# PERSONAL BIBLIOGRAPHIC DATABASES: AN INDUSTRIAL SCIENTIST'S PERSPECTIVE

by Thomas E. Wolff
AMOCO Corporation

Much has been written recently [1,2] on the use of bibliographic software for managing references and personal files. Typically, these programs have been used for preparation of bibliographies and footnotes for publications. Some programs are specialized reference managers, working with word processor documents to create ordered bibliographies by extracting reference indicator text from the documents. However, many industrial scientists' jobs involve much information streaming in, but little need or opportunity to prepare elaborately documented papers, as for technical journals. In other words, these scientists need assistance in information management, not document output management. This is the perspective that we take in evaluating and using personal bibliographic software.

> ...these scientists need assistance in information management, not document output management. This is the perspective that we take in evaluating and using personal bibliographic software.

## INFORMATION SOURCES

The sources of information are varied and becoming numerous. Documentation may be available in hardcopy format or accessible electronically or both. Typical hardcopy sources are journals and internal company documentation, although online options are often available but unused. ("Hardcopy is easier to read on the train or plane.") Compilations are frequently distributed as hardcopy bulletins. Journal titles are compiled in *Current Contents*, or abstract collections are provided on preselected subjects or customized interest areas by organizations, such as the Chemical Abstracts Service and the American Petroleum Institute. This information is generally available online, as on Dialog Information Services or STN, and is being supplied in many cases on personal computer diskette (Current Contents on diskette) or CD-ROM (Beilstein Current Facts in Chemistry) as well. Some scientists access these electronic sources themselves and are referred to as "end-users" by information specialists. On the other hand, many depend on these specialists to get the

## FIGURE 1
### A "TYPICAL" DERWENT WORLD PATENT INDEX CITATION AS DOWNLOADED IN DIALOG TAGGED FORMAT

```
 8/4/8
AX- 87-110218/16|
AX- <XRAM> C88-088483|
TI- Rhodium-carbon catalyst prepn. for use in terephthalic acid purification|
PA- (STAD )_AMOCO CORP|
AU- <INVENTORS> SCHROEDER H;  WITTMAN R L|
PN- <BASIC> EP 219288 _A_870422_8716|
PN- <EQUIVALENTS> JP 62121647_A_870602_8727; US 4728630_A_880301_8812;  CN
    8606590_A_870422_8828; EP 219288 _B_890906_8936;  DE
    3665416_G_891012_8942; ES 2011014_B_891216_9007|
AN- <PRIORITIES> US 905758 (860909); US 785055 (851007)|
AN- <APPLICATIONS> JP 86238986 (861007); EP 86307678 (861003); EP 85785055
    (851007)|
LA- English|
DS- <REGIONAL> AT;  BE;  CH;  DE;  ES;  FR;  GB;  IT;  LI;  LU;  NL;  SE|
AB- <BASIC> EP0219288
         Prepn. of a Rh-carbon catalyst (III) is effected by contacting a
    porous carbonaceous material (I), having a surface area of at least about
    600 m2/g and pH 9-11 in aq. suspension, with a pH 1-4 aq. Rh3+ salt (V)
    soln.  (II) such that (I) contains about 0.01-2 (pref. 0.5) wt.% Rh pref.
    to a depth of at least 5 microns (more pref. 10-20 micron), and pH (I)
    plus (II) is about 12-13.5.
         The pref. pH of the impregnated wet catalyst compsn. is about 6-8.
    The wt. ratio of solids to liquids is at least about 0.5 (pref. 0.5-20,
    more pref. about 0.5-2.  The contacting of (I) with (II) is by either
    soaking or spraying.
         USE/ADVANTAGE - (III) is used in the purification of terephthalic
    acid (IV) which is the starting material for prepn. of polyester fibres,
    films and resins.  (III) has enhanced selectivity for decarboxylation of
    4-carboxybenzaldehyde (the principal impurity in the prepn. of (IV)) and
    removes colour and fluorescence. (Previously notified in Week 8716) @(12pp
    Dwg.No.0/0)@|
AB- <US> 8812 US4728630
         Rhodium-on-carbon compsn. is made using an alkaline granular porous
    carbonaceous catalyst support material of surface area 600 sq. m. per g.
    or more and pH 9-11.
         Process comprises contacting support material with an aq. Rh (3+)
    salt soln. of pH 1-4 for a time to form a wetted support material contg.
    0.01-2 wt.% Rh (as metal w.r.t. dry catalyst). Support material and salt
    soln. are preselected so that sum of pH in each is 12-13.5. Pref. Rh-salt
    is rhodium trinitrate or -trichloride.
         USE - For purification of terephthalic acid under reducing
    conditions. @(7pp)@|
```

information for them. Once obtained, this information is transferred to the requesting scientist in either hardcopy or electronic document format.

## FACILITATING IDEA GENERATION

The key feature of most available information is that its delivery format generally hinders its retention and retrieval. Large reports and extensive database search output often are filed away immediately after being read (unless sent to the "circular" file). In fact, much of the value of the online database is lost when flat report documents are generated. However, bibliographic software can enhance the value of search reports, transforming them into "living" databases.

Other hardcopy documents such as company reports, journal articles and patents can be added to these databases by manual entry or by using document scanners and optical character recognition. Those skilled individuals who argue that they remember all they read miss the point. At any given time, one obtains information through the mental filters relating to one's current interests and concerns. As progress is made, one's perspective changes. Having the "old" information remain accessible offers new opportunities to learn from it. Having information from several sources in one place for easy browsing leads to the generation of new ideas. That is the principal value of creating personal bibliographic databases.

## ILLUSTRATIVE DATABASE EXAMPLES

Both the value of these databases and the complexity in creating them are best illustrated by example. Three typical cases are proposed below. The impetus for developing each, as well as the types and sources of information, are described. With each one, specific difficulties will arise and some solutions, not always optimal, may have to be implemented. These will be discussed in the rest of this article.

### CASE 1: Japanese polyolefin technology.

In the competitive market of polyolefin manufacture, many Japanese companies have strong technical positions. A possibility exists that in licensing desired technology, the licensor may also make available to us additional related technology disclosed by them during a specific multi-year period. The question arises about what these companies may have to offer us. The technology is described in hundreds of Japanese, U.S. and European patents, many of which are essentially equivalent. A database may be developed with patent citations from the following files, among others: Chemical Abstracts (CA) on STN; the Derwent World Patent Index (WPIL) and IFI/Plenum U.S. CLAIMS files on DIALOG; Patent Abstracts of Japan (JAPIO) on ORBIT Search Service.

### CASE 2: A polymer for a new electronic substrate application.

A staff team has found a new preparation and a new use for a well-known polymer in the electronic substrate area. The polymer is not currently available commercially. As their work proceeds toward market introduction, many concerns exist about the preparation of the monomers and polymer and about competing polymers for the same application. They want thorough retrospective and current awareness searches collected for all to study. Most of the files of Case 1 would be searched, as well as INSPEC and COMPENDEX engineering files on STN or DIALOG, and the RAPRA rubber and plastics file on ORBIT.

### CASE 3: Preparation and oxidation of polyalkylaromatic compounds.

For this relatively mature technology, a comprehensive database would be

developed from an alternative current awareness program. Numerous customized, biweekly Chemical Abstracts ISS profiles are distributed among the various research groups. In addition, American Petroleum Institute (API) hardcopy bulletins covering relevant literature and patents are circulated widely. These multiple sources could be supplemented and possibly replaced by a report edited by the Information Services group. A database program would be used to collect, edit, classify and format citations retrieved electronically for publication to all working in this research area. Initially, sources would be the online versions of the CA and API files, including the early access file CApreviews on STN.

## DIFFICULTIES IN DATABASE GENERATION AND USE

Each of these cases involves multiple online files, generally from two or more online systems. References are either from the "open" literature or patents. This variety and the multitude of database formats and conventions lead to considerable difficulties in importing the information into fully functional personal bibliographic databases. In developing these types of databases, we have seen problems not only in data importing, but also database structure definition, information access and handling (i.e., searching, browsing, merging and duplicate detection), and report generation.

---

**This variety and the multitude of database formats and conventions lead to considerable difficulties in importing the information into fully functional personal bibliographic databases.**

---

These areas will be discussed in relation to our experience with four bibliographic software packages: Library Master (version 1.24),

## TABLE 1

| REQUIREMENTS FOR BIBLIOGRAPHIC DATABASE PROGRAMS -- 1 | Library Master | Notebook II | Papyrus | Pro-Cite |
|---|---|---|---|---|
| *Database File Features* | | | | |
| Large number of records per database | ★ | ★ | ■ | ★ |
| Multiple record types (e.g., books, patents, etc.) | ★ | O | ■ | ★ |
| Variable and "unlimited" length fields | ★ | ■ | ■ | ★ |
| Large number of fields per record type | ■ | ■ | O | ■ |
| Field characteristics fully user-definable | ★ | ■ | O | ■ |
| Multiple field types (e.g., date, name, numerical, text) | ★ | O | ■ | ■ |
| Flexible date formats, incl. ISO (YYMMDD) | ■ | O | ■ | ■ |
| User assignable indexes for any fields | ★ | O | O | ★ |
| Compatible with patent citations | ■ | ■ | O | ■ |
| *Importing of Downloaded Data* | | | | |
| Easy to set-up, customize and use import | ■ | ■ | O | ■ |
| Multiple online systems (STN, Dialog, Orbit) | ■ | ■ | ■ | ★ |
| Multiple files | ■ | ■ | O | ★ |
| Chemical Abstracts | ■ | ■ | ■ | ★ |
| Derwent World Patent Index, IFI U.S. Claims™ | ■ | ■ | O | ■ |
| Automatic importing to proper record type | ★ | O | ■ | ★ |
| Field parsing -- standard bibliographic info. | O | O | ■ | ■ |
| Field parsing -- patent information | O | O | O | O |
| Standard input format templates available | ■ | ■ | O | ★ |
| Handles lists or tables in fields (e.g., patent info.) | ★ | O | O | O |
| Handles multiple-paragraph fields (e.g., abstract) | ★ | O | O | O |

★ Good   ■ Satisfactory   O Needs Improvement

Notebook II (version 4.05), Pro-Cite (version 1.41) with Biblio-Links, and Papyrus (version 6.0.6) [3]. Our work was carried out on an IBM-compatible PC, although Papyrus was also tested on a VAX. An updated version of Pro-Cite was recently released for the Macintosh, and the PC update is anticipated in early 1992. These four programs were evaluated *particularly with respect to our bibliographic database program requirements.* How each of these programs satisfies these requirements is indicated graphically in the accompanying tables. Detailed reviews of these programs have appeared recently [1, 4-7].

### Data Importation

Much of the data importing difficulties result from document type. Three of the four software packages allow for multiple record types within a database, but two require that the document type designator (DT in the CA file on STN) be at the beginning of the record. The restriction of Notebook II, of providing only one record type, was sufficient to eliminate it from much serious consideration. The principal complication involves dealing with *patents.* Only Papyrus includes patent record type as a default, but it is extremely limited, containing only nine information

## TABLE 2

| REQUIREMENTS FOR BIBLIOGRAPHIC DATABASE PROGRAMS -- 2 | Library Master | Notebook II | Papyrus | Pro-Cite |
|---|---|---|---|---|
| **Searching** | | | | |
| Full Boolean search logic with parentheses | ★ | O | ■ | ★ |
| All fields searchable, singly or globally | ★ | ■ | O | ★ |
| Lookup tables (indexes) for searching fields | ★ | O | ■ | ★ |
| Range searching of date and numeric fields | ★ | O | O | ■ |
| Left and right term truncation | ★ | ■ | O | ★ |
| Rapid searching | ★ | O | O | ■ |
| Proximity searching | O | O | O | O |
| Save and rerun search strategies | ■ | O | ■ | O |
| Hypertext search capability | O | O | O | O |
| **Browsing** | | | | |
| Convenient scanning of "hit" citation titles or other fields | ■ | ■ | ■ | ■ |
| Smooth browsing, including from abstract to abstract | ★ | ■ | ■ | O |
| **Editing** | | | | |
| Straightforward full screen editing | ★ | ■ | O | ★ |
| Selective locate and change | ★ | ■ | O | ★ |
| Global locate and change | O | O | O | ★ |
| Block move, copy and delete | ★ | ■ | O | ★ |

★ Good  ■ Satisfactory  O Needs Improvement

---

fields and lacking the capability to search on patent assignee. Pro-Cite provides a reasonable patent record format in a supplemental database available upon request. Library Master, having been written originally for humanities scholars, contains no patent record type, but its design makes formation of a useful patent record type straightforward.

Other problems with data importing arise from information field structure in the online databases. For DIALOG files output in tagged field format, often several fields share the two-digit field indicator, with sub-fields delineated with descriptors in less than and greater than signs. For example, in WPI records, basic and equivalent patent numbers have these field codes, respectively: **PN- <Basic>** and **PN- <Equivalents>**. Most database software requires field tags to be contained within the first several columns as defined in the importing template. Another problem is information fields that must be parsed to separate multiple types of information. The source (**SO**) field in the CA file on STN contains all the bibliographic information except date. Most patent citations combine patent number and date, or application and date within the same field. The parsing algorithms vary

among the evaluated database programs. Most of the programs can handle regular bibliographic information in non-patent citations (Library Master is still developing this functionality), but none handle patent information properly.

**Patents**

Patent citations lead to other difficulties. A typical downloaded patent citation is shown in the Figure. The most frustrating is the matter of date format. Many patent files indicate date in ISO format: YYMMDD, as 911021 for October 21, 1991. None of the database programs will understand these numbers as dates even though most handle many date formats interchangeably. Ironically, Notebook II is unsophisticated and requires that all dates be entered in a format that can be sorted as numeric data, so ISO date format would be "understood" as date input, while "normal" date input is considered text. Two additional styles of information are generally not respected by these database programs: tabular or list fields, as for equivalent patents in WPI output; or multiple-paragraph fields, as in WPI patent abstracts. Pro-Cite will word-wrap all entries, making tables or multiple paragraph fields into single, long, indecipherable paragraphs. Library Master is well designed to handle these field styles, although I have hit a limitation while trying to import a chemical structural formula from the Registry file on STN; in all cases the structure is word-wrapped. Finally, long index term fields, as in the CA file on STN, are not treated as they are in the original online database. Generally, the index terms are combined in these software programs as one field, searchable as words only. Library Master can keep multiple index term entries separate, within limits, but then the search capability to "link" terms within one index term entry is not available.

**Data Handling Within The Database**

Once the data is imported into the personal database, limitations arise in handling that information. Since browsing records and generating ideas is an important use of these databases, browsing records should

be a free-flowing process. Generally, browsing sequential records is no problem, although Pro-Cite requires multiple keystrokes to move from one record to the next and is especially cumbersome browsing abstracts in sequential records. Library Master moves smoothly between citations, and keystroke macros can be written to move automatically from abstract to abstract with one keystroke. Several of these programs try to alleviate this problem by tabulating key citation information for easing access to records. However, ideas are not usually generated from citation information. An improvement for these programs would be hypertext movement around the records, although automatic generation of useful hypertext links is not straightforward.

> **Most of the programs can handle regular bibliographic information in non-patent citations...but none handle patent information properly.**

The most useful record sets to browse are those created by searching the database on relevant criteria. Most of the evaluated software have full Boolean search logic but no proximity or "linking" operators, which would be useful in connecting CA Registry Numbers and descriptive text. Date searching is sophisticated in both Library Master and Pro-Cite, so the ability to search for patent priority information, for example, is only limited by the form of the imported information. In other words, useful information could be obtained from databases for which earliest equivalent patent dates have been created.

Duplicate detection and record merging are two additional functions of concern. Non-patent citations may come from multiple sources, such as INSPEC and COMPENDEX, and duplicate detection can generally be accomplished using the title and author fields. On the other hand, titles in the WPI file are copyrighted Derwent titles, and the patents in the CA and CLAIMS may be "equivalent" but not identical, since CA abstracts only the basic patent and CLAIMS the U.S. patent. Duplicates cannot even be detected on patent number, both because different "equivalent" patents may be cited and because patent number formats differ among online systems. The best means of duplicate detection would use earliest priority application number and date, fields that have to be created by manipulating the imported information. In this case, however, duplicates might better be merged rather than eliminated, since citations from different databases often contain different information: original patent abstract and claims are found in the CLAIMS file; copyrighted abstracts are given in the WPI and CA files. It would be particularly useful to be able to do partial record merges, such as CA abstract and U.S. claim into a WPI record. Scientists would especially want all available information on Japanese patents, which are very expensive to have translated. At this point, no database software allows such partial merging.

**TABLE 3**

| REQUIREMENTS FOR BIBLIOGRAPHIC DATABASE PROGRAMS -- 3 | Library Master | Notebook II | Papyrus | Pro-Cite |
|---|---|---|---|---|
| **Sorting and Search Set Handling** | | | | |
| Text, numeric and date sorting by field values | ★ | ■ | ■ | ★ |
| Ease of set building, subsetting and merging | ■ | ■ | O | ★ |
| Duplicate detection - user definable | O | ■ | O | ■ |
| **Reporting Output** | | | | |
| Output to screen, printer or file | ★ | ■ | ★ | ★ |
| Export in field-delimited format (e.g., STN-style) | ★ | ★ | ■ | O |
| Easy to modify output formatting | ■ | ■ | ■ | O |
| Standard output style templates available | ■ | ★ | ■ | ★ |
| **OTHER USEFUL FEATURES OF BIBLIOGRAPHIC DATABASE PROGRAMS** | | | | |
| Customizability (e.g., colors, keystroke macros) | ★ | O | ■ | O |
| Handles/uses journal or author abbreviations | O | ■ | ★ | ■ |
| Security from alteration or read-only versions | ★ | O | ★ | ★ |
| Duplicate-equivalence detection for patents | O | O | O | O |
| Partial merge of "duplicate" citations from different sources -- e.g., Derwent abstract with CA indexing and IFI claim) | O | O | O | O |
| Text highlighting - bold, underline, italics, case, sub/superscript | ★ | O | ★ | O |
| MS Windows program or PIF supplied | ■ | O | O | O |

★ Good ■ Satisfactory O Needs Improvement

## Report Output

As indicated earlier, many of these software packages excel at bibliography output. However, reports with significant text fields are not always handled well. The problems of word-wrapping lists and multiple paragraph fields generally apply on output too. A particular frustration with Pro-Cite is its inability to output citations in the same field-delimited format in which it was imported. Many of our users are used to working with STN-style citations and they want their reports to look the same. Library Master has few limitations on output formats. For Case 3 above, the report would be output with text highlighting in WordPerfect format to be read directly into a WordPerfect document for publishing.

> **A solution to many of these problems is to edit the import information with a text editor.**

## ONE SOLUTION: TEXT EDITING

A solution to many of these problems is to edit the import information with a text editor. We use KEDIT (from Mansfield Software Group, Inc.), which uses powerful text macro language REXX and its subset KEXX [8]. Routines have been written for data from several source files and for use with the database programs evaluated. Some conversions these KEXX macros carry out are summarized below:

• Move the document type (**DT**) field to top of each CA (STN) record.

• Translate ISO Dates (YYMMDD) to STN format (DD MMM YY), e.g.:
  910415 ——> 15 Apr 1991

• Edit field descriptors and create unique fields names, e.g., (Derwent on DIALOG):
  **PN- <Basic> ——> BA-**
  **PN- <Equivalents> ——> EQ-**

• Parse original patent field(s) into patent number and date fields:
  Derwent and IFI — parse the **PN** fields
  CA — rearrange **PI** and **SO** fields

• Translate application numbers in CA (STN) to Derwent (DIALOG) format:
  STN standard — CC YY-NNNNNN
  Derwent — CC YYNNNNNNN

• Change Japanese examined patent specification numbers from CA file from "Emperor year" to "Western year" for consistency with Derwent practice, e.g.:
  JP "Emperor" — JP 03052957 (Heisei)
  JP "Western" — JP 91052957

• Prepare date-sorted lists of equivalents, applications (1 entry/line, Derwent in tagged DIALOG format)

→

• Determine the earliest application priority and make separate earliest application number and date fields:

    Derwent and IFI — select from **PR** field entries

    CA — select from **AI** and/or **PRAI** fields

• Delete, add, change or move separator characters, e.g., virgule, double slashes in author, index term fields.

## CONCLUSIONS

Database software programs have generally targeted individuals involved in article writing that requires substantial footnote and bibliography generation. It appears that patent users and even the STN community are not significant constituencies of most bibliographic database program developers. Such bibliographic programs are ill-designed for patent citations. Pro-Cite is an exception because its market is "everybody." A specific STN importing program (Biblio-Links) and patent record type are available. Pro-Cite is strong in data importing and bibliography creation but its value in handling patent citations is limited by some global database features, which inhibit development of a well-designed patent workform containing all desired field types and characteristics. Library Master is a worthy text management tool with flexibility for developing most importing and bibliography-creating routines. Although the developers undoubtedly never considered patents as a potential record type, Library Master handles patents admirably. Both Papyrus and Notebook II are principally bibliographic programs with limited text management function to satisfy the needs of the industrial scientist as described in this article.

The browsing capability of most programs needs to be enhanced to facilitate idea generation. Some implementation of hypertext, with useful, automatic linking generation, should prove valuable.

A powerful text editor or word processor with "macro" language is very important in most cases for optimizing the structure of the personal bibliographic database for information retrieval. KEDIT with its macro capability has proven very effective.

Development of local area network (LAN) versions of bibliographic software will become critical for work groups to share information effectively. Library Master has recently released a LAN-version. Pro-Cite for LAN operation is in development now, according to the producer.

Copyright is an important issue to be considered in developing personal databases or published compilations of citations. Most database publishers require licenses for personal database use and for multiple copy distribution of copyrighted material.

Users must be aware that the personal bibliographic database is not a substitute for the online files when comprehensive searching is required. The searching functions of online databases are far more powerful, as they may offer additional search operators, such as linking or proximity, or database coding beyond typical index terms. On the other hand, the personal database can be a dynamic information file and a powerful browsing and idea generation tool.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] Lundeen, Gerald W. "Bibliographic Software Update." *DATABASE* 14, No. 6 (Dec. 1991): pp. 57-67.

[2] Stigleman, Sue. "Bibliography Formatting Software: A Buying Guide." *DATABASE* 15, No. 1 (Feb. 1992): pp. 15-27.

[3] Wolff, Thomas E. "Personal Bibliographic Databases: Adding Value to Search Reports." Poster session presented at the 201st American Chemical Society National Meeting, Division of Chemical Information, April 15, 1991, Atlanta.

[4] Lundeen, Gerald W. "Software for Managing Personal Files." *DATABASE* 13, No. 3 (June 1989): pp. 36-48.

[5] Kebbekus, Barbara B. "Papyrus. Version 6.0." *Journal of the American Chemical Society* 112, No. 10 (1990): p. 4092.

[6] Raeder, Aggi. "Library Master for Databases and Bibliographies." *DATABASE* 14, No. 5 (Oct. 1991): pp. 67-72.

[7] Wolff, Thomas E. "Library Master. Version 1.24." *Journal of the American Chemical Society* 114, No. 2 (1992): pp. 796-797.

[8] Wolff, Thomas E. "KEDIT: Text Editor for Post-processing Searches." *DATABASE* 15, No. 3 (June 1992): in press.

# THE AUTHOR

**THOMAS E. WOLFF** is a Staff Research Scientist in the Amoco Corporation Information Research and Analysis Group. He received his B.S. (Chemistry) from the Massachusetts Institute of Technology in 1974 and Ph.D. (Inorganic Chemistry) from Stanford University in 1980. After graduation, he was a Research Chemist in Amoco Chemical Company for ten years in the areas of polyolefin catalysis and oxidation of alkylaromatics to polycarboxylic acids. Two years ago he made the auspicious move to the Amoco Corporation information area where he provides literature, patent and business information services, principally to former colleagues in Amoco Chemical.

Communications to the author should be addressed to Thomas E. Wolff, Amoco Corporation, Mail Station F-1, P.O. Box 3011, Naperville, IL 60566-7011; 708/420-4662; Internet—mhs!amoco!thomas_e_wolff@attmail.com; STNmail—1749C; DIALMAIL—27246.