

**DATA MANAGEMENT
AT
BIOLOGICAL FIELD STATIONS
AND COASTAL MARINE LABORATORIES**

January 1992

Report of an Invitational Workshop
April 22-26, 1990
W.K. Kellogg Biological Station
Michigan State University

sponsored by
Organization of Biological Field Stations
and
Southern Association of Marine Laboratories

prepared for
National Science Foundation
Division of Biotic Systems and Resources
Biological Research Resources Program

edited by
John B. Gorentz
W.K. Kellogg Biological Station

Any opinions, findings, conclusions, or recommendations expressed in this report are those of workshop participants and do not necessarily reflect the views of the National Science Foundation.



TABLE OF CONTENTS

Preface	v
Introduction	1
Executive Summary	3
Chapter I - Data Administration	4
Chapter II - Data Standards for Collaborative Research	15
Chapter III - Computer Systems for Data Management	19
Chapter IV - Summary of the Workshop Survey and Pre-workshop Demonstrations	29
Appendix A - Participant List	42
Appendix B - Geographic Information Systems/Administrative Issues	45
Appendix C - Client/Server Database Architecture, Networks, and Biological Databases	48
Appendix D - Intersite Archival and Exchange File Structure	52
Appendix E - System Selection Overview	57
Appendix F - Workshop Survey Questionnaire	60
Appendix G - 1982 Workshop Report (reprinted)	62



PREFACE

The data on the natural populations and biological processes of a biological field station's habitats are a research resource, just as are the buildings, research equipment, and habitats themselves. Or more accurately, they are a potential resource, a potential that is realized only when the data are organized, documented, and cared for to make them usable and accessible.

Although this issue of data management has been given increasing attention the past several years, and much progress has been made, it may be that the task of developing these data resources lies largely ahead of us.

A workshop was held at the Kellogg Biological Station in 1982 to encourage and foster the development of data management at field stations. Since nearly a decade has passed, it seemed an appropriate time to assess the progress that has been made, to reexamine our goals, and to determine what can be done to encourage and lead the way to the further development of databases and their utilization.

We sought support from the National Science Foundation for a data management workshop at which representatives from field stations and coastal marine stations could examine the state of data management, share information, and propose goals and new projects to advance this important work. As terrestrial and coastal marine stations wrestle with ways to allocate their limited research resources to this need, they can and should learn from each other's successes and mistakes. Field stations as a group have some unique objectives and requirements, giving them a common interest in data management that is somewhat distinct from other data management activities and goals. We believe the essence of the workshop deliberations

held at the W.K. Kellogg Biological Station during April 22-26, 1990, has been effectively captured and documented in the report that follows.

The workshop was supported by a grant from the Biological Research Resources Program, National Science Foundation, was co-sponsored by the Organization of Biological Field Stations (OBFS) and the Southern Association of Marine Laboratories (SAML) and hosted by the Kellogg Biological Station.

Thirty-six participants were invited to the workshop, representing data managers, scientists, and administrators representing biological field stations and marine laboratories of the United States. They represented sites newly embarked on data management programs, as well as those with well-established data management facilities.

The workshop was organized into three working groups, each led by two rapporteurs. These rapporteurs compiled the findings of their respective groups, and authored the first three chapters of this report. It should be recognized, however, that each chapter contains material originally contributed by those in the other groups; the topics are interrelated and it was impossible for any one group to consider its agenda in isolation from the others. On the day prior to the workshop, to provide some background for the participants, a pre-session symposium and the results of the a pre-workshop survey were presented. These materials are summarized in the fourth chapter. Because so much of the discussion at the 1990 workshop was made in reference to the 1982 workshop, it was decided to reproduce the 1982 report in the final appendix to this report.

Co-Principal Investigators
James J. Alberts, SAML
John B. Gorentz, KBS
George H. Lauff, OBFS

RAPPORTEURS and AUTHORS

Data Administration:

William K. Michener
Ken Haddad

Data Standards:

Warren Brigham
James W. Brunt

Computer Systems for Data Management

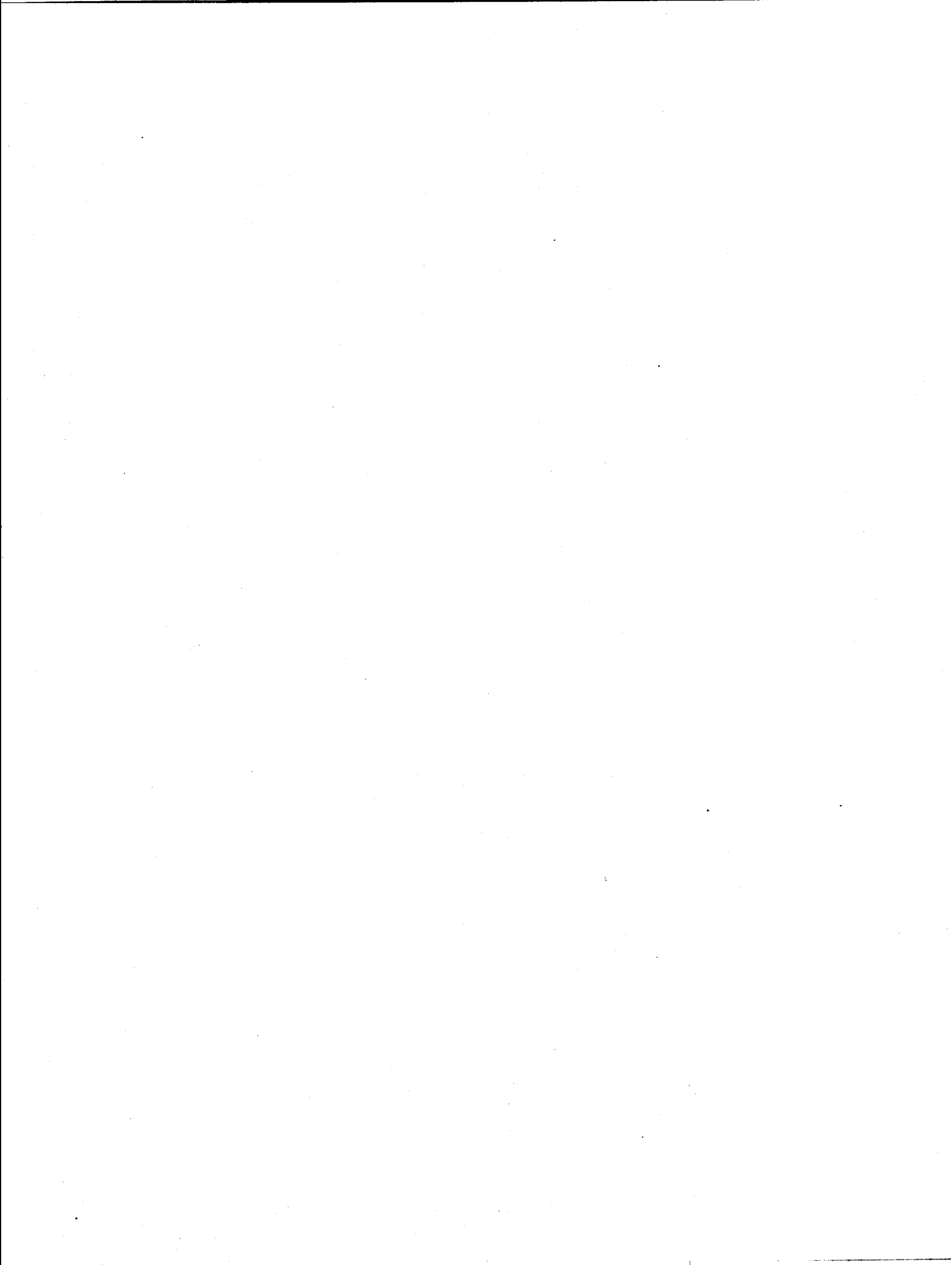
John H. Porter
Jeff Kennedy

Summary of Survey and Pre-Session

John B. Gorentz
Michael P. Hamilton

Editorial Consultant

Edie Erwin



INTRODUCTION

Science is based on the free and open exchange of information, whereby scientists can build on the work and data of those who have gone before them. Since science builds on previous work, including that represented in previous databases, scientists have a responsibility to preserve data for those who will follow after them.

In this context, the data gathered at biological field stations and marine laboratories constitute a national resource which should be preserved and made accessible for the purpose of advancing science. Long-term records of populations and biological processes in natural habitats as well as the physical and chemical environment in which they occur, are a research resource necessary to the study of ecological processes of regional and global significance.

Data sharing through the traditional system of refereed publication is not always adequate; there are unpublished data, never-to-be-published data, and raw data behind publications that need to be made available as a resource for others. Although some disagreement exists over whether available resources should be spent testing hypotheses rather than on preserving data without a clear hypothesis to be tested, it is generally agreed that the main purpose of long-term data management is to provide descriptive background data which can serve as a context for experimental studies. Research should always drive data management, rather than vice versa.

For the purpose of this publication, data management means caring for certain data so that, whatever their original purpose, they are preserved and made available for more general use, now or in the future. A field station's data management is distinct from computer management or investigator-specific data management, although it encompasses both. Comprehensive data management goals, realistic long-term planning, and solid institutional commitment are necessary for the care of data at each station. However, field stations and marine laboratories cannot manage data in isolation from each other. They need not only to collaborate and cooperate in data exchange, but also learn from each other's experiences, successes, and mistakes in developing their data management systems.

This publication is the result of deliberations by 40 representatives from stations and laboratories of all sizes. Their object was to produce a usable decision-making tool for data management planning and implementation. It is hoped that their shared wisdom will benefit the researchers, ad-

ministrators and data managers at all field stations and marine laboratories.

The first three chapters of this report represent the conclusions of each of the three working groups—Data Administration, Data Standards, and Computer Systems. The first chapter summarizes the results of a pre-workshop survey and a series of demonstrations presented by participants at a pre-workshop symposium.

The goals of the three working groups were:

DATA ADMINISTRATION

- Identify the benefits of an institutional data management program for those sites deciding whether or not to embark on one.
- Identify the types of data management that can be of use: the data management services that can be provided, the types of data to be managed, and the types of resources and staffing.
- Distinguish between that data management which is appropriately undertaken by a site and that which is best left in the hands of individual researchers.
- Identify administrative structures by which data management programs can be successful, identifying appropriate relationships between data management, research, and site administration.
- Identify realistic funding levels and methods of funding data management.
- Identify growth trajectories appropriate to field stations of both large and small size and levels of activity.
- Identify means of long-term care and storage of data.
- Consider the role of Geographic Information Systems in relationship to more traditional data management.

DATA STANDARDS

- Identify areas in which standards are needed to make data management for collaborative research more efficient, and areas in which they are best avoided because they may hinder research more than help.
- Identify the potential benefits of standards in data management.
- Identify existing protocols that might be adopted.

EXECUTIVE SUMMARY

The following is a summary of the major findings and recommendations appearing in Chapters 1-3 of this report.

Database Administration

- A data management program can benefit inland and coastal field stations by increasing scientific productivity and increasing the effectiveness of site administration.
- Those sites possessing effective data management systems remain the exception rather than the norm.
- Each field station and marine laboratory should perform a needs assessment to determine where data management fits into its overall mission, and should establish policies and directives accordingly.
- General guidelines for developing data management systems are 1) start small, 2) learn from other related institutions, and 3) find the right persons. Data management plans should allow for incremental growth.
- Training, though expensive, is likely to provide long-term benefits in productivity.
- Close communication between investigators and data managers is essential, but a site's data manager(s) should report directly to the site administrator, rather than to an individual investigator. Investigators and other site users should be involved in continual evaluation and review of data management.
- Site policies should reconcile the conflict between investigators' proprietary rights and general accessibility to data. A data ethic should be encouraged, which maintains that it is unacceptable for data sets with general utility or long-term value to remain permanently inaccessible.
- Data management should be viewed as an appropriate and necessary expense in research budgets.
- Those persons who evaluate research proposals or perform site reviews should examine how data resources are being cared for. However, funding agencies should not enforce unreasonable standard data formats.
- A mechanism is needed by which small investments, perhaps in the \$5,000-\$15,000 range, are available to get data management programs started, especially at new or small sites. These programs should be focused on specific general-use databases.

Data Standards for Collaborative Research

- Long-term studies and research on regional or global phenomena require the development and use of standards for documentation and ex-

change, so that data gathered at different times and places can be brought together for comparative analysis.

- Standards should be developed only for specific needs, with full consideration and involvement of the people who are intended to benefit from them, and should not be arbitrary or overly restrictive.
- The test of adequate documentation is that it should contain sufficient information for a future investigator who did not participate in collecting the data to be able to use it for some specific purpose.
- The Intersite Archives File Structure (Appendix D) is a recommended protocol that can be used by field stations to store and exchange data and documentation.
- A series of workshops should be funded to provide training, help field stations exchange information on data handling, and produce shareable databases. In the process, standards will be developed or adopted as needed.
- Multi-site, network-accessible databases should be funded as pilot projects.

Computer Systems for Data Management

- The single most important component of a computer system for data management is dedicated staffing to implement and operate it.
- No single computer system will be appropriate for all stations and laboratories. Systems must be tailored to achieve specific levels of data management and fit within resource constraints.
- A "top down" approach should be used in selecting computer hardware and software. The selection process should focus on data management and research tasks and the software and hardware needed to address them.
- Connection of a field station or marine laboratory to one or more wide-area networks can greatly enhance opportunities for scientific collaboration and help reduce the isolation that researchers at field stations often experience.
- Rapid changes in technology make good communication (electronic or otherwise) between data managers at different field stations critical.
- The best protections against loss of archived data are continuity of management and a strong data archiving policy. Technological backwaters and deterioration of media can be avoided by data managers who remain alert to changes in their computing environment and are aware of media limitations.
- An expansive definition of a computer system for data management can include facilities for visiting researchers. In some cases computers and computer access by visiting scientists to resident data bases are critical to the success of scientific investigations.

- Identify mechanisms by which researchers and data managers can communicate with each other to develop such standards as are needed.

COMPUTER SYSTEMS

- Provide guidelines for choosing system capabilities that can aid data management.
- Identify computer systems appropriate to both large-scale and small field stations and marine laboratories.

- Discuss the impact of new technologies on data management, including not only computers and software but also local-area and wide-area networks.

- Identify costs of networking, both initial and recurring, to assist preparing budgets.

efficiency of site administrative activities. Additional indirect benefits such as expansion of the field station's financial resource base may accrue as funding agencies continue or expand support in relation to the increasing value of that site's data as a resource.

Increased Scientific Productivity

A data management system which reduces duplication of efforts, facilitates awareness and communication of a site's data resource, and leads to better coordination of research efforts can significantly increase scientific productivity.

When data are made more freely accessible, use of data is expanded, reinterpretation of previous studies is possible (perhaps with the help of new types of analyses), an historical record for research and site use is established, duplication of effort is reduced, data are incorporated into the literature more rapidly, loss of data is prevented, and misuse of data is more easily discovered. For data sets with general utility or long term value, permanent inaccessibility is unacceptable.

Every site can benefit from a "data ethic" based on a self-evaluation of its treatment of data resources. Such an evaluation can lead to greater awareness of the current and potential value of a site's database and a recognition that specific data management activities may preserve and even enhance the value of that resource.

Increased awareness of data availability at a site through the production of a catalog of data, site bibliography, and data archive can often reduce the need to perform pilot studies and may facilitate experimental design and implementation.

Many data sets (e.g. meteorology, water quality, habitat characteristics, species lists, etc.) are of general interest to a large number of scientists. However, each scientist cannot always justify the costs of individually collecting, storing, documenting and performing the data management activities necessary to maintain the complete variety of data sets which may have relevance to his or her specific research interests. Even when scientists do have the resources to compile such data sets, they usually do not have the resources to provide the long term care necessary to make them available to a wider audience.

Sites may choose to fund, collect and store some data sets as a site activity. Relevant examples are the locations of field sites (past and ongoing), meteorological data, and other data sets which are site-specific, but of general interest and long term value. Well documented and archived meteorological and habitat data can facilitate the planning of experiments and sampling regimes by providing details about seasonal weather patterns and historical sampling locations. These activities must be carefully selected on the basis of general-

ized needs of the site's users. Sites may also wish to act as custodians or archivists for individual researchers' databases.

Where long-term data exist, it may be possible to place short term experiments into a broader temporal context. New studies may be more efficiently designed and implemented when they can be coordinated with ongoing research projects.

When data are managed as a long term resource, new investigators are often attracted to a site, and the potential exists for participation of that site in larger scale (intersite, regional, and global) comparative studies. Research sites appreciate in value as their historical databases grow.

Service to Researchers

A data management system which has the appropriate support staff can increase the efficiency of individual scientists by taking over responsibility for routine data management activities. In addition, data management consulting services provided to on-site investigators and visiting scientists regarding design and implementation of data sets, analytical tools available for interpretation of data, and hardware/software training can greatly increase scientific productivity.

Investigators working without the assistance of an organized data management system may not realize how much of their time is spent on data management. Having a site-sponsored data management system in place will not eliminate the need for investigators to spend time on data management activities; however, their productivity should increase as less time is required for more routine data management tasks.

Development and implementation of quality assurance and quality control procedures can facilitate scientific research through detection of data corrupted by human and machine errors as well as by media degradation. Other tasks, including data documentation support, data archival, and translation of data from one format to another, can often be more efficiently performed by experienced data management personnel. However, the investigator should always be involved in the process whereby the data are merged into a site's long term database system if quality is to be assured and documentation maintained.

EFFICIENT SITE MANAGEMENT

Data management can provide the means to document the project and site output that forms the basis of financial support for a field station's resources and activities. It can also enhance communication with off-site investigators, funding agencies, and institutions. Maintenance of ongoing and historical data sets can facilitate monitoring of the biological integrity of the site and provide data necessary for site impact assessment studies. The

CHAPTER I—DATABASE ADMINISTRATION

William K. Michener
Baruch Institute
University of South Carolina

and

Ken Haddad
Florida Marine Research Institute

I.1.0 INTRODUCTION

In the scientific process, answers to questions about the real world are coaxed out of data sets containing observations of patterns and processes. Various methods may be employed, but all rely on the availability of high quality, well documented data. All scientists participate in data management activities to varying degrees. Data management may therefore be viewed as a critical component of the scientific process.

Science builds on past knowledge which serves as a basis for future advances. The research community associated with field stations collects environmental data that represent a national resource which should be conserved for posterity. These data can, in many cases, be used to examine the effects of global change, loss of biodiversity, and habitat degradation. Scientists working on site-related or ecosystem-related research have a responsibility to future scientific efforts. Through improved preservation, access, and management of data, scientific research can be enhanced.

Ideally, all field research sites, stations and laboratories would have a data management system to serve the needs of current and future research. A data management system consists of both physical and functional attributes. Physical attributes include the people, hardware and software that are necessary to manage a site's database. Basic data management functions that are typically implemented to varying degrees at inland and coastal field stations include:

- a. Record keeping of ongoing research (who, what, where, and when)
- b. Organization of historical information (history of research and land use activities at the site, facilities development, site personnel, institutional support, etc.)
- c. Facilities support
- d. Individualized project support (data entry, file maintenance, security, documentation)
- e. Acquisition and maintenance of basic databases for use by multiple investigators (specimens, maps, species lists, meteorological and hydrological data, etc.)

f. Data archiving

g. Communication of data (maintain public database, network with multi-site projects)

Data exist in two primary forms at field stations and marine laboratories, site information and researcher specific data. These site-information data sets include:

1. Data on the user base

Lists of researchers and projects
Mailing lists
User statistics

2. Bibliographic data

Library catalogs
Published papers about the site
Theses and dissertations and reprints

3. Site characterization data

Meteorological data
Hydrographic records
Notes on land use

4. Inventories

Species lists
Collections
Maps and photos

Researcher specific data are generated by individual research projects and may or may not be of interest to subsequent researchers at the site.

In the following sections, we examine the benefits of data management; its current status at field stations throughout the country including obstacles to implementation; a blueprint for planning a system; suggestions for implementation; and a discussion of costs and evaluation. Since many administrators are exploring ways to store, retrieve, and analyze spatial data relevant to their sites, a separate section (Appendix B) is devoted to discussion of geographic information systems (GIS).

BENEFITS OF DATA MANAGEMENT

An effective data management program can directly benefit a site in two ways: (1) it can increase scientific productivity and (2) it can increase the

What databases not currently available can potentially be recovered and made available? What databases are anticipated in the future? Do the databases relate geographically/biologically?

Are the databases of a short-term or long-term nature?

Are databases in analog or digital form?

Do important data sources consist of photos, imagery, video, etc.

What levels of scale and/or scope are represented by the various databases (subcellular to landscape)?

3. Volume of activity:

What is the current and projected number of researchers/students?

What is the size of historical and current databases?

How many potential and actual users exist?

4. Sophistication of data generating, processing and managing activities:

What computerized storage facilities exist?

Is there access to off-site resources?

What is the potential for storage and access?

What kind of processing services and equipment are available?

What levels of expertise do the on-site personnel possess?

5. Infrastructure:

What are current and potential sources of support?

Is there an on-site library and what are its capabilities?

Can the library be used as a service node for data access?

What kind of data acquisition equipment is available?

How many support personnel are on-site?

Is it a seasonal or year-round operation?

Planning

In setting priorities, a station should identify the potential level of data usage, determine common needs, and identify potentially valuable long term data sets, including historic, current, and future data sets. Data management priorities, like research priorities, can be viewed as a compromise between what can be done and what should be done. Addressing the following questions in the light of research priorities may help set priorities for data management: (1) What do I, as a scientist, wish I knew about the history of a site? (2) If I could go back 50, 100, or 1000 years, what would I record for the future? (3) What present conditions are im-

portant enough to record for posterity? and (4) If I were presented with an historic data set, what ancillary information would I need in order to effectively make use of the data? These may be difficult questions to answer but may suggest actions to be taken.

Investigators should be consulted before a site establishes guidelines. Speculation and contemplation of future needs and priorities should be encouraged. Agreement among the on-site researchers and the external scientific community should be sought. By addressing the needs of the research community through an assessment process, one can avoid forcing unnecessary or unreasonable standards on investigators for such things as data storage and data transfer formats.

Priorities will also be affected by changes in the goals of the station and the parent institution. Funding sources can affect priorities, but they should not drive the process.

Data Archives

Many stations, after assessing needs, will conclude that they need to archive data for subsequent retrieval. Research at a site will be greatly enhanced when other data sets from that site are available. Many data sets have broad or long-term significance and should not be lost. Funding and infrastructure will be needed to support them. Stations that take on this responsibility need to ensure that important data are appropriately deposited in a system that is secure, yet allows reliable retrieval. This can be done on-site or off-site. In either case, the issues of access, longevity and quality should be addressed:

1. Access

- volume of data
- volume of requests
- level of interest
- documentation
- data formats
- cataloging
- ownership of data
- remote/on-site

2. Longevity

- primary storage media
- changing formats
- physical stability of media
- redundant storage

3. Quality

- multiple versions
- documentation
- expertise for monitoring quality
- standards

data necessary for balancing the selection of new research sites with the need to preserve the integrity of historic research sites can be cared for.

Many activities, such as visitorship, laboratory space management, and vehicle and equipment scheduling, are not usually perceived as data management, yet most field station managers perform this kind of data management on a day-to-day basis. This information is lost when there is no policy or mechanism for retention, and the data are discarded after use. The loss of these data results in lost opportunities for long range planning and improving the economies of site management.

CURRENT STATUS AND OBSTACLES TO IMPLEMENTATION

Despite the potential for increased scientific productivity, expansion of a site's financial resource base, and facilitation of site administrative activities, sites with effective data management systems remain the exception rather than the norm.

The reason for the slow and sporadic development of data management systems is sometimes attributed to the lack of an adequate staff and sustained funding. Understaffing is a problem on all operational fronts (see Chapter 4, Administration and Personnel), and data management is a time-consuming task whose needs are often underestimated.

However, effective data management systems may also be slow to develop for a number of other reasons related to: (1) a lack of recognition that, in addition to habitats, physical facilities and personnel, data are the most valuable resource that a site possesses; (2) an unrealistic or inadequate assessment of site-specific needs; (3) a lack of agreement on goals and priorities; (4) a lack of integration of data management into the overall site administrative scheme; and (5) a lack of communication among site administrators, researchers, and data managers.

BLUEPRINT FOR DATA ADMINISTRATION

The basic administrative tasks involved in establishing a data management system can be briefly stated as:

1. Identifying the user community, inventorying data and assessing their importance in light of the field station's mission.
2. Developing a data management policy appropriate to the mission and user/data profile.
3. Developing a list of data management priorities and assessing the methods and hardware/software options necessary to address those priorities.

4. Developing a justification for enhanced allocation of staff and budgeting resources devoted to data management needs, based on the preceding analyses.

Without the support of site administration, a viable data management system cannot be realized. Site administration, in conjunction with the research community, must be responsible for design, implementation, and continued support of data management. The design phase requires adequately addressing the data management needs of the present and future community of researchers likely to use the field station. Performance of a needs assessment will help determine where data management fits into the overall site mission.

Implementation of a data management system requires that considerable attention be paid to staffing, incorporation of data management into the administrative hierarchy, and funding. After initial implementation of a data management system, continuing support activities (including evaluation and management of incremental growth) must be performed. The design and implementation phases are discussed in further detail below.

INSTITUTIONAL COMMITMENT

Without an institutional commitment there can be no guarantee of continuity, and data management activities will likely be characterized by responses to short term, project-specific requests rather than the comprehensive support which is possible with a broad and well-integrated system.

NEEDS ASSESSMENT

Each station should do an assessment of its own needs and priorities. Stations differ in their needs and their ability to support data management. Some can support higher levels and intensities of data management than others. The following list presents some questions which should be examined as part of a needs assessment.

1. Mission, goals and objectives of the site:
 - Is it a preserve?
 - Is it a teaching facility?
 - Is it a research facility?
 - Does it support its own researchers or seek to attract visitors?
2. Type of scientific data being collected:
 - Is it descriptive and of general interest to various researchers?
 - Are there ongoing projects? Projects of historical interest?
 - What databases exist?

learned on costs and benefits should be used in planning the next project.

Subsequent projects should be chosen by consensus, considering overall site needs. These projects could include more individual specialized research projects but should be prioritized on the basis of cost vs. benefits. At some point it will become apparent that the next incremental step will require people and equipment necessary to accomplish the next tier of objectives. It is generally safe to assume that data management tasks will be more complex and time-consuming than anticipated. Time schedules proposed for projects will inevitably become more realistic as both data managers, scientists, and administrators become more experienced in administration and implementation of specific projects.

Eventually a site archival facility should be established. The data management system should be designed to survive the loss of key data management personnel and changes in research emphasis (driven by both investigators and funding) at the site. Specific arrangements for archiving data either on-site or off-site should be addressed by the site administration. Options include submission of data to the recognized site data management system, the parent institution, or a regional or national data bank. In the event that a station is closed, mechanisms should be established whereby responsibility for maintaining the database passes to the parent institution or "field station community" (a "sister institution" or possibly a regional/national data bank).

A key issue for data management administration is the resolution of any conflict between the investigator's proprietary rights and the need for general accessibility to data. Permanent inaccessibility to data is unacceptable, but the investigator should be able to control access to his/her data for a reasonable period of time. Several issues must be addressed when considering the issue of proprietary rights. These include the potential for allowing others to publish findings before the investigator who collected the data does, misinterpretation of the data by someone unfamiliar with the experimental design or habitat characteristics, and using the data out of context. Questions of legality, including "Who owns the data?" and "Who is liable for misuse or misinterpretation?" must also be answered.

Many universities and funding agencies have regulations regarding the fate of a data set. However, each site should have a policy regarding proprietary rights or should negotiate with individual researchers. In either case, issues of data ownership and access to particular data sets should be clarified with researchers before the project is begun.

If release will hinder research or use of the site, or violate local, state, or federal regulations (e.g., regarding endangered species), the site may wish to restrict access to the data. However, individual sites must look into the pertinent regulations and develop a policy which incorporates them.

Oversight

Ideally, the data manager(s) should report directly to the site administrator (Figure 1a). At many small sites, a single individual may function as both the site administrator and the data manager. Close communication between the investigators and the data manager is essential, although the site administrator has the ultimate responsibility of directing the data manager while also addressing the needs of the scientific community and responding to influences external to the site. These influences may include: (1) database requests from other scientists, agencies, and institutions; (2) funding agency requirements; and (3) institutional, state, or federal obligations.

Site administrators are cautioned against implementing an administrative hierarchy whereby one of their scientists assumes control of the site's data management personnel (Figure 1b), since this has a high probability of fostering conflict within that station's scientific community. For example, the perception that the personal agenda of the scientist administering data management receives precedence over the broader station objectives frequently arises. This perception, whether based on reality or not, may serve to isolate data management from the other scientists at the site.

Some sites, particularly large ones, may wish to implement an administrative hierarchy whereby the site administrator oversees a "data management steering committee" which periodically reviews the data management system and participates in establishing and prioritizing objectives (Figure 1c). For example, although data management personnel may report directly to the site administrator on a routine basis, they may also participate in monthly or quarterly reviews by the steering committee. The steering committee would ideally be comprised of the site administrator as well as a manageable number of scientists who represent the various research programs at the site. This scheme provides an important feedback mechanism to the site administrator and facilitates communication among scientists and data management personnel.

Staffing

Successful data management systems are usually staffed by persons who have a strong interest in the scientific research conducted at the site. This not only promotes effective communication with the

Plans for data archiving should take into consideration the volume of data, projected number of requests for access to it, and the anticipated level of scientific interest. The potential use can vary from a small number of scientists interested in addressing a site-specific question to a much larger inter-disciplinary scientific community that would use the database together with those from other institutions to address regional or global issues.

Regional data storage is reasonable if the volume of data and the use levels are high. Data documentation and cataloging of the data sets are crucial for access whether on-site or off.

A site should consider initially storing data sets in a standardized generic format (ASCII). This would allow flexibility in moving data sets and in accessing them remotely. The question of data ownership should be addressed early in the design phase.

The issue of longevity requires consideration of changing formats and the physical stability of media. There may be a need to ensure access through redundant storage. Disasters can destroy data on-site and off-site. If data sets have been prioritized and the most used and most critical stored in more than one location, access is preserved. The management of redundancy must be integrated into the site's data management plan. The plan should address updating data sets to avoid multiple versions which lack adequate documentation. Data sets should be tracked to manage the numbers of copies and to allow for purging of outdated data sets.

Data quality is of special concern to scientists who use data they did not themselves gather. A station must assure the highest degree of quality control over its own data and provide full documentation of data obtained from elsewhere for its own researchers. Disclaimers should be stated where appropriate.

Quality of data is linked to the development of standards for data generation and documentation. Researchers should be encouraged to fully document data before submitting it for archival storage. Participation of the scientific community in designing and implementing data set documentation can be a very valuable step towards insuring that data sets are usable in the future.

Accessing and archiving data costs money. Ideally, cost should not be a barrier to access. To encourage shared databases, stations should strive to supply data sets free of cost to those scientists who participated in their development. Data should be accessible to others at minimal or no charge, but cost recovery may be appropriate and necessary. Legal and institutional obligations regarding data accessibility will need to be addressed at each station.

One solution for small stations may be participation in development and maintenance of regional or national data banks. Data archiving in such data banks could prove cost effective, but the concept needs additional study.

The best protections against loss of archived data are continuity of management and a strong data archiving policy. Technological backwaters and deterioration of media can be avoided by data managers who remain alert to changes in their computing environment and are aware of media limitations. A data archiving plan or policy can and should ensure that a data manager remains alert and sensitive to potential problems.

IMPLEMENTATION

Successful implementation of a data management system requires that site administration pay particular attention to: (1) development of a reasonable plan which supports incremental growth, (2) effective incorporation of data management into the administrative organization, (3) assessment of costs and procurement of necessary funds and staff, (4) proprietary rights, (5) continuity of management, and (6) continuing evaluation. The research community and site administration, along with data management and funding agencies where appropriate, should discuss and agree on goals and priorities before beginning the implementation.

A plan for data management should allow for incremental growth. The guidelines are: (1) start small, (2) network, and (3) find the right person.

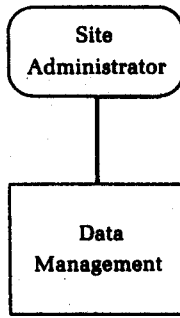
Start with a pilot project small enough to be accomplished within a reasonable time frame but important enough to have a beneficial impact. An example might be development of a reprint list, taxonomy lists, collection list, or acquisition and cataloging of aerial photography.

"Networking" means seeking advice from sister institutions with similar size, resources, and mission. Many pitfalls can be avoided by learning techniques used at other sites.

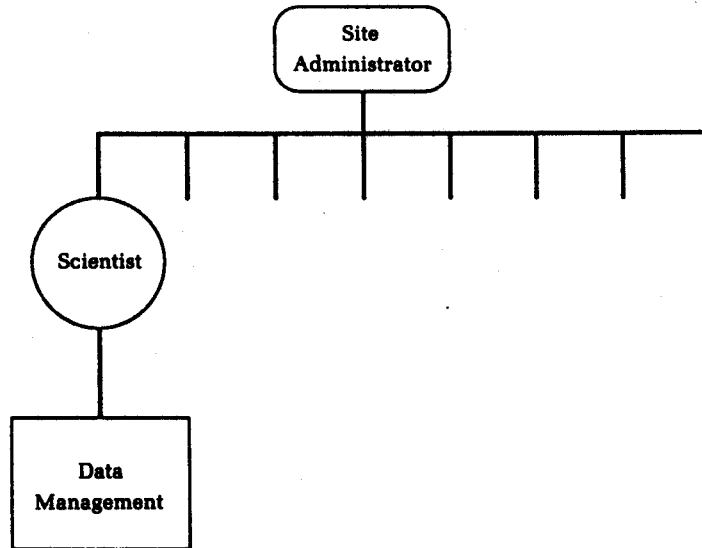
Choose a data manager who has research experience and a scientific understanding of the site's research program but also has data management skills. The person should have an interest in and enthusiasm for the data and products to ensure success. A data management system cannot necessarily be sustained over the long term as a "labor of love," but high enthusiasm and intensity are needed to get it firmly established. Good communication skills and relationships with data users are essential.

Upon completion, the pilot project should be evaluated. Maintenance costs incurred should be considered, because even when the system is established, it will not maintain itself. Lessons

A



B



C



Figure 1. Administrative Structure for Data Management.

scientific community, but often helps assure retention of data management staff at field stations which cannot offer competitive salaries.

Requisite skills and expertise needed by data management staff will be largely affected by the size of the organization and the length of time the system has been in place. Generally, the more complex the operation, the greater the need for more specialized personnel. At small sites or those with low research activity it is essential that the data manager have expertise in both science and data management and that this individual have access to appropriately qualified consultants. At some sites, primary training in biology may be appropriate, whereas at others a background in chemistry, geology, physical oceanography or other relevant disciplines may be appropriate. In any case, the initial staff member should be a scientist first.

For larger, more active sites a systems analyst/programmer should be added next. It is critical that the data management staff be able to communicate with the scientists and also have the expertise to accomplish what is needed. This means that at least some of the data management staff need to be very broadly trained. Increasing the size of the data management operation brings increasing specialization.

An alternative to making data management an adjunct to computer support is to staff data management as an adjunct to an existing library or museum. This can be appropriate at institutions where the library or museum already has strong ties to an information management and retrieval program. Links to computer support would still be needed but on a secondary basis.

The use of graduate students as a cost saving means is problematic as it may restrict continuity and could cost in additional training time. However, it does allow for student educational/financial support and with careful choice could provide talented personnel. Sites may wish to examine the possibility of developing one- to three-year undergraduate and/or graduate internships or independent study programs to accomplish specific tasks. Some tasks, such as data entry, routine quality assurance, and graphics production, may be more appropriate for temporary personnel.

The ability to communicate with a wide range of people is the most important qualification for a data manager, assuming an appropriate level of technical skills. Data managers must be able to effectively articulate the purpose and needs of the system to site administrators, researchers and the parent institution, as well as field questions and demands from on-site and off-site users. Consistent communication with other data management personnel promotes better and more creative systems.

Site specific technical skills might include organizational or curatorial expertise, experience

in data collection, storage and retrieval skills, programming knowledge, networking experience, proficiency at hardware and software maintenance, and GIS experience if appropriate. Again, the technical skills are site specific and vary depending on the size and activity level of the site. Organizational skills are basic and can include library expertise.

Strong administrative and financial support is necessary to attract and retain data management staff. Not only should personnel be adequately compensated, but data management should be provided with the necessary staff, equipment and operating budget for continual database maintenance and update. Administrative support should be demonstrated by ensuring that the data manager reports directly to the site administration. Consensus on priorities is necessary so that data managers can focus their attention on well-defined projects. Close communication with (but not supervision by) the researchers on-site will enable the data manager to be an integral part of the site's research activities. All site-related publications should acknowledge data management personnel and identify where data have been deposited, much as one does with plant and animal specimens.

Costs Associated With Implementation of Data Management

The equipment, staff, space and budgetary resources committed to data management vary widely among sites, reflecting the wide-ranging missions and academic clientele of the nation's network of field stations and marine laboratories. One site with ten scientists may require only a part-time data manager, whereas another with the same number of researchers may require two or more staff members to meet a broader range of data management duties. It is not possible to state a simple formula for the cost. However, it may be helpful in planning the implementation or expansion of data management to consider some scenarios along a continuum in which staffing is the most important limiting factor.

In the first scenario, there is no dedicated data management staff, and no formal data management, although some informal records may be kept. Documentation is spotty, if it exists at all, and backup copies of data are not maintained. Such a data management system does not require any computer and software resources. There are no long-term benefits to researchers. Without proper data administration it is not a question of whether data will be lost, but when they will be lost.

The second scenario features a part-time data manager typically capable of maintaining only one or two of the types of site data. Initial effort may be focused on identifying, acquiring, and documenting data. Except in isolated cases, data belonging to individual researchers are not managed under this arrangement. Part-time data

Annual maintenance expenses for hardware can be expected to consume approximately ten percent of the annual data management budget. Another rule of thumb is that annual recurring costs for hardware and software (i.e. updates, repairs, maintenance, license renewals) are likely to be 8-12 percent of the initial cost. Creative ways can often be found to keep these costs down, but usually at the expense of personnel time.

Training is of major significance. Hardware and software are evolving at a rapid pace and it is difficult for data management personnel to individually track these advances, learn new "tricks of the trade," or readily become proficient with new techniques and hardware/software tools. It has been demonstrated that continued expenditures in training provide a long term benefit in productivity both by the data management staff and the researchers. Training should be given a high priority even though initial expenditures might seem high. It is preferable to maintain an annual training budget and schedule.

Potential Funding Strategies

A successful data management plan must have adequate and stable funding. Potential sources of funding may be the parent institution's facilities budget or re-routing of appropriate portions of overhead costs (e.g. program management charges or indirect costs) to data management. Either mechanism requires full support of the parent institution. Any revenues generated through overhead and indirect costs associated with grant support cannot be considered completely reliable. Though probably not adequate for full support, user fees may be useful to sites catering to visiting researchers and can often result in partial cost recovery. However, this approach may tend to reduce the efficiency of data management if a "pay as you use it" approach is adopted and no mechanism is established for long term support of the facility and data management personnel.

No site/institution can rely on short term (1-2 year) grants to support ongoing data management. Funding for base level data management activities should be done with hard money. However, short term grant money can be very useful for providing start-up hardware/software, training, and other special projects. Data entry and/or conversion of previously collected data to appropriate formats, support for incremental improvements to existing systems, and the underwriting of publication costs (electronic or traditional) are possible uses of grant money. Short term grants can be used to implement specific data sets (station bibliography and data catalog, species lists, etc.). However, mechanisms should be established to cover the recurring costs of maintaining data sets after this initial investment.

Role of Funding Agencies

Data management represents a real cost for research. Therefore, it should be viewed as an appropriate and necessary expense for grant budgets. The challenge to funding agencies is to encourage ties between data management, field stations, parent institutions and the research community. Funding agencies can foster the development and support of data management systems by providing start-up funds for hardware and software, by supporting research in data management (e.g. development of more efficient database structures and quality assurance procedures, etc.), by supporting training programs, and by developing mechanisms for supporting database development.

Data management can be included as a line item in proposal budgets and as a topic to be examined during the review process. Although funding agencies should not force unreasonable standards on scientists for items such as data storage or transfer formats, they can encourage retention of data at field stations and elsewhere by asking scientists to state in their proposals what, if anything, will be done with the data when the study is completed, or as data are acquired. Referees should be encouraged to consider these factors when evaluating proposals. However, scientists should be assured of proper acknowledgment for the use of their data in any subsequent publications.

Funding mechanisms are needed for getting field stations started in data management. An initial investment of \$5,000-\$15,000 may be all that is needed to get a data management program off the ground, especially at new or small sites. A system of "mini-grants" would assist small stations in implementing a basic data management program. The possibility of internal reviews or mini-reviews for proposals of this size should be examined. The "seed project" model presently used by some funding institutions may be appropriate for initiation of data management at field stations.

Funding agencies might explore the possibility of funding the development of regional/national data banks for archiving of data from field stations. This might reduce the need for an extensive data management system at small sites.

Another important potential role for funding agencies is support for training data management personnel in the needs of field stations and providing them with the necessary skills and support group to meet these needs. This might be accomplished by sponsoring regional two or three day workshops or an exchange program whereby personnel visit sites with data management systems in operation.

managers are typically unable to provide data entry services to researchers or support comprehensive or rigorous error checking.

A data dictionary is likely to be informal, non-integrated, and not automated. A file cabinet or single microcomputer may be the only hardware used for these activities. Software should be "off-the-shelf" packages that are in wide public use and which support generic data structures, because there will be little time available for customization.

Expenditures for training and annual maintenance may be minimal, though not because of a lesser need. Part-time data managers are likely to be successful only if they have access to those who can provide training and advice, and if they take advantages of maintenance and other support services for software and hardware. Access to electronic mail networks can be used to get help and advice from other data managers, as well as to facilitate access to data by researchers.

The primary research benefit under this scenario is the creation of a persistent institutional memory, at least about a few selected types of data.

To handle the backlog of historical data sets, or the startup of a new computer system, additional resources typically will be needed.

The third scenario features a full-time data manager. This level of staffing might be appropriate to a site with ten scientists. In addition to managing site characterization and administrative data sets, a full-time staffer may also be able to manage data for a small number of individual research projects. The extent to which this is possible

depends on the size, type and complexity of the data sets, and the number of data management "clients." Provision of data entry services for a few individual research projects becomes possible, but may require contracts for service bureau data entry.

Computational environments are more variable at this level, with the more powerful personal computers, workstations, and even mainframes being used to handle the larger volumes of data.

Benefits of such a system include easy access to site information, which can in turn facilitate integration of new research projects and researchers. Because a full-time data manager is available for consultation, individual researchers can be more efficient with the time they spend on data management tasks.

The last scenario is a data management staff comprised of several individuals, typically at a more active site with a large number of scientists. The individuals may be trained in systems analysis, database programming, computational ecology, statistics, or other technical fields. There may also be a full-time data entry staff.

The computational environment is typically complex, with several different types of computers, each used for specific tasks. Computer hardware may include a microcomputer network, mini-computer, or multiple workstations. The larger the site, the greater is the need for more storage capacity and processing power. Network connections, with electronic mail, remote terminal access, and file transfer capabilities, are desirable to facilitate off-site archiving, access to external databases and transfer of data to remote researchers.

Table 1. Four scenarios representing the range of costs associated with implementing data management systems at varying levels of intensity.

Site Activity (# of scientists and users)	Personnel Required (FTE)	Hardware/ Software Costs	Annual Maintenance Expenses	Training
1-5	0.20-0.75	\$ 4,000	\$ 300-500	Self-taught
10	1	\$ 8,000	\$ 600-1,000	\$5,000
50	2.5 (plus data entry staff)	\$ 100,000	\$ 15,000 (10% of total budget)	\$25,000
100	3.5 + (plus data entry staff)	\$1,000,000	\$ 200,000 (10% of total budget)	\$25,000 +

CHAPTER II—DATA STANDARDS FOR COLLABORATIVE RESEARCH

James W. Brunt
University of New Mexico

and

Warren Brigham
Illinois Natural History Survey

RESEARCH NEEDS

Increasingly, environmental scientists are being encouraged to focus attention on regional and global issues such as biodiversity and global change. The wide geographic distribution and diversity of ecosystems encompassed by inland and coastal field stations represent a major national resource. To address large scale environmental questions, scientists will require resources such as the data generated at these facilities.

Large scale scientific issues elevate the importance of data management beyond the needs of the individual investigator. When data are regarded as "belonging to science" and, therefore, to be shared with other researchers now or in the future, data standards to enhance communication become necessary.

By using standards, researchers can save time and prevent costly mistakes in interpretation of data. The activities that suffer most from lack of standards are the arranging and organizing of data, documenting of what has been done, and sharing and exchanging of data with other researchers.

Implementation of standards is particularly important where several researchers are working on a joint project. For example, if all investigators on a project adopt standard location descriptors, all localities referred to in data sets can be communicated to other data users reliably and accurately, and the data sets can be used for comparative analysis. In addition to project-wide data standardization, site-wide standardization can result in similar benefits to both the investigators and future users of the research site and its historical data. Current data sets can be compared with future data sets.

Electronic networks are making the sharing of scientific data for comparative analysis much more feasible. Field stations, herbaria, museums and other biological information providers benefit greatly from the data communication channels provided by research networks such as the Internet. Biological databases can be made immediately accessible to ecologists, systematists and conservationists around the world (Appendix C). But the success of network accessible databases in biology depends on the ability of the disciplines and sub-

disciplines to reach consensus on elementary data models and database structures.

Data standardization at the various levels, from the raw (primary) data to structures for user access and network exchange, should have as its primary goal the advancement of the science. The emphasis should be on those standardization strategies that maximize the conduct of science and the use of the data, at any level of the information management process.

Application of standards does involve costs, however. Perhaps the greatest cost is in instances where databases must be converted to comply with "newer" standards. This implies that carefully designed standards are best applied early in the development of data management at a particular site.

Creation and implementation of data standards should not be done in an arbitrary or overly restrictive manner such that the researcher's ability to collect and process data is restricted. Proposed data standards should be examined and applied only if they enhance data management. The need is not for standards that are in some sense sophisticated or elegant, but rather, standards that active researchers will in fact use to document and archive their data.

There can be benefits to having discipline-specific standards for representing space, time, and the relevant physicochemical data associated with biological information, but the appropriate persons to develop such standards are those researchers who need them.

TYPES OF DATA STANDARDS

Organization of Data

Organization of data refers to the logical structure of data — what all the variables are, how they should be organized into different types of records, and how the variables and records should be arranged with respect to each other. Standards for organization of data make it easier for scientists to analyze and re-analyze their own data as well as share it with other researchers. The 1982 Workshop Report, Data Management at Biological Field Stations (Appendix G, Chapter 2), describes data standards with respect to design of data sets which, if

applied, would enhance data management at any site with little adverse impact on a researcher's activities. The information presented in that document is still relevant to the management of biological data.

Data Documentation

Any researcher who has tried to produce syntheses integrated over space and time using previously collected data, including data from other researchers, has probably experienced frustration due to inadequate documentation of the data.

If the documentation describing a particular data set is lost, the data become useless. While this is particularly true for archived or historical data sets, lack of proper documentation can affect any data file. Thus, for exchange and archiving, data documentation should be incorporated with the actual data as soon as practical, possibly even in the design phase of the research.

The test of adequate documentation is that it should contain sufficient information for a future investigator who did not participate in collecting the data to be able to use it for some scientific purpose.

The 1982 Workshop report describes a standard for data documentation (Appendix G, Chapter 2). Some of the information categories may not be applicable to the data at every site and some additional categories may be needed for "non-traditional" ecological data sets (i.e. remote sensing and Geographic Information System files), but the essential elements are present.

Data Exchange

Field stations and marine labs represent a heterogeneous research and computing environment. The independence and isolation of field stations has led to a tremendous variety of data management approaches usually tailored to local needs, but which make data exchange and collaboration difficult. Although inter-university and international computer networks are becoming accessible to field station user, some standards must be followed to use them for data exchange.

There is a need for a non-restrictive but powerful common-denominator structure for data sets that will encourage good practices of documentation and communication. A complete data set should be an entity that contains all of the relevant documentation, as well as a history of the data. To be complete, documentation should include comments and annotations about the data set as a whole, and also about the individual records and observations where necessary.

One generic file structure that can be used is the Intersite Archives File Structure (Conley and Brunt, Appendix D). It can facilitate an orderly approach to the design and implementation of field station and marine lab data exchange capabilities. The

structure is one that includes full documentation and comments with the data. It solves the problem of possible separation of the data from the documentation. The data itself can be of any basic type, such as statistical data, text data, graphics data (e.g. files that can be written to a graphics plotter), gene sequence data, or bit map image data.

IMPLEMENTATION

The Development and Adoption of Standards

There are several possible routes to the development and adoption of data standards. At one extreme, *de jure* standards can be put in place by fiat. At the other, *de facto* standards can be adopted by survival of the fittest.

Standards imposed from above, without full consideration and involvement of the people who are intended to benefit from and use them, are usually ignored. Equally suspect are standards promulgated for political purposes, by institutions eager to enhance their own standing, without regard for research value and technical merit.

On the other hand, standards developed through a completely *ad hoc* process tend to be developed inefficiently, with much reinventing of wheels. Standards developed in this manner tend to lack rigid definition, so that there is no way of knowing whether compliance is apparent or real. These standards, too, may have political value, in that one can easily (and truthfully) claim compliance; but very little efficiency is gained.

A relatively non-dictatorial process somewhere in between these two extremes will involve researchers and data managers in developing, testing, and using standards relevant to the full array of types and uses of ecological and environmental data.

We recommend that specific standards be developed (1) through a series of workshops at which technical resources will be examined to address specific standards and data topics, and (2) by increased use of communication networks (including both electronic and personal networks) of biological field stations and marine laboratories.

Workshops

These workshops should not be limited to the narrow issues of standards, but should include information on technology for scientific data handling in general, such as data acquisition systems, data analysis tools, data handling in general, and data exchange. Opportunities to exchange this sort of information are currently quite limited. Training programs and seminars, as well as joint efforts to create shareable databases, are needed.

Progress in data standards will be made through common consent and practice, utilizing the expertise of those with relevant experience. Training can be provided by persons knowledgeable in fundamental data management principles as they apply to scientific data.

Data managers and researchers who have dealt with issues of bringing together two or more data sets for comparative analysis have much to contribute. Librarians in the field of information science, systematists, and museum curators who are affiliated with field stations will also have relevant experience.

In order to move beyond generalities and down to practical issues, each workshop should deal with specific issues, such as electronic networking or data archiving, or with specific types of data, such as species lists, site bibliographies, data catalogs, climatological data, or spatial data (Table 2).

The result should be shareable databases that relate directly to the testing and evaluation of scientific hypotheses. In the process, standards will be proposed, developed, and tested.

Network-Accessible Databases

The development of multi-site, network-accessible

databases, such as those being developed in the plant systematics community (Appendix C) should be encouraged. Especially valuable are projects which bring people from different sites together to apply their complementary areas of expertise. Bringing data together from two or more sources will require development or adoption of standards. Even more importantly, making those data accessible on the network will immediately test the usability and usefulness of those standards and of the entire concept of shared databases.

COMMUNICATION AND EVALUATION

Any workshops or network projects should place great emphasis on communicating information about their findings and products to the community of over 200 field stations and marine laboratories. It is expected that this would encourage further testing and evaluation of the utility of standards and databases.

Table 2. A series of workshops to provide training and information exchange and produce shareable scientific databases.

Technology-Oriented Workshops

**Electronic
Networking**

Collaboration via data exchange will benefit from communication technology and expertise. Participants will learn how to use electronic mail, network file transfer, and remote access capabilities. The product of this workshop should be a plan for network access via Internet/NSFnet.

**Data
Archiving**

Methods for data storage and archiving are developed. Some attempt should be made to identify the types of data appropriate for intersite access.

Product-Oriented Workshops

**Species
Lists**

Example: Systematize lists of species at the field station and marine labs. Species inventories are basic to biodiversity studies. Strategies for data update and exchange should be studied. For certain groups, development of a central database may be appropriate.

**Site
Bibliography
Development**

Data catalogs and site bibliographies need to be developed for every site in an exchangeable manner.

**Meteorological &
Hydrologic Data**

Develop standards and methods to share these types of nearly ubiquitous data.

Spatial Data

Develop spatial data standards for geographic information systems, global positioning systems, and remote sensing, etc.

CHAPTER III—COMPUTER SYSTEMS FOR DATA MANAGEMENT

John H. Porter
University of Virginia

and
Jeff Kennedy
University of California Natural Reserve System

INTRODUCTION

Research data management is increasingly linked to computers and associated technologies. Properly integrated hardware and software systems are crucial to managing large amounts of data. This chapter examines:

- 1) hardware and software selection
- 2) typical data management computer systems
- 3) uses and implementation of networks (both local and wide-area)
- 4) technological innovations that will influence data management methodologies, and
- 5) computer systems required for archival storage with special emphasis on the needs of marine laboratories and biological field stations.

The 1982 workshop report, *Data Management at Biological Stations*, (Appendix G) provides excellent guidelines for the management of scientific data at field stations. It includes comprehensive and thoughtful discussions of software and computer systems for data management, providing a blueprint for a complete data management system. However, at a majority of field stations, the 1982 recommendations for computer systems and software remain unimplemented. This is somewhat surprising, because advances in computer and network technologies solve many of the problems identified in that report. For example, the anticipated "proliferation of microcomputers with nonstandard floppy disk formats" failed to materialize, and a few standard formats have emerged. Moreover, the extensive use of wide-area and local-area networks, which was unanticipated in the 1982 report, has reduced the need to exchange data on physical media and thus has reduced physical barriers to data exchange.

The computational environment has also become increasingly homogeneous. The distinctions between the capabilities of mainframe, mini- and microcomputers drawn in the 1982 report are rapidly diminishing. Increasing numbers of software packages run on both mainframes and microcomputers. The dominance of certain software packages in microcomputer markets has led to emergent standards for exchanging data between different brands of software, computers, and operating systems. We anticipate that the rapid improvements in price and performance will continue at an accelerated pace.

Many of the shortcomings of relational database and statistical software described in the 1982 report have also been reduced. Although there are still improvements to be made regarding data documentation (i.e., data about data, or "metadata"), many of the problems associated with the entry of textual information have been reduced or eliminated. The increasing use of Structured Query Language (SQL) by relational database packages also provides opportunities for increased standardization.

Technical advances have led to new challenges for data management. For example, the increasing use of scanners and graphical and audio data (video and remote sensing imagery, maps, photographs, and sound recordings), with their large file sizes and specialized formats, creates problems regarding data storage requirements and file exchange compatibility.

Given that the technical barriers to achieving a functional data management system have decreased, the general lack of success in fully implementing the systems envisioned in the 1982 report seems paradoxical. The consensus among workshop participants and questionnaire respondents was that these recommendations remain unimplemented in significant part because the single most important component to a successful data management system is dedicated staffing to implement and operate it. Even the most "user friendly" software interfaces and most powerful computer systems are useless for data management without dedicated individuals (possessing the requisite computer expertise and interpersonal skills) to run them.

MANAGEMENT STRUCTURE

At many biological stations and marine laboratories, data managers are expected to provide computer system support (e.g., configuring hardware, managing and installing software, and administering networks) as well as perform data management functions. It is our recommendation that data management and computer system support duties be separated whenever feasible. The advent of complex multitasking operating systems on personal computers (UNIX, OS/2, Multifinder) and sophisticated networking software can cause a significant increase in the time taken by system support tasks and concomitant decrease in the time available for data management.

SELECTING COMPUTER HARDWARE AND SOFTWARE

It is a striking comment on the rapid innovations in computer technology over the past few years that choice of a computer system is increasingly arbitrary. The distinctions between the general types of tasks that can be performed by mainframe, mini- and micro-computers has all but vanished (although there still may be significant speed differences between different size computers in performing large tasks). The convergence of mainframe, mini- and microcomputer hardware and software means that it is increasingly necessary to take a "top down" approach, centered on data management and research tasks and the software needed to address them, rather than a "bottom up" approach that starts with the selection of computer hardware and software (Figure 2).

The first step is to identify the tasks to be performed by the system: What sorts of data will be managed? What sorts of manipulations or analyses will be needed? How large are the data sets likely to be? How many users need to be supported? Based on the answers to these and similar questions software can be identified that are capable of supporting these tasks.

Software

Typically, there will be several products with similar capabilities (e.g., all relational database programs share certain basic features). Deciding between them will depend on the particular characteristics of each product: functional strengths and weaknesses, cost, speed, ease of use, prior familiarity, compatibility with other software, availability of adequate user support (from the vendor, from a knowledgeable user or users group, or from the department or parent institution), and the financial stability of the vendor (i.e., will the company still be in business five or six years from now).

Most popular software packages have their share of proponents and detractors, whose differing opinions depend largely on their familiarity with the package in question. For this reason, it is good to query several sources regarding each package. Preferably, you should test the software yourself, using your own data.

Hardware

Once potential software packages have been identified, it is time to select the type of computer to purchase. In some cases, software will only run on a computer produced by a particular manufacturer and the choice is easy. More likely, there will be a variety of computer options, each capable of running the desired software. Obvious factors to consider are cost and processing power, but equally important are capabilities for expansion and obtaining a good fit with your institutional and user environments.

Constraints in the selection of computer hardware may dictate that some software choices be re-examined to improve integration with the hardware

and to maximize how well the components function as a system. The evaluation and selection process thus becomes an iterative process. Appendix E provides more detailed guidance for hardware selection and a list of software products which were popular with workshop participants. Some products that emerged in the year following the workshop are also listed.

Selecting GIS Systems

Selection of software and hardware to implement a Geographic Information System (GIS) is an extremely challenging task. The task is complicated by the diversity of software products and approaches and by the complexity of GIS software. The term GIS covers a wide range of activities, ranging from computerized cartography to spatial analysis. No one system is best at all types of GIS work. As with selection of a general purpose data management computer system, a "top-down" approach is recommended. The first step, task identification and needs assessment, is covered in Appendix B and will not be addressed here. Selection of software must be based on the ability of specific packages to perform the required tasks, the potential for expansion, the ability to interface with existing digital data sources, ease of use (which has a major impact on the amount of training required), the initial cost of the software package and the cost of continuing support and licensing.

Once a package has been selected, decisions must be made about the hardware configuration of the system. Will the GIS be directly accessible by station researchers or only by station personnel? How many "seats" will be provided and how will they be implemented? For some systems single-user workstations (typically personal computers) may be an economical choice. In order to balance computationally intensive tasks (e.g., producing polygon overlays or network analyses) with display intensive tasks (e.g., digitizing and editing data layers) it is recommended that single-user workstations be provided with sufficient memory and software to support multitasking. For other stations, multiple single-user personal computers sharing peripherals over a LAN, or multiuser computers and terminal workstations may make more sense and facilitate centralized administration of data and computer systems. Local area networks may also play a critical role in permitting sharing of large data layers, reducing redundancy and simplifying system administration.

Obtaining sufficient online storage for large data layers is often a problem. The large size of GIS files, the number of intermediate files generated by GIS processing, and the cumulative nature of GIS data acquisition combine to strain the resources of all but the largest systems. For this reason, provision for very large data storage and backup capacity is a firm requirement for GIS computers. Inclusion of high-capacity off-line storage for backup and archival storage is strongly recommended.

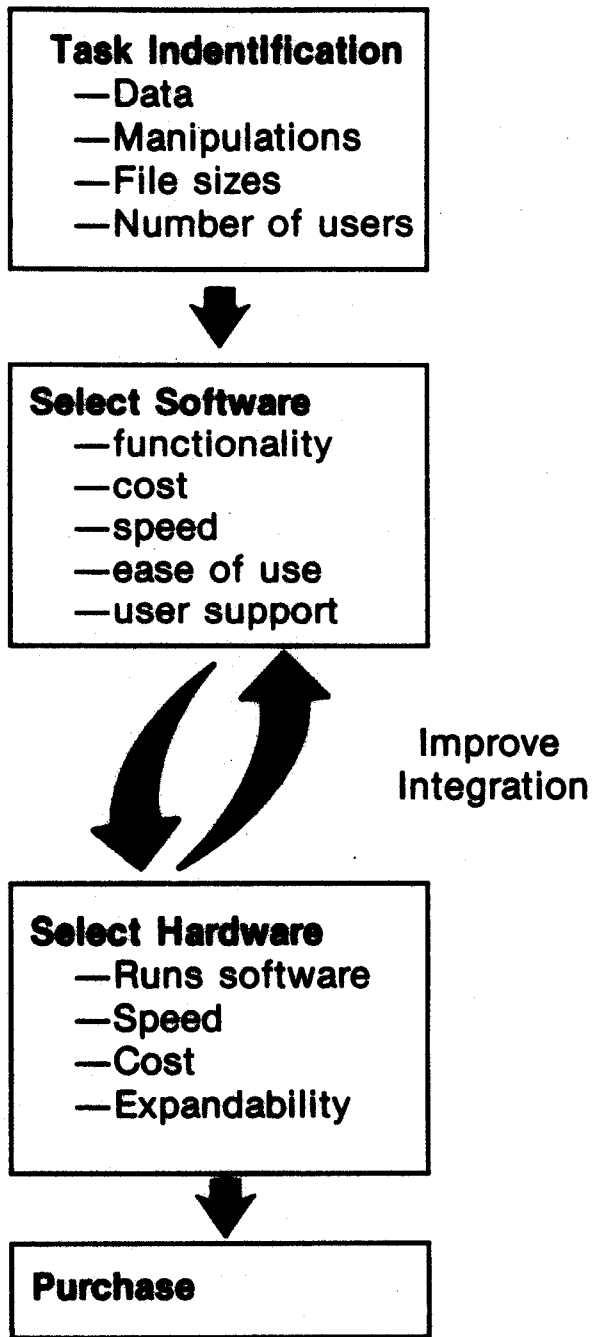


Figure 2. Selecting a computer system.

In selecting peripherals (e.g., digitizing tablets, scanners, video frame grabbers, plotters and film recorders) it is necessary to make sure that they are supported by both the hardware and the software vendor. In some cases the software vendor will bundle a computer and peripherals with the software. However, given educational discounts that are available to the research community (but not to the GIS vendor), it is often less expensive to purchase the computer and the peripherals separately.

NETWORKS

Connection to one or more networks can greatly enhance opportunities for scientific collaboration and help reduce the isolation that field stations often experience.

Local-Area Networks

Local-area networks (LAN's) can facilitate efficient use of computer resources by permitting sharing of data sets, software and computer peripherals, such as printers, disk drives and plotters (Figure 3). In a university environment, LAN's are often integrated with campus networks that are in turn connected to the wide-area networks.

A LAN can take two forms. In its most basic form, it links individual computers. Using Telnet (a program which allows you to log onto computers across a network) and FTP (a program which allows you to rapidly and accurately transfer files between computers), the LAN can be used as the avenue for accessing multiuser computers on the network or for transferring files between computers at speeds orders of magnitude faster than with a modem. In its more advanced form, "server" computers running networking software are added. This permits direct sharing of peripherals, programs and data in a way that is virtually transparent to the user. Most LAN programs support add-ons for electronic mail and automated backups. Sharing of disk drives across a LAN permits the sharing of data files (subject to security restrictions) and greatly facilitates keeping current backup copies of all data on the network.

LAN's can be used to eliminate redundancy of software and data at facilities with large numbers of individual computers. A single copy of a database or software product can be used by all the computers on the network, eliminating unnecessary duplication. Use of shared software and databases also reduces time spent installing updated software or modifying databases because only a single copy need be altered.

LAN's can take a variety of forms, but the most common consists of an Ethernet (a cabling system and electronic protocol capable of 10 Mb/s data transfer rates) running one or more types of networking software (e.g., NFS, 3Com, Appletalk, Novell, or TOPS). Network software for personal computers is usually designed so that a user (although not the network administrator) need know almost nothing about how a network operates. He or she simply operates as though

using his or her own stand-alone computer, but with the benefits of larger disk capacities, better backups and a larger variety of peripherals accessible through the LAN. An additional advantage of microcomputer LAN software is that it typically supports add-ons that make checking electronic mail as easy as turning on a computer.

Although some LAN software is specific to particular types of computers and operating systems, other software supports many different types of computers. This capability can be used to fully integrate different data management activities across computers. For example, on a network running TOPS (Transcendental Operating system) networking software, Macintosh, IBM-PC and UNIX computers can share data files regardless of which machine the data actually reside on.

Wide-Area Networks

In the past decade, the availability of personal computers has put data processing power on the desk tops and in the briefcases of many researchers. The availability of affordable data processing capability has led to a decentralization of research-related data processing. At the same time, there has been extensive growth of wide-area electronic networks that connect computers on a national and international scale.

Wide-area computer networks exist in a variety of forms with many different capabilities, including:

- Easy to access, reliable, and fast electronic mail
- Rapid and reliable transfer of text and graphics (e.g. proposals, manuscripts)
- Rapid and reliable long-distance data transfer
- Archival storage of data on distant university computers
- Better access to researchers at other institutions
- Access to mainframe computers
- Access to supercomputers
- Access to files, programs, printers and similar resources on other networks
- Access to national information and software repositories
- Access to mailing lists and mail forwarding systems

The most widely used network that supports all the functions listed above is the Internet (an association of high-speed, high-capacity, wide-area networks, including NSFnet, a network established and funded by the National Science Foundation). According to the pre-workshop survey, 22 percent of the stations who responded to the survey presently have access to the Internet. Connections to the Internet can take two basic forms. In the first form, a local area network and its computers are linked to a node on the Internet via a high-speed telephone connection. Such a link fully supports high-speed file transfers (depending on the number of network links traversed and amount of message traffic transfer speeds can range from 1,000 to 20,000 characters per second). In its second

File and Print Servers for the LAN

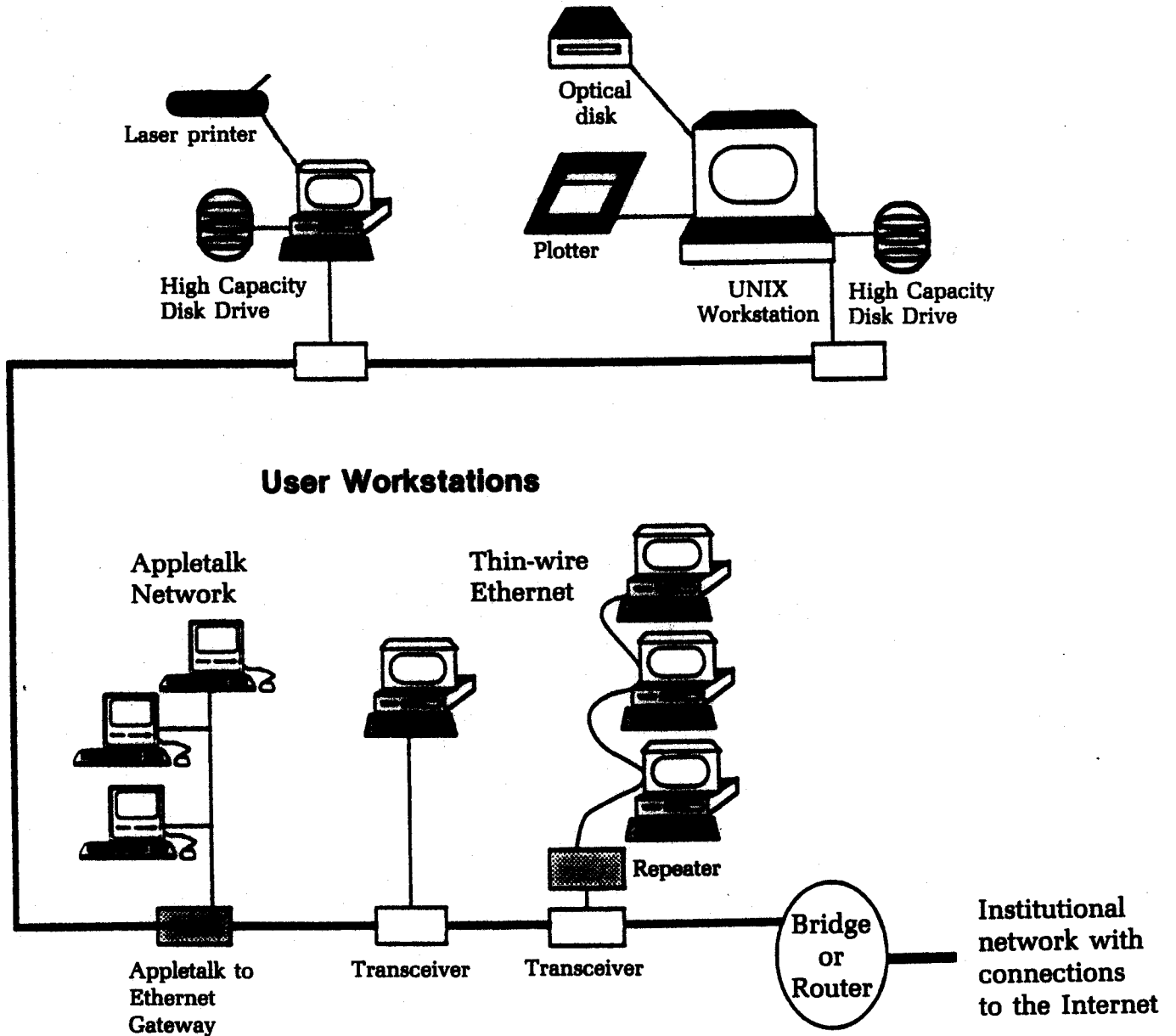


Figure 3. Typical LAN configuration

form, a modem is used to connect to a computer that is in turn attached to a LAN on the Internet (Figure 4). The speed of transfers is limited to the throughput capacity of the modem connection.

Commercial networks support subsets of the capabilities listed above. Typically these include electronic mail, access to bulletin boards and (limited) file transfer capabilities. Connections are made via modem and thus are limited in speed to the capacity of the modem. Thirty-eight percent of the field stations and marine labs surveyed support access to various electronic mail services (such as Bitnet, Omnet, and MCI mail). Unfortunately, sending mail between different electronic mail services is often difficult (electronic mail addresses become long and cumbersome) and occasionally impossible. Forty percent of the stations surveyed have no access to any kind of wide-area electronic network.

Establishing and maintaining network access entails installation costs, recurring costs (for operation and maintenance), and personnel time. Each of these costs varies widely depending on the network chosen and the location and facilities of the field station or marine lab. For example, installation costs for an Omnet account for an existing microcomputer might cost only \$300 (for modem, communications software, and Omnet fee), while installation of a full Internet connection might cost \$30,000 or more (for routing computers, cabling, network software, and network charges). An Omnet account is easily managed, whereas maintenance of a full Internet connection requires substantial time on the part of a networking expert.

A major concern for field stations is the recurring costs. The phone charge is a large part of such costs, because in many cases the remote location of the field station requires a long distance phone call or a dedicated phone line. Table 3 outlines approximate recurring costs for various network connections, each providing different levels of service. (Actual costs will be site-specific.)

The majority of the stations that do not have access to a full Internet connection through an existing institutional affiliation will find the cost of establishing their own full Internet connection outside the range of their budgets (particularly with regard to the high recurring costs). Provided that phone service is available, those stations can acquire an electronic mail box with one of the commercial mail services, the most common of which are Telemail, Omnet, MCI and Compuserve.

The process of deciding to which network a field station or marine laboratory should be connected must be guided by consideration of the needed capabilities and the cost of providing them. Identification of user needs is a critical first step. It does little good to provide a mail connection to one network when the majority of potential electronic mail correspondents are on another network. Similarly, a network that supports only limited data transfer capabilities is of little

use when large data sets are to be transferred. In selecting an electronic mail system it is important to find out what "gateways" exist for transferring messages to other networks and how difficult they are to use. The next step is to compare capabilities of individual networks to user requirements. This will yield one or more candidate networks for which cost analyses may be performed. A final network may then be selected.

The utility of electronic mail to the field biological community could be significantly enhanced by access to a mail forwarding system for field stations similar to the one recently implemented by the LTER network. This system solves several practical problems by:

- Creating simple, uniform network addresses for all users
- Sending and receiving group mailings
- Disseminating information on request by automatic reply
- Routing mail between different networks (acting as a mail gateway)
- Integrating mail and bulletin board services

This sort of service could perhaps be provided by LTER for a wider set of field stations, given appropriate funding.

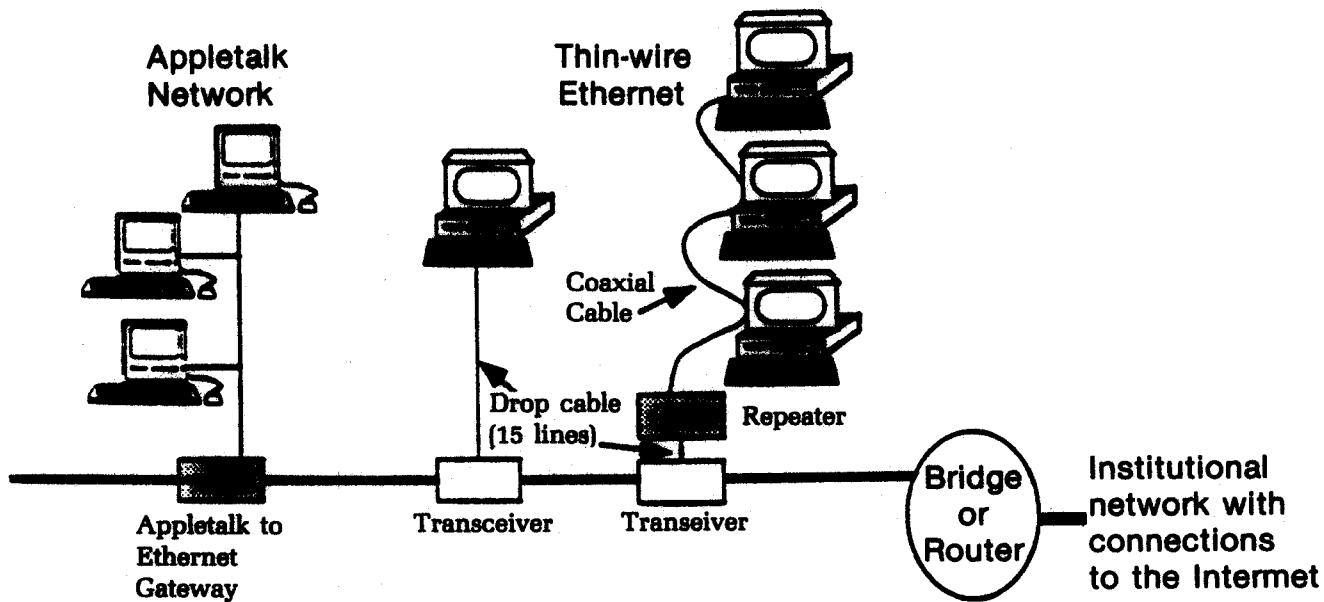
Archival Storage

A major objective of data management at field stations and marine labs is the archival storage of data. Data sets are prone to a variety of mishaps that can result in damage or loss. Data stored in printed form can deteriorate if exposed to excess moisture or heat, or spontaneously deteriorate if they are recorded on paper with a high acid content. Data stored on magnetic media can be lost because of equipment failures, power surges, extreme temperature fluctuations or simple deterioration of media over time. It is important to note that the "standard" magnetic tape or floppy disk has a recommended lifetime of only five years. It is said that there are no valuable data stored on 25 year old tapes simply because there are no readable data on 25 year old tapes.

A crucial component of archival data storage is making sure that backup copies of data are maintained. These should be kept current and stored in a location physically separated from the original data so that location specific calamities (e.g., floods, fires, hurricanes) are unlikely to damage all copies of the data. Off-site backups of data are facilitated by access to computer networks. By using transfers across a network, data can be backed up on another computer or device at a distant location without needing to transport physical media.

Data may also be lost through technical obsolescence. Optical storage devices are capable of retaining recorded data for many decades. Although optical storage reduces some of the problems associated with deterioration of media, access to data

A) Direct connections to LAN (ethernet)



B) Indirect connection to LAN via mainframe

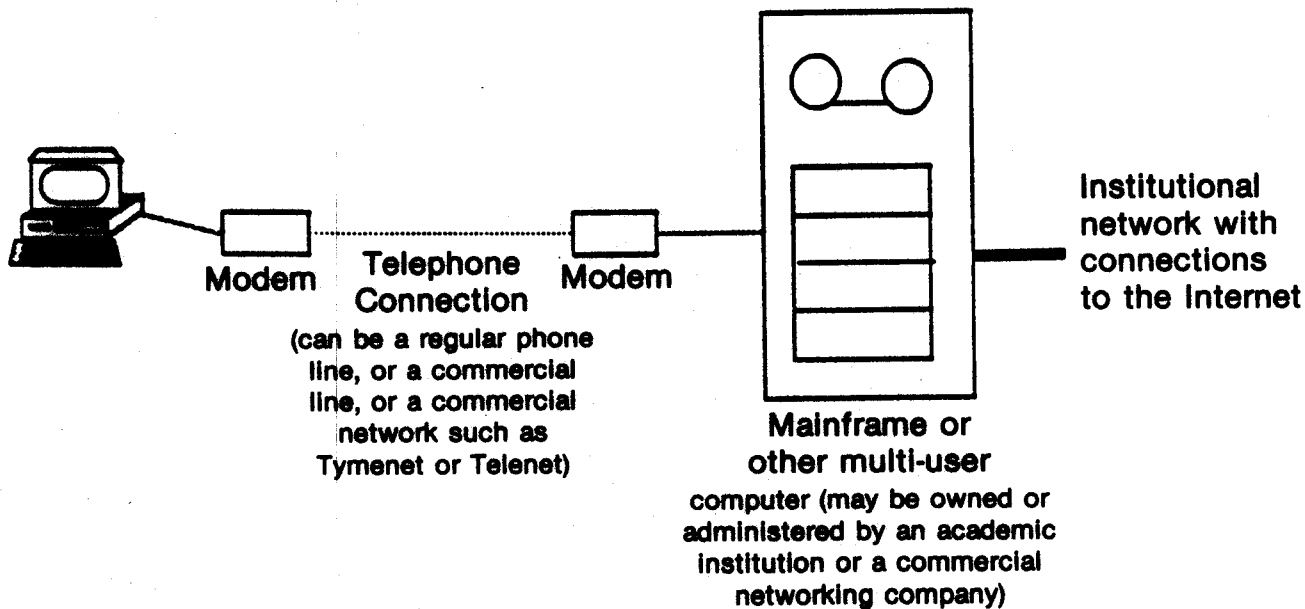


Figure 4. Connecting personal computers to the Internet.

Table 3: Approximate costs for network connections broken down by functionality. Units of cost are \$/s/y = dollars per site per year and \$/u/y = dollars per user per year.

Electronic mail only

Type of Connection	Recurring Cost	Cost of installing Network connection	Cost of on-site Equipment
Commercial (Telemail, Omnet, MCI)	600/u/y	(Phone connection)	2,000 (PC + modem)
Institutional (through a university, etc., Bitnet*, Internet mail, ...)		(Phone connection)	2,000 (PC + modem)

(*) Bitnet

Bitnet "membership" is no longer free. Organizations can join Bitnet for a fixed annual fee (\$750 - 10,000, depending on the size of the institution). This fee is usually passed through to individual users in the form of administrative costs (such as overhead costs). Bitnet is now operated together with CSNET by an organization called CREN.

Full Internet connections (electronic mail, file transfers, remote login capabilities), do-it-yourself (no network administration services provided)

Type of Connection	Recurring Cost (*)	Cost of installing Network connection	Cost of on-site Equipment
Internet direct	8,400/s/y	30,000	10,000 direct
Internet dial-up IP (SLIP)**	5,000/s/y	10,000	10,000

(*) recurring costs exclude the cost of on-site personnel

(**) IP refers to the Internet Protocols used to transfer data. SLIP is Serial-Line Internet Protocol and is a subset of IP that can operate over low-speed connections

may be lost due to rapid technological changes that render storage media obsolete and unreadable long before the end of its service life. Optical disk systems depend not only on the disk itself, but also on the disk drive which has a much shorter service life. Without a suitable drive to read it, a disk is useless even though it still retains the data.

Specialized data formats can also result in a loss of data through obsolescence, but of software rather than of hardware. Data stored in format that is readable only by a single software package can be lost if that package becomes unavailable. One way to avoid this problem is to store archival data in a simple standard format such as ASCII (American Standard Code for Information Interchange). Although this solution works well for numerical and character data, is not adequate to protect binary data, such as images. This is because binary data typically are stored in any one of a number of specialized formats. For such data it is crucial that documentation on the format be kept with the data.

Technological innovations

Computers and software are among the most rapidly changing technologies. Innovations having direct impacts on data management include automated data capture technologies, computer networks, improved data storage media, portable data entry systems and improved operating systems and software.

Automated data capture includes the use of optical scanners, image processing techniques, satellites and automated data loggers. Global positioning systems (GPS), which use radio signals from satellites to accurately calculate their current position, can also be used to automatically enter locational information.

Both local and wide-area computer networks will continue to increase in speed and can be used to integrate different types of computers into a single system. Additionally, groupware (sometimes taken to mean software which allows multiple individuals to see and edit the same data simultaneously) can facilitate collaborative projects over networked computers.

Improved data storage media take the forms of optical disks and digital audio tapes (DAT). Optical disks come in a variety of sizes and capacities, all of them capable of storing hundreds of millions of characters. DAT permits the storage of gigabytes of data on small magnetic cassettes and is extremely useful for making backup copies.

Portable data entry systems, in which a researcher types data directly into a small portable computer in the field, are increasingly popular and, when properly programmed, can help to reduce data errors by alerting researchers to apparent problems while they are still in a position to correct them.

Operating systems for microcomputers are becoming increasingly sophisticated and are converging with

those on mini- and mainframe computers. For the user, this will result in an increasingly "transparent" computing environment where it will not be possible to tell what type of computer is being used. Reusable computer program modules generated by object oriented languages can help to simplify data management tasks that require specialized attention. Despite significant increases in capabilities, software packages are growing increasingly easy to use. The trend towards sophisticated (yet easy to use) graphical user interfaces is playing an important role in this process.

Keeping pace with technology requires the obvious investment in hardware, but more importantly, it requires a significant commitment to personnel, planning, and training. The speed with which the technology used for data management evolves makes it difficult for individual data managers to keep abreast of potentially important developments. These difficulties can be ameliorated by enhanced communication among data managers. Several opportunities exist for facilitating exchanges on technology-related issues. A periodic newsletter addressing data management issues at field stations would disseminate information. An electronic bulletin board or electronic mail group lists could serve as a forum for exchanges of information between data managers and would permit a more rapid response to questions.

Facilities for Visiting Researchers

The ultimate purpose of research data management is to facilitate and improve research. For a data management system to be successful, it must be used by researchers. Field stations exist in a variety of settings and circumstances and with a diversity of missions. Computing equipment and services are as varied as are locations, and a visitor does not always have free access to facilities. Some stations provide no common-use software or hardware, whereas other stations provide aid in every facet of research activity, from data entry to data analysis.

In some cases computers and computer access by visitors to resident data bases are critical to the success of scientific investigations. Many stations have data bases that represent the only historical information available. In the absence of replicates, these data are the only way to validate many models. Expensive duplication of previous work can be avoided by identifying and using extant work.

The utility of a station or laboratory environment to visiting researchers can be enhanced by:

1. A common pool of hardware and software of a type that is currently in wide use (e.g., Word Perfect or SAS in the DOS environment), or that is very easy to learn (e.g., MacWrite or Delta Graph in the Macintosh environment).
2. A variety of materials to orient researchers to database facilities. Person-to-person interaction is the preferred mode for starting the educational

process. Electronic, video and audio material can be made available to help visitors learn more on their own in a self-paced mode. Teaching and demonstration programs that come with software are useful for this purpose. Short interactive tutorials can also be produced locally.

3. Access to electronic mail. Electronic mail is useful for both administrative and research purposes. It is widely used for pre-arrival arrangements, for access to data bases at other facilities, and to keep in touch with colleagues.

In contrast to the campus data management environment, seasonal field station users have a relatively short period of time in which to enter data. Where staffing and circumstances permit, a discussion of the proposed work between the researcher and the data manager can result in a data catalog entry that anticipates the integration of the data into the site database managed by the station. Assistance can extend to suggested data entry forms, quality assurance, portable computers and even appropriate analytic procedures.

CHAPTER IV—SUMMARY OF THE WORKSHOP SURVEY QUESTIONNAIRE AND PRE-WORKSHOP DEMONSTRATIONS

John B. Gorentz
W.K. Kellogg Biological Station
Michigan State University

and

Michael P. Hamilton
James San Jacinto Mountain Reserve
University of California, Davis

INTRODUCTION

Several months before the workshop, in November 1989, the workshop planners sent a survey questionnaire (Appendix F) to about 200 inland field stations and coastal marine laboratories. The purpose was to assess the state of data management at field stations and gather information to use for selecting representative sites to invite to the workshop.

On April 22, 1990, in a pre-workshop symposium, a series of demonstrations was presented in order to inform workshop attendees of some of the data management systems and technology in current use.

Purpose of the Questionnaire

The sponsors of the workshop (the Organization of Biological Field Stations and the Southern Association of Marine Laboratories), although having many common interests, are a diverse group with a correspondingly diverse array of data management systems. Their activities range from seasonal summer school sessions to year-round programs in research and education with large resident faculties and staffs. Their computer systems range from the personal computers of few individual investigators to networked systems and large mainframes. Their data management systems range from non-existent to just getting started to large systems with a separate staff and budget.

It was important to have representatives of the various types of field stations and their data management systems at the workshop, so some of the survey questions were designed to elicit information on relevant site characteristics. A balance was sought between experienced participants and those whose data management systems were in early stages or limited by modest resources.

In assessing the state of data management, we were less interested in examining the technology in use than in discovering what kinds of data sites have seen fit to manage, and how these data are being managed. Assuming that the best recommendations the workshop could produce would facilitate goals and objectives already adopted by

these sites, we asked questions designed to find out what was important to researchers. We also wanted to learn about common concerns and problems.

Background Issues

Although there is widespread agreement on many data management issues, perhaps much more so now than at the time of the 1982 workshop, there are also unresolved issues which influenced our choice of questions, and our interpretation of the responses. These background issues relate to the best use of limited resources, acceptable degrees of centralization, and the relative importance of technology vs. human resources.

Use of the term "data management" is usually accompanied by some unstated presuppositions. For some people, data management is whatever must be done with data, usually using computers, in order to analyze them for publication. Another view, perhaps less common now than at the time of the 1982 workshop, is that data management includes almost anything that has to do with computers and technology. For still others, data management means caring for certain data so that, whatever their original purpose, they are preserved and made available for more general use, now or in the future. This latter view was the premise of the 1990 workshop.

This workshop was based on the assumption that at field stations and marine laboratories there are historical data records worth preserving to enhance the value of the habitats for research, to provide background data, and to make long term studies possible. Some of these data sets are gathered for general use, others are the fortuitous by-product of specific research.

Without proper care, these data resources will be lost. This care entails a cost, and although there is widespread agreement that efforts to preserve data are worthwhile, there is not universal agreement that already scarce resources should be spent on data management.

Nor is it certain that the data sets compiled or otherwise preserved for general use have been used, or will be used, to advance science. Science builds upon previous work, including that represented in previous databases, and scientists have a responsibility to preserve data for those who will follow after them. There is, however, some disagreement over whether resources should be spent testing hypotheses rather than on preserving data without a clear hypothesis to be tested.

Even those who maintain that data need to be managed as a research resource will acknowledge that historical data are not yet being utilized as they could be. There are several possible barriers limiting the availability and use of existing data:

- The existence of these data sets is not commonly known. The scientific literature may serve as a partial index, but additional means are needed to make these data sets known.
- Physical access to data is difficult. Better means of electronic communication would make data sets more widely used.
- Some data sets are known to exist, and are accessible, but are too poorly documented to be useful. Better systems of documentation are needed.
- Data sets are sometimes not worth bringing together for comparison because they are too dissimilar in format, representation of data, methods, and meaning. Standards are sometimes proposed to resolve these problems, but there may be resistance to standards as being too restrictive for open-ended inquiry.

These barriers are not mutually exclusive, but disagreement as to their relative importance leads to disagreement over funding priorities.

Many data management solutions assume a certain amount of centralization. Long term care of data often implies a responsibility for data that goes beyond that of an individual investigator, and this in turn implies a degree of relinquishing control which is sometimes in conflict with the basic tradition of independent inquiry.

Perhaps most controversial is the issue of standards. A lack of standards makes comparative analysis of data sets from different sources very difficult. Researchers currently use and benefit from standards at many levels (e.g. ISO units of measurement). But there is a fear that any push toward standardized computer systems or data formats at any level will be too restrictive or too unwieldy and will interfere with research.

Finally, there is a question as to whether the greatest data management need at this time is for technology or for human resources. It is often easier to obtain funding for computers and software than for personnel, but it is possible that at the current stage of development, personnel are the limiting factor.

SURVEY METHODS

In November, 1989, questionnaires were mailed to about 200 sites that constituted the membership of the Organization of Biological Field Stations and the Southern Association of Marine Laboratories. Questionnaires were also sent to additional sites in the Long Term Ecological Research program, and to several National Marine Laboratories.

103 responses were received. Several of them were received too late to be used in selecting invitees to the workshop, but those responses are included in the results presented here.

In the questions, we did not ask about facilities and technology so much as about goals, priorities, and personnel. We tried to get the respondents to distinguish between institutional operations and those of individual research programs. Some respondents were more sensitive to the distinction than others.

The questions were open ended, because we felt that the most useful information would not necessarily fit into neat categories. In analyzing the responses, we did attempt to categorize the responses, and in the process made subjective judgments. Even though some information is presented quantitatively, the tallies were a matter of considerable interpretation on our part.

WHAT DATA ARE BEING MANAGED?

The first set of questions (Question 3a-3c) was designed to find out about the data that sites are managing as a general resource, or which could be made more available if resources permitted.

In doing this, we wanted to distinguish between those databases managed as part of a single research project, and those that are being managed for long term general use as a site responsibility. We also wanted the respondents to take a broad view of the term "database," including those data managed without sophisticated database tools or without computers, as well as non-traditional forms of data such as audio and video recordings.

The hope was that the responses to these questions would help define the subject-matter of the survey and workshop and give some idea of sites' goals and priorities.

We asked three questions about databases.

- 3a. Does your site have databases that have been compiled specifically for general use (e.g. species lists, meteorological data)? If so, please list some examples.**
- 3b. Does your site have databases originating in individual research programs, that are or could be developed into general use resources. If so, give a few examples.**
- 3c. Does your site have computerized records consisting of non-traditional forms of data, e.g. acoustic records, maps, visual images.**

In response to question 3a, 90 percent of the respondents listed one or more databases; only 10 sites said they did not manage any general use databases at all. We categorized the responses into a few arbitrary, non-discreet, non-orthogonal groups, which we tallied as follows:

Climate data: A large, clear cut category was climate databases, with a little over half of the 103 respondents listing this among their general use databases. These databases appear to include everything from records kept on paper, to automated data collection systems and electronically networked databases.

Species lists: Almost as many sites listed one or more types of species list among their general use databases. They were variously described as species lists, species inventories, and species checklists. They were usually compiled for specific taxonomic groups, such as birds, mammals, vascular plants. Two sites pointed out that they have developed taxonomic keys for their lists. One site has made its species lists a key part of its data management system, by setting up a system of standard species codes to be used in data sets.

Hydrography and hydrology: Twenty-seven sites listed one or more types of hydrography or hydrology database. This category includes databases variously described as hydrological measurements, tide measurements, stream flow, water quality, water level data, bathymetry, sedimentology, physical and chemical limnology, groundwater levels, pond levels, sea levels, hydrology, stage/discharge data, seawater temperature, salinity, stream chemistry. (We did not include precipitation databases in this category.) These databases range in scope from (for example) a modest database of stream level measurements, to a large scale information center for Galveston Bay.

Bibliographies and project lists: Sixteen sites listed some sort of bibliography or project list among their general use databases. These appear to range from simple lists to elaborately indexed computer databases. Some of them appear to be stand-alone databases. Others appear to be integrated into a larger system, serving as an index or access point to the site's other data resources. That is, a data user can search the database for a particular subject or organism, and find not only the pertinent literature, but be directed to other databases as well. Four respondents described the contents and organization in more detail. Their databases are sorted or indexed by one or more of the following categories of information: author, location, research topic, funding source, species, sampling dates, and keywords.

Miscellaneous long term monitoring: Almost all of the databases listed by the respondents could be categorized as containing long term monitoring

records. But we also noted some miscellaneous other types of long term databases, listed by 18 sites, which do not fit the above categories. They include records on flowering phenology, secondary succession, fish capture, sequences of plant surveys, vegetation on permanent plots, annual bird counts, bird migration, nesting, land use history, photo monitoring, and fire history. According to the responses to question 3c, a few sites keep databases of 35 mm slide photographs taken on a calendar schedule.

Maps and geographic data: Because of the current high level of interest in geographic information systems, we created a separate category for geographic data and map-type data. By including some responses to question 3c about non-traditional forms of data, we counted 32 sites that either have GIS systems, or have map-type databases now represented on hardcopy maps and aerial photos that could utilize GIS software or other spatial data management systems. It could be argued that this category is the largest of all, if one considers that all data that reference particular spatial locations on the earth's surface are potential GIS data. Most of the general interest data at field stations and marine labs fit this description.

In response to question 3b, which asked whether there are other data that could potentially be managed as a general resource, 76 respondents said yes, and 74 gave examples. These examples consisted of additional long term records of the types listed above, as well as point-in-time data sets. These include those resulting from (for example) three year projects, but are distinguished from continuous long term projects.

Based on these responses, it can be seen that most field stations and marine labs are in the business of data management. Even among those 10 sites that said they do not manage any general use databases, some plan to do so soon. These include field stations which are relatively new, or which have only recently adopted any data management goals.

However, among the ten are sites that have no intention of managing general use databases. These sites responded that their databases were all investigator-specific, and/or they do not think managing general use databases is an appropriate undertaking for their sites. In fact, some questioned the validity site-sponsored data management. This issue is discussed further under "Site Self-Evaluations and Recommendations."

It could be argued that our tallies under-represent the number of sites managing general use databases and the number of databases. The questionnaire asked only for examples, not a complete list. Although we intended the term database to be used in a general sense, not just applying to those data being managed in some formal DBMS, it is possible that some sites did not consider their more casual,

uncomputerized collections of data to be databases, and omitted them. But in general, it appears that the respondents used a broad definition of the term database, as we had intended.

It is more likely that the survey overestimated the amount of data management. The 50 percent of sites responding to the survey probably represented sites more active in data management than the 50 percent that did not respond. Field stations naturally want to show themselves in the best light, and some may have included investigator-specific databases in the general use category. But most respondents appear to have understood exactly what was meant by general use, and made the distinction that we wanted them to make. We suspect that many data sets which are not sufficiently managed to be readily accessible and useable were also included in the count. The databases listed should not be considered accomplishments requiring no further attention or resources.

GOALS, PRIORITIES, AND ACCOMPLISHMENTS

On the assumption that a useful set of recommendations for field stations must be consistent with existing goals, we asked in Question 6 about priorities and accomplishments as they relate to six specific objectives.

The six objectives are listed in Table 6. We asked sites to describe these by circling one or more of the following status codes for each:

- ACF—An accomplished fact
- HPG—High priority goal
- MPG—Medium priority goal
- LPG—Low priority goal
- WIP—Work is in progress
- SKF—Seeking funding
- NPL—No plans to do this

The respondents were invited to choose more than one, if necessary. Some sites circled both ACF and HPG, for example. This is a plausible answer for an objective on which an important phase has been completed even though more remains to be done. Likewise, "work in progress" and "seeking funding" often went together. One respondent circled both his personal goals and the institutional goals, his own being more ambitious.

We calculated a weighted positive score for each of the six objectives, using the formula $2 \text{ ACF} + 2 \text{ HPG} + \text{MPG} + 2 \text{ WIP} + \text{SKF}$. These scores are shown in the final column of Table 6. No attempt was made to adjust for the fact that those sites that checked more than one item are represented more heavily in this score.

Manage Selected Databases

The goal of managing selected databases had the highest overall score and the lowest negative ratings (i.e., the lowest NPL tally.) This is consis-

tent with the responses to question 3, which showed that the great majority of sites are engaged in managing some general purpose databases. (It is not certain why fourteen respondents said they had no plans to do this, when responses to question 3 when the responses to question 3 indicated that only ten percent do not manage any databases.)

In a sense, this objective is more modest than any of the others, in that it deals with only a portion of a site's data; all the others are more comprehensive, dealing with all general use data. The widespread adoption of this objective may mean sites are taking the route recommended in Chapters 1 and 3, of starting small and doing one thing at a time. Or perhaps it is unrealistic, with limited resources, for them to consider doing otherwise.

Central Catalog or Directory

The objective of implementing a central catalog or directory also scored high. More respondents indicated work in progress on this goal than on any of the others. It received many positive responses, but also a significant number of negatives (i.e. NPLs). In general, though, the responses seem to indicate that many sites believe an important first step in managing data is to take inventory.

On-line Catalog or Directory

This objective deals with a high-tech version of a data catalog, in contrast to the catalogs discussed in 6a, which are not necessarily even computerized. It concerns on-line retrieval of information about data, which presumes a high degree of automation and electronic accessibility. This high-tech objective did not score nearly as high as its low-tech counterpart. Possibly the respondents felt that, although access to data sets is important, the volume of queries does not require a high-tech, highly interactive catalog. Instant availability may not be so important. One can also infer that completeness of the catalog is more important than high-tech access, since the choices represent a real tradeoff in expenditure of human resources.

Archive or Repository

The responses to item 6f, regarding implementation of an archive or repository for all historical data on natural habitats, are somewhat puzzling. This item scored about evenly with item 6a (central catalogs), even though it is presumably much more challenging and expensive. It would seem that developing a directory to data would be only a subset of the task of implementing a complete archive.

The responses to this item are inconsistent with most of the others in the survey; the more modest tasks usually scored higher than the more elaborate ones. It is possible respondents misunderstood the question or the task, or had another objective in mind.

Table 6. Responses to Questions 6A-6F. The number of respondents selecting each status code for each objective, and a weighted total of the positive selections is shown. The status codes are explained in the text.

Objective	ACF	HPG	MP- G	LPG	WIP	SKF	NPL	Weighted
								Total
6a. Implement a central catalog or directory describing all data sets on natural habitats (i.e. data about data, computerized or not).	10	26	21	10	32	8	24	165
6b. Implement a central catalog or directory of data about data that is electronically searchable, "on-line."	9	17	18	11	21	8	37	120
6c. Manage selected databases for general use as a site/institutional responsibility.	34	24	14	8	23	8	14	184
6d. Implement a standard format for all research data.	11	6	17	16	14	3	46	82
6e. Manage working copies of all data in a unified, on-line database. (This does not necessarily mean a "centralized" database.)	8	16	17	16	18	3	34	104
6f. Implement an archive or repository for all historical data on natural habitats.	16	24	22	12	29	6	18	166

Unified, On-line Databases

Item 6e, regarding a unified on-line database, scored relatively low, and was the item least often cited as an accomplished fact. In a sense, this objective is a high-tech version of that in item 6f. As with data catalogs, the high-tech version was deemed less important.

Standard Format

The low-tech choices did not always score higher than their high-tech counterparts. Item 6d, regarding a standard format for all data, could be regarded as a low-tech subset of item 6e, a unified on-line database. But this item received the highest number of negative responses. Nearly half the respondents said they had no plans to implement it, and only six said it was a high priority goal.

This is not particularly surprising, in that it is a daunting objective. Also, there is a great deal of resistance among researchers to standards on anything that would constrain their work. What is surprising is that this item scored lower than 6e,

since a unified on-line database would seem to imply a high degree of standardization as well. It is interesting to speculate whether including the word "standards" in question 6e, without changing the essential meaning of the objective, would have resulted in a more negative reaction.

Funding

The number of respondents who are seeking funding was low for all objectives. This could be taken to mean that additional funding is not greatly needed and that objectives are being met well with existing resources. However, the question did not ask whether existing resources are sufficient, and the responses could also reflect the fact that very few sources of funds are available.

PROPRIETARY RIGHTS

Those who have been given the task of making data more accessible sooner or later run into the issue of proprietary rights. In principle, scientists are willing to share data. Science is based on the free and open exchange of information, whereby

people can build on the work and data of those who have gone before them. Scientists are often engaged in breaking down technical and political barriers that limit collaboration with others. But scientists also guard their data in order to ensure proper recognition of their work through the publication process.

Many field stations are involved in data management on the assumption that data sharing through the traditional system of publications is not adequate, and that there are unpublished data, never-to-be-published data, and raw data behind publications that need to be made available as a resource for others. They therefore need to reconcile legitimate proprietary rights with the goal of greater accessibility.

In question 5b we asked,

"How do you weigh investigators' proprietary rights to data against the goal of wider availability? Is there security against unauthorized use of data?"

Ninety-four of the 103 respondents addressed the question. The responses represented two fundamentally different attitudes. Many sites view proprietary rights as a necessary evil, while others perceive the protection of proprietary rights as an important objective. This was perhaps expressed most strongly by a site which reported, "Proprietary rights are protected to the fullest."

Twenty sites reported that proprietary rights are not an issue, or at least are not yet. Some reasons given were that proprietary data are not involved, or that there is not a centralized system. Four reported that they have no policies yet, but that it is an issue that needs addressing and is being addressed.

Thirty sites reported that they have no policy, that the issue is left to the investigator, and that the data can be accessed only through the investigator anyway. One of these respondents simply said, "Data is uninterpretable to non-investigators." This is probably a common state of affairs. Some reported that data are indeed shared by these bilateral arrangements. Six sites had systems of central access to data, but left the issue of outside access up to the investigator. In some cases the investigators exercise control by deciding whether or not their data are to be added to the central database. At others sites they exercise control over data residing in a central database through a security and authorization mechanism.

Eleven sites reported some sort of policy to limit proprietary rights, but did not have a centralized database. Most commonly the policy consisted of a limit on the time during which investigators have complete control of data; after this time they are required to make their data more accessible. These policies were usually related to site-use requirements and responsibilities for visiting re-

searchers. It would be interesting to know how effectively these policies are enforced, or whether any enforcement mechanisms are necessary.

Seven sites had policies in place that emphasized security and confidentiality, while complying with any regulations regarding open access. These tended to be some of the larger marine labs with highly centralized systems, at which researchers' rights are subservient to other purposes. These sites reflected a strong sense of ownership of and responsibility for data, with policies in place and mechanisms to enforce them.

Three sites emphasized the protection, rather than limitation, of proprietary rights.

Thirteen sites emphasized central, general purpose databases that are open to all. However, the data they contained may not have included much that was investigator-specific.

ADMINISTRATION AND PERSONNEL

Administrative and personnel factors are possibly more limiting to progress in data management than are technology and equipment. We wanted to explore the magnitude of data management tasks by determining the level of commitment to data management, including personnel resources committed.

We also wanted to learn the degree to which a field station's data management is a distinct activity, distinguished from related areas such as computer management or investigator-specific data management. We assumed that clear data management goals would be reflected in distinct data management budgets and personnel assignments.

We asked the following four questions:

- 4a. **Where does the impetus for data management arise (e.g. site administrators, interested faculty members, research programs, technical staff)?**
- 4b. **Does your site have a data manager, or other person(s) with designated responsibility for data management?**
- 4c. **What personnel are involved in data management (number of persons, positions, training, experience, fraction of time)?**
- 4d. **How is data management funded? Is there a specific budget for data management? Is it funded at the site/institution level, or on individual grants?**

Impetus for Data Management

We asked the first question, about who is pushing data management, to determine the extent to which it is research driven or technology driven. We also wanted to learn whether there was top-level administrative commitment.

We grouped the responses into categories and tallied them as follows: Researchers (68 sites), Administration (61), Technical staff (21), Long Term Ecological Research Program—LTER (9), and External (5). Two sites did not respond to this question. Many sites fit into more than one of these categories.

It does appear from the responses that data management is largely research driven. Two thirds of the sites cited researchers as the driving force, and almost as many said that administrators, who presumably have research interests foremost, were the impetus. Of course, a researcher or administrator can be overly enamored of technology for its own sake, but presumably most have not fallen into that trap. Of the 21 sites listing the technical staff as a driving force, only one listed it as the only group leading data management, and sixteen of those 21 also listed researchers.

The responses are probably a good sign, indicating that sites have their priorities in order. Research is driving data management, rather than vice versa. Data management has a supporting role, albeit an important one. As such it is not likely to take on a life of its own, unresponsive to research needs.

The number of sites listing site administrators as a driving force indicates that at a majority of sites, there is active support at the top level, and not just passive tolerance.

The five sites that indicated an external impetus were mostly government labs whose supervising agencies mandate their data management activities. The nine sites that cited the LTER program also are responding to an external impetus.

Designated Data Manager

It may be that people rather than technology are the limiting factor in successful data management. But money to fund personnel is often harder to come by than money to purchase equipment. Before making recommendations on personnel for data management, we needed a clear picture of the current personnel situation.

In response to question 4b, about a designated data manager, 41 sites reported having none. Twenty-three have a full time data manager. Thirty-six have someone doing data management part-time, including six sites at which the director is the data manager, four at which the data manager is the librarian, and two at which a laboratory manager performs this function. Two sites were unclear in their responses and are not included in the tallies.

Several of the 23 sites with full-time data managers reported that other computer-related

duties besides data management are included in the manager's workload. If the issue is data management in a strict sense, the count of 23 full-time data managers is misleadingly high. Even many of the 36 part-time data managers also do computer management as well, with data management getting a fraction of the person's attention.

The six sites with site administrators serving as data managers are generally small sites with appropriately modest goals. The four sites with a librarian-data manager suggest a route for sites to follow when it is not desirable or necessary to develop a computer and data management infrastructure: data management can be made an adjunct to library rather than computer operations.

Staff Qualifications and Background

In response to question 4c about the number of persons involved in data management and their backgrounds, one site stated, "too irregular to tabulate." This telling comment is a good summary of the overall situation. Even though most respondents did attempt to provide numbers and descriptions of those in data management, the responses taken as a whole were too irregular for us to tabulate.

This is partly because of a confusion between data management and computer management. The two types of work are often confounded, and even where distinct, are often done by the same persons. Given this situation, we could not tell which qualifications listed by the respondents were relevant to data management.

Another barrier to tabulation was the lack of comparable functions for the data management portion of the work. Combinations of staff, duties, organization, and infrastructure varied greatly. Data management is done by site administrators, faculty members, graduate students, secretaries, statisticians, librarians, and sometimes even by specially designated data managers. Various "coordinator" positions (e.g. research coordinator, scientific coordinator, site coordinator, data coordinator) have responsibility for data management among their duties. Educational levels of data managers range from high school degrees to the Ph.D. level, with many in between.

Data management is commonly done by people who are self-taught. A few data managers have backgrounds in computer science, but data managers with backgrounds and training specifically in data management are perhaps non-existent. Those whose training is primarily in computer science are uncommon.

It is not possible to determine from the responses which personnel configurations are the most successful.

Table 7. Cross-tabulation of responses to Question 4D. Responses regarding a specific budget for data management are arranged horizontally, and those regarding funding at the site/institutional level are arranged vertically.

		Specific Budget?		
		No	Yes	Total
Funding at site/institutional level?	No	48	5	53
	Yes	37	12	49
	Total	85	17	102

Data Management Funding

We asked question 4d about funding to evaluate sites' commitment to data management. We wanted to know whether data management per se is an objective distinct enough to have its own budget (whether it gets at least some funding at the institutional level, rather than exclusively from individual grants), so we could judge whether it is getting support from the top institutional level.

Of 102 sites responding to this question, only 17 had a specific data management budget. However, nearly half (49 of 102) did have at least some funding from their institution. The responses on these two issues of a specific budget and of institutional support are cross-tabulated in Table 7. But even among the 17 sites with a specific data management budget, in many cases the budget appears to be more of a computer budget than a data management budget.

Similarly, the level of institutional support is probably not as high as it might seem from the raw numbers. In some cases, the amount of support is small, as small as a bit of funding for a weather station. The fact that only one fourth of the sites with some institutional support have a specific budget is an indication that such support does not involve serious money.

SITE SELF-EVALUATIONS AND RECOMMENDATIONS

To sum up, we asked sites to evaluate their accomplishments, resources, and needs. We asked a series of five questions, starting with:

8a. What have been your most important data management accomplishments?

The data management accomplishments that the respondents listed fell mostly into three categories: data, computer systems, and administration.

The data-oriented accomplishments included: assembling historical data sets (cited by 7 sites); setting up continuous, long-term databases (3 sites); other computerization of large databases, with em-

phasis more on the data than on computerization (5); establishing systems of baseline, site-characterization, or geographic data (16); the development of site bibliographies (7); and specimen databases (3 sites).

The administrative accomplishments included: coming to grips with the need and identifying the problem (cited by 4 sites); developing an overall plan (2 sites); getting started (2); getting organized (8); establishing policies regarding the responsibilities of investigators (3); compiling data catalogs and indexes (14); a methods manual (1); establishing data archives (2); establishing standards for data entry, documentation, and format (8); establishing data quality protocols (3); development of a data management staff (4); obtaining high level support for data management (1); getting funds (5); and establishing a training program (1 site).

The computer-oriented accomplishments included: implementing database management and geographic information system software (cited by 6 sites); developing database management software (2 sites); and data entry systems (2). Five sites developed computerized databases, with an emphasis more on the computer systems than on the data. New or improved computer systems, including storage systems and networks, were cited by 13 sites. While three of these sites decentralized their computer systems, moving from mainframes to microcomputers, one site centralized its database system. Five sites installed instrumentation for automated data acquisition.

8b. What things would you now do differently, if you had them to do over? What suggestions would you give to other sites?

Not all sites responded to the above question, and some of those who did stated that they were not far enough along to answer it. But those who responded listed the following types of items:

- Take time to plan, instead of just letting things happen.
- Implement policies regarding researchers' responsibilities.

- Make sure of researchers' support, and involve them in oversight.
- Get organized sooner; catching up is hard to do.
- Start baseline data collection sooner.
- Link all data sets by location.
- Do quality control.
- Set and enforce standards to ensure consistency; set standards earlier in the process.
- Keep it simple; do not try to do everything at once; do one data set at a time.
- Spend more time on documentation of everything.
- Consult with outside experts.
- Provide training.
- Avoid mainframes.
- Use networks to keep decentralization from going too far.
- Buy commercial database software rather than developing it in-house.
- Use relational database software technology.

The last three questions of the series asked about resources for data management:

8c. What personnel resources do you think are needed to meet your data management goals? Are these resources available?

Some sites said they had adequate personnel resources. These were mostly sites that apparently had just recently received funding for new positions. Those who stated a need for additional personnel listed everything from data entry personnel to skilled professionals. Many sites with no data manager stated the need for a part-time data manager "dedicated to data management." Some who had part-time data managers emphasized the need for a full-time person. Some that had a full-time person needed more persons.

Although some sites lacked highly-trained personnel with specialized technical skills, more of them pointed to the sheer amount of time rather than skill that was needed to do the work.

A few respondents also emphasized the need for researchers to take part in data management or exercise an oversight role, or to take an interest in the sometimes mundane gathering of baseline data.

8d. What additional facilities crucial to your goals (hardware, software, etc.) are lacking?

The following were listed:

- Data collecting instrumentation, e.g., for data loggers
- Computers and computer equipment
- New or upgraded mainframe computers

- More, better, or upgraded microcomputers dedicated to data management
- Computer systems or upgrades for GIS
- Computer systems and equipment for video analysis
- Database management software
- Personnel
- Bricks and mortar, e.g., office space, physical storage space
- Local area networks, network upgrades, or communications, including links from remote sites to university campuses
- Equipment located at the field sites to reduce the need to use equipment at distant campuses
- Equipment for long-term, reliable archives

8e. Where do you think additional funding is most needed?

This question was intended to elicit the most important priorities among all the items mentioned. The need for personnel topped the list of concerns, with 52 respondents citing it, as opposed to 21 who listed computers and hardware, and 13 who listed software. The raw count understates the strong emphasis that was placed on personnel, as well as on the strong concern, expressed by 13 sites, for the stability of long-term funding for recurring costs for personnel and for the maintenance of computers, software, and data. Other needs were buildings (cited by 1 site), instrumentation (3 sites), computer network links (2), and training (2 sites).

It should also be noted that a few persons stated here and in their additional comments that their top priorities were outside the realm of data management. Some were frankly skeptical about the feasibility of managing data for general use, or the appropriateness of diverting research resources to data management, preferring to focus attention on the immediate needs of individual researchers.

Some of the skeptical comments were as follows:

"To do it right at each lab might have prohibitive costs."

"Given the extremely diverse nature of the research and the individual approach (30-50 basically unrelated research projects/yr)...I have serious questions about the potential utility of centralized data bases."

"...We often wish we had much better base line data, but given our mission, it would be difficult to justify the diversion of resources from other goals."

"...A useful topic for...discussion might be 'How do we maximize the benefit, or judge the eventual benefit, of data we collect now for future use?'"

"...I have yet to see a data mgt. system (for ecological labs) that really worked and was actually used by scientists publishing papers based on the data..."

"What is the purpose? Most of our researchers believe that maintaining long term records without specific research goals is a waste of resources. Once they answer a question their data is useless and just takes up space in a filing cabinet (after publication)."

"...keep it basic...let the researcher who wants the data do all the work."

CONCLUSIONS AND SPECULATIONS

Diversity

Although it might seem almost too obvious to mention, one of the most significant characteristics of field stations and marine laboratories is their diversity. They do have some common interests and objectives, but there is such a myriad of missions, institutional arrangements, facilities, and types of data that great care is needed in developing standards, guidelines, and recommendations for them. It is important to analyze every assumption and conclusion from the perspectives of the full range of sites. Unlike other scientific disciplines in which, for example, the issue is how to deal with huge quantities of satellite imagery, the challenge for field stations is in dealing with the great diversity of data.

Long Term Data

Most of the data that need to be made available for wider use are long term records. Managing these data requires a sustained effort that is not likely to be funded by project-specific grants.

Descriptive Data

Sites' initial efforts at data management are in the area of descriptive data, such as climatological data and species lists, rather than in the area of experimental data. Possibly this is because organized data management is like many scientific disciplines, which need to start with descriptive work before moving on to the experimental. An alternate explanation is that the main purpose of long term data management is to provide descriptive background data which can serve as a context for experimental studies, and that this will always be the focus.

Access to Data

In addition to managing descriptive data, many sites have in the past decade embarked on the development of catalogs and directories to data,

sometimes in the form of publication lists. These ventures will not only serve to make data more accessible to others, but help sites take inventory and evaluate priorities.

Dissenting Views

Although there are those who question its value and feasibility, most of the respondents took a positive view of the necessity and possibilities of data management, as indicated by their accomplishments, plans, and commitments. But it would be good for those who are committed to data management to keep the skeptics' comments in mind, because they lay bare the criteria by which data management should and will be evaluated.

Commitment to Data Management

Using the number of sites managing general use databases and those developing access mechanisms as a measure, it would appear there is great enthusiasm for data management. But judging from personnel, budgets, and other comments, data management might seem an indistinct activity, commonly confounded with computer management and short term exigencies. However, the situation has greatly improved since the time of the 1982 workshop. The survey results show a much greater agreement and understanding of the possibilities and needs than would have been found earlier.

PRE-WORKSHOP DEMONSTRATIONS AND PRESENTATIONS

By way of information and introduction to the major concept of the workshop, a day long pre-session symposium was held on April 22, 1990, highlighting examples and demonstrations of data management systems by 20 of the workshop's participants. Ten demonstrations of laptop, PC and Macintosh-based systems were presented, and discussions of other station-based capabilities were described.

Demonstrations had been pre-selected to provide a sample of diverse approaches in use at marine and inland field stations as of April 1990. Examples of the following categories of data management were demonstrated:

- field entry of data using portable computers
- automated acquisition of environmental data
- geographic information systems (microcomputers and workstations)
- relational database management for research project management
- microcomputer access to large SQL relational database
- hypertext (Hypercard) and interaction multimedia databases

- multimedia networking over Internet
- access to databases over networks using electronic mail

The participants listed below provided informal overviews of the status of their stations' data management approaches:

- John Briggs, Konza Prairie, Oracle SQL database demo
- Vic Chow, Bodega Marine Lab, MOMS, Paradox demo
- Steve McNeil, UC NRS, FileVision Database/GIS demo
- Robert Moeller, Pocono Comparative Lakes, Reflex, Paradox demo
- Jim Brunt, Sevilleta LTER, Overview of programs
- Mike Hamilton, UC James Reserve, Hypermedia GIS demo
- Paul Montagna, Marine Science Institute, Overview of program
- Grady Cantrell, Hancock Biological Station, Overview of program, Dbase III
- Fred Lohrer, Archbold Biological Station, Overview of program
- Lance Risley, Institute of Marine and Coastal Sciences, Overview of program
- Rudolph Nottrott, LTER Network Office, "ANDREW" Internet System
- Warren Brigham, Illinois Natural History Survey, Overview of GIS applications
- Bill Seitz, Texas A&M, Galveston Bay, Macintosh-based demo
- David Nebert, Institute of Marine Science, Overview of programs
- Craig Staude, Friday Harbor Labs, Mac-based demo
- Deborah Clark, La Selva Biological Station /OTS, Overview of programs
- John Heuer, Savanna River Ecology Lab, PROGRESS Database Language
- Bill Michener, Baruch Institute, Easy Entry demo
- John Porter, Virginia Coast LTER, database entry demo
- Jim Beach, Michigan State University, Demonstration of network for herbarium label data exchange

The following brief descriptions of the workshop pre-session demonstrations do not necessarily represent the range of data management approaches undertaken at marine and inland biological field stations. They do, however, reflect the diversity of ways in which scientific data management can proceed and is successfully being implemented at marine and inland biological stations.

1) ARC/INFO, Warren Brigham, Illinois Natural History Survey

The use of the ARC/INFO geographic information system running on Prime minicomputers and workstations was described. The system serves 300 users to provide a state-wide database for biodiversity, including occurrence records for distribution mapping, land use features, etc. The use of GIS to begin predicting potential habitats was demonstrated, including examples of how certain museum specimen label locations were biased by non-biological parameters such as road access. Also demonstrated was a study using GIS to increase the spatial accuracy of museum records by determining the probability surface for location descriptions on museum specimens.

2) Research projects database, John Porter, Virginia Coast Preserve LTER

A DBASE IV relational database of research project descriptions for the Virginia Coast LTER was demonstrated. The database could be sorted by date, place, location, investigator, and topic, and provided text descriptions of each project (historical and on-going).

3) Stebbins Cold Canyon Reserve GIS, Steve McNeil, University of California, Davis

A Macintosh GIS based on the program Business FileVision was demonstrated. The database links relational files about land use and environmental features to graphical display of points, lines and polygons. Query of the relational files can generate unique maps for visual display and hard copy. This program is used as an alternative to a written management plan for the University of California Stebbins Cold Canyon Reserve.

4) The "Andrew" multimedia bulletin board system, Rudolph Nottrott, University of Washington, Network Coordinator for LTER

An overview was given of the electronic network structure connecting the 17 ecological research sites which, along with the Network Coordination Office comprise the Long-Term Ecological Research network (LTER). Particular consideration was given to the national Internet and its NSFnet backbone. There are communications needs resulting from the wide geographical distribution of the LTER sites (Continental U.S., Alaska and Puerto Rico) and the dispersal of the 425 LTER researchers affiliated with over thirty institutions.

Electronic networks at three different levels are providing to ecological researchers: local-area networks (LAN), campus networks and wide-area net-

works (Internet). The functions at the highest level, the wide-area network level, include instantaneous and reliable electronic mail, access to supercomputers, access to national information and software repositories (including electronic bulletin boards), access to the LTER network office information system (mailing lists, mail forwarding system, LTER core data set catalog) and rapid long-distance transfer of data and programs, as well as text and graphics.

An overview was given of the electronic information system at the LTER network office. A detailed description was given of the LTERNET electronic mail forwarding system and a prototype installation of a multimedia electronic bulletin board (ANDREW) to be integrated with the mail system. The mail forwarding system can be reached from most major networks (Internet, Bitnet, Telemail, OMNET, UUCP, DialCom, MCI and others), and forwards messages to a user's "home" mail box on any of these networks. On request, an automatic reply function will return help information and various files stored on LTERNET. (To get initial help, send any message to forQuick@lternet.washington.edu (Internet) or forQuick@lternet (Bitnet).)

Plans for further development of the LTERNET information system were outlined, including the installation of an on-line catalog of LTER core data sets and development of this catalog into a distributed database system with local maintenance, administration and access control of all catalog entries and data sets, but with network-wide access for authorized researchers. Further development of this distributed database should include information already available at the LTER network office, such as the personnel directory, and data to be acquired in the near future (satellite images and other remotely-sensed data for all LTER sites).

5) Hypercard Bibliographic Database, Bill Seitz, Texas A & M, Galveston Bay

A bibliographic information database developed using Hypercard was demonstrated. This database runs on a Macintosh and is used for indexing maps, books, and articles. An optical scanner was used to read abstracts, and an optical character recognition program to convert bit-mapped images into ASCII characters. Approximately 2,000 records were entered in a short period of time with untrained staff. The Hypercard program can be linked to an Informix (SQL) database for rapid, relational search. Also described was a new service called MacSat which allows satellite images to be accessed directly via antenna, image processed, and displayed graphically on the Macintosh in color or 8-bit grey scale.

6) The PC-based FIS Database, John Briggs, Konza Prairie LTER

The FIFE project developed with NASA to study multi-scale remote sensing was discussed. The 100 gigabyte + image database is accessed using FIS software written by NASA and accessed by PC. The data is stored on a mainframe in an ORACLE database, and is accessed over a NOVELL network using the FIS software run by a PC. The database is accessed by about 200 users/month. The database will eventually be published on CDROM.

7) Macroscopic Ecology Laserdisc Demo, Michael Hamilton, University of California, James San Jacinto Mountains Reserve

The "data" collected at biological field stations often consists of a wide variety of types and formats, ranging from paper-based tables of numbers and text, photographs, films, illustrations, and tape recordings of sounds, to many forms of machine readable information. Computerized techniques which allow multiple forms of information to be integrated and accessed from a single microcomputer require the use of a class of tools loosely called "interactive multimedia" or "hypermedia." Hypermedia systems generally consist of 32-bit microprocessors, hard disk mass storage, videodisc or optical disk storage, digital signal processors for audio files, and appropriate software. The most widely used hypermedia platform is the Macintosh computer running software called HyperCard (tm).

The James Reserve data management program uses a hypermedia approach as an index and database integration tool to many of the Reserve's information resources. A Macintosh hypermedia database was demonstrated using Hypercard to control access to laserdisc images, record and retrieve sound files, access GIS software and display text fields which can be queried using words or phrases. This database is used to organize a time-series photomonitoring study of plant succession and vertebrate census records. Spatial fields are calculated using an ARC/INFO and displayed through Hypercard. The database is used primarily as an ecological inventory system for the field station and for teaching at the station and campus.

8) HyperCard Demo, Craig Staude, Friday Harbor Labs

HyperCard (a programming environment for Apple Macintosh) is suited for many tasks at field stations that require a short learning time, ease of use, and flexibility. Several examples were offered

demonstrating these merits, and one which currently falls short of expectations. The Friday Harbor Labs Information Program was originally developed for public relations to a general audience (e.g., open house and county fair booth). It was subsequently adapted to advertise the facilities of the station at a scientific meeting. It is a series of screens of graphic-rich information, including scanned images and simple animation, which are linked by mouse activated buttons. The demo startup stack (Macintosh jargon for "program") is used to alert new users to the peculiarities and capabilities of our Mac IIci. It is automatically displayed whenever the machine is restarted, by means of the "Set StartUp" feature of the Mac operating system. The Research Sites stack is essentially a mini-GIS. Invisible buttons overlay symbols or features on a map of the local region. When clicked, these buttons call up additional, small-area maps or text fields that describe each site in greater detail. Craig's Amazing Crustacean Database is a prototype database for storing species-specific taxonomic and collection information. Craig's Amazing Tab Inserter is a utility program that edits an imported comma-separated or space-separated ASCII file (e.g., modem accessed temperature data from a NOAA/NOS sensor) and converts it into a tab-separated ASCII file that can be printed or exported to other applications (spreadsheets, databases, etc.). The most ambitious project to date is the FHL Housing program, which finds

vacancies in housing units for visiting researchers, but it has not been implemented due to its slow response and the large number of query exceptions in the search arguments. Future versions might utilize streamlined script or add XCMDs to speed up the search process.

9) Easy Entry, Bill Michener, Baruch Institute

A database generation program called EASY ENTRY was demonstrated which is used to format data entry forms for inputting data while in the field. This system allows for the rapid uploading of data into other relational database programs running under MS-DOS.

10) SAS for database management, Paul Montagna, University of Texas, Marine Science Institute

Most users are familiar with the statistical features of SAS software (Statistical Analysis System version 6.03). However, SAS is an entire system with surprising data handling features. FSP can be used for database entry, checking and reporting. Base SAS has powerful manipulation features. Where data are maintained primarily for users who are familiar with and use SAS, it may be easiest for them to enter data directly into SAS. This eliminates the need for additional training and porting of data.

APPENDIX A—WORKSHOP PARTICIPANT LIST

The letters in brackets indicate working group participation:

[A] = Data Administration

[B] = Data Standards for Collaborative Research

[C] = Computer Systems for Data Management

Representing Field Stations and Marine Labs

Gary F. Anderson, Virginia Institute of Marine Science, College of William and Mary, Waterman's Hall, Gloucester Point, VA 23062 (Internet: gary@ches.cs.vims.edu) [A]

Michael A. Bowers, Blandy Experimental Farm, University of Virginia, P.O. Box 175, Boyce, VA 22620 [A]

Mary Bythell, West Indies Laboratory, Fairleigh Dickinson University, Teague Bay, Christiansted, USVI 00820 [A]

Barbara A. Carlson, Motte Rimrock Reserve, University of California, Riverside, Biology Department, Riverside, CA 92521 [A]

Victor Chow, Bodega Marine Laboratory and Reserve, University of California, Davis, P.O. Box 247, Bodega Bay, CA 94923 (Bitnet: UCDBML@UCDAVIS) [A]

Deborah A. Clark, La Selva Biological Station, Organization for Tropical Studies, Apartado 676, 2050 San Pedro, COSTA RICA (Internet: 3279995@mcimail.com) [A]

Philippe S. Cohen, Granite Mountains Reserve, University of California, Riverside, P.O. Box 101, Kelso, CA 92351 [A]

Robert W. Hastings, Turtle Cove Biological Research Station, Southeastern Louisiana University, P.O. Box 814, Hammond, LA 70402 [A]

Dean Kettle, Kansas Ecological Reserves, The University of Kansas, 2041 Constant Avenue, Campus West, Lawrence, KS 66047-2906 [A]

Robert Moeller, Pocono Comparative Lakes Program, Dept. of Biology, Lehigh University, Bethlehem, PA 18015 [A]

David L. Nebert, Institute of Marine Science, University of Alaska Fairbanks, AK 99775-1080 (Omnet: d.nebert) [A]

Janet Webster, Hatfield Marine Science Center, Oregon State University, 2030 S. Marine Science Drive, Newport, OR 97365 (Bitnet: hmsc@orstate) [A]

John M. Briggs, Konza Prairie Research Natural Area, Kansas State University, Division of Biology/Ackert Hall, Manhattan, KS 66502 (Bitnet: Konza@ksuvm) [B]

Walt Conley, Department of Biology, Box 3AF, Foster Hall, New Mexico State University, Las Cruces, NM 88003-0001 (Internet: wconley@nmsu.edu) [B]

Kenneth W. Cummins, Pymatuning Laboratory of Ecology, Department of Biological Sciences, University of Pittsburgh, Pittsburgh, PA 15260 [B]

John H. Heuer, Savannah River Ecology Lab, University of Georgia, Drawer E, Aiken, SC 29801 (Internet: heuer@srel.edu) [B]

Fred E. Lohrer, Archbold Biological Station, P.O. Box 2057, Lake Placid, FL 33852 [B]

Lance S. Risley, Institute of Marine and Coastal Sciences, Rutgers Pinelands Field Station, Rutgers University, P.O. Box 206, New Lisbon, NJ 08064 [B]

Bill Seitz, Moody College of Marine Technology, Texas A&M University at Galveston, P.O. Box 1675, Galveston, TX 77553 [B]

Emery R. Boose, Harvard Forest, Harvard University, Petersham, MA 01366 (Internet: eboose@lternet.washington.edu) [C]

Grady Cantrell, Hancock Biological Station & Center for Reservoir Research, Murray State University, Murray, KY 42071 [C]

Michelle Georgi, Forbes Biological Station, Illinois Natural History Survey, Box 599, Havana, IL 62644 [C]

Linda May, Horn Point Environmental Laboratories, University of Maryland P.O. Box 775, Cambridge, MD 21613 (Internet: may@umdc.umd.edu) [C]

Steve McNeil, Stebbins Cold Canyon Reserve, University of California, Davis, 144 Walker Hall, Davis, CA 95616 [C]

Paul A. Montagna, Marine Science Institute, University of Texas at Austin, Box 1267, Port Aransas, TX 78373 (Bitnet: paul@utmsi) [C]

Rudolf Nottrott, LTER network office, College of Forest Resources AR-10, University of Washington, Seattle, WA 98195 (Internet: rnottrott@lternet.washington.edu) [C]

Bob Vande Kopple, University of Michigan Biological Station, Pellston, MI 49769 (Bitnet: userhcyb@umichub) [C]

Craig Staude, Friday Harbor Labs, University of Washington, 620 University Road, Friday Harbor, WA 98250 (Bitnet: 98680@uwacdc) [C]

Nicholas Wolfe, NOAA/NMFS Beaufort Laboratory, Beaufort Laboratory Beaufort, NC 28516-9722 [C]

Workshop Pre-session Coordinator

Michael P. Hamilton, James San Jacinto Mountains Reserve, University of California, Davis, P.O. Box 1775, Idyllwild, CA 92349 [C]

Rapporteurs—Data Administration:

William K. Michener, Baruch Institute, University of South Carolina, Columbia, SC 29208 (Bitnet: a299050@univscvm, Internet: bmichener@lternet.washington.edu) [A]

Ken Haddad, Florida Marine Research Institute, 100 8th Avenue SE, St. Petersburg, FL 33701 (Omnet: r.burkhart) [A]

Rapporteurs—Data Standards for Collat

Warren Brigham, Illinois Natural History Survey, 607 E. Peab
(Bitnet: brigham@uiucdenr) [B]

James W. Brunt, Sevilleta LTER, Department of Biology, Univ
NM 87131 (Internet: jbrunt@sevilleta.unm.edu) [B]

Rapporteurs—Computer Systems for D.

John H. Porter, Virginia Coast Reserve LTER, Department of
University of Virginia, Charlottesville, VA 22903 (Internet:
jporter@lternet) [C]

Jeff Kennedy, University of California Natural Reserve System
Oakland, CA 94618 [C]

Organization of Biological Field Station:

George Lauff, Kellogg Biological Station, Michigan State Univ
Hickory Corners, MI 49060 (Internet: lauff%kbs.decnet@clv

Southern Association of Marine Labs—

James Alberts, University of Georgia Marine Institute, Sapelo

Host Site Representative

John B. Gorentz, W.K. Kellogg Biological Station, Michigan St
Drive, Hickory Corners, MI 49060 (Bitnet: gorentz@msukbs
gorentz%kbs.decnet@clvax1.cl.msu.edu)

Stephan Ozminski, W.K. Kellogg Biological Station, Michigan
Drive, Hickory Corners, MI 49060 (Bitnet: ozminski@msukt
ozminski%kbs.decnet@clvax1.cl.msu.edu) [B]

Lolita Krievs, W.K. Kellogg Biological Station, Michigan State
Hickory Corners, MI 49060 (Internet: krievs%kbs.decnet@cl

Workshop Consultant

Patricia Rich, Patricia Rich Associates, 115 Lake Forest, St. Louis

Michigan State University—Speakers &

James H. Beach, MSU Beal-Darlington Herbarium, Michigan S
Botany and Plant Pathology, East Lansing, MI (currently: H
Biology, Harvard University, 22 Divinity Avenue, Cambridg
beach@huh.harvard.edu) [B]

Paul M. Hunt, Academic Computing and Technology, Michiga
Center, East Lansing, MI 48824 (Bitnet: pmhunt@msu)

National Science Foundation - Speaker

Robert Robbins, National Science Foundation, Room 312, 1800
ternet: rrobbins@note.nsf.gov)

James L. Edwards, National Science Foundation, Division of B
G St. NW, Washington, DC 20550 (Internet: jledwards@note

APPENDIX B - GEOGRAPHIC INFORMATION SYSTEMS/ADMINISTRATIVE ISSUES

William K. Michener
Baruch Institute
University of South Carolina
and

Ken Haddad
Florida Marine Research Institute

INTRODUCTION

Management of spatial data relevant to a site is an important component of that site's data management system. A Geographic Information System (GIS) is a data management system that allows the capture, synthesis, generation, retrieval, analysis, and output of spatial data and, by some definitions, non-spatial data. Although this particular definition of GIS can be argued, there is general agreement that it is a rapidly evolving technology which is revolutionizing geographical analysis and has applications in many fields of science and resource management.

Parker (1988) and Cowen (1988) attempt to put into perspective the definitions and characteristics of a GIS as well as some of the fundamental operations. Some additional references which deal with all facets of GIS include: Burrough, 1986; Goodchild and Gopal, 1989; GIS/LIS'89; ASPRS, 1986; PE&RS, 1988; Michener, et al., 1989. In addition, almost every field of science and resource management now includes published articles and workshops related to GIS technology.

The applications of GIS technology at biological and marine field sites are numerous but can be approached through two broad and interrelated perspectives: (1) accomplishment of site management goals, and (2) accomplishment of research goals.

GIS related site management goals can range, for example, from cataloging and maintaining information generated at the research level, to conducting an integrated analysis of data collected by individual researchers, supplemented by data considered generic to a site, for management of the site's natural resources.

GIS related research goals can range from the use of spatial data for choosing research sites and the visual presentation of a researcher's data, to the use of GIS as an analytical tool for drawing scientific conclusions. In reality, the use of this technology as a research tool has only minimally been explored and in limited fields of science.

The interrelated applications of GIS technology as a tool for management and research can provide both opportunities and conflict at a field site. All aspects of GIS development at a site should be considered prior to implementation.

BLUEPRINT FOR A GIS

It is likely that, if not already implemented, many inland and coastal biological field stations are or will be considering the design and implementation of a GIS. At the site level, GIS should be considered a structured form of data management. The decision to design and implement GIS is an immediate step into a sophisticated level of data management. A site immediately goes beyond information cataloging and archiving and must be concerned with all aspects of data management and administration. All of the discussions on data management in previous chapters are relevant to GIS. In fact, particularly at the smaller sites, GIS may be the core for data management implementation.

Depending on the site and its functions, the individual researcher can have varying influences on GIS development. A concern often voiced at the research level, when administrative structure is imposed, is that science is being stifled. It is important to include the researcher in GIS design and implementation to assure that a rational data administration structure is applied and the user support base developed.

It should be recognized that GIS implementation at the site level may not be of benefit to all sites. Addressing other aspects of data management may better meet a site's needs. A given site should determine the need for a GIS from an administrative and research perspective and not assume its benefits. Individual researchers may provide the impetus for implementing a single user GIS as part of a research program. That is a site-specific issue. These observations are directed primarily at the site-initiated GIS.

There are avenues for GIS implementation outside existing data management operations. Successful GIS development and administration can occur as a parallel entity connected to data management efforts but not governed by the "data center." In fact, traditional data management administration can conflict with GIS evolution even though the principles of data management need to be applied.

If design and implementation are to occur at the site level, a GIS needs assessment should be conducted (Guptill, 1988). A GIS needs assessment is

not a trivial process, even for small or low activity sites. If the knowledge base to conduct a proper GIS needs assessment is not on site, then off-site expertise must be consulted. Because a needs assessment requires time and resources, it is often considered an impediment and consequently ignored. However, proper understanding and design are critical for long term data applications and research support, and should not be construed as an impediment to GIS implementation.

GIS IMPLEMENTATION CONSIDERATIONS

While a needs assessment should be a prerequisite to GIS implementation, elucidation of some of the important management considerations can provide a site administrator with some useful insight. Understanding the people, data, and cost considerations can facilitate successful GIS implementation.

PEOPLE AND TRAINING

Implementation at the site level should relate to the intensity of staff use and needs. This is a people consideration and will have major impact on successful implementation. Access by both the managers and researchers should drive the entire GIS development process. Technically, hardware and software play a role in implementation, but planning for longterm success must focus on the user, GIS operators, and their interactions.

Training should be considered key to GIS use and accessibility. Accessibility to the GIS can be accomplished by the availability of a skilled translator who can work with the investigators to build their understanding of the capabilities of the GIS and assist in operation of the applications software. This person should have site knowledge, a science (includes geography) background, and be well trained in the GIS applications software.

The researcher may prefer to be the analyst with hands-on skills. In this case the researcher must be trained not only in the applications software, but also in GIS concepts and principles. As with any technology, improper use and lack of understanding of the equipment can lead to error. A combination of the availability of a skilled translator and researcher training may be the best solution for optimum accessibility and effective utilization.

DATA

The use of GIS technology is dependent on the availability of data. Acquisition of hardware and software does not mean that a site has, or will ever have, a functional GIS. Selection and prioritization of data sets for entry and general access should be determined by the site administration in consultation with the site researchers. The needs of outside users should also be a consideration.

Although user needs drive the prioritization of data acquisition and maintenance, some additional issues which impact prioritization and determine successful GIS implementation are:

Spatial Resolution: Spatial resolution is probably one of the most important aspects of a GIS and one of the least understood. Spatial resolution can be divided into two components. The first component consists of the positional accuracy of an entity in the database. For example, if the location of a bald eagle nest is not accurately located it may show up in the middle of a lake, when compared to a database depicting land cover. The second component of spatial resolution is related to the user's need for detail and contains the elements of positional accuracy. Does the user need an accurate location and description of each tree in a forest, or will the location and description of the forest suffice

The subject of resolution is complex, and, if not properly addressed, could lead to wasted effort and be a major source of error in GIS analyses. Goodchild and Gopal (1989) put the question of the accuracy of spatial data, relative to GIS technology, in perspective. They suggest that the statistics do not even exist to define the error when spatial data are analyzed. It should not be concluded that GIS implementation is error bound, but that it is necessary to proceed with caution and with knowledgeable planning.

Coordinate System: Selection of an earth coordinate system is important. The three common coordinate systems are Latitude/Longitude, Universal Transverse Mercator (UTM), and State Plane Coordinates. Most coordinate systems are interconvertible, but commonality at a site may be advantageous for general communication of the data.

Quality and Documentation: Data quality and data documentation are two major issues in the development of a GIS. Variations in data quality can be amplified when analyzed in relation to other data. For example, when analyzing the relationship of soil data types (90 percent accurate) to the location of earthworm colonies (50 percent accurate), the resulting data may be only 45 percent accurate relative to hypotheses being tested. It becomes extremely important to have adequate documentation defining the source and lineage of the data and an assessment of the quality and accuracy of that data. The individual researcher or user can then determine the utility of that data set relative to the analyses they wish to conduct.

Proprietary Rights: Proprietary rights to data can often be the first controversial issue to arise when a site-initiated GIS is implemented. This issue should be anticipated and settled prior to implementation.

COST

As with any data management effort, the cost of the process and ability to support that effort should be a deciding factor in implementation. From an administrative perspective, can the site bear the long-term costs? Can a site finance common database generation and data maintenance and updating, and do so at the spatial resolution and update frequency necessary to make it useful to the researchers and other users? These are tough questions that are often ignored. The alternative to dealing with these questions in the planning process is to buy the hardware and software and hope that funds will become available for development and implementation. However, this approach has a high rate of failure.

Depending on the site, costs can be partitioned into the following:

Hardware: In addition to the initial purchase, the need for hardware evolution and expansion may be more accelerated for GIS needs than for traditional, non-spatial data management. Needs for computer hardware peripheral devices go beyond traditional printer and hard drives.

Software: Both operating and applications software can be a significant expense.

Maintenance: Hardware and software (operating and applications) often require maintenance and upgrades if a fully functional GIS is to remain operational. Routine expenses to support daily operations and replenish depleted supplies are far more than those associated with standard data processing.

Personnel: Depending on the site, GIS personnel functions can require the attention of more than one full time person for each function. The primary personnel functions are: (1) Data administration and coordination. This involves the development and maintenance of support for the GIS and all the basic administrative functions associated with data management. (2) Data capture. This refers to the identification, evaluation, and preparation of data to be entered into the GIS. This process is critical to the success of the GIS, particularly in understanding quality and accuracy of the data. (3) Data entry. This can be one of the more expensive functions, particularly for the common databases to be supported by the site. (4) Data analysis and output. This function requires technical skills and is the critical measure of successful implementation.

The costs associated with the above functions are variable and are often linked to the magnitude of the site operations.

Training: Training should include not only the technical aspects of GIS software and operations but also the development of an understanding of

the theory and algorithms applied during GIS analyses. Frequently, the GIS is treated as a black box and one is led through a process which, if not understood, leads to erroneous conclusions. Training must be a continual and scheduled process.

Data: This is frequently the most costly portion of the GIS operation. Costs include data acquisition, quality control, data maintenance and updating, data analysis and output, and archiving and security. These operations can easily cost at least four times as much as hardware and software. For example, it may be possible to acquire hardware and software for \$10-20,000, but the data necessary for implementation of an operational GIS could cost an additional \$40-80,000.

LITERATURE CITED

ASPRS (American Society for Photogrammetry and Remote Sensing). 1986. Proceedings of Geographic Information Systems Workshop. American Society for Photogrammetry and Remote Sensing, Falls Church, Virginia. 264 pp.

Burrough, P.A. 1986. Principles of Geographic Information Systems for Land Resources Assessment. (Oxford: Clarendon).

Cowen, D.J. 1988. GIS versus CAD versus DBMS: What are the differences? Photogrammetric Engineering and Remote Sensing 54(11) 1551-1555.

GIS/LIS '89. 1989. GIS/LIS '89 Proceedings. American Congress of Surveying and Mapping, American Society of Photogrammetry and Remote Sensing, Association of American Geographers, Urban and Regional Information Systems Association, AM/FM International. Volumes 1 and 2. 836 pp.

Goodchild, M. and S. Gopal. 1989. The Accuracy of Spatial Databases. Taylor and Francis, Bristol, PA. 290 pp.

Guptill, S.C. (ed.). 1988. A process for evaluating geographic information systems. U.S. Geological Survey Open-File Report 88-105. 136 pp.

Michener, W.K., D.J. Cowen, W.L. Shirley. 1989. Geographic Information Systems for Coastal Research. Proc. of Sixth Symp. on Coastal and Ocean Management/ASCE: 4791-4805.

Parker, H.D. 1988. The unique qualities of a Geographic Information System: a commentary. Photogrammetry Engineering and Remote Sensing 54(11): 1547-1549.

PE&RS (Photogrammetric Engineering and Remote Sensing). 1988. Special GIS Issue. 54(11): 1-167.

APPENDIX C—CLIENT/SERVER DATABASE ARCHITECTURE, NETWORKS, AND BIOLOGICAL DATABASES

James H. Beach
Herbaria and Museum of Comparative Biology
Harvard University
22 Divinity Avenue
Cambridge, MA 02138

INTRODUCTION

The development of the academic research networks, in particular, the NSFnet or Internet and the forthcoming National Research and Education Network (NREN), will provide the potential to make myriad biological data resources available to scientists and students around the world. Although electronic mail, interactive sessions with remote applications, and high-speed file transfer are now integral to many research programs, the development of database systems which will bring biological data resources to the networks is in its infancy.

DATABASE ARCHITECTURE

The NSFnet/NREN permits several types of long-distance access to biological data sets. The traditional and still commonplace form of communication with remote databases is one where users connect over the network to establish a terminal session with a remote host. Remote users log onto the computer and operate database application programs in the same way a local terminal user would. This is an example of "host/terminal" database architecture.

A technical characteristic of host/terminal database systems is the logical cohesion between the database manager software, which stores and manages user data files, and the application programs interacting with it (Figure C-1). A major benefit of the logical integration of the layers is ease of database system development; application programs can be tailored to fit like a glove around the features of the database manager. Remote access to data in host/terminal systems is exclusively through the host's application programs.

"Client/server" database architecture in contrast, uncouples the application programs from the database server software (Figure C-1). Client/server databases sandwich an additional logical layer to handle communication between the server and the applications, through the use of a go-between, standard query language.

The importance of the client/server model, in the context of network access to information, is that it allows the application layer programs and the

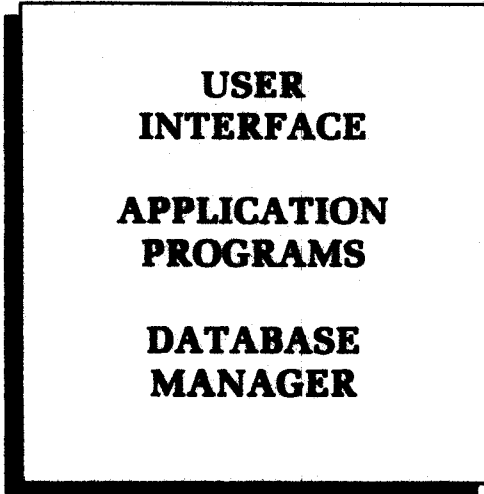
database server software to reside on different machines. Because the two layers communicate through discrete, structured messages, the conversation can be carried out between machines connected across the room, across the country, or across the globe. The development of the high-speed, high-capacity research networks strengthens the importance of client/server systems for biological databases, because institutional data servers could be queried at any time over the network by any number of applications at remote sites. A particular application could rapidly access multiple institutional servers over the network channel.

A functional difference between the client/server and host/terminal database architectures has far reaching implications for access to scientific information. In the host/terminal model, a remote network user running a (virtual) terminal session from a local computer, e.g., a desktop PC, only receives screen images; information is visually presented but there is no mechanism to download data records for local use. Capturing data from a formatted screen display, one screen at a time, is usually an imperfect process at best. As a result, access to remote information is essentially limited to the duration of the virtual terminal session. A client/server database, in contrast, transmits actual data records to the remote user's system. The records (in a standard exchange format) are then available to local programs for further processing or formatting. Note that with client/server architecture, remote users are not constrained by the application interface or program logic of the server system, but work with a familiar local application to query and obtain records from network database servers. An additional distinction of the client/server approach, in a computing environment characterized by autonomous institutions in a collaborative enterprise, is that it allows organizations to control the ongoing development of their hardware, database, and application software, while at the same time presenting a standard and stable network interface for remote client access.

STANDARDS

The scientific disciplines will need to resolve various types of data format, application and data

HOST / TERMINAL



CLIENT / SERVER



Figure C-1. Database System Architecture

communication standards in order to establish network client/server systems. Database systems developed in isolation, on small, single-user computers or on large un-networked machines may be elegantly customized for local needs, but biological databases intended to inter-operate with remote applications will need to be specified, designed and implemented on much technical common ground. The most important computerization standards for network client/server systems are:

A common set of core data definitions

Discipline or community-wide standards for core data type definitions, coding, and cataloging rules are essential for the biological information stored in systems designed for network access. These standards comprise formal descriptions and specifications for data types currently in use in non-computerized or non-networked databases. Ecologists will have an especially difficult task, due to the breadth of ecological research, but certain ecological data types are already fairly well standardized.

There are several ongoing ecology, systematics, and museum community efforts in this area, including: the LTER data catalog project, NSF-sponsored, discipline-based data workshops, and various projects of the International Working Group on Taxonomic Databases in the Plant Sciences, the Association of Systematics Collections, the Museum Computer Network, as well as several additional society and institutional efforts.

A standard exchange record format

The results of a search on a remote data server must be returned to the requesting application in a standard record format. Without such a format, result sets would not be understood by the client process, and client/server data exchange would be impossible. Included here are standards for data representation, syntax and structure specifications for records and fields. Data definition and encoding standards (above) can be applied as a part of the exchange record format definition.

The library community has standardized the definition of data elements and data record exchange formats in its highly successful MARC record format. The MARC format standards are only applied to records intended for exchange and not to database design. They have provided tremendous stability and have greatly facilitated information interchange between diverse library database systems.

Some of the organizations mentioned above have begun to investigate a MARC approach for museum data and there is growing interest in MARC formatting of biological information for record exchange purposes.

Standard network request and response protocols

For client applications to communicate with data servers, there must be a well-defined language and syntax for the interaction. Such standard protocols specify the structure and to some extent the content of the messages passed between client and server machines as part of a data request/response dialogue. They also specify how control and state information will be communicated and under what conditions diagnostic messages and result sets will be transmitted to the originator of a query.

The best example of standard protocols for the retrieval of information in a client/server architecture again comes from the libraries. That community sponsored the development of the ANSI/NISO standard Z39.50 (NISO, 1988) which specifies network session protocols for library information retrieval. The Z39.50 protocols are being used for library data exchange as part of the multi-institutional "Linked Systems Project" (Fenly and Wiggins, 1988). A thorough examination of the libraries' computing and standards infrastructure would assuredly be profitable for nascent data standards efforts in biology.

IMPLEMENTATION

Biological client/server database systems can be implemented over networks today, and they will become increasingly common as discipline, national, and international communication standards are completed. There are numerous engineering options for implementing client/server systems, but an overriding design objective is ultimate compliance to network communication standards, particularly those of the International Organization for Standardization (ISO), the U.S. National Information Standards Organization (NISO), and to the data definition standards developed within the scientific disciplines.

As a prototype example of client/server architecture, a biological client/server database system using two networked computers was demonstrated at the Data Management Workshop. A Digital VAX/VMS system functioning as the client was located at Kellogg Biological Station (KBS), while the server, a Sun Microsystems workstation, was about 70 miles away on the Michigan State University campus in East Lansing. An Ingres client application at KBS, using a query-by-form screen, recorded a query specification based on user selections and then mailed the query to an Ingres herbarium specimen data server in East Lansing. That computer parsed the contents of the network mail message and applied it as a query against the database. The result set was stored in a file, then mailed back to the KBS client application within a few minutes. The client process reported the

arrival of the result set to the user and imported the records into a local database table for further processing.

A mail-based client/server system, although in some ways a "low-tech," approach, uses a universally supported network application and is relatively easy to implement. Limitations include delays caused by network mail routing, limits on the content and length of mail messages imposed by network mail programs and the inherent difficulties of managing state information and a request/reply process with network mail. Due to short-term exigencies, the Ingres QUEL query language was used, but SQL (Structured Query Language), which is the industry standard query language for client/data server communications in relational database systems (Date, 1990; Tucker, 1990), could have been employed.

A more standard and sophisticated client/server design is a "connection-oriented" approach, whereby client and server processes enter into a real-time network dialogue. In this case, a precisely defined protocol is required for the client/server communication (e.g. ANSI/NISO Z39.50) which specifies a predictable sequence of back-and-forth control and data messages while a client is requesting or receiving information from a server. Connection-oriented, client/server, network database protocols function on top of lower-level network communication standards to form an integrated, layered stack of protocols. In contrast to the "connectionless" mail-based, client/server approach, a connection-oriented system requires sophisticated system and network-level programming to implement, but, in addition to speed, it offers numerous technical advantages for monitoring client/server sessions and for returning useful status information to the user.

SUMMARY

Client/server database systems can provide a direct and powerful method for biological database access over the NSFnet, the NREN and the international extensions of those networks. When implemented for open access, they have several advantages over host/terminal systems. The most notable are:

1. Users would not need to obtain an account on each target system they wish to query, and they would not need to learn the logic and design of each institutional database application.
2. Data records can be acquired and processed locally in the client/server model, as the data server actually copies data records and not just a refreshed screen image to the remote user. Only result sets meeting the user's query criteria are returned over the network.
3. Institutions could provide open, read-only, server-level, access to their biological data resources with limited risk or loss of administrative control.
4. Once network, client/server, interface standards are in place, institutional database system hardware, server software and applications can evolve independently and still provide open, long-term, network access.

REFERENCES CITED

- Date, C. J. 1990. An Introduction to Database Systems. Vol. 1. 5th Ed. Addison-Wesley, xxv + 854 pp.
- Fenly, J. G. and B. Wiggins. 1988. The Linked Systems Project: a networking tool for libraries. Online Computer Library Center, Dublin, Ohio. xii + 138 pp.
- National Information Standards Organization. 1988. Information Retrieval Service Definition and Protocol Specification for Library Applications [Z39.50-1988]. Transaction Publishers, xii + 52 pp. (Available from: Transaction Publishers, Rutgers University, New Brunswick, N.J., 08903, USA)
- Tucker, J. T. 1990. The inevitable merging of SQL. Unixworld 7(2): 68-70, 72, 74.

APPENDIX D—INTERSITE ARCHIVAL AND EXCHANGE FILE STRUCTURE

(Excerpt from an article submitted to Coenoses)

Walt Conley
New Mexico State University

and

James W. Brunt
University of New Mexico

An Intersite Archives File structure has been defined in order to facilitate the need for an orderly approach to the design and implementation of a data manipulation capability. The manipulations to be done are alterations on the shape and/or the content of original (archived) data files, and communication of original or descendent files to remote sites.

The Intersite Archives File (Figure D-1) is a generalized data structure that contains full documentation and comments. It is intended that the test of adequate documentation is that these files should stand alone, and that the file itself should contain sufficient information so that a future investigator who did not participate in collecting the data can use the information for some scientific purpose. The Intersite Archives File structure is intended to be used across cooperating research sites that, taken together, represent the ultimate heterogeneous computing environment. The intent is to define a generic data structure that can be useful on any hardware and software system, and that can be sent on any electronic network or file transfer system. A companion effort provides an Intersite Toolkit for obtaining information from the files; there are also tools for manipulating and screening the files. Manipulations include stripping an Archives File of various categories of information to produce descendent files that can be read by any application package.

The basic Intersite data structure is a generic ASCII flat file that contains categories of information that define the data, as well as the data itself. Intersite Archives Files can be of any basic structural type, including statistical data, text data, graphics data (e.g. files that you can write to a graphics plotter), gene sequence data, or bit map image data. Other file types will no doubt be required. Note that file type refers to the general nature of the data in the file, and not to data typing such as floating point, integer, or character. All of the data in the Intersite Archives Files are ASCII characters, and provision is made in the Intersite Toolkit for handling files containing non-printing ASCII characters which make file transfer difficult on some networks and impossible on many of the file transfer protocols.

The general categories of data in an Intersite Archives File is as follows.

\ log: A record of the history of the file; when it was initiated, updating, changes entered, locations and dates of copies of the file. Any ASCII characters with any format may be included.

\ doc: Documentation—as detailed a description as is necessary of the data contained in the file. Any ASCII characters with any format may be included. An ABSTRACT may be included here to allow automatic extensions of data dictionaries from Archives Directories. The abstract is simply a paragraph beginning with ABSTRACT and ending with a blank line; it may appear anywhere in the documentation section.

\ type: File type refers to the basic nature of the data. Statistical files are typically rows by columns tables of numeric or character data. Text files include bibliographic data, abstracts, or any prose. Graphics data refers to files which can be written to a plotter or a printer. Genome data refers to long sequences of base pairs that require line delimiters and other embedded information. Image data refers to bit map images.

File typing currently includes statistical, text, graphics, genome, and image. Other file types are possible and can be added as necessary. The only operation anticipated on file type is identification for sorting.

\ header: Header refers to labels for the columns of data in a statistical data file, or a list format text file. This allows for automatic building of data dictionaries from Archives Directories. For files of other types, the header can contain keywords that describe the data. Labels or keywords in the header are automatically retrieved for the development of data dictionaries in Intersite Archives data directories. The Intersite Toolkit provides tools that do this work.

\ data: Data refers to the actual data of the archives file—the numbers, text, etc. The data section may contain embedded comments that further describe individual records of the data.

\ log

A log of activity for the datafile including names, dates, etc.

\ doc

All the documentation needed to accompany the datafile in free format

ABSTRACT

Includes the option for an extractable abstract

\ type

A one word descriptor of the data ie., statistical, image, list, etc.

\ header

A description of the attributes for statistical data

\ data

The Data

(Includes comments)

Figure D-1 Intersite Archive File Structure

The Intersite Toolkit contains programs that manipulate the Archives File data structure, making the files ready for applications programs such as relational data management systems, statistical or graphics packages, and reporting systems such as text formatters. Any combination of categories of information in the Archives Data Files can be extracted for further use. Thus in a statistical file it is possible, for example, to quickly extract only the column labels and the table of numbers, only the ABSTRACT, only the documentation section. The Toolkit also contains compression and decompression filters (useful for disk maintenance and some communication applications), an encryption and decryption algorithm (useful for converting files with non-printing characters to files that can be sent over networks that do not handle binary data or via dial-out modem transfers), and a suite of programs that automatically build and reference a data dictionary that contains various presentations of labels, keywords, and abstracts.

For statistical and text file types, there are 2 additional formats that are useful to consider. "Table" format is the typical row X column format of statistical data with a label at the top of the column. "List" format is a transposed table, where the labels are on the left margin, giving unlimited category width but with a single column of data. List format is useful for text data such as keyworded bibliographic citations, or any similar kind of text. Note embedded comments can be included anywhere in the \ data section simply by enclosing the comment in curly brackets. The only restriction is that comments and other data cannot be mixed on the same line. (This preserves the positioning of tabular data, and serves the goal of keeping these files "readable" by humans.)

The general structure of an Intersite Archives File (type is "statistical") in Table format is shown in Fig. D-2. Note that the category indicators (\ log, \ doc, \ type, \ header, \ data) occupy a separate line but do not need to begin in any particular column. The suggested categories are optional, although deletion of any category limits the usefulness of the file and the use of the Intersite Toolkit for manipulating the files. The structure of an actual Intersite Archives File in Table format is shown in Fig. D-3. The general structure of an Intersite Archives File in List format is shown in Fig. D-4.

In the log and doc sections, there are no format requirements, and free-form text can be entered as you choose. In the header section and the data section, some structure is necessary. In the Table format, the header labels provide searching tags for

the data file manipulations (and serve as handy reminders), and the dashed lines indicate the maximum width of each column of data (which is used for subsequent manipulation of the data columns. The dashed lines are not necessary for many applications; they are useful for providing information for manipulation routines. To include them requires little, and adds considerably to the potential for cross-site data manipulation. In the List format labels appear at the left of the field, and the dashed-lines indicator for column width is not necessary. In Tables, data columns conform to the labels in that they are in the same order, and in the Table format, the data must fit within the number of columns indicated by the dashed lines.

A Table format has one or more columns, and a List format has only a single column. Columns may be of arbitrary width. The labels in each case provide for data abstraction in good applications packages in that the researcher may refer to variables by name (i.e. labels) rather than, for example, as column 3. Archives Files are specifically intended to be browsed by human researchers who want to become familiar with the data and the circumstances involved in the collection of the data. Once converted to the descendent files that will be manipulated via available relational operators (etc.), data files are not designed to be read by humans, and will be confusing to look at.

In practice, any numerical data set can be put into a rows by columns table format, and the only restriction is that the columns have some white space between them. This is the format that is typically used when recording data in the field, or when reporting data, and the Intersite data structure simply provides a computerized version of what you probably do anyway. There is a utility in the Intersite Toolkit called "extract" that can subset the standard Intersite Archives File structure (Figure D-5). This utility can create a new file with any combination of the various elements of an Archives File stripped from the original; the original is, of course, preserved intact. Other programs in the Intersite Toolkit provide manipulation and screening of the Intersite Archives Files, building of a data dictionary based on labels and keywords, extracting and sorting Abstracts, and generally obtaining information from the Archives directories.

Once the documentation has been stripped from the chosen archive files, and the files are ready for some serious work, the descendant files can be read into any applications package of your choice. A next obvious choice is entering the filtered data into a database system for further manipulation. If you use a relational database system is being used the


```

\log
***** BEGIN CHANGE LOG *****
23 December 1987. Data entered and documentation established. MAUhl
***** END CHANGE LOG *****

```

```

\doc
ABSTRACT Ant Total Density on the Jornada. This file, ant/ __ total.density, is monthly mean densities
of new colonies grouped into zones, pooled for all species. The last 5 columns represent the monthly den-
sities by year, and the first column describes the area ("zone") where the colonies were located. Data were
collected by Marsha R. Conley 1982-86.

```

These 5 species were pooled to create the file:

Code:	Scientific name:	Common name:
PODE	Pogonomyrmex desertorum	Desert Harvester Ant
PORU	Pogonomyrmex rugosus	Red Harvester Ant
MYDE	Myrmecocystus depilis	Honey-pot Ant
MYMI	Myrmecocystus mimicus	Honey-pot Ant
NOCO	Aphaenogaster cockerelli	

```

\header

```

```

Zone 1982 1983 1984 1985 1986
-----

```

```

\data

```

Playa	0.0	0.0	0.0	0.0	0.0
Mesquite Fringe	2.7	3.0	3.3	3.3	3.3
Basin Slope	6.7	8.1	8.8	7.9	6.8
Bahada	0.5	0.6	0.8	1.3	1.5
Lower Piedmont	2.8	3.1	3.2	2.4	2.6
Upper Piedmont	0.8	0.9	1.1	1.1	0.9

(Only Pogonomyrmex were found in the Upper Piedmont)

Figure D-3: Structure of an Intersite Archives File in Table format.

```

\log
Records of the history of the datafile. When it was initiated, changes entered, locations and dates of copies
of the file. Any ASCII characters with any format may be included.

```

```

\doc

```

Documentation: As detailed a description as necessary of the data contained in the file. Any characters with any format may be included. An ABSTRACT of 1 paragraph may be included anywhere in this section.

```

\type

```

Typically List files are of type text.

```

\header

```

Nothing needed here for List format. Note that the \intersite data dictionary tools will pick up the Labels at the left margin of the first record and will automatically treat them similarly to the column labels from the Table format.

```

\data

```

This is a comment. Note that the new line below is required to automatically identify the List format.

```

\begin (verbatim)

```

```

label1      line of text      ....
label2      line of text      ....
label3      line of text      ....
              :
              :
              :
labeln      line of text      ....
label1      line of text      ....
label2      line of text      ....
label3      line of text      ....
              :
              :
              :
labeln      line of text      ....

```

Figure D-4. General structure of an Intersite Archives File in List format. Note that the Labels are simply the first unbroken string of characters in each line.

APPENDIX E—SYSTEM SELECTION OVERVIEW

John H. Porter
University of Virginia

and

Jeff Kennedy
University of California Natural Reserve System

Advice on choosing a computer and software is always short-lived. Changes in systems and prices occur almost daily. Nonetheless, such advice is valuable to a person setting up a new data management system. The following sections attempt to provide needed information to new data managers on how to choose a PC computer (running MS-DOS) or a Macintosh (running the Apple operating system). What is not included is guidelines on whether to choose a PC or a "Mac." This is because the general capabilities of the two computers overlap so greatly. Choosing between them will depend on relative costs, the computing environment and the preferences of users.

SELECTING AN MS-DOS COMPUTER

The type of computer that is "best" for you depends entirely on what you want to do with it. Critical questions to ask are:

1) What sorts of activities do you want to use the computer for? Different uses have different requirements. Here is a brief table of uses and minimum desirable configurations for each.

Use	Processor	Memory	Numeric Coprocessor	Hard Disk
Word Processing	8086	640K	N	20 MB
Spreadsheets	80386SX	>1 MB	Y	30 MB
Statistics	80386SX	640K	Y	40 MB
Database	80386	640K	N	40 MB
Programming	8086	640K	Y	20 MB
Communications	8086	640K	N	20 MB
Graphics	80386	>1 MB	Y	80 MB
Multitasking	80386	4 MB	Y	40 MB

Because data management activities tend to be both computationally intensive and storage intensive, a minimum configuration for a primary data management computer would be an 80386, 80486, or 80586 central processor, with a numeric coprocessor and a large disk drive (>40 MB). Some form of high-capacity backup system (tape cartridges, Syquest or Bernoulli removable hard disks, or DAT tape cartridges) should also be added. Everything listed re-

quires a hard disk and at least 640K of memory (RAM), which will let you run 98% of all MS-DOS programs. Skimping on memory reduces costs in the short term, but increases frustration in the long run.

The 80286 processor is not listed, because 80386SX-based machines approach the price of 80286 machines, and they have the potential for expansion and for support of the OS/2 and UNIX operating systems. 80286 machines lack these capabilities. Not listed in the table is the "clock" speed of the machine. The venerable IBM-PC used 4.77 MHz, but you do not want anything that runs slower than 8 MHz. For really intensive tasks (such as graphics or multitasking) higher speeds (33 MHz and above) may be desirable. Keep in mind that disk-intensive tasks, such as using databases and statistics packages, benefit much less from a higher clock speed than from a fast disk drive or a RAM disk. The "width" and speed of the data bus will also affect the effective speed of the computer for data management.

2) Where do you want to do your computing? If you do your work in a fixed location, a desktop machine with video monitor (preferably VGA color) is a better value. If you need to compute in the field, it may be worthwhile to pay the 30% extra for a portable computer.

3) How long will it be before you buy a new computer, and how much do you plan to spend on software until then? If you plan on keeping your new computer for several years, adding new software as it becomes available, purchasing an 80386-based machine may be important. The next several years will see increasing numbers of programs that require the 80386 chip. Most of these programs will be for specialized applications (spreadsheets, graphics and multitasking) rather than word processing.

4) How much help will you need in setting up your computer and how much "down time" can you tolerate? This really affects where you buy your computer and what brand of computer you buy more than what type of machine you buy. If you feel comfortable installing boards and disk-drives, mail order can be the cheapest place to buy. If you need someone in

town to help with system setup and maintenance, it makes sense to pay a little extra to establish a relationship with a local dealer.

The brand of computer is important in determining how long it will take for computer repair. Most major domestic computer companies make their own computers with standard main boards. However, some cheaper imported computers actually come from a large number of different sources, each with variant main boards. Getting main boards for such computers can take a long time (even domestic computers may take a month or more). On the positive side, hardware failures are rare and are usually confined to individual add-on boards (not the main board), making replacement easy on all brands.

Choosing software is an art in itself that is highly dependent on the scope and difficulty of the computing tasks in question. Surveying the computer magazines for software reviews and consulting with user groups is the best source of detailed, current information affecting software selection. These software packages were recommended by attendees at the Data Management Workshop.

Word Processing: WordPerfect, Microsoft Word
 Statistics: SAS, SYSTAT, STAGraphics, PSS-PC
 Database: SAS, DBASE (III and IV), Paradox, Foxbase
 Graphics: Sigmaplot, SAS, Harvard Graphics
 Spreadsheets: Lotus 1-2-3, Excel, Quattro
 Utilities: 386Max, Norton Advanced Utilities, XTREE

SELECTING A MACINTOSH SYSTEM

As with MS-DOS machines, the type of Macintosh you need depends on your computing needs and your working environment. Critical questions include:

1) What tasks will you be using your computer for? Different uses have different requirements (suggested minimums are shown):

Use/task	Processor	Memory	Co-processor	Disk
Word Processing	68000	1-2 MB	N	20 MB
Desktop Publishing	68030	2 MB	N	20-40 MB
Spreadsheets	68030	2 MB	Y	20-30 MB
Statistics	68030	2 MB	Y	40 MB
Database	68030	2 MB	N	40 MB
Programming	68030	2 MB		20-30 MB
Communications	68000	1-2 MB	N	20-30 MB

Graphics	68030	1-2 MB	Y	40-80 MB
Multitasking	68030	2-5 MB	Y	40-80 MB
Image processing	68030	4-8 MB	Y	>80 MB
GIS	68030	4-8 MB	Y	>80 MB

Apple's release (at the time of this publication) of its System 7 operating system will require a minimum of 2 MB of random access memory. Upgrade to System 7 is not essential for simple computing, but if you have two or more Macintosh computers connected to a LAN, all must have the same System 7.0 printer drivers. Multitasking is possible using System 6.0X with Multifinder and 1 MB of RAM, but 2 MB is the practical minimum. System 7 has multitasking built-in. Both operating systems may coexist on machines connected to the same local area network.

Given that data management tasks at field stations tend to be computationally and storage intensive, the recommended minimum configuration for a primary data management computer would be a 68030 machine, such as a MacSE30, a Mac IIsi, or Mac IICI, with 4-5 MB of RAM, a built-in numeric coprocessor, and a 40 MB hard drive. RAM costs have dropped to the point where 4 MB of RAM from a mail order house can be added to a 1 MB machine for approximately \$175, installed, resulting in a 5 MB machine. The extra RAM can significantly speed processing by reducing hard disk read/write cycles. Forty-five MB Syquest removable cartridge drives are ideal for backing up and archiving files. Image processing or GIS work will require a 25 MHz Mac IICI, and preferably a 40 MHz Mac IIfx. Accelerated video display boards will vastly increase the speed of display and data analysis.

The MacClassic and MacSE, with their 68000 and 68020 processors may be fine for word processing, student use, or data entry—as opposed to data analysis—but the 68030 machines will enjoy a longer time to obsolescence. As with MS-DOS machines, disk-intensive tasks, such as data base and statistical analyses, will benefit from a fast disk drive.

2) Where do you want to do your computing? If you work in a fixed location, or you need a larger monitor than the 9-inch built-in a MacClassic or a MacSE, you will need a desktop machine with a video card and external monitor. The selection of Macintosh portables is much smaller and the costs much higher than in the MS-DOS world. For simple word processing, spreadsheeting or data logging

in the field, consider an inexpensive MS-DOS clone portable, a Radio Shack portable, or a Z88 used in conjunction with Laplink or MacLink Plus file transfer and cable packages. Data and graphics analysis can then be done on your office machine with the uploaded data.

3) How long do you plan to keep your computer before upgrading and how much do you plan to spend on software until then? In general it is cheaper in the long run to buy a more sophisticated machine initially than to upgrade at a later date. Buying a 68030 machine will give you a longer usable lifetime for the machine. An SE30 is the cheapest 68030 machine, but it has only one slot for an add-in board such as an external video board. The Mac IIsi currently provides the greatest combination of low cost, expandability, functionality and ease of access and repair. The 68030 is also compatible with Apple's version of the UNIX operating system, AUX 2.0.

4) How much reliability, service, and support do you need for your system? Buying your CPU

(Central Processing Unit) and peripherals from an Apple authorized dealer gives you one stop shopping, and subsequent service, but Apple limits its warranty to one year, and the quality of post-purchase service and support varies significantly from dealer to dealer. Buying your peripherals from third-party vendors can earn you 2-5 year warranties and often improved support, but at the expense of having to deal with multiple manufacturers and/or dealers. Research the support and repair programs of each purchase carefully. Local Macintosh user groups and bulletin boards are good sources for this information.

Software selection is highly dependent on the scope and difficulty of the computing tasks in question. Again, user groups, bulletin boards (such as ZMAG on the Compuserve Information Network) and magazine reviews are excellent sources of current information. The following software packages were recommended by attendees at the Data Management Workshop and/or were given high rankings among 1000 Macintosh products evaluated in MacUser 7(8):135-220.

TASK

SOFTWARE

Word Processing:	Word, MacWrite II, WordPerfect, TeachText, WriteNow
Page Layout & Desktop Publishing:	PageMaker, Framemaker, QuarkXPress, Fast Forms
Desktop Presentation:	Persuasion, PowerPoint, More
Multimedia:	MacroMind Director, Media Tracks, MacRecorder Sound System, Audiomedia
Hypermedia:	HyperCard, Reports
Spreadsheets:	Excel, WingZ, Works, Parameter Manager Plus
Statistics:	SYSTAT, DataDesk, SPSS, Statview II, JMP
Graphing & Charting:	DeltaGraph, KaleidaGraph, Igor, MacSpin (see also statistics & spreadsheet programs, above)
Mathematical Equation Writing/Solving/Modelling:	Mathematica, Theorist, Expressionist, Stella, Extend
Data Acquisition & Lab Instrument Interface:	LabVIEW 2, MacADIOS, MacLab
Flatfile Database:	FileMaker Pro (quasi-relational), Borland Reflex Plus (quasi-relational), DAtabase
Relational Database:	4th Dimension, FoxBASE +/Mac, Omnis, Panorama, Double Helix
Bibliographic Database:	EndNote & EndNote Plus, EndLink
Communications & Multiplatform Connectivity:	MicroPhone II, White Knight (Red Ryder), SmartCom II, VersaTerm Pro, Kermit, TinCan, ZTerm, MacTerminal, Timbuktu, MacLinkPlus/PC, LapLink Mac III
Networking & E-Mail:	AppleShare, MacTOPS (small networks, primarily), Novell Netware, Microsoft Mail, QuickMail
Multitasking:	System 7.0, MultiFinder (System 6.0x)
Paint/Draw Graphics:	Canvas, Illustrator, Freehand, MacDraft, Mac Draw, MacPaint, Studio/1 & Studio/32, Super 3D, Swivel 3D
Image Processing:	Pixel Paint, Photoshop, Image (National Institutes of Health shareware), Digital Darkroom, Spyglass View/Transform/Dicer
CAD:	Claris CAD, MiniCad+ 3.0, VersaCad, Ashlar Vellum
GIS:	MacGIS (U. Oregon), MacGIS (Cornell), Map II, ESRI ArcView (for download & display of ArcInfo data & images on a Mac), Business File Vision/File Vision IV (a poor-man's quasi-GIS)

APPENDIX F WORKSHOP SURVEY QUESTIONNAIRE

Workshop on Data Management for Inland and Coastal Field Stations
April 1990

Pre-Workshop Survey Questionnaire
November 1989

3. Data bases
- a. Does your site have databases that have been compiled specifically for general use (e.g. species lists, meteorological data)? If so, please list some examples.
 - b. Does your site have databases originating in individual research programs, that are or could be developed into general-use resources. If so, give a few examples.
 - c. Does your site have computerized records consisting of non-traditional forms of data, e.g. acoustic records, maps, visual images.
4. Administration and Personnel
- a. Where does the impetus for data management arise (e.g. site administrators, interested faculty members, research programs, technical staff)?
 - b. Does your site have a data manager, or other person(s) with designated responsibility for data management?
 - c. What personnel are involved in data management (number of persons, positions, training, experience, fraction of time)?
 - d. How is data management funded? Is there a specific budget for data management? Is it funded at the site/institution level, or on individual grants?
5. Availability of data
- a. How can a person find out what data are available at your site? Is there a catalog or directory of data? If so, what information is kept, and how is it organized?
 - b. How do you weigh investigators' "proprietary" rights to data against the goal of wider availability? Is there security against unauthorized use of data?
 - c. Does your site have standardized quality control procedures for data?

Use the enclosed envelope to send your responses to:

Data Management Workshop
W.K. Kellogg Biological Station
Michigan State University
3700 East Gull Lake Drive
Hickory Corners, MI 49060

Please reply by December 11, so we can make your responses available to the planning group which expects to meet later in the month.

The questions are somewhat open ended, based on the assumption that the most useful information you can give us won't fall into neat categories. Please feel free to add explanatory comments, using additional sheets of paper if necessary.

If you have questions, please contact John Gorentz at the above address, or at 616 671-2221, or by electronic mail at [gorentz@msu.kbs \(Bitnet\)](mailto:gorentz@msu.kbs (Bitnet)), [jgorentz@internet.cfr.washington.edu \(Internet\)](mailto:jgorentz@internet.cfr.washington.edu (Internet)) or [J.GORENTZ \(Usenet\)](mailto:J.GORENTZ (Usenet)).

1. Questionnaire respondent(s)

Institution:

Date:

Name(s) and position(s):

Mailing address:

Phone:

Electronic mail:

2. Your site and institution

To inform workshop planners who are unfamiliar with your site, please include copies of any brochures or materials describing your site, its facilities, habitats and ecosystems, types of research, level of activity, and other programs.

6. Goals and Objectives

Below are items that could represent data management goals for your institute or field station. Following each item, circle the status code(s) that best describe your site in relation to the goal.

Status Codes

ACF - An accomplished fact
VIP - Work is in progress
HPC - High priority goal
SKF - Seeking funding for this goal
MPC - Medium priority goal
NPL - No plans to do this
LPC - Low priority goal

a. Implement a central catalog or directory describing all data sets on natural habitats (i.e. data about data, computerized or not).

Status: ACF HPC MPC LPC VIP SKF NPL

b. Implement a central catalog or directory of data about data that is electronically searchable, "on-line".

Status: ACF HPC MPC LPC VIP SKF NPL

c. Manage selected databases for general use as a site/institutional responsibility.

Status: ACF HPC MPC LPC VIP SKF NPL

d. Implement a standard format for all research data.

Status: ACF HPC MPC LPC VIP SKF NPL

e. Manage working copies of all data in a unified, on-line database. (This doesn't necessarily mean a "centralized" database.)

Status: ACF HPC MPC LPC VIP SKF NPL

f. Implement an archive or repository for all historical data on natural habitats.

Status: ACF HPC MPC LPC VIP SKF NPL

7. Facilities

a. What facilities (computers, hardware, software) are the most important to your data management system?

b. List any electronic mail or other network links your site has to the outside world.

8. Evaluation

a. What have been your most important data management accomplishments?

b. What things would you now do differently, if you had them to do over? What suggestions would you give to other sites?

c. What personnel resources do you think are needed to meet your data management goals? Are these resources now available?

d. What additional facilities crucial to your goals (hardware, software, etc.) are lacking?

e. Where do you think additional funding is most needed?

9. Other comments on data management not covered in the foregoing: