OUTLINE OF ISSUES RELATING TO ELECTRONIC RECORDS
Brian Nelson Burford
NH State Archives

| General Issues | Sub-Issue | Specifics or Examples |
|---|---|---|
| Definitions | | |
| | What is a "record"? | Durable information showing that an event happened, or a decision was made, and created at or soon after the event or decision by a person who witnesses the event or decision. |
| | What is an electronic file? | Information stored as a combination of machines that manipulate physical characteristics of media in such a way as to store meaning for later use. Electronic Records require a combination of hardware and software to interpret this information. The media is organized in such a way as to inform the machine if the "switch" is on or off (also expressed as 0s and 1s, or binary data). Each "switch" is called a bit, and they are organized in groups of 8 called a byte. A file may contain instructions (platform or application) or data to be manipulated according to the application instructions. |
| | What is Digital Preservation? | The ability to maintain data (zeros and ones) in a structure over time in such a way to allow one to reinterpret those zeros and ones with the same meaning as when the data was created. The purpose is the preservation of the information preserved in the "zeros and ones" for later use by humans. |
| Sources of Data | | |
| | Capture born digital from sources inside organization (government or corporate department, division, or other agency) | files created by users on corporate network |
| | Capture web-sites | documenting information available for the public to utilize over time |
| | On-line forms | information supplied by users of a form |

| | | |
|---|---|---|
| | | located on a web-site |
| | Import from outside sources | files created by users outside firewall, and imported |
| | conversion of magnetic to digital | e.g., VHS to digital |
| | Scan existing fixed documents (e.g., paper, microfilm, sound recordings, motion picture film) to raster image, sound or video files | |

| | | |
|---|---|---|
| Hardware | | |
| | CPU speeds | |
| | internal or external storage | |
| | devises | floppy drives, CD readers and writers; DVD readers and writers |
| | documentation | |
| Software | | |
| | Programming Languages | Java<br>C++<br>Perl<br>Python (an *interpreted, interactive, object-oriented* programming language)<br>HTML<br>XML  (universal language)<br>[many, many others] |
| | Platform | Microsoft (MS) DOS[2]<br>Unix[2]<br>Linux[1]<br>FreeBSD[1]<br>NetBSD[1]<br>Debian GNU/Linux[1] |
| | Desktop | Windows XP, NT, 2000[2]<br>Gnome[1]<br>KDE[1] |
| | Server | MS Server[2]<br>Apache[1] |
| | Applications (examples) | MS Word  (text)[2]<br>MS Excel  (spreadsheet)[2]<br>MS Outlook  (email, a form of text)[2]<br>Mozilla (email program)[1]<br>MS Access  (database)[2] |

| | | |
|---|---|---|
| | | MSSQL $(database)^2$<br>Oracle $(database)^2$<br>MySQL $(database)^1$<br>PostgreSQL $(database)^1$<br>DSpace $(database)^1$<br>Greenstone (digital library) $^1$<br>JHove (metadata extractor)$^1$<br>Manifold $(GIS)^1$<br>AutoCAD $(CAD)^2$<br>Ascent Capture (Electronic Record Management [ERM])$^2$<br>Laserfiche $(ERM)^2$<br>Alchemy (Object database)$^2$<br>Veritas Electronic Vault$^2$ (electronic record archiving software)<br>iLumin (email archiving software)$^2$<br>Zantaz EAS$^2$ (email archiving software)<br>University of Michigan Digital Library eXtension Service (DLXS) (digital library)$^1$<br>ePrints (Institutional Repository software)$^1$<br><br>*In the list of applications above, the trailing superscript number indicates Open$^1$ or Proprietary$^2$* |
| | Open$^1$ vs Proprietary$^2$ | Open Source: the source code is made publicly available and a general license to use software is granted for future use.<br><br>Proprietary: the computer code is owned by a private entity, and future rights of use and changes are controlled by that private entity. |
| | Versions & backwards compatibility | |
| | | new features which interpret data differently |
| | | could effect "presentation" or appearance of document – and hence our interpretation of the document (e.g., formatting: columns, paragraph indentations, some characters, etc) |

| | | documents containing inserted objects (photos, graphs) |
|---|---|---|
| | documentation | instructions on how software works, required hardware environment, and how it interprets data |
| Data | | |
| | File Formats | *Notice that this list is very similar to the "application" list above, because specific applications usually create a data file in a specific open or proprietary structure or format. The preservation issue is, in fact, whether data files created by one application can be properly interpreted by another application program. A FORMAT is data arranged in a specific way following systematic rules, so the data interpreter will understand the rules and know how to retrieve the knowledge contained therein.*<br><br>ASCII (text)<br>txt (text)<br>RTF (rich text format)<br>DOC (MS Word text)<br>PDF/A (image) PDF version 1.4<br>    [ISO 19005-1.]<br>    security hole?<br>    PDF files after 1.4?<br>TIFF (raster image)<br>JPEG2000 (raster image)<br>GIF (raster image)<br>BMP (bitmap – raster image)<br>PNG (personal network graphics – raster?)<br>AI (Adobe Illustrator) (vector image)<br>CDR (CorelDRAW) (vector image)<br>CMX (Corel Exchange) (vector image)<br>CGM Computer Graphics Metafile (vector)<br>DXF (AutoCAD) (vector)<br>WMF Windows Metafile (vector)<br>MPEG (motion image)<br>RealAudio (sound)<br>WAVE (sound)<br>MP3 (sound)<br>ABC Amber Converter (Conversion) |

| | | Xena (XML intermediate)<br><br>raster is "resolution dependent" – number of pixels. The more pixels, the clearer the image.<br><br>Vector defines shapes and combinations of shapes mathematically. |
|---|---|---|
| | Multi-formatted file structures | For example, TIFF files may have IPTC or XMP metadata formatted within it, meaning one file, but it needs two or three different software "translators" to make sense of it. Also, GIS ESRI 3 files have three separate bit-streams |
| | Optical Character Recognition (OCR) | ABBYY<br>Finereader<br>Prime Recognition |
| | Compression | Loss-less such as Group 4 *(i.e. TIFF4)*<br>Lousy such as the LZW rules<br>Mixed |
| | Aggregation | Zip files |
| | Authenticity & Certification | malicious alteration;<br>unintended alteration *(e.g., during migration);*<br>data corruption<br>checksums *(sum of zeros and ones in the original data)*<br>hash Values *(value derived from the relative positions of the zeros and ones in the original data)* |
| | File Naming Conventions | to ease recovering info;<br>complexity begins when multiple users are creating files |
| | Dynamic Records | indexing changes to file to show earlier versions of data |
| | Organization & Indexing | Object database<br>Folder hierarchy<br>ERM (electronic record management program) |
| | Corruption | natural decay of magnetic medium;<br>chemical reactions in CDs or DVDs;<br>physical damage to medium |
| | "bit-level" preservation | keeping the digital structure of the original electronic document |

|  |  |  |
|---|---|---|
| Metadata |  |  |
|  | Standards | Dublin Core; <br> modified Dublin Core <br> Australian Govt. |
|  | Capture method | automated; <br> user determined; <br> mixed |
|  | Method to tie metadata to data file | XML wrapper |
|  | Thesaurus | used to assist indexing or searching on metadata <br> Protégé[1] (Stanford) <br> TheW32[2] (Tim Craven) – assists user to create a Thesaurus. |
| Security |  |  |
|  | Firewalls | control of access |
|  | Encryption | control of access |
|  | Read-Only | limitation of access |
|  | Malware | Viruses; <br> Worms; <br> Trojans; <br> Spyware; <br> sometimes embedded or inserted in applications or data files |
|  | File Sharing |  |
|  | Electronic Signatures |  |
| Preservation Storage Media |  |  |
|  | Magnetic | hard drive(s) <br> CDs (not all are magnetic) <br> *there are even differences between magnetic (writable) CDs—gold or silver CDs have longer "life" expectancy* <br> DVDs (not all are magnetic) <br> floppy disks *(very short expectancy – 2 years from manufacture?)* <br> tape <br> storage sticks |
|  | Optical Disks | (also called WORM – "Write Once Read Many" – but this is not the same as Worms listed under Security). Optical storage works on the technology of bounced light off uneven reflective surface – more durable than magnetic |

| | | |
|---|---|---|
| | Fixed | Computer Output to Microfilm *(COM)* print to paper |
| | Mixed | Information Life-cycle Management *(ILM) – tiered according to value; uses paper, film, electronic according to need* |
| | | |
| Preservation Strategy | | |
| | Retention Schedules | NH Municipal Records Board RSA |
| | Migration | open the file in a newer version, allowing that program to change the data-structure to the new requirements |
| | Encapsulation | storing both the data and the application files together, perhaps allowing the future use of the application program to open and interpret the data |
| | Emulation | writing an application program that translates former application instructions into current instruction for data interpretation (similar to "drivers" which translate instructions for hardware devises) |
| | "Computer Museum" | maintaining working examples of all hardware and software in order to run the programs and data files into the future – requires being able to obtain "antiquated" parts which wear out. And requires knowledge of how to use program (so save the documentation). Requires storage space. |
| | Hybrid Technology | Computer Output on Microfilm (COM) or scanning film images [computer (text, raster, vector, spreadsheet) to microfilm to computer (raster)] |
| | LOCKS | "Lots of copies keeps it safe" – many copies distributed to various institutions improves chances of survival ("playing the odds") |
| | Periodic "inspection" for data loss | quality control. Whether microfilm, magnetic, or optical. |

|  | Periodic medium refresh |  |
|---|---|---|
|  | Checksums; error-correcting codes | audits of the bit-streams to ensure the data is transferred to new medium intact. |
|  |  |  |
| Access |  |  |
|  | Right-To-Know (a.k.a., Freedom of Information) | RSA 91A |
|  | Privacy | This is for several reasons including identity theft or to prevent abuse of personal information (e.g., HIPPA, Patriot Act.) |
|  | Intellectual Property Rights | ownership of the original medium; copyrights to information/interpretation |
|  | Discovery | court-ordered during legal suits; requirement to preserve and share ALL information described in writ |
|  | Forensics | methods used to recover data thought to be deleted or lost<br><br>Restorer2000<br>WordFIX (Cimaware)<br>AccessFIX (Cimaware)<br>ExcelFIX (Cimaware)<br>OfficeFIX (Cimaware)<br>OfficeRecovery (http://www.officerecovery.com/)<br>OnTrack |
|  |  |  |
| Transfer of Data |  |  |
|  | Compression | Lossless vs. lousy |
|  | Security | encryption<br>checksums & hash values |
|  | Methods | CDs or DVDs<br>File Transfer Protocol (FTP) |
|  |  |  |
| Disaster Recovery |  |  |
|  | back-ups | kept at a different site |
|  |  |  |
| Disposal |  |  |
|  | Delete key | erases pointer, but not data until overwritten |
|  | Reformat disc | Formatting a magnetic medium |

| | | prepares the medium for use by the operating system. The surface of the disk is checked for physical and magnetic defects, and then an address structure is added (e.g., FAT, NTFS, etc). The addresses may be made up of tracks, sectors, clusters, cylinders, etc. Included is a root directory, which is a list of addresses. |
|---|---|---|
| | Software scrub | overwrite with random data[3]– usually multiple overwrites. One theorist (Peter Gutman in Australia) recommends up to 35 overwrites.<br><br>SDelete |
| | Degaussing | erase all magnetic info on entire disc with a powerful magnetic |
| | Destruction of medium | physically destroy the hard drive, floppy disc, CD or DVD – most secure |

created June 8, 2005
revised July 12, 2005; July 13, 2005
minor addition Feb 25, 2006; March 14, 2006, March 20, 2006
Revisions April 4, 2006: in medium, file transfer, applications, XML

Brian Nelson Burford
NH State Archives

---

[3] "random" sequences are not technically possible in computers, and are really only QUASI-random. Hence, because patterns may arise, "randomization" is not an absolutely secure form of data destruction.