
Using New Technologies to Provide Easy Access to Research Databases

by Andy Covell¹

Manager, Research Data Center
Syracuse University

The Research Data Center, a unit within Syracuse University's central computer services organization, provides fee-based programming and data management services for researchers engaged in data-intensive research. Research Data Center analysts routinely develop strategies for managing, processing, and analyzing research databases. Mainframe based access to magnetic tape and disk is the prevailing strategy for providing access to large research databases. With recent technological developments a number of alternatives can now be realistically considered, and the Research Data Center is investigating these alternatives. This paper summarizes the information obtained over the last 6 months of that investigation.

INTRODUCTION

Most research databases are created through the collection, entry and organization of data specifically for research (e.g. survey research) or are derived from data collected outside the research enterprise for reasons other than academic research (e.g. government databases).

There are three basic types of research databases:

- *Raw files* are electronically readable files which are not formatted for any particular software package so they cannot be analyzed directly. Raw files are accessed very infrequently, and are usually stored off-line once they are read into a master file.
- *Master files* are data files which have been formatted for some software package (e.g. SAS) so they can be easily accessed for direct analysis or to obtain extracts. They are usually static, and they are typically accessed on a regular basis by some community of researchers over an extended period. Storing master files off-line is common, although on-line access is clearly preferable.
- *Analysis files* are data files created for a specific research analysis. They are formatted for some software package and contain the derived variables needed to answer a specific research question. Analysis files, the working data sets of a research project, are created, modified, and analyzed with great frequency during analysis. They are rarely accessed once the research is completed. Researchers always prefer on-line access to the analysis files which

support current research, while the analysis files of previous research can be stored off-line.

Access to large research databases (raw files, master files, and even analysis files) is often limited - commonly access is through mainframe attached magnetic tape drives. Recent technological developments can significantly enhance access to large databases. High speed networks, an ever-increasing array of on-line and near-line storage alternatives, and user friendly, flexible database software are rapidly evolving technologies that can provide fast easy access to large research databases.

NETWORKS

High speed networks offer fast access to data stored on remote computers. With a few simple commands researchers can retrieve data from a remote computer, even if it is from a different model computer across the country. Research databases no longer need to be written out to magnetic tape to be transferred from one computer to another.

Over the past five years high speed networks have become a fact of life at US universities. Most have spent considerable sums on high speed campus networks, and millions have been spent on regional networks which link campus nets together and a national backbone which links the regionals. One result is the Internet, a large conglomeration of interconnected campus, regional and national networks.

Basic, consistently implemented network services—remote login, file transfer, and electronic mail—are available on all Internet computers (except some desktop computers which may lack electronic mail). Researchers use the same procedure to access a file whether the file is on a computer located in another building on campus or on an Internet computer located in another part of the country.

File transfer is provided by a service called FTP (which stands for File Transfer Protocol). FTP, which is a basic service available on all Internet computers, provides broad access to network data resources. Unfortunately, FTP sometimes uses network capacity unnecessarily because FTP always transmits a complete file even if only a small extract is needed.

Another network service of interest to data-intensive researchers is Network File System (NFS) which enables transparent remote file access. With NFS, researchers can access a remote file or directory as if it were on his or her local computer. Thus, NFS overcomes the major shortcoming of FTP (NFS does not transmit an entire file). Unfortunately, NFS is not yet as widely available as FTP.

While FTP and NFS enable network access to data stored on almost all Internet computers, they do not automatically convert binary coded numeric information between different computer models. Thus, FTP does not support totally transparent data sharing. Some database software vendors have solved this problem (see section IV).

FTP is a nearly universal tool for providing network access to data resources. But do the Internet and local campus networks really work fast enough to support network access to large research databases? Or does network traffic and the existence of "weak links" (e.g. low-speed regional network connections that become a data transmission bottleneck) reduce performance to the point that transmitting large research databases is infeasible?

To get a realistic assessment of network performance, an informal test was carried out with the help of Jim Jacobs of the University of California at San Diego. A file containing roughly 10 Megabytes (an extract of the

General Social Survey) was repeatedly transferred, in roughly two hour intervals, through the national Internet (from San Diego to Syracuse) and through the Syracuse University Internet (between workstations in separate buildings) on Wednesday, May 9, 1990. The data transfer rate was recorded each time the file was transferred.

The results of the test, summarized in Table 1, indicate that campus networks which run at ethernet speeds, such as the Syracuse University Internet, are indeed suitable for large data file transmission, while long distance transmission via Internet is limited — for larger files magnetic tape is still an attractive method. Of course within a couple of years the Internet's national backbone network and many of the regional networks will be seeing a 24-fold increase in performance. The day may soon come when magnetic tapes are used to transfer data only in rare instances.

MASS STORAGE

One of the most striking developments over the past few years in mass storage technology is the proliferation of high capacity storage products. Not too many years ago, if one had a large research database there were only two realistic storage alternatives: it could go on mainframe mag tape or, with a little help from the mainframe systems folks, it might be put up on mainframe magnetic disk. Today there are a slew of other alternatives suitable for a range of large research database applications —

TABLE 1. Results of Informal Internet Data Transfer Test

<u>(EDT)</u>	<u>Transfer Rate</u>	<u>ElapsedTime</u>	<u>Transfer Rate</u>	<u>Elapsed Time</u>
6:33AM	12 kbytes/sec.	14 min.	76 kbytes/sec.	2 min.11 sec.
9:18AM	11 kbytes/sec.	15 min.	74 kbytes/sec.	2 min.15 sec.
11:27AM	10 kbytes/sec.	17 min.	77 kbytes/sec.	2 min. 9 sec.
1:21PM	8 kbytes/sec.	21 min.	71 kbytes/sec.	2 min.20 sec.
3:09PM	8.6 kbytes/sec.	19 min.	76 kbytes/sec.	2 min.11 sec.
5:03PM	8.8 kbytes/sec.	19 min.	76 kbytes/sec.	2 min.11 sec.
7:47PM	8.6 kbytes/sec.	19 min.	71 kbytes/sec.	2 min.20 sec.
9:08PM	7.8 kbytes/sec.	21 min.	75 kbytes/sec.	2 min.13 sec.
11:24PM	11 kbytes/sec.	15 min.	71 kbytes/sec.	2 min.20 sec.

from 700 megabyte hard disks which cost a few thousand dollars and can attach directly to a PC or workstation to mainframe attached, terabyte capacity (i.e. a thousand gigabytes) optical jukeboxes. See the Appendix for descriptions of a wide variety of storage products.

With networks enabling high speed access to a heterogeneous mix of computers, one is no longer constrained to a particular computer platform when evaluating mass storage alternatives. Other characteristics of the application — e.g. required capacity, expected frequency of access, and life expectancy of the data — can be matched against the storage technologies available for a variety of platforms.

While there are hundreds of options available for storing large amounts of research data, almost all storage products store data on one of three basic media: magnetic tape, magnetic disk, or optical disk.

MAGNETIC TAPE

Magnetic tape is the storage technology upon which many research data libraries were built, and it remains a dominant research data storage medium on many campuses.

Universities made the rather substantial investment to provide tape access some time ago, so researchers have ready access to central tape storage facilities and tape drives. Magnetic tapes remain an attractive medium for storing large databases because the main cost is that of new tapes, which are relatively inexpensive. A standard reel of 9-track tape, which holds up to 180 megabytes of data, can be purchased for around \$15. IBM mainframe 3480 cartridge tapes, which hold slightly more, cost under \$10.

Tapes are often used to transport databases between institutions. Tapes can be packed and shipped overnight, and standard tape formats exist which can be read and written on every major university campus in this country.

One of the major problems with magnetic tape is slow data access — the operator intervention required to mount a tape and the sequential processing of magnetic tape results in access time measured in minutes.

Another problem with magnetic tape is limited archival life. Tape is a relatively fragile storage media, with an average archival life of somewhere around five years. Those charged with maintaining access to data on tape for extended periods must periodically “exercise” each tape to ensure readability and prevent print-through.

Compact, high capacity tape products commonly used to back up workstation and minicomputer hard disks have

potential as storage media for research databases. 8mm tape store up to 2.3 gigabytes of data in compact cartridge, and an 8mm Exabyte drive (Exabyte is the only 8mm drive manufacturer although several companies sell Exabyte drives) runs \$4,000. 4mm tapes, also known as Digital Audio tapes (DAT), are compact cartridges (smaller than a pack of cigarettes) which store 1.3 gigabytes of data. Although 4mm tape drives do not have the installed base of the 8mm drives, buyers are attracted to DAT because drives are made by more than one vendor. Significant increases in the capacity of both 4mm and 8mm tapes are expected.

Magnetic Disk

The basics of magnetic disk technology have not changed significantly in twenty years, but continual improvements have resulted in steady increases in capacity and performance and a steady decrease in cost per megabyte. Magnetic disks are the obvious choice if high performance on-line access to research databases is required.

Hard disk are now available which attach directly to PCs or workstations and hold around 700 megabytes of data. They can be purchased for as low as \$2,500. For those with greater appetite for local storage, several drives can be daisy-chained together to provide access to several gigabytes of on-line storage.

Network servers are computers that are dedicated to the task of data access for network client computers. They typically have attached disks that are faster than those attached to individual desktop computers. Network server performance is likely to be boosted in the near future with the introduction of RAID (Redundant Array of Inexpensive Disks) technology. RAID servers will achieve much faster transfer rates by transmitting data in parallel, using multiple disks and read/write heads.

Mainframe disk drives, also known as Direct Access Storage Devices (DASD), currently offer the best overall performance. Mainframe DASD, which can provide on-line access to hundreds gigabytes to hundreds of mainframe users, are typically used for demanding time-sharing and transaction processing applications.

OPTICAL DISK

There are three basic optical technologies on the market today: CD-ROM which is primarily a publishing media with data disks created and distributed by information providers, WORM disks which enable a single write followed by unlimited read access, and erasable magneto-optical disks which allow unlimited read/write access. Digital paper is a newer ultra-high density optical technology which is just becoming available in commercial products.

CD-ROM, WORM, and magneto-optical are very dense optical disk storage technologies, with the disks typically available as removable cartridges or platters. Optical disks are slower than magnetic disks, but the drives are not as susceptible to head crashes and other malfunctions. Optical disks also have a lengthy archival life with some vendors claiming up to 30 years.

CD-ROM disks, developed originally for the audio industry, are 4.7" disks that hold roughly 600 megabytes of data. Data is formatted and written on a CD-ROM master disk (a process called mastering) which acts as a template for "stamping" copies. CD-ROM disks can be mastered for one to two thousand dollars; disks stamped from the master run \$2 to \$3 per disk.

CD-ROM is popular medium for distributing textual and bibliographic databases and is beginning to catch on as a medium for distributing research data files, a development boosted by the Census Bureau's decision to distribute much of the 1990 data on CD-ROM. Its main attraction is on-line access to large databases from desktop workstations.

In many instances, CD-ROM is a good alternative to dial-up access to expensive information services such as Dialog. However, the general suitability of CD-ROM for research database access is questionable. CD-ROM is relatively slow when compared to other on-line storage media (e.g. hard disks and WORM disks) so it is far from ideal for supporting multi-user on-line access to research databases (though several CD-ROM network server packages are on the market). Furthermore, the computers which many CD-ROM distributors target, often lack the computing resources to effectively handle the analysis files that are typically derived from the large master files distributed on CD-ROM.

WORM (Write Once Read Many) is a high capacity, locally written storage media with a lengthy archival life. WORM is faster than CD-ROM, although WORM drives are not nearly as fast as most magnetic drives.

One of the striking things about WORM technology is the wide range of WORM products. Unlike CD-ROM, WORM is available in several sizes and configurations — from 5-1/4" disks which hold hundreds of megabytes to 14" platters with an 8.2 gigabyte capacity. Most WORM drives act like magnetic disks, although main-frame attached drives typically emulate a tape drive. WORM is available as a single drive removable cartridge drive in some products and in jukebox configurations in others. WORM products are available for the complete range of computer platforms — from PCs to mainframes.

Magneto-optical disk, a recently introduced optical

technology, is a high capacity, fully erasable, removable storage media. It shares many of the properties of WORM drives (similar capacity to store data, comparable performance, long archival life), but it is not available in such a wide range of products — 5-1/4" disks which can store several hundred megabytes are the norm.

Digital paper is an ultra-high write-once optical media which can be produced in large sheets and reels. Only one product using digital paper is on the market right now (the Creo 1003 tape drive which stores a terabyte of data on a single reel of tape); and one product that was being planned has been dropped (a Bernoulli drive based on optical paper). The future of digital paper is unclear, but if the technology takes off, it could become the storage media of the future.

DATABASE SOFTWARE

Database software packages enhance database access by relieving the end-user from the burden of knowing the physical characteristics of each variable, for example where each variable is physically located, how long each variable is, and so forth. Once a raw file has been read into a database package, end-users can simply access variables by name, usually with some a flexible, user-friendly query language. This is an excellent approach for master files which are used by groups of researchers.

Database access can also be enhanced by using the indexing capability built into most database software. A database index is essentially a computerized lookup table that speeds data access, similar to the way the index in the back of a book works. By indexing the variables which are frequently sorted on or used to select extracts, researchers can realize significant time savings.

The ability of some database packages to provide transparent access to remote databases is another way database software can enhance access to research databases. Several database packages (e.g. Ingres and Oracle) advertise remote access to data on different platforms through high speed networks. SAS will soon offer an add-on product, called SAS Connect, which will provide the same capability for SAS datasets.

¹ Paper presented at the 16th Annual Conference of the International Association for Social Science Information Service and Technology (IASSIST), Poughkeepsie, New York, May 30-June 2, 1990.