

A New Framework for Building Digital Library Collections

George Buchanan, David Bainbridge, Katherine Don and Ian H. Witten

Waikato University, New Zealand
{g.buchanan, d.bainbridge, kdon, i.witten}@cs.waikato.ac.nz

Abstract. *This paper introduces a new framework for building digital library collections and contrasts it with existing systems. It describes a radical new step in the development of a widely-used open-source digital library system, Greenstone, which has evolved over many years. It is supported by a fresh implementation, which forced us to rethink the entire design rather than making incremental improvements. The redesign capitalizes on the best ideas from the existing system, which have been refined and developed to open new avenues through which users can tailor their collections. We demonstrate its flexibility by showing how digital library collections can be extended and altered to satisfy new requirements.*

Keywords:

1 Introduction

The Greenstone digital library software provides a wide range of tools for building digital library collections [10]. Based on our own extensive and varied experience, and that of others [9], we have designed a new framework for building digital library collections. We call it “Greenstone 3” to distinguish it from the earlier system, “Greenstone 2”. The new framework is supported by a fresh implementation that is completely independent of the existing one. It capitalizes on the best ideas from the existing system, which we have further refined and developed to open new avenues through which users can tailor their digital library collections. Several pertinent new open standards have emerged since the original design many years ago, and a key objective is to incorporate them into the new design.

The trend towards increasingly open, flexible architectures can be traced in the development of digital library protocols. The Open Archives Protocol for Metadata Harvesting has provided a simple base-line for metadata access, and subsequent work strives to base component-based, modular protocols upon it [7]. The METS document framework [4] provides an open, extensible system for representing documents in digital repositories. Greenstone 3 adopts the same approach for a range of digital library functions, as this paper demonstrates.

Set against the move towards standard protocols, today’s digital library systems must confront an increasing range of document formats and media, architectural designs for browsing and classification, indexing requirements, and user

interface techniques. This drives the demand for open architectures, but makes it hard to provide well-formed and supportive infrastructures that maximise reuse of individual components and the reliability of the system as a whole.

This paper proceeds as follows. First we introduce and describe the configuration structure and document representation used through our new building architecture. The second section describes the building process itself. This is followed by the model that we have adopted to ensure that the new architecture is extensible. Fourth, we compare this with the previous Greenstone 2 design and present an example use-case scenario, identifying the key advantages of the new open collection-building architecture. We conclude by relating the new architecture to existing collection-building systems.

2 Collection Configuration and Document Representation

Before describing the process by which collections are built, we will show how the collection-building process is configured, and how documents are represented internally within the system. In [11] we described our “plugin” mechanism and showed how it provides great flexibility for such varied tasks as ingesting both document files and metadata files, extracting text and metadata from different file formats, and expanding compressed files. It has proved so successful that we have adopted it in the new framework. However, we depart radically from our earlier system by using the METS framework (which postdates that system) throughout for internal document representation.

2.1 Configuration

Collections are designed individually, and the structure of a collection is encapsulated in an XML file called the “collection configuration file.” Its contents include build-time configuration options such as:

- The document types that the collection should recognise
- The metadata access structures or “classifiers” that are to be provided for users to browse the collection, such as by Titles A–Z
- The full-text indexes to build for searching the collection.

The configuration file also contains run-time information about the collections, for example display options.

The bulk of an example configuration file can be seen in Fig. 1—all that is omitted is the part concerning the purely runtime configuration. It defines a simple collection of HTML pages and plain text files, along with XML files that define metadata that is to be associated with these documents. A document may consist of several separate files—for example, in this collection any images associated with the HTML pages will be contained within the HTML document rather than being identified as documents in their own right. Or a document may comprise complementary pairs of files—like a Word file and matching PDF version.

```

<documenttypes>
  <document type="HTML">
  <document type="Text">
  <document type="MetadataXML">
</doctypes>

<search type="mg">
  <index name="dtx" level="document" field="text">
    <displayItem name="name" lang="en">entire documents</displayItem>
    <displayItem name="name" lang="fr">documents entiers</displayItem>
    <displayItem name="name" lang="es">documentos enteros</displayItem>
  </index>
  <index name="stx" level="section" field="text">
    <displayItem name="name" lang="en">chapters</displayItem>
    <displayItem name="name" lang="fr">chapitres</displayItem>
    <displayItem name="name" lang="es">cap?tulos</displayItem>
  </index>
  <index name="stt" level="section" field="title">
    <displayItem name="name" lang="en">section titles</displayItem>
    <displayItem name="name" lang="fr">titres des sections</displayItem>
    <displayItem name="name" lang="es">t?tulos de las secciones</displayItem>
  </index>
  <format>
    <gsf:template match="documentNode"><td valign="top"><gsf:link><gsf:icon/></gsf:link>
    </td><td><gsf:metadata name="Title"/></td></gsf:template>
  </format>
</search>

<browse>
  <classifier name="CLSubject" type="Hierarchy" file="subject.xml" field="Subject"/>
  <classifier name="CLTitle" type="AZList" />
  <classifier name="CLOrganization" type="Hierarchy" file="organization.xml"
  field="Organization"/>
  <classifier name="CLKeyword" type="List" file="howto.xml" field="HowTo">
    <displayItem name="name" lang="en">HowTo</displayItem>
  </classifier>
</browse>

```

Fig. 1. A collection configuration file

The *search* block in Fig. 1 defines the full-text-searchable indexes that will be built. In this case the MG indexer[12] will be used to build document-level and section-level indexes of the main text, and a section-level index of *Title* metadata. The *browse* block defines the metadata-based browsing mechanisms through which users can access the collection. This configuration file specifies a hierarchical subject browser by defining a Hierarchy-type classifier of each document's *Subject* metadata, and the file that defines that hierarchy is named. Other browsing classifiers are also defined. Run-time configuration items can be seen in both the classification and indexation parts of the file. For example, the "displayItem" nodes give names for the indexes (in English, French and Spanish) and for one of the browsers.

The Greenstone 3 building program reads this file and configures the components to be used by the collection-building process. Many default options are implied. For instance, this collection would automatically expand Zip files—because the Zip expansion plugin is enabled by default. The process followed during building will be described later.

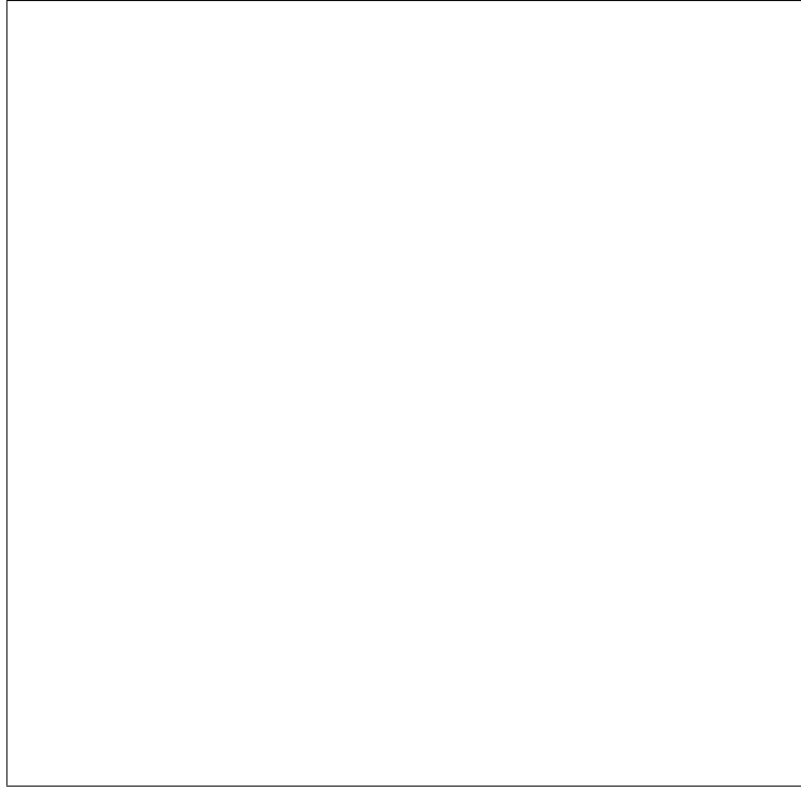


Fig. 2. An example METS Document

2.2 METS and Document Representation

We were eager to adopt an open, standard form for modelling the content and metadata of a document. The METS framework supports a free selection of metadata formats and avoids duplicating data that is already stored in constituent files [4]. Using this standard, open framework should facilitate interoperation with other digital library systems.

Fig. 2 illustrates a simple METS document consisting of two chapters, each stored in an HTML file, linked by a central HTML file. Some bespoke metadata is kept on each chapter, whilst the document-level metadata is stored in Dublin Core format.

METS documents contain up to seven sections, of which five are illustrated in Fig. 2. These are the METS Header, Administrative Metadata, Descriptive Metadata, File lists and Structural lists. The structural content is the only required section, and the content of the remaining sections is mapped onto the encoded structure. For the sake of clarity, this example shows only the connections between sections and the descriptive metadata and files.

A document is encoded into METS form by a plugin that corresponds to the type of the document. Each such plugin extracts any appropriate metadata from the document files (e.g. title), catalogues the files that comprise the document, and extracts the document text in subsequent phases (extraction and indexation; see Section 3). In other words, it encapsulates the document content and acts as a standard gateway to its METS representation. This means that other phases of the building process are able to access data without having to understand the particulars of each file format—even if the data is stored as a binary object.

Document text can be provided in one of three ways, as the processing component requires. Plain Unicode is a simple baseline recognized throughout the system. At a higher level of sophistication, an XML representation is supported. Finally, indexers can use the underlying binary files if they so desire. A complete document plugin provides all these formats, though the XML representation is optional. In any case, a document may have no textual content whatsoever (in which case all forms are empty). For example, a PostScript file can be translated into an XML format. The PostScript plugin allows a plain-text reduction of this, but if an indexer wishes to treat PostScript files in a special way, it can access the original files directly.

METS documents are partitioned into several components, the most common elements being structure, files, administrative metadata and descriptive metadata. A given document may form part of a single file, or be an agglomeration of individual files and parts of files. Metadata can be associated with a document as a whole or with any part of its structure. A document may also have more than one structural form. For instance, it may have both a page based structure and a logical chapter-and-section based structure. Metadata may be provided in several different formats and standards, and this allows for the native coding of metadata on the document in different formats—for instance RFC1807 and Dublin Core.

When a document first enters the library, it is given a unique identifier. That identifier will remain with the document throughout any subsequent revisions, and is recorded within the METS framework.

Having described the configuration controls and document representation that underpin the new collection-building architecture, we now describe the building process itself.

3 Building Digital Library Collections

In our new architecture, collections are built in several distinct phases, which occur in sequence.

Expansion. Compressed files such as Zip archives are expanded, and links to web sites are expanded into lists of constituent web pages.

Recognition. All files are sent to the RecognitionManager, which identifies groups of files as documents.

Encoding. All the recognised documents are catalogued for the subsequent phases of building.

Fig. 3. The Greenstone 3 Building Architecture

Extraction. Every document is passed through extractors, which use special processing algorithms to extract information from the document (e.g. title, keyphrases) or add metadata stored in special files.

Classification. Documents are assigned to classifiers (e.g. topical classifiers, ordered list classifiers) depending on their inherent and extracted metadata.

Indexation. Documents are sent to indexers to build indexes that support later searching.

Validation. Post-building checks are carried out on the collection as a whole, and on its constituent documents.

The rest of this section describes each phase in greater detail.

3.1 Expansion

Before any individual documents are identified, any files that are compressed, or refer to external sources, are expanded. The case of a Zip archive has already been mentioned; other cases include OAI repositories, and links to web sites via a file of URLs. The purpose of this phase is to create a complete list of all files that are available for building the collection. These raw files are then passed to the next phase for grouping into documents.

3.2 Recognition

The relationship between files and documents is complex. A file of emails may be seen as one file that contains many documents. Conversely, a JPEG image and a plain text file that describes it may constitute a two-file document. Our architecture uses special “recogniser” plugins to address the problem of distinguishing documents from files.

When the building process is initialised, all plugins listed in the collection configuration file are registered by a “Recogniser Manager.” When the Expansion phase terminates, this Manager is handed each file in the collection in turn. It sends these files to each plugin. For some simple cases—like a plain text document comprising a single file—the plugins are trivial. However, in other circumstances conflicts can emerge, and the Recogniser Manager includes a scheme for resolving these.

For example, suppose a collection includes both web pages and digital photographs. A single JPEG file may be a document in its own right, or form part of a web page. By default, JPEG files that are referred to by web pages, being constituent parts of a larger whole, are not treated as documents in their own right, while JPEGs that do not appear in the web pages are identified as separate, discrete documents. This default behaviour can be adjusted, if desired, to allow *all* JPEG files to be treated as documents, even ones that also appear in a web document.

The same issue arises with any image file on a web page. Furthermore, such a file (for instance a logo) could well be part of more than one document. In this case it is treated as a subcomponent of all the documents in which it appears. A web page and associated files cannot form a strict subset of another web page, because the .html files must necessarily differ. The same strategy is used in other cases.

Once the Recogniser Manager has taken a pass through the entire inventory of files and identified the documents, it resolves any conflicts that occur and hands the output on to the Encoding phase. At this stage all that is known of a document is the list of files that constitute its content. No metadata has yet been assigned—even that found within the document itself. The content and representation of any metadata internal to the document is identified in the next phase.

3.3 Encoding

Documents appear in a wide variety of forms, textual encodings and file formats. The Encoding phase creates the document framework that is used throughout the remainder of the building process. This framework supports the retrieval of the full text of the document. It contains metadata defined within the document, and will also contain any metadata added to it. Finally, it binds the document, with its representation, into the classification and indexation structures. The METS document framework described in Section 2.2 is used to encapsulate the files, metadata and structure of each document. At the end of the Encoding phase, documents are ready to be processed by extracting further metadata, indexation and storage.

3.4 Extraction

The Extraction phase is controlled by the Extraction Manager, which filters the documents through a series of extraction plugins. These can have one of two origins:

- Metadata-only files identified in the recognition phase above
- Explicit mention in the collection configuration file

The former method is typically used for files of descriptive metadata that are provided with the documents. The latter is used to invoke automatic extraction tools such as the KEA Keyphrase Extractor [5, 13].

At the conclusion of the Extraction phase, all metadata has been collected and stored in the METS document representation.¹ The next two phases place the documents into the context of the collection as a whole.

3.5 Classification

The Classification phase uses each document's metadata to site it within the browsing structures, or "classifiers," that are supplied by the collection to access the documents it contains. Each classification is represented by a plugin specified in the collection configuration file. There may be more than one instance of any given classification type, achieved by initiating a separate copy of the plugin for each classification. Classifications are stored in a database, and when a document is matched to a node in a classification, this fact is recorded in both the document and the classification node.

3.6 Indexation

Indexing follows classification. To support the widest range of indexing systems, indexers are also treated as plugins. Again, the indexer—or indexers—used when building a collection are determined by the collection configuration file. For each

¹ George, is this true?

indexer, parameters can be supplied that control its behaviour. Indexers can build several different indexes, each with its own parameters. In the case of MG[12], for example, the structure level (e.g. document, section) and field (full text, title, etc.) can be controlled. In the example configuration file given in Fig. 1, section- and document-level indexes of the full text are built, as is a section-level index of titles.

3.7 Validation

Validation plugins are used to control the quality of the collection by checking the final form of each document, and the collection as a whole. For example, a collection may insist on a complete metadata set—that is, all compulsory fields must be present. Or the databases may be checked for empty classification nodes, which may be suppressed. Validation could be done in earlier phases, but we separate it into a final phase for architectural clarity.

3.8 Summary

Greenstone 3 uses a seven-phase build architecture. The first three phases identify and store each document. The next three phases enrich and index those documents. The final phase provides for quality control. The phases support different behaviours and activities that are required by different types of software, and a set of distinct functions that should assist programmers in extending any given implementation of the architecture.

Individual documents are recorded in the METS framework, which provides a standard form that gives access to different metadata formats, file formats and document structures. The document plugins give the indexation and classification components access to content without having to decode or transcribe documents.

4 Extendibility: Managers and Plugins

The central role of plugins in each phase should now be apparent. However, we have made little mention of the structure of plugins themselves, or how they connect to the core system. The architecture achieves extensibility through *Plugins* and *Managers*. Each phase of the building process is controlled by a Manager—e.g. the Recogniser Manager, the Indexer Manager. Managers are configured through the collection configuration file when the building process starts; and in some cases further configuration occurs when special files—like metadata-only files—are found when building. Each manager coordinates the plugins for its phase of the build cycle.

Each manager works in the same way. It initialises and configures its plugins. Then, when its phase runs, it takes each plugin in turn and tells it to do any necessary preparatory work. Next it passes each document or file in turn to that plugin. Finally it sends the plugin a termination message.

Document plugins are unique in that they apply to several different phases of the building operation, whereas other plugins appear only within one phase. Managers accept a list of plugins from the configuration file. New plugins are readily installed into the system by placing them into the plugin directory corresponding to their type, and adding them to the configuration file of one or more collections.

There is a standardised interface for all extraction, indexation and classification plugins. This includes:

getNumberOfPasses Return the number of passes required by the plugin

startPass(int) Begin a pass for this plugin

processDocument(Document) Process the given document

endPass(int) Terminate a given pass

tidyup Close the plugin

configure(Node) Pass configuration information to the plugin.

The configuration information mentioned in the last method is taken from the collection configuration file, where it is specified as an XML node. It is particular to each plugin.

The *getNumberOfPasses* method allows plugins to request multiple passes over the documents. The MG indexation plugin, for example, requires a minimum of two passes plus a further two passes for each index it builds. Running each plugin separately and serially keeps memory overheads down.

5 Advance over earlier work

The architecture that we have described capitalizes on lessons learned from the existing Greenstone digital library system, and incorporates some very significant improvements. At the time the earlier system was designed (1998) several important open standards did not exist, or were available only in draft form. For example, the METS Document Framework, a key component of the new architecture, was unborn.

The original Greenstone design adopted *ad hoc* file formats that became inconsistent with emerging norms. Collection configuration was viewed as a compilation process directed by a configuration file whose format was similar to the Makefile format [2] rather than by directives couched in a structured markup language. It used an HTML-like structure, the Greenstone Archive format, to store documents and metadata, which is less structured than modern XML. The internal formats were also inconsistent with each other—configuration files being flat files whereas document files were derived from HTML. These simple and apparently cosmetic differences magnified into large-scale divergence.

5.1 Architecture

In Greenstone 2, collections are constructed in two phases: *importing* and *building*. The first parallels the Expansion, Recognition and Encoding phases of

Greenstone 3, while the second mirrors the Extraction, Classification and Indexation phases. Both phases use the same set of plugins, which are listed in the collection’s configuration file and loaded separately in each phase—despite the fact that some plugins only pertain to one phase.

Document plugins, like *HTMLPlug*, played a rather different role from that described above. Documents were identified by a single key file—e.g. an HTML file—and complex document structures were poorly supported. A special generic plugin was later developed that recognized multi-component documents such as a Word file and its corresponding PDF version. However, the plugins had to be carefully ordered in the configuration file to recognize this compound document. Once a file was captured by a plugin, subsequent plugins could not see it. The Recognition phase described above facilitates complex document forms and arbitrates competing claims for a file by different document types.

5.2 Plugin details

Document plugins handle Encoding very differently in the new design. Originally, all documents were encoded into the standard Greenstone Archive Format as soon as they were encountered. This duplicated content, and could lose, or render inaccessible, some information in the original file. The benefit was that all documents were presented to subsequent phases in a standardised format.

In the new architecture, plugins convert document content to either plain text or well-formed XML, which allows formats such as TEI [6] to be retained verbatim. The METS Framework provides a standard container for metadata and structural information. This open public standard retains the benefits of the original Greenstone Archive Format, and goes further by allowing multiple structures within the same document.

Plugins in Greenstone 2 could be of three types—Expansion (then called “structural”), Document, and Extraction (called “metadata”). All were placed in the same pipeline. The new architecture makes a clearer distinction between phases and allows a richer variety of documents. For most tasks, it supports clearer, more compact, interfaces. In addition, monolithic parts of the old design, such as indexation, are now readily adjusted and extended in the flexible manner exhibited by other phases (such as Expansion).

5.3 Example: The Kids Digital Library

We briefly present an actual digital libraries that was difficult to accommodate within the earlier design, and show how it benefits from the new architecture. In the Kids Digital Library [8] each document can belong to several collections. Some collections are private (e.g. a child’s own documents), others public. Some documents are unchanging (accepted final essays); others are under continually revision by a restricted group of users. Students can annotate the work of others, and teachers provide feedback too. All collections support full-text searching

and metadata-based browsing; some have additional browsing facilities (see below). Collections sometimes change rapidly over a short period—e.g. during a scheduled class.

Using the new architecture, many of these requirements can be delivered far more easily. Consistent document identifiers and explicit change history allows revisions of a document, and comments made between changes, to be tracked. Incremental indexing is supported by recording accession and revision data as METS administrative metadata. The METS framework allows files to be duplicated between collections yet reside in only one location.

The Kids Digital Library featured some unusual browsing classifiers such as the “Top ten” and “Latest ten” stories. These require simple support for feature extraction. In the earlier design, this could only be done by altering the building scripts—which caused portability problems, and made it difficult to install new versions of the software. The new architecture circumvents the problem by modularizing feature extraction.

6 Comparison with other Digital Library Systems

Cheshire II [3] emphasises the construction of digital libraries from original scanned documents. The process described for collection building reveals a system of fixed metadata fields and a strict control of the format in which documents are presented to the system. Nowhere is support for feature extraction, expansion of compressed files, or novel indexes described.

The CORR² is built on the NCSTRL software [1]. The CORR documentation reports that the system requires documents to be submitted in a standard source format. No documentation on CORR or NCSTRL describes the parameters for collection configuration.

6.1 Discussion

7 TODO

NCSTRL, CORR, Dienst, Eprints, Dspace, OpenDLib, Cheshire, Citadel (Ed Fox), Mariam, Perseus, Alexandria, NSDL

8 Conclusion

In the new collection building architecture we have described, the building process is segmented into a number of distinct phases. Once documents are identified, they are encoded into a flexible, open framework (METS) and are passed in that form to the succeeding phases of the build process. Within each phase,

² <http://arxiv.org/archive/cs/intro.html>

the elements are componentised to support greater portability and simpler development. The build process is configured through a simple XML format file which is readily extensible for future components.

We have implemented this architecture in a practical digital library system. The use of a document framework rather than standardised document encoding has not created any obstacle to providing existing features. Rather, it supports our move to more complex document file structures and permits any processing of the document (e.g. by an indexer) to either refer to the original binary or to extracted content as it chooses.

We believe that the challenge of providing effective, portable solutions to building digital library collections is an important one. As file formats diversify, the issue of digital preservation grows. Similarly, if extraction and indexation technologies are not made more portable their adoption will slow. Our new architecture provides a step forward in meeting these challenges.

References

1. J. R. Davis and C. Lagoze. Ncstrl: Design and deployment of a globally distributed digital library. *Journal of the American Society for Information Science*, 51(3):273–280, 2000.
2. Free Software Foundation. *GNU make Manual (version 3.80)*, 2002.
3. R. R. Larson and C. Carson. Information access for a digital library: Cheshire ii and the berkeley environmental digital library. In *Proceedings ASIS '99*, pages 515–535. Information Today, 1999.
4. Library of Congress. *Metadata Encoding and Transmission Standard (METS)*.
5. G. W. Paynter, I. H. Witten, S. J. Cunningham, and G. Buchanan. Scalable browsing for large collections: A case study. In *Proceedings of the Fifth ACM International Conference on Digital Libraries*, pages 215–218, June 2000.
6. C. Sperberg-McQueen and L. Burnard, editors. *Guidelines for Electronic Text Encoding and Interchange*. TEI P3 Text Encoding Initiative, Oxford, 1999.
7. H. Suleman and E. A. Fox. Designing protocols in support of digital library componentization. In *Proceedings of the 6th European Conference on Research and Advanced Technology for Digital Libraries*, pages 568–582. Springer-Verlag, 2002.
8. Y. L. Theng, N. Mohd-Nasir, G. Buchanan, B. Fields, H. Thimbleby, and N. Cassidy. Dynamic digital libraries for children. In *Proceedings of the first ACM/IEEE-CS joint conference on Digital libraries*, pages 406–415. ACM Press, 2001.
9. I. H. Witten. Examples of practical digital libraries: collections built internationally using greenstone. *D-Lib Magazine*, 9(3), 2003.
10. I. H. Witten and D. Bainbridge. *How to build a digital library*. Morgan Kaufmann, San Francisco, CA., 2003.
11. I. H. Witten, D. Bainbridge, G. W. Paynter, and S. J. Boddie. Importing documents and metadata into digital libraries: Requirements analysis and an extensible architecture. In *Proceedings of the European Conference on Digital Libraries*, pages 390–405, Sept. 2002.
12. I. H. Witten, A. Moffat, and T. C. Bell. *Managing gigabytes: compressing and indexing documents and images. (second edition)*. Morgan Kaufmann, San Francisco, CA., 1999.

13. I. H. Witten, G. W. Paynter, E. Frank, C. Gutwin, and C. G. Nevill-Manning. KEA: Practical automatic keyphrase extraction. In *ACM DL*, pages 254–255, 1999.