [30, 12]. Identifying the fundamental meanings of a word is a notoriously slippery problem – human lexicographers do not agree among themselves as to what constitutes a separate word sense. One pragmatic approach involves examining translation equivalents: where translators use distinct words in the translation language, we have evidence of a substantive different meaning. Thus Table 1 displays one cluster of word meanings for the polysemous Greek word *arche*, derived from comparison of a Greek source text and five separate English translations. *Arche* can mean “empire,” “government,” “political office,” and “beginning”; by grouping together the words that occur around it, we are able to use translations to identifyits intended sense. The six texts as a whole are aligned according to the standard citation scheme rather than at the word or sentence level (c. 42 Greek words per chunk), with sections themselves in four of the five English translations automatically aligned with the Greek original. The experiment thus explores what can be done with translations of canonical texts, with minimal extra tagging, that will populate large emerging digital libraries.

**Table 1.** Parallel text analysis: word clusters associated with uses of the Greek word *arche* in Thucydides (c. 150,000 words) and five English translations. Translation equivalents are underlined. The clusters capture the senses “empire,” “government,” “political office,” and “beginning.” The cluster headed “ancient” (marked in bold) captures a distinct word that happens to share the stem *arch-*.

| empire | dominion | power | government |
|---|---|---|---|
| office | government | magistrates | people |
| dominion | power | rule | Hellenes |
| magistrates | Theseus | people | council |
| **ancient** | descendants | temples | Pythian |
| whom | beginning | pits | just |
| called | Zancle | Pangaeus | originally |

The translation analysis points to four elements of text mining relevant to digital libraries. First, this function will improve as digital libraries grow larger, because we will have access to more and more translations of source texts into a range of languages. Second, the clustering of word groups is computationally intensive and the current algorithm is not suited to providing real time results. Third, while such exploratory data may not begin as part of the general digital library, the results of such analysis, once generated, may become a domain specific service available to those reading Greek (or similar languages for which this service is suitable). We may well find domain specific front ends, integrating data from several larger collections into a new hybrid information space designed for particular communities. Fourth, the output of the parallel text analysis is useful in its own right to human analysts, but this output also provides a foundation for cross-language information retrieval, machine translation and other multi-lingual services.

### 3.3 User Contributions

Every large collection contains errors and, while these are finite and may be corrected in a finite period of time, by the time original errors may be fully corrected, scholarship will have marched on, rendering bibliography, front matters, and annotations in need of revision. While some automated processes do approach perfection, even these still generate errors, and most automated processes have error rates far removed from 100%. Some processes (such as assigning a sense to a given instance of a word or analyzing the syntax of a sentence) will yield disagreement among experts.

We need to consider mechanisms to collect information from our user communities [33, 4]. In some cases, the amateurs will probably perform better than the academics: professional historians may chafe at genealogists fixated on precise identifications of people, places and organizations or antiquarians fretting about the precise buttons a particular person might have worn, but such attention to detail can be of immense value when we are refining the results of our collections.

There are two categories of contribution. First, we need annotations and reference articles that incorporate at least some full text. Wikipedia has demonstrated both immense strengths and troubling weaknesses in this area. More conservative efforts such as PlanetMath [16], based on Noosphere, have arguably produced more consistent results, but they are much more focused efforts and have created around 5,000 encyclopedia entries rather than the 1,000,000 in Wikipedia [24]. The NEH has funded "Pleaides: An Online Workspace for Ancient Geography," which will explore community created scholarly content [22]. The Perseus DL will include Wiki pages for every data object and every addressable text chunk. Our optimistic hypothesis is that community driven commentaries, created by dedicated amateurs struggling to understand the texts, may prove more useful than commentaries produced by experts: we expect many errors at first but that the churning of user responses will weed out mistakes and that Wiki commentaries will evolve into accurate instruments. We are less concerned with the potential errors than with whether such Wiki commentaries will attract a critical mass of contributors.

Second, we need to collect user feedback on propositional data: e.g., whether "London" in passage X is London, Ontario, rather than London, England; whether "saucia" is a nominative singular adjective rather than ablative; whether a certain dative is an indirect object of the verb rather than a dative of possession with a nearby noun. Propositional data does not always have a single, clearly correct answer, but we can collect alternatives and store them in a structured format.

We created two initial mechanisms to collect propositional user input, allowing users to match a particular word sense and morphological analysis appropriate to a given word in a given passage. Figure 5 illustrates the results of three processes. First, a morphological analyser generates an exhaustive list of possible analyses for the form "saucia." Second, another module examines the immediate context and the relative frequency of the forms and possible dictionary entries (if the word is lexically ambiguous), then calculates probabilities for each alternative analsysis. Accuracy of the automated disambiguation stands at 76%. Third, users can cast votes of their own, whether to reinforce the automatic analysis or to suggest an alternative.

As of March 7, 2006, users have cast 7,597 votes to disambiguate Greek and Latin words with more than one possible morphological analysis. The overall accuracy for individual voting stands at 89%, but the improvement over the performance of the automated disambiguation is substantially higher, since users overwhelmingly vote on words for which the system has assigned the wrong analysis. (While the overall accuracy of automatic disambiguation is 76%, its accuracy on words that users vote on is only 34%.) And since 43% of word tokens have only one morphological sense (and do not need therefore to be disambiguated), user voting with 89% accuracy on ambiguous forms has the potential to deliver an overall system with 93.7% accuracy if every ambiguous word

**Fig. 5.** System to register votes against machine generated choice of correct morphological analysis: heuristic driven morphological analysis has been a long-term service in the Perseus DL, but we have only recently applied machine learning to estimate the correct morphological analyses for a given word in a given passage. Users can vote for alternative analyses if they disagree with the automatic choice. This system has been adopted to evaluate unfiltered, anonymous contributions, providing a possible baseline for more demanding models.

receives one vote. We could solicit expert users to fill in the remaining accuracy gap, but a better solution may simply be to focus on enlarging the contributor base: the more individual votes per word, the more likely that all, when taken together, will be correct.

## 4   User-Centered Digital Libraries: Customization and Personalization

The three incunabular constraints of early digital libraries all act as much against as for the user. The catalogue model tells users what exists and sometimes points them to on-line versions of the source object, but then its job is done and the users must do what they can with what they find. The digital codex model may incorporate searching and convert citations to links, but the author creates fixed content and structure around which users must work. The static library can learn neither on its own nor from its users. Early digital libraries thus, not surprisingly, replicate the hegemony of library, author, and publisher.

Digital libraries can, however, shift the balance further toward the user and toward the active life of the mind. More finely grained data objects, automated processes and decentralized user contributions all should interact, with the digital library progressively growing better structured and more self-aware. Initial data structures seed automated processes which classify and mine data. Users evaluate classification results and feed their contributions back into the system. Data mining suggests new patterns, which in turn complement or revise previous schemes, leading to the discovery and classification of new structures within the same set of digital objects.

Two fundamental strategies should be available to the user. First, users should be able to customize the environment [2]. Such customization needs to reflect not only default page layouts and simple preferences but also much more elaborate models of user knowledge. Figure 6 shows a customized report for a user who has studied Latin with a particular textbook (one of about thirty Greek and Latin textbooks whose vocabularies we have modeled on a chapter by chapter basis). Multiple readers with different backgrounds can thus see in a given passage which terms are likely

**Fig. 6.** Customized knowledge profile: the digital library knows the textbook with which the user has worked and analyzes probable known and unseen vocabulary in a given passage. Because the user has specified a profile and the system has responded, this is an example of customization.

to be novel and which they have encountered before. The fundamental principle can be applied to scientific terminology – which is, of course, easier to recognize than natural language.



**Fig. 7.** Personalized knowledge profile: the digital library does not know the background of the user but has analyzed four initial questions, compared these with past question patterns and suggested which of the remaining three hundred words the current user most likely to query. User behavior clusters into distinct classes and this approach has been able to predict 67% of subsequent queries. Because the system has taken the initiative rather than the user, this is an example of personalization.

Second, personalization augments the user-initiated decisions of customization: digital libraries should be able to analyze user behavior and background and offer new automatically generated configurations of information [36, 29, 32]. Figure 7 illustrates a recommender system that compares records of questions from earlier readers who had read a particular text. By mining past behaviors, the system can quickly learn to predict most of the subsequent questions that new users will pose.

# 5 Conclusion

Google with its massive library project, more recently Microsoft, Yahoo and others in the Open Content Alliance, and potentially the EU in its i2010 initiative are poised to assemble massive, but coarse, libraries of digitized books that have the potential to reinforce usage models based on print distribution. This paper provides initial examples of a post-incunabular design strategy utilized at the Perseus Project for its Greco-Roman collection but, we hope, scalable to other domains, focused on the principles of customization and personalization built upon fine grained digital objects, automated processes and decentralized user contributions.

# References

[1] H. S. Baird, V. Govindaraju, and D. P. Lopresti. Document analysis systems for digital libraries: Challenges and opportunities. In *Document Analysis Systems*, pages 1–16, 2004.

[2] N. Beagrie. Plenty of room at the bottom? Personal digital libraries and collections. *D-Lib Magazine*, 11(6), 2005. http://dlib.anu.edu.au/dlib/june05/beagrie/06beagrie.html.

[3] J. Bradley. Documents and data: Modelling materials for humanities research in XML and relational databases. *Literary and Linguistic Computing*, 20(1), 2005.

[4] M.A.B. Burkard. Collaboration on medieval charters–Wikipedia in the humanities. In *Proceedings of the XVI International Conference of the Association for History and Computing*, pages 91–94, 2005.

[5] D. J. Cohen and R. Rosenzweig. Web of lies? Historical knowledge on the internet. *First Monday*, 10(12), Dec. 2005. http://www.firstmonday.org/issues/issue10_12/cohen/.

[6] G. Crane. Cultural heritage digital libraries: Needs and components. In *ECDL*, Rome, Italy, 16-18 Sept. 2002.

[7] G. Crane, R. F. Chavez, A. Mahoney, T. L. Milbank, J. A. Rydberg-Cox, D. A. Smith, and C. E. Wulfman. Drudgery and deep thought: Designing a digital library for the humanities. *Communications of the ACM*, 44(5):35–40, 2001.

[8] G. Crane, C. E. Wulfman, L. M. Cerrato, A. Mahoney, T. L. Milbank, D. Mimno, J. A. Rydberg-Cox, D. A. Smith, and C. York. Towards a cultural heritage digital library. In *JCDL*, pages 75–86, Houston, TX, June 2003.

[9] A. Dekhtyar, I. E. Iacob, J. W. Jaromczyk, K. Kiernan, N. Moore, and D. C. Porter. Support for XML markup of image-based electronic editions. *International Journal on Digital Libraries*, 6(1):55–69, Feb. 2006.

[10] J. S. Downie, J. Unsworth, B. Yu, D. Tcheng, G. Rockwell, and S. J. Ramsay. A revolutionary approach to humanities computing?: Tools development and the D2K data-mining framework. In *Annual Joint Conference of The Association for Computers and the Humanities & The Association for Literary and Linguistic Computing*, 2005.

[11] F. Esposito, D. Malerba, G. Semeraro, S. Ferilli, O. Altamura, T. M. A. Basile, M. Berardi, M. Ceci, and N. Di Mauro. Machine learning methods for automatically processing historical documents: From paper acquisition to XML transformation. In *DIAL*, volume 1, pages 328–335. IEEE Computer Society, 2004.

[12] F. C. Gey, N. Kando, and C. Peters. Cross-language information retrieval: the way ahead. *Information Processing and Management*, 41(3):415–431, 2005.

[13] http://docsouth.unc.edu/.

[14] http://en.wikipedia.org/wiki/Boule.

[15] http://memory.loc.gov/ammem/index.html.

[16] http://planetmath.org/.

[17] http://print.google.com/googleprint/library.html.

[18] http://seneca.stoa.org/projects/demos/article_council?page=1&greekEncoding=UnicodeC.

[19] http://www.nsdl.org.

[20] http://www.opencontentalliance.org/.

[21] http://www.theeuropeanlibrary.org/portal/index.htm.

[22] http://www.unc.edu/awmc/pleiades.html.

[23] Y. E. Ioannidis, D. Maier, S. Abiteboul, P. Buneman, S. B. Davidson, E. A. Fox, A. Y. Halevy, C. A. Knoblock, F. Rabitti, H. J. Schek, and G. Weikum. Digital library information-technology infrastructures. *International Journal on Digital Libraries*, 5(4):266–274, 2005.

[24] A. Krowne. Building a digital library the commons-based peer production way. *D-Lib Magazine*, 9(10), 2003. http://www.dlib.org/dlib/october03/krowne/10krowne.html.

[25] A. Krowne and M. Halbert. An initial evaluation of automated organization for digital library browsing. In *JCDL*, pages 246–255, New York, NY, USA, 2005. ACM Press.

[26] C. Lagoze, D. B. Krafft, S. Payette, and S. Jesuroga. What is a digital library anymore, anyway? Beyond search and access in the NSDL. *D-Lib*, 11(11), 2005. http://www.dlib.org/dlib/november05/lagoze/11lagoze.html.

[27] Carl Lagoze, Sandy Payette, Edwin Shin, and Chris Wilper. Fedora: an architecture for complex objects and their relationships. *Int. J. Digit. Libr.*, 6(2):124–138, 2006.

[28] M. Lesk. The qualitative advantages of quantities of information: Bigger is better. *J. Zhejiang Univ. Science*, 11(6A), 2005.

[29] E. J. Neuhold, C. Niederée, and A. Stewart. Personalization in digital libraries: An extended view. In *ICADL*, 2003.

[30] D. W. Oard. Language technologies for scalable digital libraries. In *International Conference on Digital Libraries*, 2004. Invited Paper.

[31] G. Pant, K. Tsioutsiouliklis, J. Johnson, and C. L. Giles. Panorama: Extending digital libraries with topical crawlers. In *JCDL*, pages 142–150, New York, NY, USA, 2004. ACM Press.

[32] M. E. Renda and U. Straccia. A personalized collaborative digital library environment: a model and an application. *Information Processing and Management*, 41(1):5–21, 2005.

[33] M. Richardson and P. Domingos. Building large knowledge bases by mass collaboration. In *K-CAP*, pages 129–137, New York, NY, USA, 2003. ACM Press.

[34] M. Riva and V. Zafrin. Extending the text: digital editions and the hypertextual paradigm. In *HYPERTEXT*, pages 205–207, New York, NY, USA, 2005. ACM Press.

[35] John Russell. Making it personal: information that adapts to the reader. In *SIGDOC '03: Proceedings of the 21st annual international conference on Documentation*, pages 160–166, New York, NY, USA, 2003. ACM Press.

[36] A. F. Smeaton and J. Callan. Personalisation and recommender systems in digital libraries. *International Journal on Digital Libraries*, 5:299–308, 2005.

[37] E. S Villamil, C. González Muñoz, and R. C. Carrasco. XMLibrary search: an XML search engine oriented to digital libraries. *Lecture Notes in Computer Science - Research and Advanced Technology for Digital Libraries*, 3652:81–91, 2005.

[38] Gabriel Weaver, Barbara Strickland, Alison Jones, and Gregory Crane. Quantifying the accuracy of relational statements in Wikipedia: A methodology. In *To appear in JCDL 06: Proceedings of the 6th ACM/IEEE-CS joint conference on Digital libraries*, 2006.

[39] R. Witte. An integration architecture for user-centric document creation, retrieval, and analysis. In *Proceedings of the VLDB Workshop on Information Integration on the Web*, 2004.

[40] I. H. Witten, K. J. Don, M. Dewsnip, and V. Tablan. Text mining in a digital library. *International Journal on Digital Libraries*, 4(1):56–59, 2004.

[41] Z. Zhuang, R. Wagle, and C. L. Giles. What's there and what's not?: Focused crawling for missing documents in digital libraries. In *JCDL*, pages 301–310, 2005.