# RODA:
# digital preservation
# for the portuguese
# public administration

José Carlos Ramalho
jcr@di.uminho.pt

Miguel Ferreira
mferreira@dsi.uminho.pt

Rui Castro
Rcastro@iantt.pt

Luis Faria
lfaria@iantt.pt

Francisco Barbedo
frbarbedo@iantt.pt

Cecília Henriques
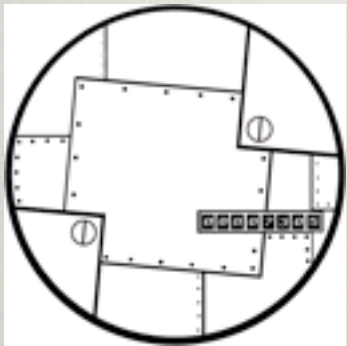chenriques@iantt.pt

Glória Santos
gloria@iantt.pt

Luis Corujo
lcorujo@iantt.pt

# Context

Digitarq (2003-now)
- metadata management (EAD based)
- digital object management (NISO MIX)

RODA (2006-2008)
- metadata management (EAD based)
- digital object management (...)
- digital preservation protocols and policies
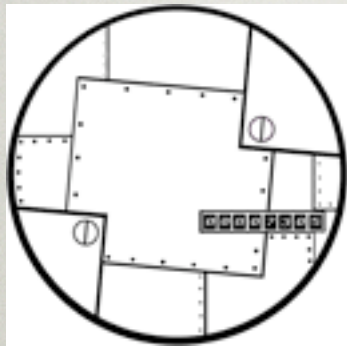
CRAV: Readers Virtual Room (2006-2007)
- request management
- document workflow

# Context



Digitarq (2003-now)
• metadata management (EAD based)
• digital object management (NISO MIX)

RODA (2006-20...)
• metadata mana...
• digital object m...
• digital preserva...

CRAV: Readers V...
• request manage...
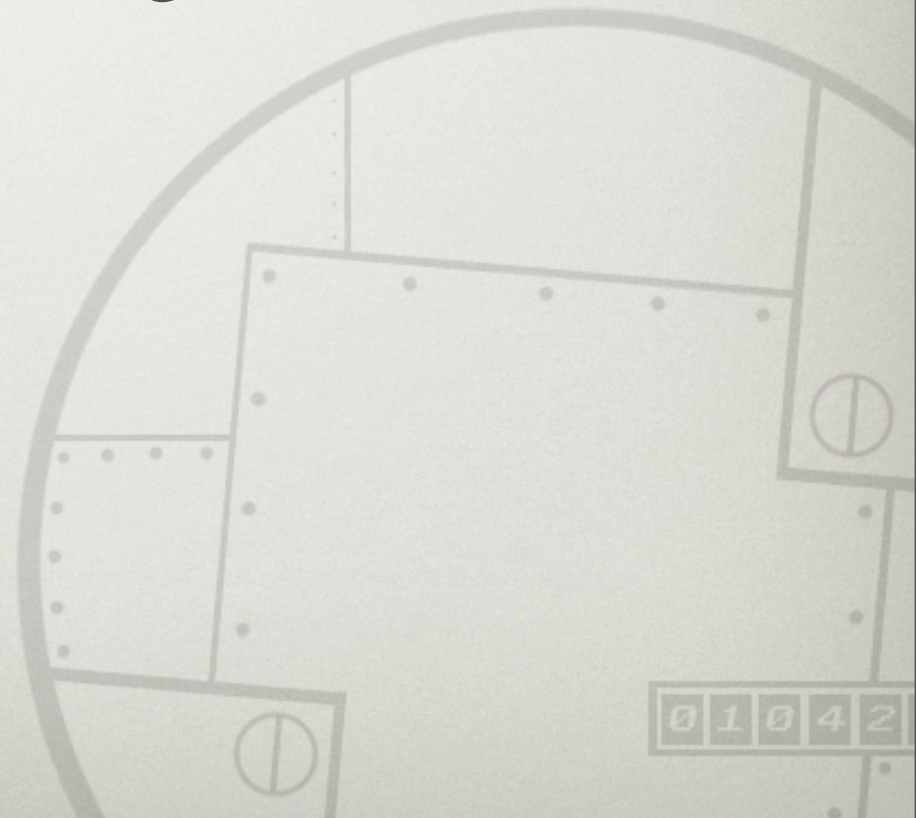• document wor...

**Partners/Contracters:**

• National Directory Board of Archives

• Photography National Archive

• Oporto's county Archive

• Some city hall archives (can grow exponencially)
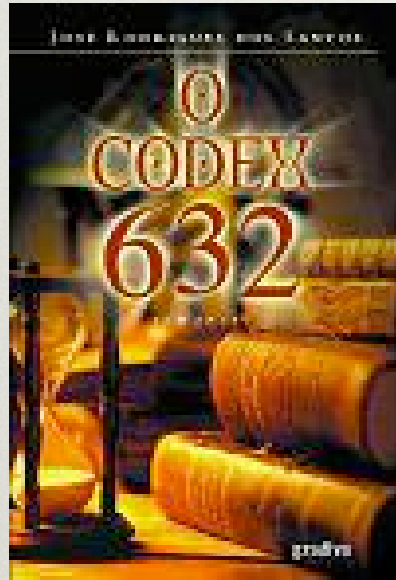
# RODA: Motivation

- Today History is being made in the digital world;

- Digital Object production grows everyday;

- There are no structures to support incorporation, management and long-term preservation of digital objects;

- We have to preserve the digital memory, heritage and testimonials of public organizations.

  - Example: SGU work

# Some Requisites/Questions?

- How do we achieve Authenticity?

- How do we describe and classify DO?

- How can we implement digital preservation?

# Authenticity



"O Codex 632" by José Rodrigues dos Santos

Subject: Who really was Cristophoros Colombus?

Was he italian? Spanish? Or a portuguese belonging to a jewish family?

# Authenticity

We must trust our sources: in ancient History there are no direct speech or evidence.
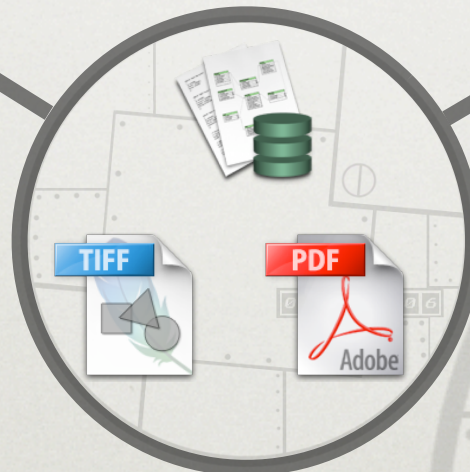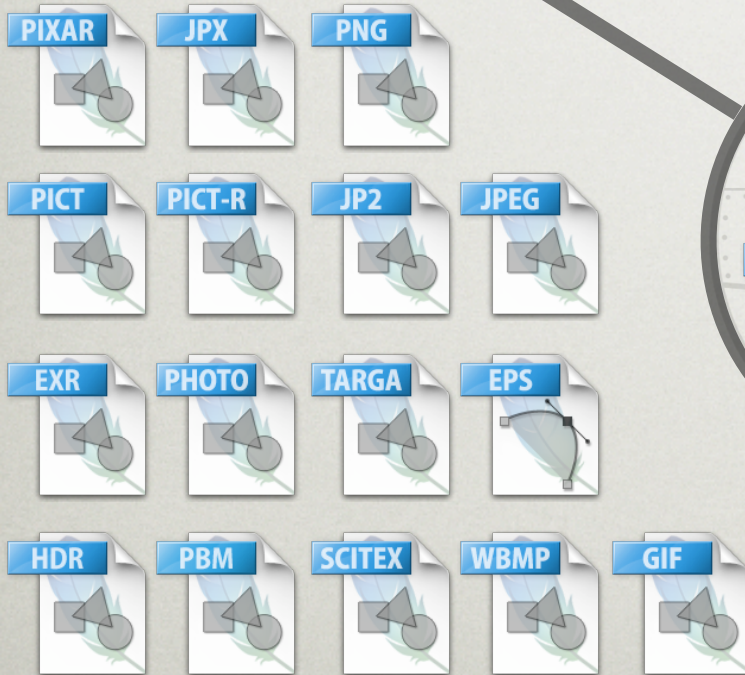
EX: the bible

# Authenticity

We must trust our sources: in ancient History there are no direct speech or evidence.

EX: the bible

How do we become trustful?

# Authenticity

We must trust our sources: in ancient History there are no direct speech or evidence.

EX: the bible

How do we become trustful?

- Reputation

- Documenting every action taken upon DOs

# Digital Object Classes

# DO Anatomy

Conceptual
level

Logical
level

Physical
level

Database
Text Doc.
Still Image

SQL S...

Ms Word Doc.

PNG image

...

**If one of these levels becomes obsolete we loose access to the DO**

# DO Preservation Strategies

- Focusing the **physical/logical object**

  o Centered in preserving information in her **logical format** or/and **physical support**

  o Uses original technology associated to these objects to ensure the access to them

  o **Technology preservation**

- Focusing the **conceptual object**

  o Centered in **preserving the object core properties in a way that is independent from hardware** and software

  o **Conceptual object preservation**

# Emulation

# Emulation

**Emulator: application** capable of reproducing the behaviour of an hardware/software platform. Ex: ZX Spectrum, GBA, ...

# Emulation

**Emulator: application** capable of reproducing the behaviour of an hardware/software platform.
Ex: ZX Spectrum, GBA, ...

- Advantages
  - o Original technological context recriation
  - o Object's *look & feel* preservation
- Disadvantages
  - o Emulators also become obsolete
  - o Users have to operate obsolete systems
  - o Creating emulators is a complex task
  - o Copyright problems
  - o To preserve a complete operating system to be able to visualize a single document may be overwhelming
  - o Information reuse in not guaranteed

# Encapsulation

# Encapsulation

Preserving the **original bit stream** together with enough metadata capable of ensuring its future interpretation and access

# Encapsulation

Preserving the **original bit stream** together with enough metadata capable of ensuring its future interpretation and access

- Advantages

  o It allows the postponement of preservation **responsibilities**

  o Targeted for objects that will be accessed in a far future

  o **Emulator and visualizer developement is delayed**

- Disadvantages

  o **Complex objects** have **complex specifications**

  o An **incomplete specification** can have nasty effects

# Conceptual object preservation

**Migration:** periodic DO transfer from one hw/sw configuration into an updated one (centered in preserving significant properties other then preserving the original bit stream).

### Advantages
– DO are disseminated in formats known to users
– No need to preserve the original hw/sw platform
– Most used strategy and the only that has worked so far

### Disadvantages
– Possible loss of information during conversion
– Continued maintenance is needed
– **In the longterm perspective costs are high**

# Conceptual object preservation

**Migration:** periodic DO transfer from one hw/sw configuration into an updated one (centered in preserving significant properties other then preserving the original bit stream).

Advantages
- DO are
- No need to preserve the original hw/sw platform
- Most used strategy and the only that has worked so far

Disadvantages
- Possible loss of information during conversion
- Continued maintenance is needed
- **In the longterm perspective costs are high**

What are the significant properties?

# Preservation Services



$t_1$   <.9, .8, .95, .1>

$t_2$   <.5, .3, .95, .6>

$t_3$   <.7, .5, .65, .1>

$t_4$   <.9, .6, .9, .7>

$t_5$   <.9, .8, .6, .1>

$t_6$   <.3, .6, .95, .1>

$t_7$   <.5, .3, .95, 1>

# Preservation Services

**CRiB project: http://crib.dsi.uminho.pt**

# Open Archival Information System



OAIS

Preservation Planning

Data Management

DM

DM

Producer

SIP

Ingest

AIP

Archival Storage

AIP

Access

DIP

Consumer

Administration

**Management**

ISO 14721

# OAIS (Functional Components)

- Ingestion

  - **Reception**, **validation**, **transformation/normalization**, description of the whole package submitted by the producer

- Storage

  - Ensures information preservation at physical/logical level (e.g. refreshing, migration, integrity checks, disaster recovery, etc.)

- Metadata management

  - Responsible for the management of stored DOs

# OAIS (Information Packages

- Submission Information Package (SIP)
  - ✳ **Digital Object**
  - ✳ **Metadata created by producer**
    - ▸ **too open...**
- Archival Information Package (AIP)
  - ✳ **Digital Object to be stored**
  - ✳ **Metadata:** enough to ensure DO's preservation and access
    - ▸ model defined by PREMIS
- Dissemination Information Package (DIP)
  - DO transformed into the **format** that will be **delivered** to the **consumer**
  - **Metadata**

# Ingestion

# Ingestion



**Submission Contract**
- SIP specification
- Ingestion workflow specification

# SIP Structure (example)



one still image

# SIP Structure (example)



one still image



criation
properties:
  - date
  - hardware
  - ...

# SIP Structure (example)



one still image

criation properties:
- date
- hardware
- ...

Technical Metadata:
- color
- dimensions
- ...

# SIP Structure (example)



one still image

criation
properties:
- date
- hardware
- ...

Technical Metadata:
- color
- dimensions
- ...

Descriptive Metadata:
- producer
- colection
- ...

# SIP Structure (example)



one still image

Manifest

criation properties:
- date
- hardware
- ...

Technical Metadata:
- color
- dimensions
- ...

Descriptive Metadata:
- producer
- colection
- ...

# SIP Structure (example)



one still image

Manifest

criation properties:
- date
- hardware
- ...

Technical Metadata:
- color
- dimensions
- ...

Descriptive Metadata:
- producer
- colection
- ...

Compressed File

18

# SIP Structure (+complex)

# SIP Structure (+complex)



```
10010010110
10100100101
```

# SIP Structure (+complex)



DO = Image+

# SIP Structure (+complex)



DO = Image+    Properties

# SIP Structure (+complex)



DO = Image+    Properties    Technical Metadata

# SIP Structure (+complex)



DO = Image+     Properties     Technical Me...     Descriptive Metadata

# SIP Structure (+complex)

Manifest

DO = Image+    Properties    Technical Metadata    Descriptive Metadata

# SIP Structure (+complex)



Manifest

DO = Image+    Properties    Technical Metadata    Descriptive Metadata

**Compressed File**

# Ingestion Workflow

## SIP ⟶ AIP

- integrity check
- virus check
- generation of preservation metada (PREMIS)
- conversion to a normalized format
- generation of technical metadata
- generation of preservation metadata (PREMIS)

# AIP Storage

# Normalization

# Data Model

# Stages

- Analysis and Planning

- Prototyping

- Testing and Dissemination

# Planning and Analysis

# Requisites

- Graphical Interface for Ingestion process

- Producer registry

- SIP production tool

- Ingestion feedback

- Partial Ingestion

- "Quarantine" zone: cache, ingestion buffer

- SIP validation

- Error reporting

- Persistent identifiers

- PREMIS event generation

- DIP digital signature

- ...

# Development framework



27

# Requisites based comparaison

# Matching data models



DSpace

# Matching data models



Fedora

# Roda Data Model

# Roda Data Model



**Description Objects**

# Roda Data Model



**Description Objects**

**Representation Objects**

# Roda Data Model



Description Objects

Representation Objects

Preservation Objects

# Architecture

# RODA Schemas

# Prototyping

# Preserving Conceptual Object

**Conceptual level**

**Logical level**

**Physical level**

Database
Text Doc.
Still Image

SQL Server

Access

PDF Doc.

Ms Word Doc.

PNG image

Hard Disc

Tape

...

# Text Documents and Still Images

# Text Documents and Still Images

- EAD elements capture most of the significant properties: provenance, producer history, context, ...

- Content is kept in a normalized format: PDF and uncompressed TIFF.

# Text Documents and Still Images

<EAD>

- EAD elements capture most of the significant properties: provenance, producer history, context, ...

- Content is kept in a normalized format: PDF and uncompressed TIFF.

# Text Documents and Still Images

<EAD>

<EAD>

- EAD elements capture most of the significant properties: provenance, producer history, context, ...

- Content is kept in a normalized format: PDF and uncompressed TIFF.

# Databases

- Data?

- Structure?

- Views?

- Reports?

- Stored Procedures?

- ...

# Databases

- Data?

- Structure?

- Views?

- Reports?

- Stored Procedures?

- ...

First prototype:

- Data

- Structure

# SIP Builder

# DBML

- Platform and RDBMS independent

- Stores the DB structure and information

- BLOBs are exported and preserved as standalone files in the representation

- Transformations to SQL and back are defined

# DBML

- Platform and RDBMS independent

- Stores

- BLOB
  standa

- Transf
  define

```
<TABLE NAME="Districts">
        <COLUMNS>
                        <COLUMN NAME="code" TYPE="int" NULL="no"/>
                ...
                </COLUMNS>
                <KEYS>
                <PKEY TYPE="simple">
                        <FIELD NAME=""/>
                </PKEY>
                <PKEY TYPE="compound">
                        <FIELD NAME=""/>
                        <FIELD NAME=""/>
                </PKEY>
                <KEY NAME="" REF=""/>

                ...
                </KEYS>
        </TABLE>
```

# DBML

- P ... <DATA>

  - <products>
    - <products-REG>
      - <code> a122 </code>
      - <description> milk </description>

      ...
    - </products-REG>
    - <products-REG>

      ...
    - </products-REG>
  - </products>

  ...
  - </DATA>

  ...

MS independent

- S

- B

s

- Transf

defined

```
                          <COLUMN NAME="code" TYPE="int" NULL="no"/>

              UMNS>

              TYPE="simple">
                          <FIELD NAME=""/>
              >
              TYPE="compound">
                          <FIELD NAME=""/>
                          <FIELD NAME=""/>
              </PKEY>
              <KEY NAME="" REF=""/>

              ...
              </KEYS>
</TABLE>
```
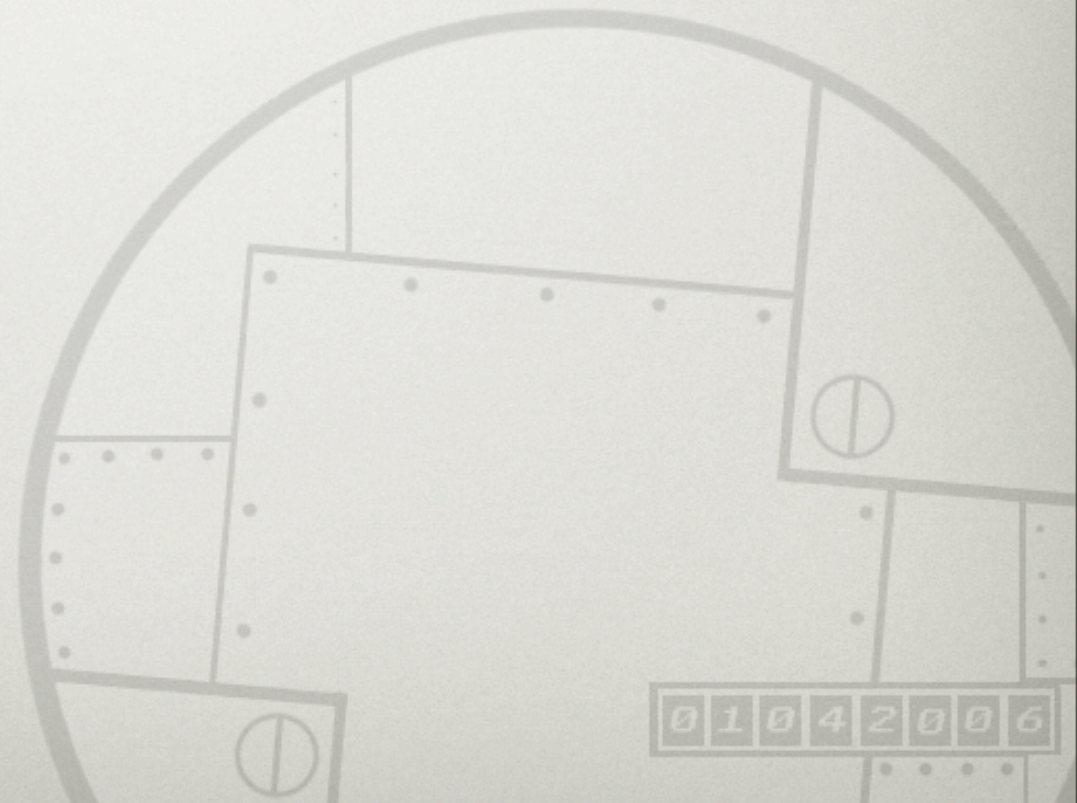
# DBML

DB SIP composition:

- METS file for packaging and organizing

- EAD file describing intellectual properties

- DBML file(s)

- DO for each found BLOB

- METS file + MIX for each DO

# SIP -> AIP

- Check and validation ...

- Generate SQL file

- Generate PREMIS

# Dissemination

- Abstract Database Creation: a database of databases... Ingests databases from DBML (DBML-->SQL$_{adb}$);

- Specific Database Creation: execute the SQL file in the selected RDMS

# Dissemination

# DB Abstract Schema

# Browser

# Browser

# Browser

# Browser

# Search Engine

# Final thoughts

*"Data Preservation is a people problem"*
*Michael Lesk*

# Final thoughts

*"Data Preservation is a people problem"*
*Michael Lesk*

- People need to be trained to save data in a proper way.
- What to preserve? Data, Structure, Semantics...
- Preservation is for future users but only today users vote on budget
- We need to make data collecting people have preservation concerns
- Preservation is fault tolerance. All systems are imperfect

# RODA Homepage

Let's Preserve Tomorrow's History...

# Questions?