# Archives and Computers: Description and Retrieval

Michael Middleton, School of Librarianship, University of New South Wales.
Baiba Irving, Mitchell Library, State Library of New South Wales.

The traditional forms of demarcation between information handling disciplines are being blurred by the increasing use of machine-based systems. The record descriptions and formats that were once easily placed within the respective domains of archivists, records managers and librarians are now less readily delineated. The increasing sophistication of data structuring techniques and the development of flexible data base management systems is forcing new perspectives on customary methods of document description and retrieval.

The computer is encroaching upon the realm of the archivist in two ways.[1] On the one hand, more and more archival material is appearing in machine-readable form; on the other, an increasing number of computer-generated finding aids is being developed to provide access to archival material.

The problem of describing machine-readable records has recently been addressed by a number of writers.[2] It is a matter of increasing concern because more and more institutions are turning to the capture and manipulation of information in machine-readable form in a wide range of areas such as correspondence, minutes, sales and marketing data, inventories, experimental and exploratory data.

In many cases, systems for these different areas are independently organised on word processors and/or computers within an organisation. The distinction between word processors and computers is, however, decreasing rapidly. This, coupled with management interest in integrating systems for the optimal use of information resources, has resulted in the development of sophisticated data base management systems,[3] whose ultimate aim is ease of access to all organisational information flow.

If such developments lead in the direction predicted by Strassman,[4] with an "information middleman" holding the key to machine-based records of all types through an online personal workstation, then the description of data (whether archival or not, machine-readable or not)

and the development of finding aids to retrieve that data, become more obviously part of the same problem.

This means that the ultimate retrieval of archival records (in online or alternative easily accessible forms) will become much more dependent upon the original description of the documents at the time of creation. The computer-based records management systems, such as those described by Butler and Nicholson[5] may then evolve into components of full-scale information systems; perhaps the "front-end" of archival retrieval systems.

These developments should put greatly increased pressure on records managers to have data efficiently organised. Archivists, in their turn, must recognise the imperative necessity of consulting with records managers so that the problems of description are addressed at the stage of records creation. (The same need for consultation applies in the area of appraisal).[6]

## Finding aids

Lytle[7] distinguishes provenance and subject content indexing as alternative methods of subject retrieval from archives. The emphasis of archival theory and practice on provenance has meant that content indexing has always been assigned a lesser role in archival finding aids although many do incorporate some form of such indexing. Computer manipulation techniques have potential application in all areas of archival description but they raise particularly interesting possibilities in regard to content indexing because of the extensive descriptor permutation and matching capabilities that the computer permits.

Projects to automate archival finding aids began in the 1960's. The early developments in the United States are reviewed by Hickerson, Winters and Beale.[8] In English-speaking countries, the best known systems that have been developed are SPINDEX in the United States and PROSPEC in the United Kingdom.[9]

SPINDEX had a very shaky beginning but in its SPINDEX II manifestation it has engendered much more support.[10] SPINDEX II, developed by the National Archives and Records Service in the United States, consists of a number of data entry, file maintenance, register and index production programs. Its capabilities and applications have been well-documented.[11] From it has evolved SPINDEX III, designed to use greatly enhanced photocomposition techniques for its listings and greater selection capabilities from data files.

PROSPEC[12] is the name used by the Public Record Office in the United Kingdom for its application of the Institute of Electrical Engineers computer-based INSPEC service. It too uses data entry based on designated elements of information for the creation of a data base from which are produced a number of retrieval aids (in addition to management tools).[13]

Content indexing has received particular attention at "collecting" archives, where subject retrieval is often difficult because control records were never created for, or did not accompany the transfer of, document collections. At the Baltimore Region Institutional Studies Centre, for example, an online application called ARCHON has been developed for retrieval at file unit level based on subject, geographic location, and date variables.[14] A complete file description includes:[15]

1. The file unit citation.
2. The file unit folder heading (if any).
3. Subject descriptions applicable to the file unit's contents, modified by qualifiers (subheadings) for form, aspect, and material. As many descriptor terms as are necessary to adequately describe the file unit are assigned.
4. Names of persons or corporate bodies.
5. Geographical descriptors as appropriate.
6. The range of years covered as appropriate.
7. New concepts for possible use as descriptors to cover subjects not already represented in the thesaurus.

In Australia, computerised archival finding aids have begun to emerge but only, as yet, on a very small scale. At the Australian National University Archives of Business and Labour, the regularly published overview of holdings is computer-produced and work is proceeding on the automated retrieval of cartographic information.[16] The Tasmanian State Archives uses computer listings to assist in the administrative control of its holdings. At James Cook University in Townsville, a sophisticated program has been designed to retrieve information on and about the interviews collected in the School of History's Oral History Project. Automated archival finding aids are on the long-term agenda of many institutions, such as Australian Archives, the Archives Office of New South Wales and the West Australian State Archives.[17] As yet, however, discussion and planning has been institution-based and there has been no regional or national consideration of matters such as standardisation of terminology and descriptive practice.

### Machine-readable formats for archival material

The production of machine-readable data bases relating to archives requires that the data be structured in a manner that facilitates computer handling. Developments in computing are such that the data base creator may be very flexible with definitions of structure, to the extent of being idiosyncratic. In particular, the development of data base management systems (dbms) gives users the option of entering data into a machine using their own definitions for data labels. The internal structuring of data within the machine is handled by the dbms software and need not concern the person using the data.

Such flexibility is not appropriate when the exchange of data between institutions is involved. The encoding of data in machine-readable form on carriers such as magnetic tape facilitates wider dissemination. If this is to be achieved effectively, then documented standards are required for the formatting of data to allow the ready use of those data by remote users with different types of computing equipment.

Two formats being developed in this area should be examined carefully, by archivists, for their potential usefulness. They are both part of the U.S. Library of Congress' Machine Readable Cataloguing (MARC) series of formats that have been produced for the interchange of bibliographic information on magnetic tape incorporating the basic structure of American National Standard ANSI Z39.2.[18] These formats are for the description of:

a. Manuscripts[19] (published in 1973)

b. Machine-readable data files — MRDF[20] (draft made available in 1980).

The *Manuscript* format provides specifications and content designators (tags, indicators, subfield codes to identify data in the machine record) for collections and for single manuscripts. The Library of Congress is not itself distributing records of this kind on magnetic tape. It has published the format as a guide to those institutions which are involved in the creation of such computer-produced retrieval aids. The format follows the structure of all the MARC formats based on ANSI Z39.2 and its International Standards Organization equivalent ISO-2709.[21] Thus, each magnetic tape record consists of LEADER — RECORD DIRECTORY — CONTROL FIELDS — VARIABLE FIELDS. The allowable control and variable fields are shown in Figure 1. A sample manuscript record is given in Figure 2.

The *Machine-Readable Data File* format is styled along the same lines as the manuscript format. The *MRDF* is based on the *Anglo-American Cataloguing Rules* definition:

A body of information coded by methods that require the use of a machine (typically not always a computer) for processing. Examples include files stored on magnetic tape, punched cards, aperture cards, disk packs, mark-sensed cards, and optical character recognition font documents. The term . . . embraces both data stored in machine-readable form and the programs used to process that data.[22]

The data fields suggested for *MRDF* are shown in Figure 3.

## Development of individualised data base systems

Though machine-readable data need to be transferred from one place to another in standardised format, once it is to be made available for access the local rules may prevail. In most cases, access is now directly online so effective searching of fields requires the structured definition of those fields according to the constraints of the retrieval program.

Figure 1
MARC Manuscript Data Fields for Magnetic Trips

| Tag or Tag Range | Description |
|---|---|
| a) Control fields | |
| 001 | Manuscript control number |
| 002 | Proposed subrecord map of directory |
| 008 | Fixed length data elements |
| | (38 character positions in which are encoded such data as date entered onto file, country or repository, form of reproduction and processing status.) |
| b) Variable fields | |
| 010 | LC card number |
| 011 | Linking LC card number |
| 035 | Local system number |
| 041 | Languages |
| 042 | Search Code |
| 043 | Geographic area code |
| 090 | Shelf location |
| 091 | Microfilm shelf location |
| 100/130 | Main entry — personal; corporate; conference; uniform title |
| 240 | Uniform title |
| 241 | Romanized title |
| 245 | Title statement |
| 260 | Imprint |
| 300 | Physical description |
| 302 | Item Count |
| 303 | Unit Count |
| 304 | Linear Footage |
| 350 | Value |
| 490 | Series statement |
| 500 | General note |
| 506 | Restrictions |
| 520 | Scope and contents note |
| 535 | Repository |
| 540 | Literary rights |
| 541 | Provenance |
| 543 | Solicitation information |
| 545 | Biographical tracings |
| 555 | Finding aids |
| 600/690 | Subject heading — personal; corporate; conference; uniform title; topical; geographic; profession or activity |
| 700/740 | Added entry — personal; corporate; conference; uniform title; variant title |

Figure 2:   Sample manuscript record in MARC format

**Leader**

| 00988 | n | ß | c | ßß | 2 | 2 | ßß2∅5 | 1ßß4 5∅∅ | ∅∅1∅∅13∅∅∅∅ | ∅∅8∅∅4∅∅∅13 | 1∅∅∅∅32∅∅53 ℱ |

∅  　　　　　　　24　　　　　　36　　　　　　48

| 245∅∅52∅∅∅85 | 3∅∅∅∅19∅∅137 | 535∅∅3∅∅∅156 | 52∅∅296∅∅186 | 541∅∅3∅∅∅482 ℱ |

60　　　　　72　　　　　84　　　　　96　　　　　1∅8

| 65∅∅∅69∅∅551 | 6∅∅∅∅35∅∅586 | 6∅∅∅∅41∅∅627 | 6∅∅∅∅31∅∅658 | 6∅∅∅∅23∅∅681 ℱ |

12∅　　　　132　　　　144　　　　156　　　　168

Control Number | Fixed Fields

| 6∅∅∅∅31∅∅712 | 6∅∅∅∅41∅∅753ℱ | ms∅67∅∅∅26∅ | ℛ | ℱ | 72∅914 | i | 1839 | 191∅ | miu | ßß ℱ |

18∅　　　　192　　　　2∅5　　　　13

Main Entry

| ßßß | ß | ßßßßßß | ß | ß | ß | ß | ß | eng | ß | ℱ | 1∅ | $aDuveen,ßAlbert,$ecollector. ℱ |

53

Title

| ℱ | ∅∅ | $aAlbertßDuveenßautographßcollection,$f1839-1910. | ℱ ℱ |

85

Physical Description | Repository |   | Scope and Contents

| ßß | $aca.ß55∅ßitems. | ℱ | ∅∅ | $aArchivesßofßAmericanßArt. | ℱ | ßß | $aCorrespondenceß ℱ |

137　　　　　　156　　　　　　　　　　186

| andßmiscellaneousßdocuments.ßUnrelatedßlettersßbyßAmericanßartistsßofßtheß ℱ |

| 19thßcentury,ßgatheredßbyßDuveenßasßanßAmericanßartßreferenceßgroup.ß ℱ |

| PersonsßrepresentedßincludeßAlbertßBierstadt,ßF.ßE.ßChurch,ßTimothyßCole,ß ℱ |

| CharlesßR.ßLeslie,ßWilliamßPage,ßandßWorthingtonßWhittredge. ℛ |

Provenance | Subject - Topical

| ℱ | ßß | $aGiftßofßMr.ßDuveen,ß1956. | ℱ | ß5 | $aArtßandßartists$xAutographß ℱ |

482

Subject - Personal Name

| collections$y19thßandß20thßcenturies. | ℱ | 1∅ | $aBierstadt,ßAlbert,$d183∅-19∅2. ℱ |

586

Subject - Personal Name | Subject - Personal Name

| ℱ | 15 | $aLeslie,ßCharlesßR. | ℱ | 1∅ | $aPage,ßWilliam,$d1811-1885. ℱ |

681　　　　　　　　　　712

Subject - Personal Name

| ℱ | 1∅ | $aWhittredge,ßWorthington,$d1820-191∅. | ℱℛ |

753

Subject - Personal Name |   | Subject - Personal Name

| ℱ | 1∅ | $aChurch,ßFrederickßEdwin,$d1826-19∅∅. | ℱ | 1∅ | $aCole,ßTimothy,$d1852-1931. ℱ |

627　　　　　　　　　　　　　　658

ß – Blank   ℱ – Field terminator   ℛ – Record terminator

Figure 3
MARC MRDF Proposed Data Fields

| Tag or Tag Range | Description |
|---|---|
| a) Control Fields | |
| 001 | Control number (probably as assigned by individual systems) |
| 007 | Physical description data elements |
| | (7 character positions in which are encoded such data as media for storage and media for distribution.) |
| 008 | Fixed length data elements |
| | (40 characters positions in which are encoded such data as date entered on file, country of production, and language.) |
| b) Variable fields | |
| 010 | Library of Congress Number |
| 017 | Copyright Registration Number |
| 035 | Local System Number |
| 036 | MRDF Number |
| 040 | Cataloguing Source |
| 041 | Languages (Text Files) |
| 042 | Centre of Responsibility |
| 050 | Library of Congress Classification Number |
| 052 | Geographic Classification Number |
| 072 | Subject Category Codes |
| 082 | Dewey Decimal Classification Number |
| 100/130 | Main Entry — personal; corporate; conference; uniform title |
| 210 | Abbreviated Title/Acronym or Short Title |
| 214 | Augmented Title |
| 241 | Romanized Title |
| 242 | Translation of Title |
| 245 | Title Statement |
| 250 | Edition Statement |
| 260 | Publication, Production, Distribution, or Generation |
| 265 | Source of Distribution |
| 300 | File Description |
| 315 | Frequency of Serially-Issued MRDF |
| 351 | Software Programming Language and Number of Source Program Statements |
| 352 | Computer Requirements |
| 353 | Peripheral requirements |
| 354 | File Structure/Sort Sequence |
| 362 | Date and numbering for Serially-Issued MRDF |
| 400/490 | Series Statement — personal; corporate; conference; title; other |
| 500/589 | Notes — field 500 is for general notes, and there are 24 other separate fields that may be used to identify: limitations; intended audience; sponsor; associated documentation; sampling procedures, etc., if applicable |
| 600/653 | Subject Added Entry — personal; corporate; conference; uniform title; topical; geographic; reversed geographic; uncontrolled |
| 700/740 | Added Entry — personal; corporate; conference; uniform title variant title |
| 800/830 | Series Added Entry — personal; corporate; conference; uniform title |
| 850 | Holdings |

A well-known example of this in Australia is the STAIRS retrieval program that provides access to data bases on AUSINET.[23] The files constituting the AUSINET data bases are mainly bibliographic in nature. Although the STAIRS program allows the data structure in each field to be described differently, the file descriptions have been standardised. Thus, most files have fields labelled AUTHOR, TITLE, IMPRINT, DESCRIPTORS, and so on. Where appropriate, individual differences can be incorporated, such as DEGREE on the HDEG (Union list of higher degree theses) file.

An example of the use of a MARC-formatted file with STAIRS is the Australian National Bibliography (ANBB on AUSINET). This file is produced by the National Library of Australia in MARC format but the large number of allowable MARC fields is reduced to a considerably smaller number of field labels searchable with STAIRS. For example, all main and added entry author fields, of whatever type, become simply AUTHOR fields.

There would be no difficulty in building an archival finding aid on this system using a similar description of data and allowing for fields labelled PROVENANCE, SERIES DESCRIPTION, SERIES DATE RANGE, and so on.

At the University of New South Wales, we created a small test file of archival material along these lines. The file was developed for the use of students studying for the Diploma of Information Management (Archives Administration) and they are able to add to it, use it for online information retrieval, and produce a variety of permuted indexes from it.

Such procedures are possible because of the processing capabilities of a program called RIQS (Remote Information Query System).[24] RIQS, although not as flexible an information retrieval program as STAIRS or commercial retrieval programs such as DIALOG, ORBIT and ELHILL, nevertheless provides a relatively simple-to-use data base creation facility and enables the production of indexes in batch processing mode. Retrieval can be performed in either batch or online mode. Figure 4 depicts the definition of the record structure provided for the Archives file at the University and shows the first record fed into the file. Figure 5 shows how the same record may be retrieved with a subject search and displayed with control over the display format.

For this file at the University, the basic information used was descriptions of record series from the *Concise Guide* (and *Supplements*) to the holdings of the N.S.W. State Archives. Since the *Concise Guide* descriptions are not standardised in terminology or content, they had to be re-written and/or extended in order to supply the information considered necessary for the file. In addition, the lack of a thesaurus or authority list of approved subject index terms made the assignation of such terms difficult and time-consuming. (The Keyword thesaurus

Figure 4

File definition and input of data for
1 complete record followed by 1 partial record

```
/Job
archive,t100,cm100000.
user,77b1672,rts.
attach,rias/un=public.
rias.
save,file=archive.
dispose(output=pr/od=math,Jn=archive)
/eor
create archive file
record definition
(1)      organisation
(2)      agency recording
(3)      series title
(4)      date range - series
(5)      date range - contents
(6)      shelf metres
(7)      unit quantity
(8)      physical nature
(9)      series description
(10)     index terms
(11)     arrangement
(12)     related series
(13)     location
(14)     location of original
(15)     reference
(16)     restriction
(17)     end1
(18)     end2
multiple (4)(5)(8)(10) thru (15)
data restrictions
types
date (4) (5)
decimal (6)
input data
(1)nsw(2)city of sydney improvement board
(3)minute books(4)28oct1884*17feb1896
(5)28oct1884*17feb1896
(6)0.1(7)2 volumes(8)volume
(9)minutes of board meetings including transcripts of evidence
taken during the hearing of references and appeals
(10)sydney-buildings*buildings-regulation*sydney municipal council
(11)chronological(13)1/2129 - 2130
(16)none
(1)nsw(2)city of sydney improvement board
(3)register of references from city building surveyor
(4)29may1890*10feb1896
(5)29may1890*23feb1896
(6)0.04
(7)1 volume(8)volume
(9)register of references by city building surveyor to board of
buildings considered to be improperly constructed or dangerous.
information shown includes reference number,date, board meeting at
which disposed of, particulars of buildings, names of owners, details of
fee payment.
```

Figure 5

Retrieval and display of records
from the file ARCHIVE

*Note:* the first record listed is part of the first record input into the file, as shown
in Fig. 4, only some of the information in each record is printed out, and
it may be reformatted with a fair amount of flexibility.

```
 enter search command or type halt

? if #10 eq'sydney-buildinss' display  zap across #2/
? tab 10 #3 space 4 #4/
? #9/tab 20 #8/'    ---' end

 searching initiated


 no. of reports on display file =     3

 do you want the display reports listed
? y

     city of sydney improvement board
        minute books    28oct1884 * 17feb1896
     minutes of board meetings including transcripts of evidence taken
     during the hearing of references and appeals
                volume
        ---
     city of sydney improvement board
        references from city building surveyor    27aus1879 *
     10feb1896.
     city of sydney improvement act, clauses 27 and 29,provided for
     the city building surveyor to notify the board of structures
     deemed unsatisfactory odangerous.  the board could then take
     action on these references.   series comprises files containing
     correspondence, plans, surveyor's reference and board notices of
     inspection, hearings and decisions.   each file has a reference
     number.
                file * plan
        ---
     city of sydney improvement board
        reports of city building surveyor    27aus1879 * 1Jul1884
     reports on buildings etc comprising transcripts of letters from
     city building surveyor together with partial minutes of evidence
     taken at board meetings.  index: names of owners and
     streets.front of volume contains list of cases heard by board
     1879-80
                volume
        ---

*  *  *  *  *
```

developed for N.S.W. government departments by the Records Management Office[25] was not, at the time, available but could provide the basis for a valuable archival thesaurus). Such difficulties need to be confronted and resolved by everyone contemplating the automation of archival finding aids. They do not, however, obscure one of the most valuable aspects of this experiment: training in the identification of the essential as well as the less necessary components of the description of archival series and items.

Lawrence McCrank, in a report on a symposium on "Archival Automation", states that:

> archival automation involves more than mechanization or storage for inventory control . . . Special attention must be paid to: a) revising traditional descriptions and the structure of finding aids; b) creating multiple avenues for retrieval in addition to provenance; c) developing indexing systems that accommodate access by a standardized vocabulary or thesaurus; and d) creating evaluation techniques.[26]

In this paper we have explored the first two of these avenues and in so doing have dealt with some of the possibilities, as well as the difficulties, that computer-based systems present for archivists in regard to the description of records and the production of finding aids.

We have tried to show that, increasingly, archivists need to communicate with people in other information disciplines and areas of practice in order to identify common areas of concern as part of the process of solving specific archival problems, in particular, those arising from automation. As Michael Cook notes, "archivists may hope to learn from the experience of their near colleagues."[27]

There are of course many other areas where computer-based systems have implications for archival theory and practice such as appraisal, storage, conservation and reference use.[28] Although we have not been able to deal with such matters here, our general observation about the need to find out about developments in information in related disciplines would equally apply. The cost of ignoring developments in information technology and practice could well be the erosion of the archivist's legitimate areas of concern.

## FOOTNOTES

1.  For the most recent overview of the archival implications of computer-based systems, see Michael Cook, *Archives and the Computer* (London, Sydney, etc.: Butterworths, 1980). The most recent bibliography on the subject is Richard M. Kesner, comp. and ed., *Automation, Machine-Readable Records, and Archival Administration: An Annotated Bibliography* (Chicago Society of American Archivists, 1980).

2.  M. Roper, "The changing face of the file: machine-readable records and the archivist", *Archives*, Vol. 14, No. 63 (1980), 145-150; L. Bell, "The archival implications of machine-readable records", *Archivum,* 26 (1979), 85-92.

3.  G. L. Wolfendale, ed., *Data Base management systems; proceedings of the joint ANU/ACS one-day seminar held at the Computer Centre of the Australian National University, 17th November, 1976* (Canberra: Australian National University, 1977); A. F. Cardenas, *Data base management systems* (Boston: Allyn and Bacon, 1979).

4. P. A. Strassman, "The office of the future: information management of the new age", *Technology Review,* Vol. 82, No. 3 (1980), 54-65.

5. D. J. Butler and W. H. Nicholson, "ARMS — a computer based records management system developed by Tyne and Wear County Council", *Journal of the Society of Archivists,* Vol. 6, No. 4 (1979), 200-208.

6. A number of records management applications are considered in Cook, *Archives and the Computer,* Ch. 2.

7. R. H. Lytle, "Intellectual access to archives I: Provenance and content indexing methods of subject retrieval", *American Archivist,* Vol. 43, No. 1 (1980), 64-75; Cook, *Archives and the Computer,* pp. 35-40, considers the general question of the indexing of archival material.

8. H. T. Hickerson, J. Winthers and V. Beale, *SPINDEX II at Cornell University and a review of archival automation in the United States* (Ithaca, N.Y.: Cornell University Libraries, 1976).

9. General accounts of proposed and operating systems are available in Ch. 3 of Cook, *Archives and the Computer* and L. Bell and M. Roper, eds., *Proceedings of an international seminar on Automatic Data Processing in Archives* (London: H.M.S.O., 1975).

10. H. T. Hickerson, ed., *SPINDEX users conference: proceedings of a meeting held at Cornell University, Ithaca, New York, March 31 and April 1, 1978* (Ithaca, N.Y.: Cornell University Libraries, 1979).

11. Hickerson, *SPINDEX II at Cornell University;* Hickerson, *SPINDEX users* conference; S. E. Hannestad, "SPINDEX II: a computerized approach to preparing guides to archives and manuscripts" in S. Lusingnan and J. S. North, eds., *Proceedings of the 3rd international conference on computing in the humanities, Waterloo, Ontario, Canada, 2-6 August 1977* (Waterloo, Ontario: the University, 1977), pp. 273-282.

12. F. McCall, *PROSPEC Manual* (London: Public Record Office, 1974).

13. Great Britain, Public Record Office, *Nineteenth report of the Advisory Council on Public Records* (London: H.M.S.O., 1970); "Inside the New Kew Repository", *American Archivist,* Vol. 42, No. 2 (1979), 223-226; Cook, *Archives and the Computer,* pp. 76-83.

14. A. M. Newburger and P. M. Rosenburg, "Automation and access: finding aids for urban archives", *Drexel Library Quarterly,* Vol. 13, No. 4 (1977), 45-59.

15. *Ibid.,* 53.

16. M. Saclier, "Computer Applications at the Australian National University Archives of Business and Labour", *ADPA,* Vol.2, No.1 (1976), 16-18; M. Saclier, "A Progress Report on the Use of the Computer in a Small Archives", unpublished paper delivered to a workshop on *Computers: Archival Problems and Applications,* arranged by the Sydney Branch of the Australian Society of Archivists, 7 June 1980.

17. See the papers by Christopher Coggin and by Peter Scott in the forthcoming publication of papers delivered to the Third Biennial Conference of the Australian Society of Archivists, Melbourne, May 1981.

18. *American National Standard for Bibliographic Information Interchange on Magnetic Tape* (ANSI Z39.2 — 1971) (New York, ANSI, 1979).

19. U.S. Library of Congress, MARC Development Office, *Manuscripts: a MARC format; specifications for magnetic tapes containing catalog records for single manuscripts or manuscript collections* (Washington, 1973).

20. U.S. Library of Congress, Network Development Office, *Machine-readable data files: a MARC format (DRAFT — revised 7.1.80)* (Washington, 1980).

21. International Standards Organization, *Documentation — format for bibliographic information interchange on magnetic tape* (ISO-2709-1973) (Geneva: ISO, 1873).

22. *Anglo-American Cataloguing Rules,* 2nd ed. (Chicago: American Library Association, 1978), pp. 202-203.

23. *An introduction to AUSINET* (Melbourne: ACI Computer Services, 1980).

24.  B. Mittman and L. Borman, *Personalized data base system* (Los Angeles: Melville, 1975).
25.  Records Management Office of N.S.W., *Principles of Keyword Classification* (March 1978) and *Manual of Keyword Classification (October 1978).*
26.  L. J. McCrank, "Archival automation: future access to the past", *Bulletin of the American Society for Information Science,* Vol. 6, No. 6 (1980), 30-31.
27.  Cook, *Archives and the Computer,* p. 9.
28.  See C. M. Dollar, "Appraising machine-readable records", *American Archivist,* Vol. 41 (1978), 423-430 and Cook, *Archives and the Computer,* Ch. 4.