

**Can document-genre metadata improve information access
to large digital collections?**

Kevin Crowston and Barbara H. Kwasnik

Syracuse University

School of Information Studies

4-206 Centre for Science and Technology

Syracuse, NY 13244-4100 USA

Telephone: +1 (315) 443-1676

Fax: +1 (315) 443-5806

Email: {bkwasnik, crowston}@syr.edu

Draft of September 21, 2003

Submitted to *Library Trends*

**Can document-genre metadata improve information access
to large digital collections?**

Abstract

We discuss the issues of resolving the information-retrieval problem in large digital collections through the identification and use of document genres. Explicit identification of genre seems particularly important for such collections because any search usually retrieves documents with a diversity of genres that are undifferentiated by obvious clues as to their identity. As well, because most genres are characterized by both form and purpose, identifying the genre of a document provides information as to the document's purpose and its fit to the user's situation, which can be otherwise difficult to assess. We begin by outlining the possible role of genre identification in the information-retrieval process. Our assumption is that genre identification would enhance searching, first because we know that topic alone is not enough to define an information problem and second because search results containing genre information would be more easily understandable. Next, we discuss how information professionals have traditionally tackled the issues of representing genre in settings where topical representation is the norm. Finally, we address the issues of studying the efficacy of identifying genre in large digital collections. Because genre is often an implicit notion, studying it in a systematic way presents many problems. We outline a research protocol that would provide guidance for identifying Web document genres, for observing how genre is used in searching and evaluating search results, and finally for representing and visualizing genres.

Can document-genre metadata improve information access to large digital collections?

Introduction

Computerized information-access systems face a fundamental limitation: they know what documents say, but not what they mean or for what purposes they might be useful. Extracting and representing the meaning of documents is difficult and time-consuming to do, and automatic systems still have significant limitations. We note though that humans rarely have to read every word of a document to understand its purpose. Instead, people take a shortcut: they start by identifying the kinds of documents they are faced with (i.e., the document's genre), and then use different types of documents in appropriate ways. For example, a grant proposal is used differently from a syllabus, a product brochure or a bank statement. Accordingly, differences in an information situation are often reflected in the *kind* of document that is considered helpful (e.g., a problem set, a lesson plan and a tutorial about mathematics are all about math but useful in different situations). Information-access systems would be more useful for many tasks if they could similarly distinguish the purpose of documents and handle them in appropriate ways.

In this paper we discuss the possibility of improving information access in large digital collections through the identification and use of document genre as a facet of document and query representation. First, we provide some historical background on the concept of genre and the approach it provides to the problem of incorporating context into information retrieval. We outline the framework of the information-retrieval problem

Document-genre

with respect to genre, and some traditional resolutions that have been attempted. Finally we outline a research agenda that addresses some of the questions and issues that investigating genre entails.

Theory: Document genre

Rhetoricians since Aristotle have attempted to classify communications with similar form or purpose into types or “genres.” Numerous definitions of genre, or discourse type, have been suggested (e.g., Longacre, 1983; Miller, 1984; Swales, 1990). In our discussion, we draw on the definition of genre proposed by Yates and Orlikowski (1992), who describe genre as “a distinctive type of communicative action, characterized by a socially recognized communicative purpose and common aspects of form” (Orlikowski & Yates, p. 543). For instance, this document is an example of the journal-article genre. It has a form familiar to most researchers and practitioners, and is monitored by the journal’s editorial policies as well as the profession’s communication practices. There are many document genres: some common, such as a report or a newsletter, and others restricted to specific domains, such as the course syllabus or a problem set in higher education. Genre is applicable to electronic as well as physical documents. For example, in a study of Web documents, Crowston and Williams (2000) were able to identify documents of many familiar genres and of a few genres that seemed to be new to the Web, such as the home page (Dillon and Gushrowski, 2000) or the hotlist.

Genre is useful because it makes documents more easily recognizable and understandable to recipients, thus reducing the cognitive load of processing them

Document-genre

(Bartlett, 1932/1967). Yates and Sumner (1997) argue that on the Web, genres help in both the production and consumption of documents because genre adds “fixity” in a medium that does not otherwise distinguish very well between text types (say, a book and a post-it). For example, since this article conforms to the journal-article genre, a reader can quickly determine the purpose of our communication, locate relevant sections, evaluate the document’s contribution and possibly to use it to prepare build further queries. In our preliminary studies of people searching the Web (Roussinov, Crowston, Nilan, Kwasnik, Liu & Cai, 2001), we observed that the genre of the document was one of the clues used in assessing document relevance, value, quality and usefulness.

The problems of information access.

To explain how genre can be useful, we will first briefly review the problem faced by an information-access system. An information-access system has three components: 1) the users, who approach the system with contextually-based information needs, 2) a store of information (e.g., the documents or databases), and 3) an intermediating mechanism to connect needs and information. The intermediary may be a person, a search algorithm, a browsing environment, or a summarizer, among others.

The basic process of matching users’ needs to potentially useful information in the system is complicated by many factors. First, problems may occur due to improper or incomplete representations of the information itself. When the information-access system is created, the documents or texts must be represented in such a way that they can be retrieved again as needed. Librarianship has occupied itself for over a century with systematic approaches to organizing and representing information in systems. In creating

Document-genre

bibliographic records we call this process cataloging; in organizing actual documents or topics for meaningful retrieval, we call it classification; in providing access to bibliographic databases we call the representation process indexing.

There are similar processes of information representation on the Web and in many other applications in which large stores of information are prepared for eventual use in the future. Many of the schemes are adaptations of traditional schemes, such as that used on Yahoo.com, the Dublin Core Project, or the GEM Metadata Project for educational materials (<http://www.thegateway.org/>). Others comprise grass-roots, emergent sorts of organization and representations, such as the evolving classification on eBay.com or on amazon.com) (Kwasnik & Liu, 2000; Kwasnik, 2002). An increasingly popular approach relies solely on the full-text of the documents.

Another problem that may arise is that the process itself of matching users' queries to the document representations may be inadequate or faulty. Much effort on the part of information scientists has been spent in developing and perfecting search strategies, including various matching algorithms, probabilistic techniques, citation mapping, and natural language processing. These efforts struggle with many obstacles—among them the difficulties of evaluating search results in real environments, as well as problems of scaling, reliability and the representation bottleneck.

On the user side, we have people in need of information. Often, though, users are unable to precisely specify what it is they need, and even if they do, the way in which humans articulate their needs produces a great variety of expression. The problem of appropriate representation of users' queries is not just a question of finding the correct

Document-genre

representation according to some absolute criteria. Because information use is situated in specific contexts, there is also the need to be able to represent the information in such a way that a match can be made not only on the level of physical description and topic, say, but also in terms of matching the information with a potential use. For example, consider a person approaching a system with the query “I want to prepare a Passover dinner.” At a certain level we can see that there is a need for concrete information in the form of actual recipes. We might even interpret this as a “known item search.” However, recipes may satisfy the need only partially, since the person may want to know much more about the rituals and meanings of a Passover dinner and not just the food itself. The information need may be either broader than what is asked for, or much narrower and specific. We know that people ask for what they expect they can get that will most closely match what they *really* want, and thus their requests are often presented in a compromised form.

Thus, we can see that topic alone is not enough to define an information problem because different users may require different solutions to seemingly similar information problems. Indeed, even the same user may require different information at different times. These different needs arise because the situation (or context) of a user determines not only what topics are requested and what strategies are invoked in searching and evaluating output, but also what types of resources are considered relevant and useful. For example, methods for learning mathematics (a topic) may be construed differently by a student, by a parent and by a classroom teacher because of their different information-use situations. While we know that it is important to understand the situation of the user, the representation of the situation and then its implementation in a system is a difficult problem. Our efforts to create user profiles, universal situation grammars, and so on

Document-genre

suffer from limitations of scope to specific domains and lack of extensibility and flexibility.

Why we think identification of genre would be useful

We suggest that enhancing document representations by incorporating non-topical characteristics of the documents that signal their purpose—that is, their genre—will enrich document (and query) representations in such a way that they resonate more truly with the information need of a user as situated in a particular context.

Because most genres are characterized by both form and purpose, identifying the genre of a document provides information as to the document's purpose and its fit to the user's situation, which can be otherwise difficult to assess. For instance, a university professor looking for information about computer database systems for the class that she teaches would most likely be interested in documents of educational genres (e.g., syllabi, assignments, class notes). On the other hand, when working on a research paper in the database area, the same professor would more likely appreciate scholarly work (e.g., papers, annotated biographies, calls for papers). The relevant documents for these two searches would be quite different, even though the topic and query keywords might be nearly the same.

Explicit identification of genre seems particularly important for large digital collections because any search of these collections usually retrieves documents with a diversity of genres, and, what is worse, these genres are undifferentiated by obvious clues to their identity. This is in contrast to non-digital information-seeking situations in which the searcher generally has an idea of what sorts of documents exist in the collection. Even

Document-genre

if he or she does not, clues of physical form and location increase the chances that a document's genre is recognized. For example, a user searching in a library can visually distinguish CDs from books, from encyclopedias, or from newspapers. Similarly, a user searching a database containing only journal articles has already implicitly restricted the search to that genre of documents. On the Web, however, a search of a large and diverse document collection will usually retrieve some documents of relevant genres along with many documents of irrelevant genres—a low precision result—even if all retrieved documents conform to search specifications regarding the topical content of the document.

Recognition of genre also has implications for automated methods of representing documents, such as automated summarization and indexing. A one-size-fits-all approach to summarizing or evaluating Web documents without regard for their form and function is likely to misrepresent many of them. For example, a newspaper article can be summarized by the first few sentences of the document, but such an approach will not work for a home page or a frequently-asked-questions document (FAQ) (Marcu, 1997). When medical information is sought, it makes a difference to the evaluation whether the document retrieved is a newsletter, a personal homepage, or a hospital's patient-information site.

How librarians have addressed the notion of genre in library information systems

We do not mean to imply that information science has never addressed the notion of genre, or that genre has not been incorporated into any information-representation schemes. Indeed, several classification systems allow some articulation of genre, and

Document-genre

many metadata standards, including the Dublin Core, include a field for genre. The treatment of genre is limited or not very well defined, however. Our understanding of the nature and role of document genre is still nascent, and so the use of this kind of information is underdeveloped in information-retrieval systems. Furthermore, it is not clear whether the extension of genre designations originally designed for physical collections will export well to digital ones.

Historically, most library information systems took genre for granted since most collections contained only a limited array of document types. The exceptions are literary genres (such as poetry) and publication types (such as almanacs or newspapers), which have had a lively existence in explicit document representation for several centuries. Aside from these, the primary facets of access to documents in traditional systems are descriptive components and subject, while genre is relatively rare. The descriptive access points derive from traditional ways of talking about books and book-like documents, and include: title, author, place and name of publisher, edition, date, series, physical description in terms of pages, size, volumes, and sometimes information about components, which are called *analytics*. The subject analysis of a document captures what a document is about—that is, its topic.

Librarians and information scientists have recognized that the topical approach is extremely important, but insufficient in some situations, and completely inappropriate in others. Not every document is necessarily about something. Sometimes the document's nature as a document represents the most important or useful aspect of it. For instance, on the one hand we can say that a book may be *about* symphonies—their history or

Document-genre

structure—but what is Beethoven’s Fifth *about*? It simply is. A symphony has a form and identifiable characteristics but it does not have a readily identifiable topic, *per se*, except that which can be attributed to it through subtle and non-consensual processes of interpretation. As the notion of *document* becomes broader and more diverse, as it does in the environment of the Web, we can see how the concept of subject does not stretch very well to cover all types of information.

In response to the need to identify a document’s form or genre in addition to its subject, librarians have created auxiliary tools in the form of tables and subdivisions to be used with existing topically based classification and subject-heading schedules. Here are a few examples:

The *Dewey Decimal Classification (DDC)* (Dewey, Mitchell, Beall, Matthews & New, 1996) provides several ways to denote a document’s form or genre. The first is to incorporate a designation in the number itself. This is used in the 800s, which cover *belles lettres*. The first part of the number designates country/language, and the final digits represent the genre—1 for poetry, 2 for drama, 3 for fiction, 5 for speeches...7 for humor and satire, and so on:

Document-genre

English poetry	821
English drama	822
English fiction	823
English speeches	825
English humor and satire	827
Bulgarian poetry	891.81
Bulgarian drama	891.82
Bulgarian fiction	891.83
Bulgarian speeches	891.85
Bulgarian humor and satire	891.87

These genre designations are limited to the genres generally accepted by Western literary scholars, and do not necessarily do a good job of describing emerging, culturally diverse, or hybrid genres. Still, it is a way of privileging genre in the organization of literary works. It is interesting to note, however, that most public libraries do not make use of this formal system for fiction and arrange such works by author, with the *ad hoc* tradition of separating out popular genres into separate sections for easy access and browseability: Mysteries, Romances, Science Fiction, and so on.

Another technique in the *DDC* is to use suffixes from the Tables. The number for the topic is established, and then suffixes from Table 1 are added to denote the form or genre. For example:

Middle Eastern Cooking	641.5956
Middle Eastern Cooking Encyclopedia	641.5956+03 (dictionaries & encyclopedias)
Middle Eastern Cooking Magazine	641.5956+05 (serial publication)

In physical collections, the suffixes serve to distinguish materials on the same topic but in different publication formats one from the other. This notion of form/genre evolved from the physical distinctions of publication and document types, and thus is grounded in publishing practices and realities. The further interpretation of how such documents will be used remains implicit in the nature of the forms themselves, but has

Document-genre

practical implications for collections. For instance, many dictionaries and encyclopedias comprise the non-circulating reference collection; magazines are indexed and stored differently than are books, and so on. In terms of digital collections, however, where the physical clues of publication format are largely absent, these suffixes might provide useful indicators for sorting and filtering search results.

Another way in which subject is indicated on the bibliographic record is through the use of subject terms from a thesaurus or list, such as the *Library of Congress Subject Headings (LCSH)*. The *LCSH* comprises an evolving list of terms used by catalogers to assign subject designations to a work. Terms can denote topics, such as “sonnets,” in which case this would be a work *about* sonnets, not the sonnets themselves. Proper names may also be subjects. For example, a document *about* William Shakespeare will be assigned Shakespeare’s name as a subject, while a work *by* William Shakespeare would not. Modern cataloging practices abound in confusions about topic, creative responsibility and genre/form, since in many documents these three are inextricably fused. This confusion extends to searchers as well, who do not realize that searching for a genre using *LCSH* is problematic at best.

This distinction of reserving subject headings for topics/subjects only, is somewhat moderated by the addition of a subdivision. There are several kinds of subdivisions that can be used to “subdivide” a subject by time, geographical location, and further topical aspects. For example:

Witchcraft—Sweden
Witchcraft—15th Century
Witchcraft—Biblical teaching

Document-genre

The subdivisions of interest here, though, are the ones from the Form Subdivisions list. This type of subdivision allows the cataloger to further describe a work by its form or literary genre. This list is limited to several hundred well-established types. The genres included have literary warrant, since every subject heading and division in the *LCSH* was developed for an existing, rather than a hypothetical, work.

- Witchcraft—Bibliography
- Witchcraft—Case studies
- Witchcraft—Dictionaries
- Witchcraft—Handbooks, manuals, etc.
- Witchcraft—Periodicals
- Witchcraft—Poetry

The fact remains, however, that form and genre are not, as a rule, an important finding aid in traditional systems. For instance the work: *Final environmental impact statement for the Green Mountain and Finger Lakes National Forests land resource management plan* is assigned the following subject headings from the *LCSH*.

- Forest reserves—Vermont Green Mountain National Forest
- Forest reserves—New York (State)—Finger Lakes National Forest
- Forest management—Vermont Green Mountain National Forest
- Forest management—New York (State)—Finger Lakes National Forest
- Green Mountain National Forest (Vt.)
- Finger Lakes National Forest (N.Y.)

Thus, this work can be retrieved by either of the two national forests covered in the report and by two topics: *forest reserves* and *forest management*. It is not possible to retrieve this work as an *environmental impact statement* except for the coincidence that the terms appear in the title and would come up on a keyword search. There are many genres, such as this one, that serve a useful purpose as templates and are of interest in their own right, aside from the specific topic, but since this work is an example of an

Document-genre

environmental impact statement, rather than about one, there is no subject heading assigned for this important aspect of the document.

Some libraries recognize that genre and form are often perceived as “topical” and have made some additional access points to accommodate this. For instance, the Rare Book, Manuscript, and Special Collections Library at Duke University (<http://scriptorium.lib.duke.edu.genre-headings.html>) has an interesting set of auxiliary tools for searching its collection. One of these is a Genre/Form list from which genre terms can be used for searching as if they were topics. Here is a sample of terms from that list:

- Accounts
- Business letters
- Manuscripts
- Official reports
- Pattern books
- Petitions
- Recipes
- Seals
- Subliterary papyri
- Tax returns
- Vouchers

It is immediately obvious how very helpful such a list might be in studying the communicative forms of the cultures represented in the collection.

How to study genre

Having presented our case for understanding more about document genres in order to enhance retrieval of information from large digital collections, we turn now to the issues of precisely how we might study this phenomenon. Because genre is often an implicit and subtle notion, studying it in a systematic way presents many problems. Our overarching question is: would identification of document genres improve information

Document-genre

access technologies in large digital collections such as digital libraries and the Web? This question cannot be answered directly, given the current understanding of genre or of genre's role in information retrieval. Thus, we envision a research agenda for investigating genre that proceeds through a series of componential studies, each of which we see as necessary for a full understanding. Thus in answering the central question with respect to genre it is necessary to investigate the following:

- The identification of Web document genres from the users' perspective and articulated in the users' own terms;
- The creation of a faceted (i.e., multidimensional) classification of these genres that can be used for controlled investigations in later stages of study;
- An investigation of how users integrate genre metadata into their own searching, evaluation, and use of documents;
- An evaluation of the degree to which incorporation of genre metadata in information-access systems makes a difference to the effectiveness of searching, sorting, ranking, and eventual use of documents; and
- An evaluation of various interfaces for visualizing and presenting genre metadata once it has been identified.

We also recognize that studying genre cannot be a once-and-for-all endeavor, since new genres are emerging all the time, and old ones are being used in ways that are different than originally conceived. Thus, we propose that any study of genre must also establish a conceptual framework from within which to design continuing investigations. That is, we need a set of working hypotheses based on what we know about genre as a

