Preservation and Management Strategies for Exceptionally Large Data Formats: 'Big Data'

Review of nature of technologies and formats

Tony Austin & Jen Mitcham 10 May 2007

This report has been produced as part of the Big Data Project. It is a technical review of each of the 'Big Data' technologies currently practised by archaeologists with a consideration of data formats for preservation and future dissemination. As well as data acquisition there will be an analysis phase to any project. Survey normally involves a series of traverses over a spatially defined area. Composite mosaics can be produced as either part of acquisition or as part of post processing. The composite can then be fed into a range of geospatial tools including 3-D visualization. Examples include Geographical Information Systems (GIS) and Computer Aided Design (CAD) software.

Discussion as indicated by the Big Data questionnaire¹ and the project case studies² focuses on the following technologies

- Sonar (single beam, bathymetry and sub bottom profiling)
- Acoustic Tracking
- 3D Laser Scanning
- Geophysics
- Geographic (eg GIS)
- LiDAR
- Digital Video

Raster (still) images and Computer Aided Design (CAD) also featured in the questionnaire but are covered more than adequately elsewhere. See, for example, the recent AHDS **Digital Image Archiving Study**³ and the **CAD: A Guide to Good Practice**⁴

¹ <u>http://ads.ahds.ac.uk/project/bigdata/survey_results/bigdata_quest_final.doc</u>

² <u>http://ads.ahds.ac.uk/project/bigdata/caseStudies.html</u>

³ <u>http://www.jisc.ac.uk/uploaded_documents/FinaldraftImagesArchivingStudy.pdf</u>

⁴ <u>http://ads.ahds.ac.uk/project/goodguides/cad/</u>

A tabular summary of Big Data formats can be found at the end of this document (table 1). This structure has been adopted rather than considering formats under each technology as the formats often span technologies. For example, SEG Y available in a number of maritime applications (see table 1) is a generic seismic survey format. This summary does not pretend to be inclusive but rather a representative flavour of the vast range of formats that seem to be associated with Big Data.

Sonar

Sonar (**SOund NA**vigation and **R**anging) is a simple technique used by maritime archaeologists to detect wrecks. It uses sound waves to detect and locate submerged objects or measure the distance to the floor of a body of water and can be combined with a Global Positioning System (GPS) and other sensors to accurately locate features of interest. A useful overview or Maritime Survey techniques can be found on the Woods Hole Science Center (part of the Unites States Geological Survey) website⁵.

Bathymetry (single beam and multibeam sonar)



Illustration: top view of the multibeam data of Hazardous, lost in November 1706, when she was run aground in Bracklesham Bay © Wessex Archaeology

Single beam scanning sends a single pulse from a transducer directly downwards and measures the time taken for the reflected energy from the seabed to return. This time is multiplied by the speed of sound in the prevalent water conditions and divided by two to give the depth of a single point.

Multibeam sonar sends sound waves across the seabed beneath and to either side of the survey vessel, producing spot heights for many thousands of points on the seabed as the vessel moves forward. This allows for the production of accurate 3D terrain models of the sea floor from which objects on the seabed can be recorded and quantified. Wessex Archaeology used multibeam bathymetry during the Wrecks on the Seabed project (Big Data

⁵ <u>http://woodshole.er.usgs.gov/operations/sfmapping/index.htm</u>

case study)⁶. As well as the raw data itself, 3D terrain models, 3D fly through movies and 2D georeferenced images were created. The 2D images were then used as a base for site plans and divers were able to use offset and triangulation to record other objects on to the plans.

The data

Why should we archive?

For future interpretation of data. Seeing anomalies in the results not seen before? For monitoring condition and erosion of wreck sites For targeting areas for future dives/fieldwork

Problems and issues

Many bathymetric systems use proprietary software. The extent to which this software supports open standards or openly published specifications is largely unknown. Data exchange between systems may also be problematic.

Specialised metadata

Metadata to be recorded alongside the data itself includes: Equipment used (make and model) Equipment settings Assessment of accuracy? Methodology Software used Processing carried out

Associated formats include

Generic Sensor Format (.gsf), HYPACK (.hsx, .hs2), MGD77 (.mgd77), eXtended Triton Format (.xtf), Fledermaus (.sd, .scene – visualisation)

⁶ http://www.wessexarch.co.uk/projects/marine/alsf/wrecks_seabed/multibeam_sonar.html

Sidescan sonar



Illustration: This image created with sidescan data clearly shows a ship wreck protruding from the seabed $\ \odot$ Wessex Archaeology

Sidescan sonar is a device used by maritime archaeologists to locate submerged structures and artefacts. The equipment consists of a 'fish' that is towed along behind the boat emitting a high frequency pulse of sound. Echoes bounce back from any feature protruding from the sea bed thus recording the location of remains. The sidescan sonar is so named because pulses are sent in a wide angle, not only straight down, but also to the sides. Each pulse records a strip of the seabed and as the boat slowly advances, a bigger picture can be obtained. As well as being a useful means of detecting undiscovered wreck sites, sidescan data can also be used to detect the extents and character of known wrecks.

The data

The data tends to be in a wide range of little known proprietary and binary formats. Although there are some open standards such as SEG Y around. The software packages associated with sidescan sonar may support ASCII or openly published binary exports.

Why should we archive?

For future-interpretation of data. Seeing anomalies in the results not seen before? For monitoring condition and erosion of wreck sites For targeting areas for future dives/fieldwork

Problems and issues

Many sidescan systems use proprietary software. The extent to which this software supports open standards or openly published specifications is largely unknown. Data exchange between systems may also be problematic.

Specialised metadata

Metadata to be recorded alongside the data itself includes: Equipment used (make and model) Equipment settings Assessment of accuracy? Methodology Software used Processing carried out

Associated formats include

eXtended Triton Format (.xtf), SEG-Y, CODA (.cod, .cda), Q-MIPS (.dat), HYPACK (.hsx, .hs2), MSTIFF (.mst)

Sub bottom profiling



Illustration: example of sub-bottom profiler data © Wessex Archaeology

Powerful low frequency echo-sounders have been developed for providing profiles of the upper layers of the ocean bottom. Specifically sub-bottom profiling is used by marine archaeologists to detect wrecks and deposits below the surface of the sea floor. The buried extents of known wreck sites can be traced using an acoustic pulse to penetrate the sediment below the sea bed. Echoes from surfaces or the horizons between different geological layers are returned and recorded by the profiler and the sequence of deposition and subsequent erosion can be recorded. The case study, Wessex Archaeology, utilised sub bottom profiling⁷ for the Wrecks on the Seabed project

The data

The data tends to be in a wide range of little known proprietary and binary formats. Although there are some open standards such as SEG Y around. The software packages associated with sub bottom profiling may support ASCII or openly published binary exports.

Why should we archive?

For future-interpretation of data. Seeing anomalies in the results not seen before? For monitoring condition and erosion of wreck sites For targeting areas for future dives/fieldwork

Problems and issues

Many systems use proprietary software. The extent to which this software supports open standards or openly published specifications is largely unknown. Data exchange between systems may also be problematic.

⁷ <u>http://www.wessexarch.co.uk/projects/marine/alsf/wrecks_seabed/sub_bottom_profiler.html</u>

Specialised metadata

Metadata to be recorded alongside the data itself includes: Equipment used (make and model) Equipment settings Assessment of accuracy? Methodology Software used Processing carried out

Associated formats include

CODA (.cod, .cda), QMIPS (.dat), SEG Y (.segy), eXtended Triton Format (.xtf)

Acoustic Tracking



llustration: diagram showing how acoustic tracking devices keep track of the divers' location at any one time $^{\odot}{\rm Wessex}$ Archaeology

Acoustic tracking can be used to keep a log of a diver's location throughout the dive. Sound signals are emitted by a beacon attached to the diver and picked up by a transceiver attached to the side of the boat. The relative position of the diver underwater can be calculated and these relative co-ordinates can be used to calculate an absolute location for the diver. Additional equipment may be needed to compensate for the motion of the vessel in the water. Acoustic Tracking was utilised for the Wrecks on the Seabed project⁸.

The data

Normal practice is to use a data logger for collection. Generally the data will be in the form of structured ASCII text. As such it will be easy to import into other packages such as a GIS or database. Wessex Archaeology supplied their Acoustic Tracking data as a Microsoft Access database

Why should we archive?

For Wessex Archaeology this data was seen as crucial to the project archive as it sets much of the other maritime archaeology project data in context. Will need to refer to this database to establish where the diver was when individual photographs were taken, segments of digital video recorded or general observations made.

Problems and issues

Possibly processed and not the raw data.

Specialised metadata

Metadata to be recorded alongside the data itself includes: Equipment used (make and model)

⁸ <u>http://www.wessexarch.co.uk/projects/marine/alsf/wrecks_seabed/acoustic_diver_tracking.html</u>

Equipment settings Assessment of accuracy Methodology Software used Processing carried out

Associated formats include

ASCII text formats

3D Laser Scanning



Illustration: Solid model created from point cloud laser scan data from stone 7 of Castlerigg Stone Circle in Cumbria - image from Breaking Through Rock Art Recording project [©] Durham University

There are a wide variety of applications of laser scanning as a tool for capturing 3D survey data within archaeology. A common application of this technology is as a tool for recording and analysing rock art, but subjects can range from a small artefact to a whole site or landscape. A 3D image of Rievaulx Abbey was recently created by Archaeoptics in 10 minutes. The benefit of this technique is that a visually appealing and reasonably accurate copy of a real world site or object can quickly be created and manipulated on screen.

When a laser scanner is directed at the subject to be scanned, a laser light is emitted and reflected back from the surface of the subject. The scanner can then calculate the distance to this surface by measuring the time it takes, and x, y and z points relative to the scanner can be recorded. Absolute co-ordinates can then be created by georeferencing the position of the scanner. Some scanners may also record colour values for each point scanned and the reflection intensity of the surface (see Trinks *et al*, 2005⁹)

Huge datasets are produced using this technology. The recent project by Wessex Archaeology and Archaeoptics to scan Stonehenge reported that each scan took "3 seconds to complete and acquiring 300,000 discrete 3D points per scan. A total of 9 million measurements were collected in just 30 minutes" (see Goskar *et al*, 2003¹⁰). It is not surprising that laser scanner data files can be many gigabytes in size.

The data

Also online at http://www.britarch.ac.uk/ba/ba73/feat1.shtml

⁹ Trinks I, Díaz-Andreu M, Hobbs R & K. Sharpe, K. 2005. 'Digital rock art recording: visualising petroglyphs using 3D laser scanner data', Rock Art Research 22, p. 131-9 Also online at <u>http://www.dur.ac.uk/m.diaz-andreu/articles/2005_RAR_Trinks_et_al.pdf</u>

¹⁰ Goskar T, with Carty A, Cripps P, Brayne C, & Vickers D. 2003. 'The Stonehenge Lasershow', British Archaeology 73

There are a number of different types of data that are created as a laser scanning project progresses:

Primary data produced through this technique is point cloud data. Point clouds essentially consist of raw XYZ data, to locate each point in space, plus if recorded, RGB data to record the colour of each point.

Firstly there are the raw observations as collected by the scanning equipment in a number of different proprietary formats.

Numerous scans may be carried out to record a complex subject, with the scanning equipment moved to a different position each time. This will create a large number of data files. All of these individual scans would then need to be stitched together in order to create a composite mosaic of the whole subject. From the point cloud data, it is possible to create a solid model of the subject, such as that illustrated above. A cut down or decimated version of the raw XYZ data may be used to create a dataset of a more manageable size for processing, viewing and analysing.

Why should we archive?

In archaeology it is thought that perhaps one of the main opportunities we will gain from storing and re-using this data in the future is that successive scans of the same sites may be used to monitor erosion or other physical changes to the site. The Durham University Fading Rock Art Landscapes project¹¹ was set up with just this in mind.

Data could also be re-processed in different ways to create new models and allow for new interpretations of the data. With technologies such as this it is very easy to create large datasets from a high resolution scan and then be hampered by a lack of storage space and processing power when attempting to view and interpret the resulting dataset.

Problems and issues

There are a fairly small range of software tools for viewing laser scan data. Huge file sizes may hamper reuse - may be more appropriate for researchers to interrogate a cut down or decimated version of the laser scanning data as this will be easier to process. No standard data format currently exists for laser scanning data. This should be addressed.

Which data do we actually need to archive? The raw data as created by the laser scanner? ie: a separate file for each scan - not yet combined to produce full composite scan of whole object. Or perhaps the composite scan is fine - will this have undergone additional processing? Processed results are also useful. If theories were reached as a result of looking at a particular version of

¹¹ <u>http://www.dur.ac.uk/prehistoric.art/</u>

the dataset, is it worth keeping this also so future researchers can see how a particular theory came about?

Specialised metadata

Re-use potential is maximised if relevant metadata exists for laser scanning data. Both technical information about the survey and more obvious information about the context of the scan are required. Lists are published by Heritage3D¹² include:

Date of capture Scanning system used Company name Monument name Weather during scanning Point density on the object Technical information relating to the scanning equipment itself - may include triangulation, timed pulse, phase comparison

Associated formats include

XYZ (.xyz), Visualisation ToolKit (.vtk - processed), LAS (.las), Riscan Pro (.3dd), National Transfer Format (.ntf), OBJ (.obj), Spatial Data Transfer Standard (various), Drawing eXchange Format (.dxf - processed)

¹² <u>http://www.ceg.ncl.ac.uk/heritage3d/downloads/TLS%20formats%20V1.pdf</u>

Geophysics



As stated in the ADS **Geophysical Data in Archaeology Guide to Good Practice**¹³, the increasing size and sampling resolution of geophysical surveys in archaeology is resulting in the accumulation of increasing quantities of data. However the most common techniques, resistivity and magnetometer surveying generally do not produce datasets that are large enough to fall under the remit of the Big Data Project. The one land-based geophysical technique that can produce exceptionally large datasets is Ground Penetrating Radar.

In a Ground Penetrating Radar survey, the instrument is dragged along the ground at a constant speed and electromagnetic pulses are sent into the ground by an antenna. As the pulses come into contact with objects and layers within the ground, they are reflected back to the instrument picked up by a receiving antenna. A single GPR transect creates a vertical 2D image of the subsurface features. If numerous transacts are carried out within a grid, these images can be combined to create a 3D depiction of the results. The amount of data being collected using Ground Penetrating Radar is likely to increase as the technology moves along. For example, Terravision¹⁴ by Geophysical Survey Systems, Inc (GSSI) is a new and very advanced piece of equipment for carrying out GPR survey. It features a 14 antenna array and 6 foot wide survey path with a data collection speed of up to 10mph. The internal storage capacity of the equipment can be up to 32GB. The quantities of data that could be produced if archaeologists started to use equipment like this would be substantial.

The data

Basically these will be x, y and z co-ordinates. A good overview of GPR data formats is available on one of the USGS websites including a number of format specifications¹⁵. GPR data from the Where Rivers Meet case study was in the proprietary and binary DZT format. This stores data as radians

¹³ http://ads.ahds.ac.uk/project/goodguides/geophys/

¹⁴ <u>http://www.geophysical.com/TerraVision.htm</u>

¹⁵ <u>http://pubs.usgs.gov/of/2002/ofr-02-0166/ofr02_166.pdf</u>

Why should we archive?

For future-interpretation of data. Seeing anomalies in the results not seen before?

For monitoring condition and erosion of the archaeology For targeting areas for future fieldwork

Interpretation of GPR data can be fairly subjective, users wishing to re-use the GPR data may wish to go back to raw results rather than use someone elses subjective interpretations

Problems and issues

Metadata! If the raw GPR data is to be archived and re-used it has limited value without the field notebooks used to record location of each transect. This metadata will most probably be in a paper format and would require some time and effort to digitise and rationalise for general re-use. Specialised metadata Metadata to be recorded alongside the data itself includes: Equipment used (make and model) Equipment settings Assessment of accuracy? Methodology Software used Processing carried out Check this against G2GP, pres handbook, proc doc Data formats/file formats

Associated formats include

DEM (.dem), DLG (.dlg), RADAN[™] DZT (.dzt), NTF (.ntf), RAMAC – RD3/RAD (.rd3, .rad), Fledermaus (.sd, .scene – processed), SDTS (various plus .ddf), SEG 2 (.dat, .sg2), SEG Y (.segy)

Geographic (eg GIS)



Illustration: High Resolution raster images such as this aerial photograph of Aberystwyth can be used within a Geographic Information System. Image taken from **Mapping Medieval Townscapes: a digital atlas of the new towns of Edward I** project¹⁶

Though many archaeologists use GIS without creating exceptionally large datasets, large file sizes can be an issue for some GIS projects. High resolution raster layers such as scanned aerial photographs, satellite imagery and sometimes digital elevation models (DEM) within a GIS can be very large. An high quality scanned colour aerial photo at a 'ground-resolution' of, 20 cm per pixel, could be as large as 250 MB. Hundreds of these photos could be needed to give coverage of a large study area such a county. This would have huge implications for the future archiving and re-use of the project data. See the AHDS GIS Guide to Good Practice¹⁷

The data

Data can consist of a georeferenced high resolution raster images such as a geotiff or vector x,y,z data. Data is generally collected using other technologies and processed within a GIS to generate project outcomes.

The established preservation strategy for GIS vector data is to migrate to ESRI Shape (SHP) and Export (E00) formats which are generally seen as *de facto* standards in light of alternatives¹⁸. However, the development of the Open Source GDAL Geospatial Data Abstraction Library¹⁹ is allowing popular GIS formats to be abstracted to GML, an XML based language. This is being adopted as a preservation strategy by the ADS. The process can be reverse engineered if required.

Why should we archive?

For future interpretation of data. Seeing anomalies in the results not seen before?

¹⁶ http://ads.ahds.ac.uk/catalogue/specColl/atlas_ahrb_2005/

¹⁷ <u>http://ads.ahds.ac.uk/project/goodguides/gis/</u>

¹⁸ http://www.jisc.ac.uk/uploaded_documents/FinaldraftImagesArchivingStudy.pdf 5.4

¹⁹ <u>http://www.gdal.org/</u>

For monitoring condition and erosion of landscapes For targeting areas for future fieldwork

Problems and issues

Copyright?! Satellite imagery and aerial photographs will often have been obtained from other organisations and are not necessarily owned by the project. Specialised metadata How data acquired/created Copyright Processing carried out

Associated formats include

GEOTIFF (.tiff), TIFF World file (.tfw), JPEG World file (.jgw), DEM (.dem), ESRI Shape file (.shp), ESRI export (.e00), MOSS export (.exp), GML (.gml), GRASS (various), MapInfo (.tab), MapInfo interchange Format (.mif), IDRISI raster (.rst, rdc), IDRISI vector (.vct and others)

Lidar



Illustration: LiDAR image of Hambledon Hill showing the hillfort with an internal long barrow © Environment Agency.

The name LiDAR comes from 'Light Direction And Ranging'. The technique²⁰ ²¹ involves scanning a pulsed laser along the ground from an aircraft. By monitoring the direction and speed of the incoming reflected pulse, the layout of the landscape can be recorded with some accuracy. Once xyz co-ordinates have been collected in this way, they can be used to generate a digital elevation model of the study area that can be viewed within a GIS package. The benefits to archaeologists are obvious²². Any undulations in the ground surface will be recorded, thus even very slight rises and depressions of archaeological earthworks will be visible within the elevation model. The researcher is able to alter the light direction and intensity within the model to look for features which would otherwise be hard to pick out. English Heritage²³ reports that LiDAR can measure between 20,000 to 100,000 points per second. It is therefore not surprising that exceptionally large datasets are rapidly created.

Lidar data is normally supplied by specialist organisations; usually commercial, although it is worth noting that some infoterra²⁴, a leading supplier of 'geoinformation solutions' data is available to academic users through MIMAS (Manchester Information & Associated Services) under a CHEST agreement²⁵. Also some material is available from the Environment Agency²⁶ and the same for academic use through MIMAS (Manchester Information & Associated

²⁰ http://aarg.univie.ac.at/events/aarg2000/abstracts/abstracts2000.html#72

²¹ <u>http://www.english-heritage.org.uk/server/show/nav.8730</u>

²² http://aarg.univie.ac.at/events/aarg2000/abstracts/abstracts2000.html#73

²³ <u>http://www.english-heritage.org.uk/server/show/nav.8730</u>

²⁴ <u>http://www.infoterra.co.uk/</u>

²⁵ <u>http://landmap.ac.uk/download/infoterra_download.htm</u>

²⁶ http://www.environment-agency.gov.uk/science/monitoring/131047/?lang= e

Services)²⁷ but this turns out to be mostly processed data either as raster images or in GIS formats.

The data

Raw data as recorded from the aeroplane is slightly different, containing time of flight of the laser and position of the plane in WGS84²⁸ (World Geodetic System) coordinates. This is normally converted into a local system such as OSGB36²⁹ xyz which is much more useful. It is not believed that there is any benefit in archiving the actual raw data.

XYZ ASCII files (point clouds) are referenced to a coordinate system such as the British National Grid. Data may also include an i value to record the intensity of the returned signal. Two data files will be produced; first pulse and last pulse. The first pulse data records the highest point or the first pulse returned to the aircraft. This could be the top of a tree for example. Last pulse data records the last pulse back to the sensor and this would represent the height of the ground underneath the tree. Last pulse data could therefore be used to create Digital Terrain Models (DTM = DEM or Digital Elevation Model) of the landscape.

Why should we archive?

For future-interpretation of data. Seeing anomalies in the results not seen before? For monitoring condition and erosion of landscapes For targeting areas for future fieldwork

Problems and issues

Copyright of data. Archaeologists are unlikely to be able to afford the equipment to carry out their own LiDAR survey so will either purchase existing data or get contractors to carry out the survey for them. This was the case with one of the Big Data case studies, Where Rivers Meet, where the data was supplied by infoterra. Clearly this may prevent reuse by or archiving with a third party

Specialised metadata

Metadata to be recorded alongside the data itself includes: Most important metadata to accompany LiDAR data Instrument make and model (thus wavelength of laser) Altitude of flight Effective point spacing (ie: 1 point per metre)

²⁷ <u>http://landmap.mimas.ac.uk/lidar/lidar.html</u>

²⁸ <u>http://en.wikipedia.org/wiki/WGS84</u>

²⁹ http://en.wikipedia.org/wiki/British national grid reference system

Copyright Other metadata that it would be useful to have Time of year of flight Weather conditions/ground condition report (weather for previous week too as need to assess ground water levels) Instrument settings QA tolerance (for example ± 15cm)

Associated formats include

LAS (.las), XYZ (.xyz), DML (.dml), Blue Sky (.txt), NTF (.ntf), Doppler Markup Language (.dml)

Digital Video/Audio



Illustration: Screenshot taken from a digital video extract mounted by Wessex Archaeology on their 'Wrecks on the Seabed' project web pages © Wessex Archaeology

Digital video, notorious for producing large file sizes, consists of a series of digital images that when viewed in succession, create an impression of movement. It may or may not be associated with digital audio. Digital video is becoming more and more popular as a means of recording archaeology, particularly amongst maritime archaeologists where sites are less easily accessible than terrestrial sites. If a whole dive is recorded in this way by means of a hat-mounted camera, the generated file will be substantial in size. Digital video can also be used to record terrestrial archaeology, for example to record excavations in progress, condition surveys, experimental archaeology and interviews. English Heritage have been very involved in using video to benefit Archaeology including the production of video diaries of fieldwork³⁰. TheBamburgh Research Project³¹ has more recently made extensive use of digital video to record both the archaeological processes and the social context of their training excavations. Though Internet Archaeology noted in 1997 (EVA Conference paper) that few archaeologists have access to the technology to create digital video³², it has moved along guickly. With digital cameras that shoot video becoming cheaper and more accessible to a wider range of users plus the easy availability of video editing software, it is likely that use of digital video within archaeology will increase in popularity. Similarly, digital audio can be used by archaeologists. Perhaps not as useful as video to record excavations, artefacts and sites, but oral history projects and interviews with archaeologists are increasingly using digital audio as a medium.

The recent AHDS **Digital Moving Images and Sound Archiving Study**³³ provides informative guidance.

The data

³⁰ <u>http://www.english-heritage.org.uk/upload/pdf/Conservation Bulletin 27.pdf</u>

³¹ <u>http://www.bamburghresearchproject.co.uk/media.htm</u>

³² <u>http://intarch.ac.uk/news/eva97.html</u>

³³ <u>http://ahds.ac.uk/about/projects/archiving-studies/moving-images-sound-archiving-final.pdf</u>

Data has often originally recorded onto DV tape. This is the most economical way to store it but it should be noted that tape degrades relatively quickly. Also video tape is rapidly being superseded by disc based technologies both in terms of recording and viewing. Can be transferred to disk based storage where large scale storage devices are increasingly affordable.

Round 1 of Wessex Archaeology's Wrecks on the Seabed project (a case study – appendix D) has produced somewhere in the region of 75 gigabytes worth of dive footage. Even if all the data is worth retaining something of this size should be unproblematic.

Why should we archive?

In the case of the Wessex data associated with the track log (see Acoustic Tracking above) digital video in a maritime context provides an important record of what was seen by the diver at what particular points on the underwater site. Though not many future users would wish to view full unedited footage of the dive, it is important to preserve this information. Digital video could be utilised as a tool to assess the condition of a wreck site and monitor damage over time. In short it pulls together components of a project just as the traditional paper site diary and more recent video diaries do. Such videos also become a source for historiography.

Problems and issues

Can be very long - Wrecks in the Seabed project has produced around 40 hours of digital video, much of which would be unclear and murky with sections where not very much is happening! Not sure if someone would want to sit down and view the whole lot. Most users would probably be happy with a cleaner edited version showing some of the highlights of the dive.

DV tape has an increasingly short lifespan and should be migrated to disc based storage. DVDs also have a finite lifespan with hard drives (internal or external) providing the securest medium for storage.

Specialised metadata

Digital video in an underwater context may be associated with some record of where the diver and thus the camera was at any one time. Also, the metadata we would expect from any deposit of digital video: Software, version and platform Name and version of video codec (where appropriate) video dimension (in pixels) frame rate per second (fps) bit rate Name and version of audio codec including sample frequency, bit-rate and channel information Length (hours, minutes, seconds) of file File size

Associated formats include

MPEG 1 (.mpg, .mpeg), MPEG 3 (.mpg, .mpeg), MPEG 4 (.mpg4) , MXF (.mxf)

Table 1: Formats review

File extension, name	Description	Properties	Comments
.3dd	Proprietary, binary format used	Proprietary	Export to ASCII
RiSCAN Pro ³⁴	by Riegl's Laser (3D) scanning software; RiSCAN Pro. This can export as a variety of other	Binary Raw data or can be	based format for preservation along with
Laser scanning	formats including ASCII, DXF,		suitable
Point cloud	apparently supports 'Smooth		metadata.
Mesh	data transfer by using the well		
	documented RiSCAN PRO		
	XML-project format ³³³ which		
	metadata although we were		
	unable to locate the schema		
	definition.		
.cod	Difficult to locate information	Proprietary	As a general
.cda	about CODA formats. A number	Binary Bow data or	guideline export
CODA ³⁶	processing sonar data state that	can be	standards if
	they don't accept data in these		possible and
Seismic survey	formats; however, products		then to ASCII.
including	produced by CODA		Will need
Sidescan sonar	technologies generally support		supporting
Sub-bollon proming	SEG Y.		melauala
.CSV	Delimited as the name suggests	Open	The archival
.dat	is structured (usually ASCII)	standard ³⁷	dream for the
.txt	text. Comma Separated Values	ASCII	long term
.XYZ	(CSV) is pernaps the best	Raw data or	preservation of
	associated particularly with	Call De	information about
Delimited text	spreadsheets. Other popular		data collection,
	delimiters include tab and pipe.		etc held as
	The .txt extension can reference		metadata
	structured as well as		
	unstructured data. While .xyz is		
	data.		

³⁴ <u>http://www.3dlasermapping.com/uk/3d/software/RiSCAN%20PRO.pdf</u>

³⁷ <u>http://en.wikipedia.org/wiki/ASCII</u> for an introduction

³⁵<u>http://www.riegl.com/terrestrial_scanners/3d_software_selection_guide_/software_package_/riegl_so</u><u>ftware_index.htm</u>

³⁶ <u>http://www.codaoctopus.com/survey/coda/index.asp</u>

.dat (see .csv)			
.dat	Appears to be generated by data loggers such as the Triton	Proprietary? Binary	Appears to be migrated as a
QMIPS	Technology ISIS Data Logger	Raw data	matter of course
	system. Such data is often		to other formats
Sidescan sonar	converted to other formats such		(see entries for
Sub-bottom profiling	as SEG Y as, for example, by		these).
	the US Geological Survey		
	(USGS) using an in-nouse script		
	processing ³⁸ Information about		
	the OMIPS format can be		
	viewed on the USGS website		
	including the file header		
	details ³⁹ .		
.ddf	See SDTS		
DDF: Data			
Description File			
.dem	ASCII based format developed	Published	Suited to data
	by the United States Geological	standard ⁴⁰	exchange and
DEM: Digital	Survey (USGS). Sources note	ASCII	preservation but
Elevation Model	that these are raster images.	Raw data or	is essentially
	Described as largely	can be	
Elevation Models	standard (see below) but older		15 II - 566 SD I S).
Mesh	datasets still in DEM formats		
Pointcloud	Still supported by many		
	geospatial processing		
	packages. There is a freely		
	available viewer for many		
	USGS formats which is a		
	limited-feature version of		
	commercial software called		
	Global Mapper [®] . Possibly US		
dla		Published	Suited to data
.uig	(USGS) digital line graph (DLG)	standard ⁴³	exchange and
DLG: Digital Line	files are digital vector	ASCII	preservation but
Graph	representations of cartographic	Raw data or	is essentially
	information. Data files of	can be	deprecated (but
DEM: Digital	topographic and planimetric		is it – see SDTS).
Elevation Models	map features are derived from		
	either aerial photographs or		
	from cartographic source		
	materials using manual and		
	There is a freely available		
	viewer for many USGS formats		
	which is a limited-feature		

³⁸ <u>http://woodshole.er.usgs.gov/operations/sfmapping/SUchirp.htm</u>

³⁹ <u>http://woodshole.er.usgs.gov/operations/sfmapping/qmips.htm</u>

⁴⁰ <u>http://rockyweb.cr.usgs.gov/nmpstds/demstds.html</u>

⁴¹ <u>http://edc.usgs.gov/guides/dlg.html</u>

	version of commercial software called Global Mapper ⁴² . Like DEM more recent data is in SDTS format. Possibly US centric?		
.dxf DXF: Drawing eXchange Format 3D including Point cloud CAD Mesh	Published and maintained by AutoDesk vendors of AutoCAD. Was seen for a long time as a <i>de facto</i> standard for the exchange of CAD files ⁴⁴ but then Autodesk stopped publishing (after v. 12) for DXF associated with new versions of AutoCAD. They have; however, recently published the standard for AutoCAD 2008 and several previous versions ⁴⁵ .	Proprietary but Published (currently) ASCII or Binary Raw and processed	Until recently Version migration was seen as the only real way of securing the long term preservation of CAD material; however, use of GDAL/OGR is a possible (as yet untested) strategy (see GML below). Also see the emergence of OpenDWG, IGES and STEP as described in the recent Digital Image Archiving Study ⁴⁶
.dzt RADAN™ DZT GPR: Ground Penetrating Radar	Proprietary, binary format in use with Geophysical Survey Systems, Inc. (GSSI) applications ⁴⁷ . A DZT limited functionality viewer and a RADAN to ASCII converter are available from GSSI ⁴⁸ .	Published ⁴⁹ (currently) Binary Raw data	The ASCII export may be suitable for preservation with supporting metadata
.e00 ESRI Export file GIS	'The ESRI E00 interchange data format combines spatial and descriptive information for vector and raster images in a single ASCII file. It is mainly used to exchange files between different versions of ARC/INFO, but can also be read by many other GIS	Proprietary Not published ASCII Processed (usually)	Usable as an exchange format. No longer seen as the best option for preservation.

⁴³ <u>http://rockyweb.cr.usgs.gov/nmpstds/dlgstds.html</u>

⁴² <u>http://mcmcweb.er.usgs.gov/drc/dlgv32pro/index.html</u>

⁴⁴ Walker, R. (ed.) 1993. *AGI Standards Committee GIS Dictionary*. Association for Geographic Information

- ⁴⁵ <u>http://usa.autodesk.com/adsk/servlet/item?siteID=123112&id=8446698</u>
- ⁴⁶ <u>http://www.jisc.ac.uk/uploaded_documents/FinaldraftImagesArchivingStudy.pdf</u>
- ⁴⁷ <u>http://www.geophysical.com/</u>

⁴⁸ <u>http://www.geophysical.com/softwareutilities.htm</u>

⁴⁹ <u>http://pubs.usgs.gov/of/2002/ofr-02-0166/ofr02_166.pdf</u>

	programs. It is a common format for GIS data found on the Internet ⁵⁰ . This format is proprietary and not in the public domain. An informal analysis of the format is available ⁵¹ . Despite this it has been recognised for many years as the best option for exchange and preservation purposes in being ASCII based and having wide vendor support. A better option is now available in the form of migrating ESRI Shape files directly to GML using GDAL libraries (see below).		
.exp	Use of this original open source	Published	Possible
MOSS: Map Overlay	extent that it proved quite	Processed	exchange format
Statistical System	difficult to track down its source	(usually)	
GIS vector	code ² ; however, its export		
	exchange due to support by		
	other packages.		
various ⁵⁴ including	'MOSS export files contain polygon data extracted from the U.S. Department of Interior's MOSS public domain GIS. These consist of points, lines, or closed polygon loops (possibly with islands), and a 30- character attribute field referred to as the subject value ^{,53} .	Proprietary	Not suited for
.fla	raster graphics, a scripting	but published	preservation
.sfw	language called ActionScript	under licence	
.as	and bi-directional streaming of	Binary	
.asc	audio and video'~. SFW files	(deliverables)	
.TIV	are binaries compiled from	Processed	
Flash [®]	contain simple ActionScript		
	source code (which can also be		
2D, 3D animation	embedded in a SFW file) whilst		
	ASC files contain server-side		
	ActionScript. FLV files represent		
	Flash video clips. Adobe [⊯] took		

⁵⁰ <u>http://www.jisc.ac.uk/uploaded_documents/FinaldraftImagesArchivingStudy.pdf</u> 5.4

⁵¹ <u>http://avce00.maptools.org/docs/v7_e00_cover.html#OVERALL</u>

53 http://ahds.ac.uk/preservation/gis-preservation-handbook.pdf

⁵⁵ <u>http://en.wikipedia.org/wiki/Adobe_Flash</u>

⁵² <u>ftp://ftp.blm.gov/pub/gis/</u>

⁵⁴ <u>http://en.wikipedia.org/wiki/Adobe Flash#Related file formats and extensions</u>

	over Macromedia and hence Flash and now maintain the specification which is currently freely available under licence ⁵⁶		
.gml GML: Geography Markup Language Geospatial data Including GIS CAD	XML (and hence ASCII) based standard for geospatially referenced data. This encoding specification was developed and is maintained by the Open Geospatial Consortium (OGC). The Ordnance Survey (OS) supply MasterMap [®] mapping data as GML ⁵⁷ . Many GIS packages including ESRI and MapInfo products now support GML. The emergence of the Geospatial Data Abstraction Library (GDAL/OGR) is starting to provide the means to easily migrate geospatial data into formats such as GML for	Published standard ⁵⁹ ASCII Processed	GML is very suited for preservation and data exchange of geospatial data.
	exchange ⁵⁸		
various	Grass is an Open Source package ⁶⁰ . Like other GIS a	Openly published (see	Not a preservation
GRASS	versions of GRASS is	Binary and	GML if possible
GIS	represented by a number of files grouped in a directory; coor, topo, cidx (all binary) and head, dbln, hist (ASCII) ⁶¹ . Attribute data is stored in an associated database. GRASS also has GDAL libraries built in so exports to other formats including GML should be an option.	ASCII Processed (usually)	
.gst	(GSF) is described as 'for use	Binary	Possible use as
Generic Sensor Format	as an exchange format in the Department of Defense Bathymetric Library (DoDBL)'.	Raw data	format if widely supported.
Bathymetric data	The specification is currently openly published ⁶² . As well as the generic it allows attributes specific to a wide range of		

⁵⁶ <u>http://www.adobe.com/licensing/developer/</u>

⁵⁷ <u>http://www.ordnancesurvey.co.uk/oswebsite/products/osmastermap/information/technical/gml2.html</u>

⁵⁸ <u>http://www.gdal.org/index.html</u>

60

⁶¹ <u>http://grass.itc.it/grass63/manuals/grass63vectorlib_2007_04_21_refman.pdf</u>

⁶² http://www.ldeo.columbia.edu/res/pi/MB-System/formatdoc/gsf_spec.pdf

⁵⁹ <u>http://www.opengis.net/gml/</u>

	bathymetric surveying systems		
	to be included.		
.hsx .hs2	To quote from the Hypack 'HYSWEEP [®] survey has a Text	Proprietary ASCII (.hsx)	Text logging provides the
HYPACK Inc.	allowing raw data to be stored in a format that can be inspected	Raw data or can be	preservation
Sidescan sonar	and modified by most editing		
Bathymetric data	program (Windows Wordpad for		
(single beam?)	example). Easy inspection of files is the advantage of text		
	logging - the disadvantage is		
	larger files and slower load time.		
	important to you, it is best to		
	choose the HYSWEEP [®] binary		
	format (HS2) ⁶³ . The manual		
	also contains format		
iaw	Identical to TIFF World files		
.)9**	(see below)		
JGW: JPEG World			
file			
.las	The LAS format is described as	Published ⁶⁹	Specifically
.las	The LAS format is described as 'a public file format for the	Published ⁶⁹ Binary	Specifically designed for the
.las LAS	The LAS format is described as 'a public file format for the interchange of LIDAR data	Published ⁶⁹ Binary Raw data or	Specifically designed for the exchange of
.las LAS	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ .	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along with industry stakeholders ⁶⁵ .	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along with industry stakeholders ⁶⁵ . Discussions of extending LAS to additionally handle terrestrial	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along with industry stakeholders ⁶⁵ . Discussions of extending LAS to additionally handle terrestrial laser scanning data are actively	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along with industry stakeholders ⁶⁵ . Discussions of extending LAS to additionally handle terrestrial laser scanning data are actively taking place ⁶⁶ . A recent	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along with industry stakeholders ⁶⁵ . Discussions of extending LAS to additionally handle terrestrial laser scanning data are actively taking place ⁶⁶ . A recent addendum to the English	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along with industry stakeholders ⁶⁵ . Discussions of extending LAS to additionally handle terrestrial laser scanning data are actively taking place ⁶⁶ . A recent addendum to the English Heritage Metric Survey Specification covering laser	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along with industry stakeholders ⁶⁵ . Discussions of extending LAS to additionally handle terrestrial laser scanning data are actively taking place ⁶⁶ . A recent addendum to the English Heritage Metric Survey Specification covering laser scanning supports LAS as a	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.
.las LAS Lidar Laser scanning	The LAS format is described as 'a public file format for the interchange of LIDAR data between vendors and customers. This binary file format is an alternative to proprietary systems or a generic ASCII file interchange system used by many companies' ⁶⁴ . The American Society for Photogrammetry & Remote Sensing (ASPRS) endorses and supports the use of LAS along with industry stakeholders ⁶⁵ . Discussions of extending LAS to additionally handle terrestrial laser scanning data are actively taking place ⁶⁶ . A recent addendum to the English Heritage Metric Survey Specification covering laser scanning supports LAS as a data exchange and archival	Published ⁶⁹ Binary Raw data or can be.	Specifically designed for the exchange of data; a role for which it has strong support. In being a binary format would not be seen as suited to a long term preservation role as ASCII text alternatives exist.

⁶³ <u>http://www.hypack.com/documentation.asp</u>

⁶⁴ <u>http://www.lasformat.org/</u>

⁶⁵ <u>http://www.asprs.org/society/committees/lidar/lidar_format.html</u>

⁶⁶ www.ceg.ncl.ac.uk/heritage3d/downloads/TLS%20formats%20V1.pdf

⁶⁷ <u>http://www.ceg.ncl.ac.uk/heritage3d/downloads%5Caddendum2006.pdf</u>

	Such usage has not been formalised as vet.		
.map .tab .dat .id .ind MapInfo TAB GIS	MapInfo's native format. It is regulated by MapInfo as a proprietary format and is not openly published. It comprises of a number of related files ⁷⁰ . These are a mixture of binary and ASCII. The best option for preservation is to export to GML using GDAL libraries (see	Proprietary Not published Binary and ASCII files Processed usually	Not suited for data exchange or preservation. Export to GML or MIF with support metadata
.mad77	above). Described as a 'format for the	Published	Possible
MGD77 Geophysical data Including Bathymetric Magnetic Gravity	exchange of digital underway (?) geophysics data'. It was developed by the US National Geophysical Data Center (NGDC) following an international workshop in 1977 'Workshop for Marine Geophysical Data Formats' ⁷¹ . Has been revised relatively recently. UNESCO note of MGD77 that the 'format has experienced much success over the last 20 years. It has been sanctioned by the Intergovernmental Oceanographic Commission (IOC) as an accepted standard for international data exchange, and it has been translated by IOC into French, Japanese and Russian. Most contributors of data to NGDC now send their	ASCII Raw	exchange format Possible preservation format
.mif	Like the interchange format of	Proprietary	Possible
.mia MIF: MapInfo Interchange format GIS	Its competitor, ESRI (Export .e00), MIF files are ASCII and support is widespread in GIS applications. The MIF file contains geometric data whilst the optional MID file has header and attribute data as delimited text. The format specification is currently	ASCII (currently) ⁷³ ASCII Processed (usually)	format because of widespread support

⁶⁸ http://www.english-heritage.org.uk/server/show/nav.001002003003007001

⁶⁹ <u>http://www.lasformat.org/documents/ASPRS%20LAS%20Format%20Documentation%20-%20V1.1%20-%2003.07.05.pdf</u>

⁷⁰ http://en.wikipedia.org/wiki/MapInfo TAB format

⁷¹ <u>http://www.ngdc.noaa.gov/seg/gravity/document/html/mgd77.shtml#general</u>

72 http://ioc.unesco.org/iocweb/iocpub/iocpdf/tc045.pdf

73 http://extranet.mapinfo.com/common/library/interchange_file.pdf

	available from MapInfo.		
.mpg	An International ISO/IEC	Published	Suitable for
.mpeg	(11172) developed by the	open	preservation. The
1 0	Moving Picture Experts Group	standard ⁷⁴	AHDS currently
MPEG-1:	(MPEG) for Video CD (VCD)	Binary	recommend de-
	and less commonly DVD-Video.	Processed	multiplexing of
Video	Provides reasonable quality	usually	video and audio
Audio	audio/video plavback	,	channels.
	comparable to VHS tape. The		
	MPEG-1 Audio Laver III		
	equates to MP3 audio.		
.mpg	As MPEG-1, an ISO/IEC	Published	Suitable for
mpeg	(13818) standard but for DVD	open	preservation. The
	as well as various flavours of	standard ⁷⁶	AHDS currently
MPEG-2	TV 'MPEG-2 video is not	Binary	recommend de-
	optimized for low bit-rates (less	Processed	multiplexing of
Video	than 1 Mbit/s) but outperforms	usually	video and audio
Audio	MPEG-1 at 3 Mbit/s and	actionly	channels
	above ⁷⁵ and hence much		ondiminioion
	higher quality		
.mpg4	Another MPEG ISO/IEC	Published	In being an
	(14496) standard concerned	open	online streaming
MPEG-4	with 'web (streaming media)	standard ⁷⁸	standard could
	and CD distribution	Binary	be used for
Video	conversation (videophone), and	Processed	dissemination
Audio	broadcast television. all of which		
	benefit from compressing the		
	AV stream ⁷⁷ .		
.mst	Based on TIFF format v. 5. A	Proprietary	As a general
	format specification is currently	Published	quideline export
Marine Sonic	accessible ⁷⁹ . Notes that	(currently)	to more open
Technology MSTIFF	'Although TIFFs allow for	Binary	standards if
	customization of the format,	Raw data or	possible.
Sidescan sonar	MSTL decided it was better to	can be	•
	use the basic structure and		
	create our own MSTL specific		
	tags instead of trying to fit all of		
	our proprietary information into		
	the TIFF'.		
.mxf	A generic wrapper or container	Published	Generally seen
	for moving images. Developed	open	as an emerging
MXF: Material	as an open standard by the US	standard ⁸¹	standard that
Exchange Format	Society of Motion Picture and	Binary and	may have the
-	Television Engineers (SMPTE).	ASCI	potential to
Video	The recent Digital Moving	Processed	become a
Audio	Images and Sound Archiving	usually	preservation

⁷⁴ <u>http://www.iso.ch/iso/en/CatalogueDetailPage.CatalogueDetail?CSNUMBER=25371</u>

⁷⁵ http://en.wikipedia.org/wiki/MPEG-2

⁷⁶<u>http://www.iso.ch/iso/en/CatalogueDetailPage.CatalogueDetail?CSNUMBER=37679&ICS1=35&IC S2=40&ICS3=</u>

⁷⁷ http://en.wikipedia.org/wiki/MPEG-4

⁷⁸ <u>http://www.iso.ch/iso/en/CatalogueDetailPage.CatalogueDetail?CSNUMBER=38559</u>

⁷⁹ <u>http://www.marinesonic.com/documents/mstiff.pdf</u>

	Study undertaken by the AHDS notes that 'MXF is related to the AAF format (see above). It was specifically developed for optimised interchange and archiving of multimedia content. Although currently (2006) too new to be widely used it is emerging as a standard' ⁸⁰ .		standard.
.nc NetCDF: Network Common Data Form Scientific data including Bathymetric Lidar and others?	NetCDF 'is a set of software libraries and machine- independent data formats that support the creation, access, and sharing of array-oriented scientific data ⁸² .Openly published ⁸³ . Libraries freely available under licence. Tools include ncgen and ncdump which respectively generate from and dump to ASCII. Also supports the sub-setting of datasets. Appears widely used for scientific including bathymetric data, for example, the NERC British Oceanographic Data Centre (BODC) ⁸⁴ .	Published Binary Raw or can be	This could provide an ideal mechanism for preservation and data sharing through storing once and generating binary or ASCII as requested
.ntf NTF: National Transfer Format Geospatial data including Point cloud CAD Digital Elevation Models (DEM) Lidar	Complex ASCII based storage and transfer format for vector and raster images (same extension). Largely used by the OS for distributing pre- MasterMap data (see GML). It is a British Standard BS 7567 'Electronic Transfer of Geographic Information'. A wide range of NTF converters are available to, for example, popular GIS formats. Lidar data as supplied has often been processed in terms of coordinate transformation and decimation.	Published standard ⁸⁵ ASCII Raw and processed	In being ASCII based and published it should be suited for both transfer and preservation. Unclear; however, as to how wide its usage is outside of the OS where it is being superseded by GML
.obj OBJ	A simple ASCII based format for representing 3D geometry. Initially developed by Wavefront	Published ASCII Raw data or	Wide support suggests a possible data

⁸¹ http://www.smpte.org/standards/

⁸⁰ <u>http://ahds.ac.uk/about/projects/archiving-studies/moving-images-sound-archiving-final.pdf</u> 5.2.3

⁸³ <u>http://www.unidata.ucar.edu/software/netcdf/docs/netcdf/File-Format-Specification.html#File-Format-Specification</u>

84 http://www.bodc.ac.uk/data/online_delivery/gebco/

⁸⁵ <u>http://www.bsistandards.co.uk/shop/products_view.php?prod=6536</u>

⁸² <u>http://www.unidata.ucar.edu/software/netcdf/</u>

3D including Laser scanning Mesh Point cloud	Technologies. The format is apparently open and has wide support amongst both software vendors and open source community. Whilst the format specification is available on numerous websites ⁸⁶ we were unable to identify a format maintainer. There are numerous converters available for OBJ files.	can be	exchange format. In being ASCII based it could act as a preservation format
.rd3 .rad RAMAC – RD3/RAD GPR: Ground Penetrating Radar	Used natively by Malå GeoScience equipment. Data is in a binary RD3 file whilst header information is in an ASCII text file. Apart from the file extension they share the same name. A number of open source tools exist for manipulating RD3 files. For example, GPR IDL tools ⁸⁷ can be down loaded from Source Forge and used to convert files into x, y, z.	Published ⁸⁸ (currently) Binary (ASCII header) Raw	Move to ASCII text for preservation purposes
.rst .rdc IDRISI raster GIS	Native raster image format for Clark Labs ⁸⁹ IDRISI GIS software. Associated file types include Raster Documentation (RDC) files. Can be viewed as ASCII.	Proprietary Published ⁹⁰ (currently) Processed (usually)	
.sd .scene Fledermaus Visualisation of a range of spatially referenced data including Multibeam sonar Digital Elevation Models (DEM = DTM) GIS (various) CAD (.dxf)	Visualisation toolkit for 2D, 3D and movies. Supports the import and export of a large range of spatially referenced data types ⁹¹ . A reduced functionality viewer is available for download ⁹² .	Proprietary Binary (part ASCII) Processed	Can represent a project outcome (i.e. a presentation format). Not suited for preservation but can export as a range of formats

⁸⁶ <u>http://people.scs.fsu.edu/~burkardt/txt/obj_format.txt</u>

- ⁸⁸ <u>http://pubs.usgs.gov/of/2002/ofr-02-0166/ofr02_166.pdf</u>
- ⁸⁹ http://www.clarklabs.org/about/
- ⁹⁰ http://www.clarklabs.org/about/

⁹¹ <u>http://www.ivs3d.com/products/fledermaus/</u>

92 http://www.ivs3d.com/download/iview3d_download.html

⁸⁷ <u>http://gpr-idl-tools.sourceforge.net/tutorial.html</u>

Various	An Earth Science standard	Published	Well supported
including .ddf	developed by the USGS for	standard ⁹⁵	as a data
	data exchange. Downloaded	Binary	exchange
SDTS: Spatial Data	files are a tarred (zipped)	Raw data or	standard but
Transfer Standard	directory which as well as data	can be	probably US
	contains numbers of DDF or		centric.
Geospatial data	data description files.		
	Compliance with SDTS is a		
Terrain	requirement for federal		
Image	agencies in the US. Supports		
	Raster and vector data. There		
	are large numbers of tools and		
	translators for extracting data		
	from SDTS to various formats.		
	In some cases this involves		
	extraction to earlier standards		
	such as DLG (see above)		
	which suggests SDTS is a		
	CDAL (ass CML shave) support		
	GDAL (see GiviL above) support		
	An undete to verious SEC	Openly	Doosiblo
.syz	formate including SEG V by the	Dpeniy	Pussible
.uai	Society of Exploration	Binory	Export to ASCII
SEC 2	Goophysicists (SEG) Pathor	Dinary Dow data	with suitable
3662	strangely there sooms to be	Naw uala	motodoto if
	numbers of SEC 2 to SEC V		nossible
including	converters available. Does this		possible
GPR: Ground	mean SEG V is still better		
Penetrating Radar	supported? Seismic Univ is a		
	popular freeware package for		
	working with SEG and other		
	seismic formats ⁹⁶		
segv	An openly published format ⁹⁸ by	Published	Can be
	the Society of Exploration	Binary	converted to
SEG Y	Geophysicists (SEG), Originally	Raw data	ASCII for
	(rev. 0) developed in 1973 for		preservation
Seismic survey	use with IBM 9 track tapes and		purposes.
including	mainframe computers and using		Possibly useful
Sub-bottom profiling	EBCDIC (an alternative to		as a data
Sidescan sonar	ASCII encoding rarely used		exchange format
GPR: Ground	today) descriptive headers. The		as it appears
Penetrating Radar	standard was updated (rev. 1) in		widely supported.
U U	2001 to accommodate ASCII		
	textual file headers and the use		
	of a wider range of media. It		
	should be noted that in the		

93 http://www.fws.gov/data/gisconv/sdts2av.html

⁹⁴ http://home.gdal.org/projects/sdts/

⁹⁵ <u>http://mcmcweb.er.usgs.gov/sdts/standard.html</u>

⁹⁶ <u>http://www.cwp.mines.edu/cwpcodes/index.html</u>

97 http://www.seg.org/publications/tech-stand/seg_2.pdf

98 http://www.seg.org/publications/tech-stand/

	interim between revisions a number of flavours of SEG Y appeared trying to overcome the limitations of rev. 0. SEG Y to ASCII converters exist as, for example, made available by the USGS ⁹⁹ . A limited functionality SEG Y viewer can be downloaded from Phoenix Data Solutions ¹⁰⁰		
.shp and lots of	Well documented ¹⁰² and supported by other GIS vendors	Proprietary Published	Because of
formats ¹⁰¹	such as MapInfo ¹⁰³ . It has been	(mostly?)	can usually be
	described as 'developed and	Binary	used as an
ESRI SHAPE file	regulated by ESRI as a (mostly)	Processed	exchange format
GIS	interoperability among ESRI and other software products ¹⁰⁴ . ESRI Export format (E00) and more recently GML (see entries above) are seen as the best preservation ontions	usuany	
.svg	XML (and hence ASCII) based	Published	Suited for both
	format for 2D vector graphics.	open standard	dissemination
SVG: Scalable	Specification developed and	XML (ASCII)	and the long term
vector Graphics	Web Consortium (W3C) ¹⁰⁵ The	data	simpler 2D vector
2D vector images	specification notes that 'For accessibility reasons, if there is an original source document containing higher-level structure and semantics, it is recommended that the higher- level information be made available somehow, either by making the original source document available, or making an alternative version available in an alternative format which conveys the higher-level information, or by using SVG's facilities to include the higher-	uata	graphics. It should be noted that Adobe has recently dropped its support for SVG ¹⁰⁶ but support is still widespread.

99 http://pubs.usgs.gov/of/2005/1311/of2005-1311.pdf

¹⁰⁰ http://www.phoenixdatasolutions.co.uk/seisvu.htm

- ¹⁰¹ http://ahds.ac.uk/preservation/gis-preservation-handbook.pdf
- ¹⁰² http://www.esri.com/library/whitepapers/pdfs/shapefile.pdf

¹⁰⁴ http://en.wikipedia.org/wiki/Shapefile

- ¹⁰⁵ <u>http://www.w3.org/TR/SVG/</u>
- ¹⁰⁶ <u>http://www.adobe.com/svg/eol.html</u>

¹⁰³<u>http://extranet.mapinfo.com/products/Features.cfm?ProductID=1044&productcategoryid=1#Importing%20and%20Exporting%20Data</u>

.tfw TFW: TIFF World file ESRI GIS products (others?)	level information within the SVG content. This suggests that for archival purposes use might need restricting to simpler models. A mechanism for georeferencing images developed by ESRI (GIS software vendor). As such similar to GEOTIFF (see above) but in this case the metadata is held in a separate ASCII text file ¹⁰⁷ . TIFF World files will be small in themselves but may be associated with large images	Proprietary ASCII (but associated image will be binary) Processed	That the metadata (spatial information is in ASCII could be seen as good for preservation.
.tiff GEOTIFF TIFF GIS and other image processing packages	The GEOTIFF standard is in the public domain. It allows metadata, specifically georeferencing to be embedded within a TIFF image. There is complete conformance to the current TIFF 6.0 specification. As the recent Digital Image Archiving Study notes 'The use of uncompressed TIFF version 6 <as format="" preservation=""> is the best strategy at the current time, but a watching brief should be maintained on JPEG2000 as an emerging preservation format'¹⁰⁸. TIFF is also a public domain format currently maintained by Adobe^{® 109}. It should be noted that the size of a TIFF file is limited to 4GB¹¹⁰.</as>	Public domain ¹¹¹ Binary Processed	Despite being a binary format TIFF has long been recognised as a <i>de facto</i> preservation standard for raster images. Binary is currently the only real option for the bitstream encodings of raster images.
.txt (see .csv)		<u>^</u>	
.txt Blue Sky Lidar Point cloud Mesh	Blue Sky are a leading European are described as 'at the forefront of imaging technology and geospatial data ¹¹² . Straight forward XYZ data (see below) ¹¹³ . Lidar data as supplied has often been processed in terms of	ASCII Raw (ish)	Suited for preservation with suitable metadata

¹⁰⁷ <u>http://support.esri.com/index.cfm?fa=knowledgebase.techArticles.articleShow&d=17489</u>

¹⁰⁸ <u>http://www.jisc.ac.uk/uploaded_documents/FinaldraftImagesArchivingStudy.pdf</u> 1.4.i

- ¹⁰⁹ http://partners.adobe.com/public/developer/tiff/index.html
- ¹¹⁰ http://www.awaresystems.be/imaging/tiff/faq.html#q8
- ¹¹¹ http://remotesensing.org/geotiff/spec/geotiffhome.html
- 112 http://www.bluesky-world.com/index.html
- 113 http://www.bluesky-world.com/dem-lidar.html

	coordinate transformation and decimation.		
.vct .vdc .mdb .adc IDRISI vector GIS	Native vector graphics format for Clark Labs ¹¹⁴ IDRISI GIS software. Associated file types include Vector Document (VDC) files, additional attributes within Microsoft Access database (MDB) files and Attribute Documentation (ADC) files. See RST for IDRISI raster format. Can be viewed as ASCII	Proprietary Published ¹¹⁵ (currently) Binary (ASCII export) Processed (usually)	Export to ASCII along with suitable metadata for preservation.
.vpf VPF: Vector Product Format Vector graphics	Developed by the U.S. Defense Mapping Agency. 'The Vector Product Format (VPF) is a standard format, structure, and organization for large geographic databases that are based on a georelational data model and are intended for direct use' ¹¹⁶ . It may be difficult to work with 'Don't try to read VPF unless absolutely necessary. It's a dog of a format' ¹¹⁷	Open standard ¹¹⁸ ASCII Raw or processed	Possible transfer format but may be complex.
.vtk Visualization Toolkit (VTK) ¹¹⁹ 3D computer graphics, image processing, and visualization and thus Laser scanning	X, Y and Z co-ordinates plus the file header and tail which make the data readable in programs such as Paraview ¹²⁰ which uses the Visualization Toolkit and is a freely downloadable viewer for the visualization of large data sets such as point clouds. It should be noted that Paraview and VTK support a very wide range of technologies and formats and that recent implementations support XML based formats with older ASCII/binary formats now considered as legacy but still supported. VTK has not been considered here beyond point	Published (i.e. the toolkit is open source) ASCII or Binary Processed (probably)	With an open source viewer the .vtk files supplied by one of the case studies seem a very reasonable way of disseminating point cloud data. That this is now seen as a legacy format might contradict this; however. In being ASCII text suited for preservation of

¹¹⁴ http://www.clarklabs.org/about/

¹¹⁵ <u>http://geography.rutgers.edu/courses/06fall/321/Public/Idrisi_Guide/Andes%20Manual.pdf</u> ch. 5

¹¹⁶ <u>http://www.nga.mil/portal/site/nga01/index.jsp?epi-</u> content=GENERIC&itemID=a2986591e1b3af00VgnVCMServer23727a95RCRD&beanID=16296300 80&viewID=Article

¹¹⁷ http://www.vterrain.org/Culture/vector.html

¹¹⁸ http://www.nga.mil/NGASiteContent/StaticFiles/OCR/vpf_main.pdf

¹¹⁹ http://public.kitware.com/VTK/

120 http://www.paraview.org/HTML/Download.html

	cloud data supplied by the Breaking through Rock Art case study.		research outcomes.
.wav			
WaveForm			
.wmv .asf	Proprietary video and audio codec		
WMV: Windows Media Video			
Video Audio			
.wrl VRML: Virtual Reality Modelling Language 3D graphics	As VRML 97 a published ISO (14772-1) standard for 3D vector graphics. Designed with the internet in mind. As such requires a plug-in or viewer ¹²¹ . Apparently still popular especially for the exchange of CAD drawings but is slowly being superseded by other standards such as X3D (below)	Published open standard ¹²² ASCII Processed	Possible exchange format. In being ASCII based has the potential to act as a preservation format but aging.
various X3D 3D graphics	Developed as a replacement for VRML (above) by the web3D consortium ¹²³ this ISO (19775) standard is XML based although a binary specification has been more recently released as an ISO (19776-3) standard. It is backwardly compatible with VRML. It is noted as being compatible with the MPEG-4 (above) specification. Like VRML requires a plug-in or viewer.	Published open standard ¹²⁴ ASCII and binary flavours Processed usually	With XML being ASCII based this has archival possibilities.
.xml	XML ¹²⁵ is a general-purpose markup language geared	Published open	Ideal for exchange and
XML: eXtensible Markup Language	towards facilitating the sharing of data. An XML document is said to be 'well formed' and	standard ¹²⁷ ASCII RAW or	preservation if an established schema exists
increasing range of technologies	when it conforms to XMLs syntactical rules. It is described as valid when it conforms to semantic rules defined in a published schema. Many XML documents use a different file	processed	

¹²¹ http://vads.ahds.ac.uk/guides/vr_guide/sect37.html

122 http://www.web3d.org/x3d/specifications/vrml/

¹²³ http://www.web3d.org/x3d/

¹²⁴http://www.web3d.org/x3d/specifications

¹²⁵ <u>http://en.whttp://www.web3d.org/x3d/specificationsikipedia.org/wiki/Xml</u>

	extension, for example .gml (see above). Others such as MIDAS XML developed by the Forum on Information Standards in Heritage (FISH) ¹²⁶ are explicit in having the .xml extension.		
.xtf	As described by the Triton Imaging Inc 'The XTF file format	Proprietary Binary	Possibly very suited for data
eXtended Triton Format ¹²⁸	was created to answer the need for saving many different types of sonar, navigation, telemetry and bathymotry information	Raw data or can be	exchange if industry support is widespread.
Sidescan sonar Sub-bottom profiling Bathymetric data	and bathymetry information. The format can easily be extended to include various types of data that may be encountered in the future'. Currently a Publicly Available Specification. Also described as an 'industry standard' for sonar. Some packages supporting XTF provide for ASCII text exports		Where possible ASCII text exports with suitable metadata would provide the best long term preservation environment
.xyz (see .csv)			
.xyz .xyzrgb	Point cloud data - simply the X, Y and Z coordinates of each scanned point, sometimes with	ASCII (can be binary) Raw(ish)	ASCII text is seen as the best option for long
XYZ	Red, Green and Blue colour values also. XYZ data is often		term preservation along with
Laser scanning	decimated to make dataset more manageable. Depending on purpose this can often be done without discernable loss of detail. Lidar data as supplied has often been processed in terms of coordinate transformation and decimation.		suitable metadata

¹²⁷ http://www.w3.org/XML/

¹²⁶ <u>http://www.heritage-standards.org/</u>

¹²⁸<u>http://www.tritonimaginginc.com/site/content/public/downloads/FileFormatInfo/Xtf%20File%</u> 20Format_X21.pdf