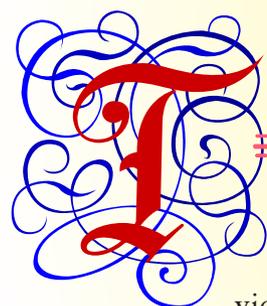Ian H. Witten, Craig Nevill-Manning, Rodger McNab,
and Sally Jo Cunningham

Full-text retrieval can substitute
for bibliographic metadata as a means of
accessing large library collections.

New Zealand

# A Public Library Based on
# Full-text Retrieval

HE NEW ZEALAND DIGITAL LIBRARY AIMS TO IMPOSE STRUCTURE ON ANAR-
chic and uncataloged repositories of information providing information
consumers with effective tools to locate and peruse what they need. Our
goal is to produce an easy-to-use digital library system that runs on inex-
pensive computers at the information providers' own sites and offers a ser-
vice that providers maintain. New Zealand's geographical isolation amplifies
the benefits of networked digital libraries, in terms of both cost and timeliness of access to infor-
mation. We are collaborating with the MeDoc project (see Endres and Fuhr, this issue) in Ger-
many to provide local indexes to German language technical reports, and with the U.S.-based
*Journal of Biological Chemistry* to field-test novel browsing techniques.

The project rests on five basic planks:

- We avoid manual processing of source material, and avoid making assumptions about the document repositories from which it is collected; for example, we do not require bibliographic metadata.
- Access is via a full-text index of the entire contents of each document, rather than document surrogates.
- We are concerned with user interface aspects and the real needs of library users.

- Our systems must operate in geographically remote locations with high Internet costs—an environment in which the benefits of networked library technology is especially striking.
- We aim to produce a library scheme that operates on small, inexpensive servers.

Full-text indexes are provided to several substantial collections of information. These collections serve as case studies. They drive our research by providing technical challenges for indexing, and human interface challenges for retrieval. This arti-
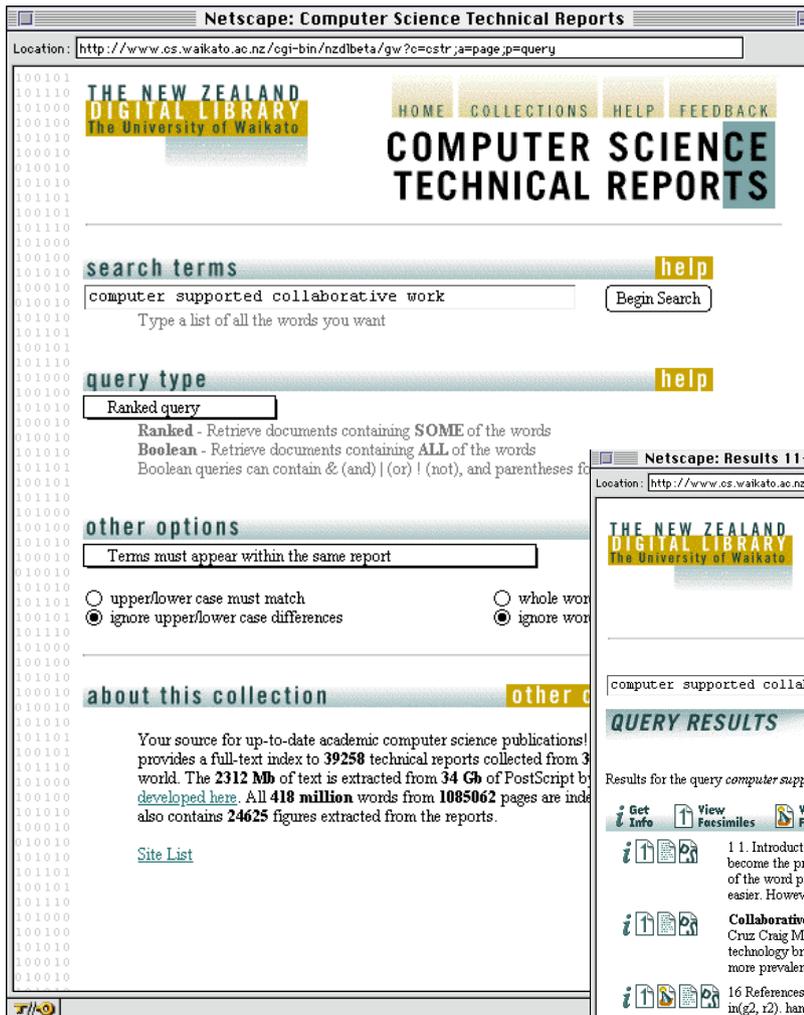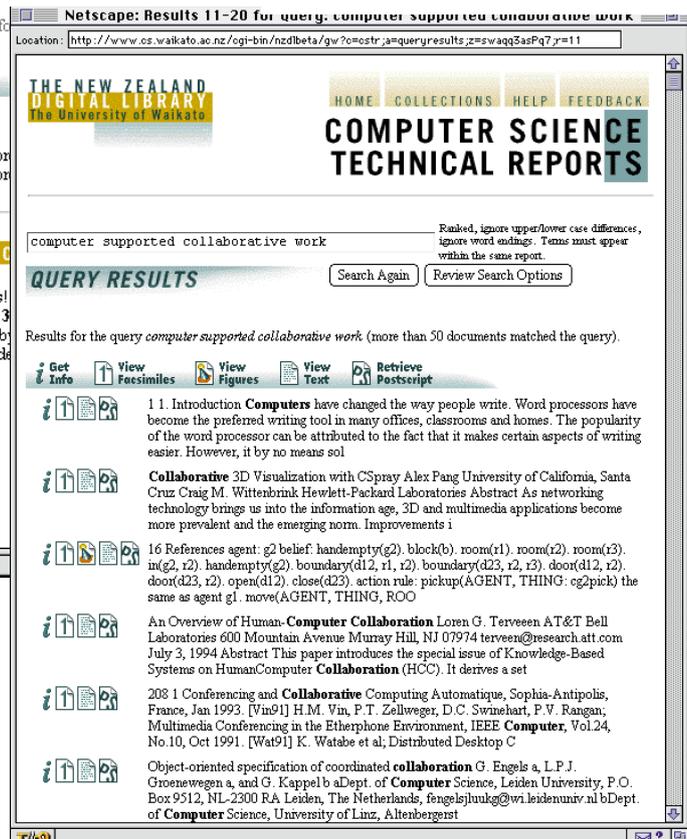
**Figure 1.** Computer Science Technical Report query page and a typical response.

cle focuses on the collections: technical details of mechanisms [7], protocols [3], and novel prototype interfaces [5] are available elsewhere.

### CS Technical Collections

We have indexed 40,000 computer science technical reports—one million pages, 400 million words—from 300 sites internationally. This represents 2.3Gb of text, extracted automatically from 34Gb of PostScript/PDF source.

*Querying and retrieval.* Figure 1 shows a typical query and response. Searching facilities include ranked and Boolean queries; the ability to stem individual terms and/or render them case-insensitive; phrase searching; document, page, or first-page-only searches. Since the collection is not formally cataloged users cannot perform conventional author/title/date searches, but most reports include this sort of bibliographic information in

the first page and fielded search can be approximated by restricting the query's scope appropriately.

The response shows the first few words of retrieved documents. Buttons allow users to retrieve the original document's URL, size, creation date, and download date; the PostScript itself; a facsimile image of the first page; the document text with query terms highlighted; and (where possible) the figures it contains. The extracted text and figures enable a quick scan to determine whether the document is worth downloading.
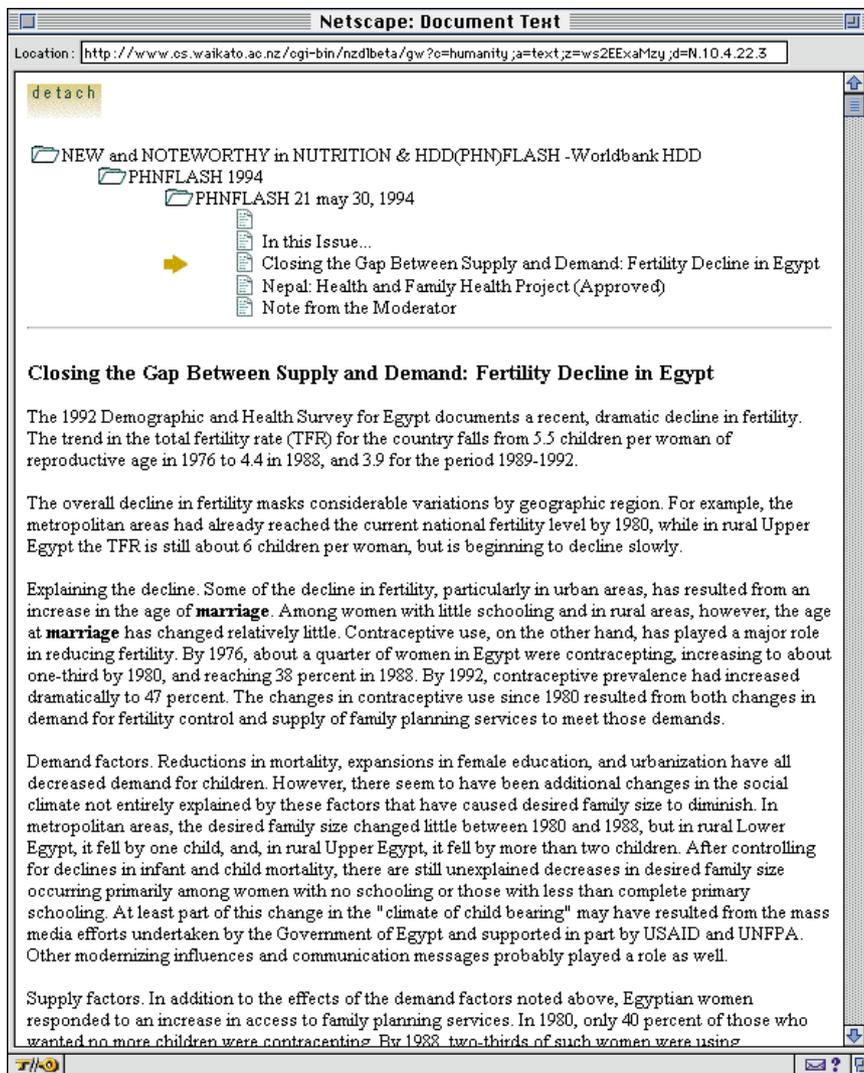
**Figure 2.** Browsing the result of a query for "marriage" in the Humanity Development Library.

total of 2.2Gb (coincidentally, just the same size as the original text), and 6% of the total PostScript source.

All sizes increase linearly with the volume of text. Retrieval is independent of database size, taking two disk seeks per query term and two per document retrieved. Database inversion is a potential limiting factor, but a recently developed algorithm can invert an estimated 5Gb of raw text in 12 hours with only 40Mb of main memory [4]. Collections of tens of gigabytes seem quite feasible. However, extrapolation on the basis of raw size is a gross oversimplification. Other factors, such as number of documents or richness of vocabulary, can send things awry. Our actual experience extends to rich-vocabulary (library catalog) collections of 10 million short documents, and 3.5Gb of raw text; inversion times are a few hours.

As collections grow, it will inevitably become necessary to partition the database. Partitioning may be desirable in any case. Indeed the separation between collections in our current system is dictated by user needs rather than technical limitations. Scalable full-text indexes exist: experiments indicate that the results of retrieval performed in parallel on the segments of a partitioned database can be combined with little degradation in effectiveness. We will rely on such techniques to provide full scalability.

## Other Collections

To demonstrate how digital libraries can benefit diverse groups of users we have built many other collections of publicly available information.

*The Computists' Communique and TidBITS magazines:* Searchable on individual news item, its title, and individual paragraph.

*FAQ Archive:* Searchable within entire FAQ list, by subject heading, or individual paragraph.

*Collection maintenance.* Other technical report servers suffer frequent maintenance problems caused by changes in the bibliography file format, and inconsistencies in the information, in the repositories that they index. This is why we do not use cataloging information stored with the repository itself. The collection is maintained automatically by examining the technical report repositories periodically for changes and updates. New sites are detected by various means (mostly manual) and added to the index.

*Size and scalability.* The public-domain MG search engine is tailored for highly efficient storage of full-text databases [6]. The text in this collection is compressed to 840Mb; several indexes are added, totaling 700Mb; and first-page facsimiles and extracted figures add another 680Mb, for a
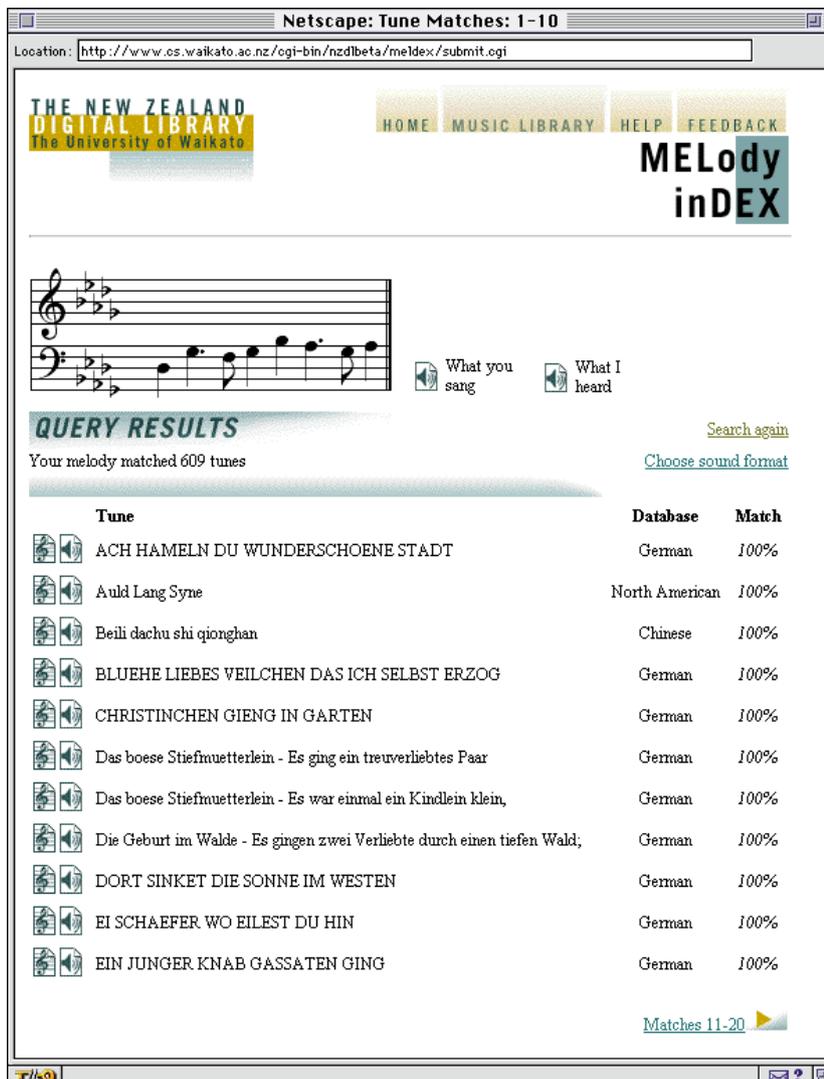
**Figure 3.** Using the melody index

*The HCI Bibliography:* Searchable by reference entry, with or without abstract.

*Humanity Development Library:* Humanitarian and development information collected by the Global Help Project.

*Indigenous Peoples:* Position papers, resolutions, treaties, U.N. documents, speeches and declarations on social, political, strategic, economic, and human rights issues.

*Oxford Text Archive and Project Gutenberg:* Public-domain collections of English text.

Major differences between these collections include the source and format of information, updating policy, granularity of searching (document, page, paragraph, and so on), different kinds of index (titles, references, abstracts), structure and format in which output is displayed, and provision of summaries of the collection's contents, for example, lists of titles and hierarchical, Web-accessible browsers as shown in Figure 2. The challenge is to enable information providers to tailor the system to new document sets without programming.

Our melody index retrieves music on the basis of notes that are sung, hummed, or played [2]. Users can literally sing a few bars and have melodies containing that motif retrieved. Such systems will enable researchers to analyze music for recurring themes or duplicated phrases, and musicians and casual users alike will retrieve compositions based on remembered (even imperfectly remembered) passages.

The system transcribes melodies automatically from microphone input, searches a database for tunes containing similar melodic sequences, and ranks matches using features such as melodic contour, musical intervals and rhythm. Figure 3 shows the response when a user sang the first eight notes of *Auld Lang Syne*. The transcribed input appears at the top; titles of similar items, ranked according to matching score, appear below. Any of the tunes may be selected for audio replay or visual display.

The database comprises 9,400 international folk tunes (half a million notes) stored in musical notation, an amalgam of the Digital Tradition database of North American folk songs with the Essen database of European and Chinese melodies. Though few machine-readable score collections are presently available, optical music recognition technology will change this situation. We have built a prototype service that accepts images of music and returns corresponding audio or MIDI files (based on [1]). Music submitted for processing can be automatically added to the database of indexed tunes.

## Creating, Distributing, and Integrating Collections

It takes a skilled technician a few hours to build a simple collection and integrate it into the digital library framework, excluding any time needed to gather the information. There is remarkable variety in the requirements for different collections. We have already created a dozen and experience shows that each poses new problems. Differences occur in the format of the source information, interface structure, indexing needs, and maintenance/update regime. We estimate the number of collections must triple before we reach a state where we can be reasonably confident that new ones can be accommodated within the existing framework. At that point we plan to design an end-user interface for collection building.

Subsets of the collections can be written to CD-ROM. Coupled with a local Web server, this provides a powerful means whereby users without convenient access to the Web, or for whom access is unreliable or expensive, can benefit from the same interface and facilities.

The NZDL is organized as several independent collections, each showing what the underlying technology can do in a different area. Joint searches can occur by simply constructing a combined index. But interesting issues arise when cross-searching different kinds of information, for instance computer science technical reports, bibliography collections, and paper titles and abstracts—because automatic relevance ranking must be compared for different kinds of text, and because user interface issues are nontrivial. For example, a paragraph-level search on a document collection might be initiated if a title search fails, or a search on a name might be extended from a bibliography to a report collection if few items match.

Value can be added to collections by mining the full text and correlating the contents with different sources. Examples include matching bibliography collections against technical report collections to infer metadata; eliminating duplication by self-matching bibliography entries; identifying reference lists and inserting hyperlinks to cited documents in the collection (and in bibliographies); analyzing document sets to infer topic clusters. Phrases suitable for a browsable hierarchical index can be inferred directly from the text, and help users build intuition about the content of a collection [5].

The NZDL will enable institutional end-users to create focused information collections in particular areas—in contrast to the less selective approach taken by Internet search engines—and make them publicly available. We have paid particular attention to the specific needs of a variety of document sources, providing flexible index granularity as well as browsing and searching interfaces tailored to content.

Digital libraries can be constructed whenever repositories of suitable text exist. Extracting all search information from the documents themselves is feasible if full-text indexing is used, and eliminates the manual cataloging involved in library creation and maintenance. Although New Zealand's geographical isolation provided the initial impetus for making access to information in electronic form simpler and more efficient, the techniques we have devised are applicable internationally. **C**

### References

1. Bainbridge, D. and Bell, T.C. An extensible optical music recognition system. In *Proceedings of the 1996 Australasian Computer Science Conference.* (Melbourne, Australia), 308–317.
2. McNab, R.J., Smith, L.A., Witten, I.H., Henderson, C.L. and Cunningham, S.J. (1996) Toward the digital music library: Tune retrieval from acoustic input. In *Proceedings of ACM Digital Libraries.* (1996), 11–18.
3. McNab, R.J., Witten, I.H., and Boddie, S.J. A distributed digital library architecture incorporating different index styles. In *Proceedings of the IEEE Conference on Advances in Digital Libraries.* (April 1998, Santa Barbara, Calif).
4. Moffat, A. and Bell, T.A.H. In situ generation of compressed inverted files. *J. Amer. Soc. Info. Sci. 46*, 7, 537–550.
5. Nevill-Manning, C.G., Witten, I.H., and Paynter, G.W. Browsing in digital libraries: A phrase-based approach. In *Proceedings of ACM Digital Libraries.* (1997), 230–236.
6. Witten, I.H., Moffat, A., and Bell, T.C. *Managing Gigabytes: Compressing and Indexing Documents and Images.* Van Nostrand Reinhold, 1994.
7. Witten, I.H., Cunningham, S.J., and Apperley, M.D. The New Zealand digital library project. D-*Lib. Mag.*, Nov. 1996; www.dllib.org.

**IAN H. WITTEN** (ihw@cs.waikato.ac.nz) is a professor of computer science at the University of Waikato, Hamilton, NZ.
**CRAIG G. NEVILL-MANNING** (cnevill@stanford.edu) is a post-doctoral fellow in the Department of Biochemistry at Stanford University, Stanford, CA.
**RODGER MCNAB** (rjmcnab@cs.waikato.ac.nz) is a research assistant in the Computer Science Department at the University of Waikato, Hamilton, NZ.
**SALLY JO CUNNINGHAM** (sallyjo@cs.waikato.ac.nz) is a senior lecturer in the Computer Science Department at the University of Waikato, Hamilton, NZ.