

The digital preservation technological context

Michael Day,
Digital Curation Centre
UKOLN, University of Bath
m.day@ukoln.ac.uk

La preservación del patrimonio digital: conceptos básicos y principales iniciativas, Madrid, 14-16 March 2006

<http://www.ukoln.ac.uk/>



Session overview

- Introductory comments
- Technical issues
- Preservation strategies
- Preservation metadata and shared infrastructure

<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Introductory comments

<http://www.ukoln.ac.uk/>



Digital preservation (1)

- Concerns continued access (and use)
- Digital preservation is NOT just about technology
- Unites a range of interrelated issues:
 - "... the planning, resource allocation, and application of preservation methods and technologies to ensure that digital information of continuing value remains accessible and usable" - Margaret Hedstrom (1998)

<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Digital preservation (2)

- Is sometimes now characterised as 'digital stewardship' or 'digital curation'
 - The concept of data curation originated in data-rich scientific domains like bioinformatics
 - Curation - "The activity of managing and promoting the use of data from its point of creation, to ensure it is fit for contemporary purpose, and available for discovery and reuse" - Philip Lord, *et al.* (2004)
 - "Maintaining and adding value to a trusted body of information for current and future use" -- DCC presentation at CNI (2005)

<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

The fragility of digital content

The main technical issues

<http://www.ukoln.ac.uk/>



General comments

- Digital information is dependent on its technical environment
- Physical objects are subject to:
 - Physical deterioration
 - Technology obsolescence
- Relatively short timescales



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Storage media (1)

- A major focus of concern in the 1970s and 1980s
- Current media types
 - Typically, magnetic or optical tape and disks, various devices (e.g., memory sticks)
 - Examples include: CD-ROM, DVD (optical), DAT, DLT (magnetic)
- Unknown lifetimes
 - Subject to differences in quality or storage conditions
 - But relatively short lifetimes compared to paper or good quality microform



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Storage media (2)

- Technical solutions:
 - Periodic copying of data bits on to new media or types of media (refreshing)
 - Longer lasting media
 - Migrating to good-quality microform or paper (!)
- In an organised preservation system, regular routines (quality checking, backup, replication, refreshing, etc.) will help solve the media longevity issue



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Technology obsolescence (1)

- A set of much bigger problems
- Software dependence
 - Digital content is, at least in part, dependent on the configurations of hardware and software (applications and operating systems) that were originally used to interpret or display them
- Hardware and software obsolescence
 - Application software and operating systems are upgraded regularly
 - Hardware becomes obsolete or needs repair



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Technology obsolescence (2)

- Technical solutions
 - Various preservation strategies have been developed to cope with the obsolescence problem
 - For the most part, these depend on the existence of a continual programme of active management (life cycle management)
 - Supported by systems that implement the various functional entities identified by the Reference Model for an Open Archival Information System (OAIS)
 - Preservation strategies can only be seen in this wider context



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Layers of meaning (1)

- Digital objects are logical entities not fixed to any one particular physical carrier
- Three layers (Thibodeau, 2002):
 - Physical objects: the actual bits stored on a particular medium
 - Logical objects: defines how these bits are used by application software, based on data types (e.g. ASCII); in order to understand (or preserve) the byte-streams, we need to know how to process them
 - Conceptual objects: what humans deal with in the real world, meaningful units of information



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Layers of meaning (2)

- On which of these layers should preservation activities focus?
 - We need to preserve the ability to reproduce the objects, not just the bits
 - In fact, we could change the bits and logical representation and still reproduce an authentic conceptual object



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Authenticity and integrity

- Digital information can easily be changed (e.g., by design or accident)
- How can we trust that an object is what it claims to be?
- Mechanisms are available at the bit level (e.g. checksums), but will this be sufficient?



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Problems of scale

- An increasing flood of 'born-digital' data
 - Data deluge in science and engineering
 - » Petabytes generated by high throughput instruments, streamed from sensors and satellites, etc.
 - The World Wide Web
 - » Comprises billions of pages + "deep Web"
 - » Internet Archive = >1 petabyte, and growing @ 20 Tb. per month (<http://www.archive.org/>)
 - 5 exabytes of new information created in 2002:
 - » <http://www.sims.berkeley.edu/research/projects/how-much-info-2003/>



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Some general principles (1)

- Most of the technical problems associated with long-term digital preservation can be solved if a life-cycle management approach is adopted
 - i.e. a continual programme of active management
 - Ideally, combines both managerial and technical processes, e.g., as in the OAIS Model
 - Many current systems (e.g. repository software) are attempting to support this approach
 - Preservation strategies need to be seen in this wider context
- Preservation needs to be considered at a very early stage in an object's life-cycle



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Some general principles (2)

- Need to identify and understand the 'significant properties' of an object
 - Focuses on the essential
 - Helps with choosing an acceptable preservation strategy
- Encapsulation may have some benefits
 - Surrounding the digital object - at least conceptually - with all of the information needed to decode and understand it (including software)
 - Produces autonomous 'self-describing' objects, reduces external dependencies; linked to the Information Package concept in the OAIS Reference Model
- Keep the original byte-stream in any case



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Digital preservation strategies



<http://www.ukoln.ac.uk/>



Preservation strategies

- Three main families:
 - Technology preservation
 - Technology emulation
 - Information migration
- Also:
 - Digital archaeology (rescue)



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Technology preservation

- The preservation of an information object together with all of the hardware and software needed to interpret it
 - Successfully preserves the look, feel and behaviour of the whole system (at least while the hardware and software still functions)
 - May have a role for historically important hardware
 - Problems with storage and ongoing maintenance, missing documentation
 - Would inevitably lead to 'museums' of "ageing and incompatible computer hardware" -- Mary Feeney
 - May have a short-term role for supporting the rescue of digital objects (digital archaeology)



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Technology emulation (1)

- Preserving the original bit-streams and application software; running this on emulator programs that mimic the behaviour of obsolete hardware
- Emulators change over time
 - Chaining, rehosting
 - Emulation Virtual Machines
 - » Running emulators on simplified 'virtual machines' that can be run on a range of different platforms
 - » Virtual machines are migrated so the original bit-streams do not have to be



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Technology emulation (2)

- Benefits:
 - Technique already widely used, e.g. for emulating different hardware, computer games
 - Preserves the original bits
 - Reduces the need for regular object transformations (but emulators and virtual machines may themselves need to be migrated)
 - Retains 'look-and-feel'
 - May be the only approach possible where objects are complex or dependent on executable code
 - Less 'understanding' of formats is needed; little incremental cost in keeping additional formats



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Technology emulation (3)

- Issues
 - Which organisations have the technical skills necessary to implement the strategy?
 - Preserving 'look and feel' may not be needed for all objects
 - It will be difficult to *know* definitively whether user experience has been accurately preserved
- Conclusions
 - Promising family of approaches
 - Needs further practical application and research



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Information migration (1)

- Managed transformations
 - A set of organised tasks designed to achieve the periodic transfer of digital information from one hardware and software configuration to another, or from one generation of computer technology to a subsequent one - CPA/RLG report (1996)
 - Abandons attempts to keep old technology (or substitutes for it) working
 - A 'known' solution used by data archives and software vendors (e.g., a linear migration strategy is used by software vendors for some data types, e.g. Microsoft Office files)
 - Focuses on the *content* of objects



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Information migration (2)

- Main types (from OAIS Model)
 - Refreshment
 - Replication
 - Repackaging
 - Transformation
- Issues
 - Labour intensive
 - There can be problems with ensuring the 'integrity and authenticity' of objects
 - Transformations need to be documented (part of the preservation metadata)



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Information migration (3)

- Uses
 - Seems to be most suitable for dealing with large collections of similar objects
 - Migration can often be combined with some form of standardisation process, e.g., on ingest
 - » ASCII
 - » Bit-mapped-page images
 - » Well-defined XML formats
 - Migration on Request (CAMiLEON project)
 - » Keep original bits, migrate the rendering tools



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Digital archaeology

- Not so much a preservation strategy, but the default situation if we fail to adopt one
- Using various techniques to recover digital content from obsolete or damaged physical objects (media, hardware, etc.)
 - A time consuming process, needs specialised equipment and (in most cases) adequate documentation
 - Considered to be expensive (and risky)
 - Remains an option for content deemed to be of value



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Choosing a strategy (1)

- Preservation strategies are not in competition (different strategies will work together)
 - A suggestion that we should keep the original bits (with documentation) in any case
- But the strategy chosen has implications for:
 - The technical infrastructure required (and metadata)
 - Collection management priorities
 - Rights management
 - » e.g. Owning the rights to re-engineer software
 - Costs



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Choosing a strategy (2)

- Tools for supporting preservation decisions, e.g.
 - Preservation strategies
 - Target formats for transformations
- Nationaal Archief (Netherlands) testbed project
- Vienna University of Technology utility analysis tool
- Both developed further by the Digital Preservation cluster of the DELOS Network of Excellence



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DICIC

Case study

Rescue of content from BBC
Domesday videodiscs



<http://www.ukoln.ac.uk/>



DICIC

Rescue of BBC Domesday (1)

- BBC Domesday project (1986)
 - To commemorate the 900th Anniversary of the original Domesday survey
 - Two interactive videodiscs (12")
 - Mixture of textual material (some produced by schools), maps, statistical data, images and video
 - Technical basis:
 - Hardware: BBC Master Series microcomputer and Philips Laservision (LV-ROM) player
 - Some software in ROM chip, others on the discs
 - System obsolete by end of 1990s; working hardware becoming more difficult to find



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Rescue of BBC Domesday (2)

- CAMILEON project
 - Proof of concept for the emulation approach
 - Converted data into media-neutral form
 - Adapted an existing emulator for the BBC microcomputer to render Domesday content
- The National Archives (and partners)
 - Reengineered the whole system for use on Windows PCs
 - Digital versions of images and video converted from original master tapes (still held by BBC)
 - Developed an improved interface
 - Web version: <http://domesday1986.com/>



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

The screenshot shows the 'Domesday 1986 - County Data' page. It features a map of the South West of England with a yellow box highlighting a specific area. Below the map is a list of parishes with their respective grid coordinates. The interface includes navigation buttons like 'Zoom Out' and 'Zoom In', and a 'Darken' button to toggle the map's appearance. Information about the current map and parish is displayed at the bottom.



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

This screenshot shows a different view of the 'Domesday 1986 - County Data' page. The map highlights a different area, and the list of parishes is updated accordingly. A photograph of a landscape, identified as 'Kimmeridge Bay & Torrishay Cap', is displayed in the bottom right corner. The interface includes navigation buttons and information about the current map and parish.



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

This screenshot shows another view of the 'Domesday 1986 - County Data' page. The map highlights a different area, and the list of parishes is updated. A photograph of a group of people, identified as 'Llansawney Cornmill, 1985', is displayed in the bottom right corner. The interface includes navigation buttons and information about the current map and parish.



<http://www.ukoln.ac.uk/>



La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006

Preservation metadata and shared infrastructures



<http://www.ukoln.ac.uk/>



Preservation metadata (1)

- All digital preservation strategies depend - to a greater or lesser extent - on the creation, capture and maintenance of metadata
 - Preservation metadata:
 - The "information a repository uses to support the digital preservation process," specifically "the functions of maintaining viability, renderability, understandability, authenticity, and identity in a preservation context" (PREMIS Data Dictionary, 2005)
 - Cuts across older categorisations of metadata (descriptive, administrative, structural)



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DIC|C

Preservation metadata (2)

- PREMIS Working Group
 - Preservation Metadata: Implementation Strategies
 - Working Group sponsored by OCLC and RLG
 - Reviewed earlier Metadata Framework document and existing practice
 - Focused on implementation and definition of 'core' metadata
 - PREMIS Data Dictionary (May 2005)



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DIC|C

Preservation metadata (3)

- PREMIS Data Dictionary
 - Less explicitly based on OAIS Information Model structure than older OCLC/RLG Framework
 - Based on own data model
 - Defines some of the semantic units for: Objects, Events, Agents, Rights
 - Supports automatic capture, where possible
- PREMIS also provides:
 - An XML implementation, e.g. for use in a packaging format like METS (Metadata Encoding and Transmission Standard)
 - Maintenance activity (Library of Congress)



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DIC|C

Shared infrastructures

- For example: registries for sharing information about, or for identifying or validating formats, etc.
 - There is "... a pressing need to establish reliable, sustained repositories of file format specifications, documentation, and related software" (Lawrence, *et al.*, 2000)
 - DSpace 'bitstream format registry'
 - Global Digital Format Registry (GDFR)
 - » Some components exist, e.g. Typed Object Model, JHOVE tool
 - DCC Representation Information registry



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DIC|C

Some final comments

- The technical issues of digital preservation are only one part of a multidimensional problem
- Progress has been made on addressing technical problems
- Need for sustainability and co-operation
- Need for people with the appropriate skills



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DIC|C

Acknowledgements

UKOLN is funded by the Museums, Libraries and Archives Council, the Joint Information Systems Committee (JISC) of the UK higher and further education funding councils, as well as by project funding from the JISC, the European Union and other sources. UKOLN also receives support from the University of Bath, where it is based: <http://www.ukoln.ac.uk/>



The Digital Curation Centre is funded by the JISC and the UK e-Science Programme: <http://www.dcc.ac.uk/>



<http://www.ukoln.ac.uk/>

La preservación del patrimonio digital, Madrid, 14 al 16 marzo 2006



DIC|C