

Hybrid Archives

Survey of existing data harvesting/transfer techniques and tools

Document Details

Document type: Project Report
Author: Tony Austin
Draft / Version: Version 1.0
Date of edition completion: 17/06/2004

Contents

Contents	2
Introduction.....	3
Security and Risks	3
Validation.....	4
Physical devices	5
Floppy disk ⁴	5
CD ⁵ including DVD ⁶	5
Zip disk ⁷	6
Tape	6
Hard drives	6
Network technologies	7
Electronic mail.....	7
Shared file systems	7
File Transfer Protocol (FTP).....	7
HyperText Transfer Protocol (HTTP)	8
Sample tools.....	8
Webdrive ²¹	8
Samba ²²	8
rsync ²³	8
wget ²⁴	9
Mirror ²⁵	9
Summary	9
References	11
URLs	11

Introduction

There are essentially two broad methods for moving digital data; the use of physical media and the employment of network technologies. Physical devices are currently the favoured method of data transfer to archives, and organisations such as the Arts & Humanities Data Service and the UK Data Archive dedicate sections of their web site to a description of the media formats that are accepted for deposit. Although exact figures on the use of such transfer methods is difficult to find, if we take the Archaeology Data Service as a typical archive within the arts and humanities sector, a significant majority of deposits (80%) arrive on physical media (figure 1).

mediumType	number	%	physical/network
3.5" floppy disc	21	17	P
5.25" floppy disc	0	0	P
CDROM	77	61	P
DVD	0	0	P
Zip Disk	2	2	P
DAT tape	0	0	P
Hard drive(USB, etc)	0	0	P
FTP	3	2	N
email	22	17	N
HTTP download	1	1	N
Drive mapping	0	0	N

Fig. 1 ADS: Depositor delivery media type (as per Collections Management Database)

The AHDS place emphasis upon removable media formats as the primary method of deposit. It is therefore unsurprising that it comprises 80% of total deposits. Removable media formats deposited with the AHDS typically parallel those in widespread use within the computing industry. The 3.5 inch floppy disk was the prevalent format for deposit five years ago. However, the Archaeology Data Service, has not received a deposit on floppy disk format since 2000. In contrast, the number of deposits provided on CD-ROM media format has grown exponentially. It may be speculated the prevalence of such media is linked to two factors: the falling cost of creation devices and the increase in the amount and complexity of data produced by research projects.

Network technology is currently used in 20% of cases as a transfer method. This statistic is primarily composed of transfers via e-mail. However, this medium is only practical for very small archives as a result of restrictions on the size of mail boxes. It is also common for institutions to restrict the type and size of files transferred using this method. The absence of other network methodologies in a significant number is a possible area of concern. However, the Hybrid Archive model is to be applied to depositors that would find it difficult to deposit their data using other methods. As a result, some degree of familiarity with these transfer methods is to be expected.

Security and Risks

The main threats when transferring data on physical media are theft or loss. Media can also be damaged in transit or affected by adverse environments (e.g. Jones & Beagrie, 2001). There is also the possibility that the data receiver will not have a suitable device for reading the data. Files may have been created using a different operating system such as Mac OS 9

used by the Apple Mac although utilities such as TransMac¹ can help here. There is also variance within device classes and the formats that a device within a class might support.

There are also potential problems when using networks for data transfer. Except for catastrophic events data loss and corruption are largely handled today by the almost ubiquitous use of TCP/IP (Transmission Control Protocol/Internet Protocol) suite of protocols. TCP/IP uses what is known as the 'three-way handshake' to ensure unambiguous communication (see Comer, 1997) in the form of acknowledgements. Transmissions that are not acknowledged because of packet loss, delay, etc are rebroadcast. A further problem concerns bandwidth. Large data transfers may seriously affect low capacity network links and hence other users of that link. This should be very much kept in mind when deciding how to transfer data. Before data transfer can occur a connection has to be made which can be problematic. Destination servers can be down or firewall configuration and passwords changed. A firewall can also enforce a maximum file size that is allowed through. A major problem with some technologies is that unencrypted passwords may be passed through the network. Methods that support encryption should always be used.

Whichever method of transfer is employed copies of any data being transferred should be maintained until data integrity is agreed by both ends of transfer (see below)

Validation

Minimal validation would be to open files following transfer to check that they render what might be expected, for example, a file with a jpg extension should produce an image. Although proprietary (a minimal cost) Quick View Plus² from Stellant™ claims to provide for the viewing of virtually any business document in more than 225 Windows, UNIX, Macintosh, DOS and Internet file formats and is a useful tool for any digital archive. It also provides an option to view files in hexadecimal which provides a hex/byte view of a file. A more sophisticated approach will involve the examination of file header information to confirm the format of a digital object. The JHOVE³ (JSTOR/Harvard Object Validation Environment) project has developed a number of open source tools and modules for the identification of specific formats including versioning as automated processes. As yet the project does not cover all the formats likely to be encountered by a digital archive.

More formalised validation will also involve the generation of checksum or fixity values at either end of the transfer for comparative purposes. There are a number of algorithms available for generating checksums. Many produce checksums that are either a 16 or 32 bit value which may not be enough to ensure a necessarily unique value for a file. MD5 (Message Digest number 5) generates a 128 bit value which effectively supplies a unique value for a file (see Rivest, 1992). MD5 has seen a wide take up. It is in use with many software suppliers so that users can verify downloads. Moving a file should not affect its size or dates associated with it so applying the MD5 algorithm at both ends of a transfer should produce the same value. MD5 is available for most software platforms as Open Source software. The basic tool deals with individual files. To be useful for validating significant data transfers it needs to be built into a utility that can deal with multiple files and directories with MD5 values automatically generated with other metadata such as file name and an organisation identification (i.e. which end of the transfer). This simple metadata from both ends of a transfer can be fed into something like a database for comparison. Any MD5 value that does not occur twice (once for each end of the transfer) is indicative that something has gone wrong with the transfer of a particular file.

A possible alternative is to use file synchronisation software. This is normally used for backup or mirroring web sites and may not be suited to general archival usage in that it can validate a transfer but does not generate something like a checksum value which has future use in determining the ongoing status of a file.

Physical devices

Floppy disk⁴

Floppy disks first became available in the early seventies. Originally 8 inch they were progressively reduced in size to as small as 2.5 inches until finally stabilising on the 3.5 inch disk we know today. Capacity climbed from 250 kb for the first read/write capable drives to a standard 1.44 mb on a double sided and high density formatted disk; a six fold increase. It is fair to say that floppy disks were the mainstay for physically moving data until relatively recently. If packaged properly they are unlikely to sustain physical damage in transit unless exposed to magnetic fields (they are a magnetic media). The main drawback to floppy disks is their capacity. For example, one deposit with the ADS, the Newham archive, consisted of over 250 floppy disks containing several thousand files (see Austin et al, 2001). Clearly the work moving data on and off of the disks at each ends of the transfer was substantial. The development of compression algorithms and disk spanning has helped to reduce such laborious operations in that data is 'zipped ' into a single file which could be spread over multiple but fewer disks.

The 3.5 inch floppy sees less and less usage as high capacity alternatives become a standard PC feature and is likely to follow the once popular 5.25 inch into obscurity with computer manufacturers such as Dell and Apple having already phased out 3.5 inch drives as standard on some models.

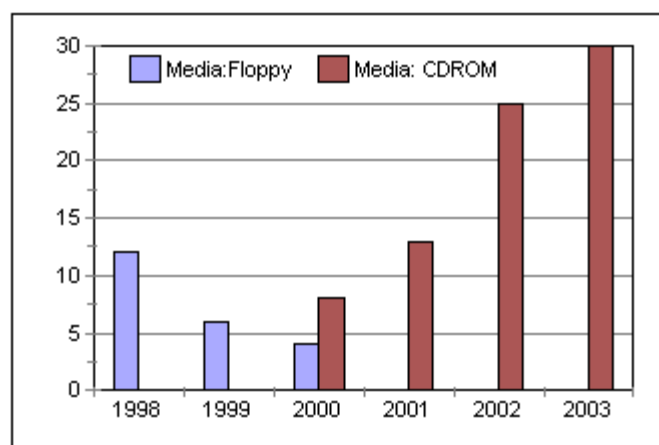


Fig. 2 Demise of the floppy. ADS: Depositor delivery media type by year (as per Collections Management Database)

CD⁵ including DVD⁶

With CD-R (write once read many) and CD-RW (read/write) drives becoming much more popular the CD is becoming the default media for physical data transfer. The first read only CDROMs appeared in the mid eighties. In 1993 CD-R devices appeared followed by CD-RW in 1997. At this time it appeared that DVD would supersede CD technologies but disagreements over standards have held back DVD read/write. Compared to floppy disks CDs and DVDs are much more robust in having no moving parts. Both are optical media in that lasers are used in read/write operations. Hence they do not have the floppy disc susceptibility to damage from magnetic anomalies although some recordable disks use magneto-optical, CD-MO, technology. However, the main advantage is a massively increased capacity. Typically CDROMs can hold 650 mb of data and DVDs around 5gb although DVD-18 which is dual-sided and dual-layered has a capacity of 17gb.

Early CD drives often had manufacturer specific read/ write operations but a core standard for audio known as the 'Red Book' was agreed. There are various extensions to this core including one for computer data known as the 'Yellow Book'. Drives that support these extensions, known as MultiRead capable CD-ROM drives, and current and all future generations of DVD-ROM drive will read discs made by CD-R and CD-RW devices. As already noted this is not the case with DVD drives where disks created using DVD-RAM, DVD+RW (DVD Plus RW) or DVD/RW have incompatible formats. Thus, for the latter, it must be established whether drives at each end of a transfer are compatible.

Zip disk⁷

The Iomega⁸ Zip disk is the leading product from a number of high capacity or 'Super Floppies'. It requires an Iomega drive to read/write to the disk. Iomega claim to have sold over 50 million drives. A Zip Disk has a much higher quality magnetic coating on the disk. Coupled with a read/write head in a Zip drive which is similar to that used in a hard drive means many thousands of tracks can be written per inch compared to just over 100 on a standard floppy. A Zip drive also allows a variable number of sectors per track compared to a fixed number with a conventional floppy and thus a better use of disk space.

Zip disks come with 100, 250 and recently 750 mb capacity. At first glance they would appear to be competitive when compared to CD-R and CD-RW technologies; however, disks are expensive compared to CDROMs. This is perhaps reflected in fig. 1 above where Zip disks account for only 2% of deposits.

Tape

DAT⁹ or Digital Audio Tape is another magnetic media. Magnetic tape preceded disks but has survived largely as a high capacity backup media as defined in 1998 in the DDS (Digital Data Storage) protocol which is a suite of standards (there is a second less common standard; DataDAT). The latest, DDS-4 supports a capacity of 20/40gb and data transfer rates of 2.4/4.8 mbps and is backwardly compatible with earlier specifications. As with all magnetic media it is susceptible to magnetic anomalies. It can also be victim to tape stretch. DAT systems, however, have built in error checking in the form of cyclical redundancy checks and error correction codes for each data packet.

DAT technology uses helical scan recording which is inherently slower than linear recording. The non-proprietary Linear Tape Open or LTO¹⁰ standard was developed by Hewlett-Packard, IBM and Seagate, LTO technology supports linear multi-channel and bi-directional formats. It can provide data transfer rates of up to 40mbps.

As well as backup there is no inherent reason why this media shouldn't be used for moving data into an archive although it would require the latter to have a suitable reader. Tape can be in the form of a cassette or cartridge which eases transfer.

Hard drives

External hard drives have been available for many years with the possibility to effect data transfer by simply moving data on to it on one machine and then transporting it and connecting it to a new machine.

Modern externals can 'plug and play' through the USB¹¹ (Universal Serial Bus) port on Windows 2000 and XP and some Linux installations. Currently these drives are reaching a terabyte in capacity, for example, the laCie Bigger Drive¹² which is physically no bigger than an external CDROM drive. The laCie drive also supports FireWire¹³, a cross-platform, high speed serial data bus.

A new class of miniaturised external USB drives have appeared in recent years that use flash memory and thus are solid state storage devices. Various described as pen drives, thumb

drives, jump drives and key chain drives which gives a good indication of size and portability. They can have a capacity of several gigabytes. These devices could be a useful way for moving data around. Between these and the laCie drive are a whole range of USB drives such as the palm sized Amacom IOdisk¹⁴ with a capacity of up to 80gb. The latter could manage all but the biggest datasets.

Network technologies

Electronic mail

Data can be moved around as email attachments. This, however, is only suited to small digital objects as mail boxes tend to be small. Also firewalls can restrict the size of files and the type of file allowed through.

Shared file systems

Many tools for networked data exchange are based of sharing file systems across a network.

NFS (Network File System) was developed in the late 1980's and is defined in RFC 1094¹⁵. It allows the mounting of a remote file system locally in such a way that it is transparent to users. It is described as machine, operating system, network architecture, and transport protocol independent. It operates within a client/server model with the server said to export a file system. This is often described as a shared file system. Its main use is for supplying file systems to diskless clients or clustering machines but it could be used for efficient data transfer. An archive could make an area on a server available for export. A depositor could mount this locally and move data onto the NFS and then unmount. With the data transferred the archive removes NFS availability. NFS is normally used within a technical environment and may not be suited to general use.

A similar network technology supports CIFS or Common Internet File Systems protocol. CIFS is developed from Microsoft's native file sharing protocol; SMB or Server Message Block protocol (e.g. Leach & Perry, 1996). CIFS piggybacks on TCP/IP and hence a client/server model. It is further described as cross-Platform. More generally IFS or Internet File System tools use other Internet protocols such as FTP, HTTP and WebDav (see below) to mount remote file systems locally. CIFS and IFS-based tools are generally much more user friendly. Some examples are given below.

Synchronisation software will use shared file systems. In mounting a remote file system a full range of commands become available. Tools can be built using copy and compare commands for example. These generally exist as command line tools.

Peer-to-Peer or P2P technology is a large-scale collaborative type of file sharing. Napster the Internet music share site used an hybrid-P2P model. P2P has been linked to digital archiving (Cooper and Garcia-Molina, 2001). Many see inherent weaknesses in such systems¹⁶. Firstly, large scale collaborative creates major security problems. Data integrity is also difficult to ascertain and maintain. Thirdly data components may be offline when required.

File Transfer Protocol (FTP)

FTP is as its full name suggests a set of rules for the movement of files over a network. It predates TCP/IP and hence the Internet and had to be rewritten to fit this client/server architecture. Originally, a command line utility, today there are many GUI (Graphical User Interface) clients available for all Operating Systems (see below). FTP is very useful for moving data around including batch processing but is not used as often as it could be. As already suggested this may be because it sounds technical but could also be because it is not specifically encouraged by archival organisations. There have been concerns over security as

FTP does not encrypt login information. This can be overcome by using FTP under SSH (Secure SHell) which does encrypt. Thus a user uses SSH to connect to a remote host and then FTPs back within this secure session. Care should also be taken to use the right mode; bitstream or ascii for text, or files can be corrupted. Unix FTP has a major problem in that it cannot handle directory structures. A usual way around this is to tar or zip a collection of directories and files and FTP this. This does; however, require intervention at the other end of a transfer to expand the zip or tar file.

HyperText Transfer Protocol (HTTP)

HTTP can be used to transfer data by simply setting up a web page with embedded links for files to be downloaded although this is somewhat cumbersome.

A recent development; WebDav¹⁹ (Web Distributed Authoring and Versioning), uses HTTP extensions to enable distributed web authoring tools and is described as a 'writable, collaborative medium' and a 'network file system'. As such it can be used for data transfer and in using HTTP it has advantages over FTP in supporting strong authentication, encryption, proxy support, and caching. Collaborative environments have their own security implications in that by definition there is group access. The environment also has to be set up. It should be noted that WebDav is currently a proposed standard referenced by RFC 3744²⁰

Sample tools

Webdrive²¹

Webdrive is proprietary but of minimal cost. It is IFS based and includes support for FTP, WebDav, SSH and SSL. Its main selling point is user friendliness in that is very easy to configure and that it maps a remote drive into a familiar Windows Explorer interface alongside other drives available to the user (hard drive, floppy, CDROM, etc). The user can then drag and drop files and directories (including sub directories) to and from the remote drive. In a use scenario a data provider would be given temporary access via SSH to a designated and controlled area on a server at a remote archive. There would be the need for a checksum generator tool

Samba²²

Samba is an example of a CIFS tool allowing mainly Microsoft Windows users to access Unix file systems and print services using TCP/IP. Samba can in fact be used by any SMB/CIFS-enabled client. It can be fairly user friendly at the Client or PC end through the provision of GUI interfaces but is relatively complex to configure and administer at the server end. Samba is open source freely available under the GNU General Public License. It has a similar use scenario to WebDrive.

rsync²³

rsync is a command line file transfer program for Unix (including Linux) systems. It is open source software and is freely available under the GNU General Public License version 2. It is primarily aimed at backup and mirroring in that an algorithm only sends the differences in files across the network. Thus it can be described as providing a very fast method for these purposes. Comparing files normally requires both files to be locally accessible. rsync mounts remote file systems locally. Thus both copies of a file appear local and the various diff commands available under Unix and called by rsync can be used for comparisons. rsync can also be used for simply moving data as a copy of a file does not have to be present at both ends of the transfer.

Thus rsync can be fast; however, it is also expensive in terms of set up in that it requires a C (CC not GCC) compiler to make the executable and a familiarity with compiling in options and using 'make' commands. Also normally used with secure shell protocols such as SSH. Thus it may well require a raft of software installation.

wget²⁴

Again a command line utility which runs on most Unix/Linux operating systems as well as Microsoft Windows. It is open source and distributed under the GNU General Public License. As its name implies data movement is one way rather than the put and get supported by, for example, FTP. Documentation suggests it supports the Internet protocols HTTP, HTTPS and FTP (with limited options). This would imply that a source would have to open up access to an archival organisation which could be problematic. wget has; however, some major advantages. It can harvest data recursively through directory structures which standard Unix FTP can't. Further it is described as 'a utility for non-interactive download of files from the Web' in that it can easily be called from scripts, cron (scheduled) jobs, etc. Also wget was specifically developed to support slow network connections and to retry if a file download fails in its primary role of mirroring web and ftp sites. It supports proxy servers which can lighten network load and speed up retrieval.

Mirror²⁵

Appears to be a similar tool to wget for mirroring websites but with less functionality. Described as 'Very easy to use, good for simple backups or laptop synchronization'.

Summary

Clearly one method of data transfer will not meet archival requirements. The main dictate will be finding common ground between source and destination of the transfer.

Media	Plus	Minus	Usage
Floppy disk		Limited capacity Support dwindling	Probably none
CD	Standardised Commonplace Capacity of 650 mb	Only suited for medium sized transfers	Good for general, middling transfers
DVD	High capacity up to 17gb	Not standardised	Specific where source/destination share standard
Zip	Capacity up to 750mb	Disks expensive compared to CD technology. Not ubiquitous	General, middling transfers where source/destination share technology
Tape	High capacity.	Not well supported beyond backup regimes	Specific where source/destination share standard
Hard drive	Modern 'plug and play' USB drives are both portable and high capacity.	Not 'throwaway' media like the CD	Good for high capacity transfers as USB ports become ubiquitous
Email	Network transfer. Ease of use.	Capacity generally restricted by service providers	Good for very small transfers
Shared drive	Once set up can appear seamless to users	Generally technical (but see tools below). Security implications. Takes up bandwidth	Specific transfers where source/ destination share standards/tools
FTP	Well established. Lots of GUI tools available	Can be technical (command line but see tools). Takes up bandwidth. Possible security issues unless used with SSH. Unix transfers can't transcend directory structures	Good for transfers. With wealth of GUI tools could be used more.
HTTP	Better support for authentication, encryption, proxy services, and caching than FTP	Tools geared toward mirroring or collaborative environments. Takes up bandwidth	Largely untested in an archival environment but could be used for transfers
WebDrive	User friendly PC GUI interface. Supports FTP, SSH, WebDav	Proprietary (minimal cost).	Very good for data transfer.
Samba	Open Source CIFS with GUI client end	relatively complex to configure and administer at the server end	Very good for data transfer
rsync	Open source. Potential for very high speed data transfer (will lose this in simple data transfer)	Technical. Command line tool. Unix to Unix. Primarily aimed at backup and mirroring data	Specific
wget	Open Source Supports FTP (partly), HTTP, HTTPS.Can harvest recursively. Can automate. Download retry Proxy support	One way through use of get from destination. Primary role of mirroring	Specific
Mirror	Easy to use	Primary role of mirroring and synchronisation	Specific

References

- Austin, T., Robinson, D. & Westcott, K. 2001. 'A digital future for our excavated past' in Z Stancic and T Veljanovski (eds) *Computing Archaeology for Understanding the Past: CAA 2000*, BAR International Series 931, ArcheoPress, Oxford pp 289-296
- Comer, D. 1997. *Computer Networks and Internets*, Prentice Hall (USA)
- Cooper, B & Garcia-Molina, H. 2001. 'Creating trading networks of digital archives', 1st ACM/IEEE Joint Conference on Digital Libraries (also online at <http://dbpubs.stanford.edu:8090/pub/2001-23> - downloaded 15 July 2004)
- Hertel, C. 2003. *Implementing CIFS: The Common Internet File System*, Prentice Hall (USA)
- Jones, M. & Beagrie, N. 2001. *Preservation Management of Digital Materials A Handbook*, The British Library
- Leach, P. & Perry, D. 1996. 'CIFS: A Common Internet File System', Microsoft Internet Developer November 1996 (also online at <http://www.microsoft.com/mind/1196/cifs.asp> - downloaded 15 July 2004)
- Rivest, R. 1992. 'The MD5 Message-Digest Algorithm', Network Working Group RFC 1321 (online at <http://theory.lcs.mit.edu/~rivest/Rivest-MD5.txt> - downloaded 5 July 2004)

URLs

- 1 TransMac <http://www.asy.com/scrtm.htm>
- 2 Quick View Plus
http://www.stellent.com/stellent3/idcplg?IdcService=SS_GET_PAGE&nodeId=66
- 3 JHOVE <http://hul.harvard.edu/jhove/>
- 4 Floppy disks <http://www accurite.com/FloppyPrimer.html>
- 5 CDROM <http://www.pctechguide.com/08cd-rom.htm>
- 6 DVD <http://www.pctechguide.com/10dvd.htm>
- 7 Zip drive/disks <http://computer.howstuffworks.com/question277.htm>
- 8 Iomega <http://www.iomega.com/global/>
- 9 DAT tape <http://www.pctechguide.com/15tape.htm#DAT>
- 10 LTO tape http://www.pctechguide.com/15tape2.htm#Linear_Tape_Open
- 11 USB <http://www.usb.org/faq>
- 12 FireWire <http://developer.apple.com/firewire/>
- 13 LaCie Bigger Disk http://www.lacie.com/download/datasheets/BiggerDisk_triple_en.pdf
- 14 Amacom IOdisk http://www.amacom-tech.com/google_iodisk.html

- 15 RFC 1094 <http://www.cse.ohio-state.edu/cgi-bin/rfc/rfc1094.html>
- 16 P2P critique http://www.solution6.com/ne_newslst_detail_nbf.asp?id=1390
- 17 FTP <http://www.w3.org/Protocols/rfc959/>
- 18 SSH <http://aset.its.psu.edu/internet/ssh/>
- 19 WebDav <http://www.webdav.org/>
- 20 RFC 3744 <http://www.webdav.org/specs/rfc3744.html>
- 21 WebDrive <http://www.webdrive.com/>
- 22 Samba <http://www.samba.org/>
- 23 rsync <http://samba.anu.edu.au/rsync/>
- 24 wget <http://www.gnu.org/software/wget/>
- 25 Mirror <http://home.in.tum.de/~jain/index.html>