# Practical Database Preservation

Download resources

**Wifi: database preservation workshop**
**Pass: database**

http://provided_address

http://192.168.2.1

Luis Faria lfaria@keep.pt
Bruno Ferreira bferreira@keep.pt

DLM 8th Triennial Conference, September 15th 2017, Brighton, UK

KEEPSOLUTIONS
University of Minho SPIN-OFF

Why do we archive databases?

# Databases hold a lot of important information

They support organisations' internal information systems

A record of the interactions between the public sector and tax payers

May be used to store scientific data

# They hold the records that do not exist in any other form

They keep record day to day business activities

# Legal or patrimonial reasons

Legislation in some countries mandates that databases are periodically archived

# Financial reasons

Maintenance costs, license fees

# Efficient way to preserve a large number of records

Little manual effort is involved when compared to archiving record by record

**Nº 10**

Aos _vinte e cinco_ dias do mez de _Fevereiro_ do anno de mil novecentos _____ n'esta Egreja parochial de _São Bartholomeu de Campello, concelho de Baião_, Diocese do Porto; eu o Presbytero Joaquim Catharino Pinto Caturão, parocho da mesma freguezia baptisei solemnemente _____ um individuo do sexo feminino a quem dei o nome de _Maria Glória_ e que nasceu n'esta freguezia lugar do Jugil de _____

ás nove horas da noite do dia _dezeseis_ do mez de _Fevereiro_ do anno de mil novecentos _____ filho legitimo de _Arthur d'Arnado Pinto e Anna Pessoa_, proprietarios naturaes, recebidos, parochianos e moradores no lugar do Jugil de _____ d'esta referida freguezia

neto paterno _de João d'Arnado Pinto e Josefa Candida_

materno _de Manoel Pinto e Maria Joaquina de Jesus_

Foi padrinho _____ Ferreira ____ madrinha Joaquina _____ madrinha Baltha... de Jesus proprietarios do lugar do Jugil de _____ d'este mencionada freguezia

os quaes todos _____ serem os proprios. E para constar se lavrou em duplicado este assento, que, depois de ser lido e conferido perante os padrinhos, _____ só o parocho os padrinhos _____ por não saberem escrever. Era ut supra

O parocho Joaquim Catharino Pinto Caturão

*[Left margin notes, partly legible:]*
Maria Gloria Nª ... Cerisa ... tinha viuva por ter fallecido seu marido José Mon... Dias no dia 17 de Dezembro do anno corrente. Registo do obito nº 394 p'la conservatoria. Baião, 17 de Dezembro de 1956. O conservador ...

... José Monteiro ... paroquia do registo civil de Baião ... em 3 Julho de 1921 Assento nº 70 de 1921. Em 23 de Outubro de 1972 o conservador ...

Nº 3 - faleceu em 6 de Fevereiro de 1950, na freguezia de Gove, registo nº ... Assento nº 37 Ano 1950 Em 6 de Fevereiro de 1950.

---

**Nº 11**

Aos _quatro_ dias do mez de _Março_ do anno de mil novecentos _____ n'esta Egreja parochial de _São Bartholomeu de Campello, concelho de Baião_, Diocese do Porto; eu o Presbytero Joaquim Catharino Pinto Caturão, parocho da mesma freguezia baptisei solemnemente _____ um individuo do sexo masculino a quem dei o nome de _Joaquim_ e que nasceu n'esta freguezia lugar de _Amarelhe_

ás sete horas da manhã do dia _vinte e cinco_ do mez de _Fevereiro_ do anno de mil novecentos _____ filho legitimo de _Manoel Pinto Carneiro e Maria da Graça_ proprietarios do lugar de Amarelhe, naturaes, recebidos e parochianos d'esta d'esta freguezia

neto paterno _de José Pinto Carneiro e de Maria Rosa de Jesus_

materno _de Bernardo Carneiro e Leonor e Marina de Jesus_

Foi padrinho _Joaquim Nicolau e Oliveira e madrinha Emilia de Jesus proprietarios do lugar de Villares d'esta mesma freguezia_

os quaes todos _____ serem os proprios. E para constar, se lavrou em duplicado este assento, que, depois de ser lido e conferido perante os padrinhos, _____ só com o padrinho por que a madrinha por não saber escrever. Era ut supra

Joaquim Nicolau de Oliveira

O parocho Joaquim Catharino Pinto Caturão

*[Left margin:]*
Joaquim
Faleceu em 24 Agosto de 1903
...
Casou ...
... Villares
*[Right margin/side note, printed:]* Porto — Typ. a vapor da Officina de S. José

*[Margin note, printed/stamped:]* Porto—Typ. a vapor da Officina de S. José

6
Carvalho

**N.º 10**

Aos vinte e cinco dias do mez de Janeiro do anno de mil novecentos ... n'esta Egreja parochial de São Bartholomeu de Campello, concelho de Baião, diocese do Porto; eu Presbytero Joaquim Catharino Pinto de Teixeira, parocho da mesma freguezia baptisei solemnemente um individuo do sexo feminino a quem dei o nome de Maria Glória e que nasceu n'esta freguezia lugar e Juzi de ... ás nove horas da noite do dia dezeseis do mez de Fevereiro do anno de mil novecentos ... filho legitima de Arthur d'Arrada Pinto e Anna Pereira, proprietarios naturaes, recebidos, parochianos e moradores no lugar do Juzi de esta referida freguezia ...

neto paterno de João d'Arrada Pinto e Joaquina Candida materno de Manuel Pinto de Teixeira e Joaquina de Jesus

Foi padrinho Teixeira Cabral Ferreira e madrinha Joaquina e Junça, madrinha Bath... ... de Jesus proprietarios do lugar do Juzi de desta mencionada freguezia.

os quaes todos s... serem os proprios. E para constar se lavrou em duplicado este assento, que, depois de ser lido e conferido perante os padrinhos, assigno só porque os padrinhos ... em mão de bem ... Era ut supra O parocho Joaquim Catharino Pinto de Teixeira

*(left margin annotations, partly illegible)*
Maria Glória
N.º 10 baptisei...
... viuva
... falleceu...
marido José Mon...
... 17 ...
bro do anno...
... Registo
do obito n.º 39...
... curia
17 de dezembro
de 1956...
...
José Montei...
...
... civil...
... em 8
de Dezembro de
1921 Assento
n.º 70 de 1921
Em 23 de Out...
bro de 1972...
Conservador...
...
N.º 3 - falleceu em
6 de fevereiro de
1950, na fregue...
... e obito...
Concelho. Assento
n.º 37 Ano 1950
Em 6 de Fevereiro
de 1950.

**N.º 11**

Aos quatro dias do mez de Março do anno de mil novecentos ... n'esta Egreja parochial de São Bartholomeu de Campello, concelho de Baião, diocese do Porto; eu Presbytero Joaquim Catharino Pinto de Teixeira, parocho da mesma freguezia baptisei solemnemente um individuo do sexo masculino a quem dei o nome de Joaquim e que nasceu n'esta freguezia lugar de Amarelhe ás sete horas da manhã do dia vinte e cinco do mez de Fevereiro do anno de mil novecentos ... filho legitimo de Manuel Pinto Carneiro e Maria da Graça proprietarios do lugar de Amarelhe, naturaes, recebidos e parochianos n'esta e desta freguezia ...

neto paterno de José Pinto Carneiro e Anna Maria de Jesus materno de Bernardo Carneiro e Joanna e Anna de Jesus

Foi padrinho Joaquim Nicolau Oliveira e madrinha ... Julia de Jesus proprietarios do lugar de Vittorezes desta mesma freguezia

os quaes todos s... serem os proprios. E para constar, se lavrou em duplicado este assento, que, depois de ser lido e conferido perante os padrinhos, assigno só com o padrinho por que a madrinha ... não sabe escrever. Era ut supra Joaquim Nicolau de Oliveira O parocho Joaquim Catharino Pinto de Teixeira

*(margin annotations)*
Joaquim
Falleceu a 24 de ...
... n.º 903
Casou
a obito

BObjectType=N'UVITSQ', BShowInternalTable=N'Yes', BOrderBy=N'T', BUpdateUsage=1

Results | Messages

| | Object Name | Ty... | Rows | Total(MB) | - | Unused(MB) | == | Used(MB) | = | Index(MB) | + |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 23 | sys.xml_index_nodes_162099618_32000 (Production.ProductModel) | IT | 350 | 0.133 | - | 0.078 | == | 0.055 | = | 0.055 | + |
| 24 | sys.xml_index_nodes_162099618_32001 (Production.ProductModel) | IT | 650 | 0.195 | - | 0.086 | == | 0.109 | = | 0.109 | + |
| 25 | sys.xml_index_nodes_270624007_32000 (Sales.Store) | IT | 9113 | 0.570 | - | 0.055 | == | 0.516 | = | 0.516 | + |
| 26 | dbo.AWBuildVersion | U | 1 | 0.016 | - | 0.000 | == | 0.016 | = | 0.008 | + |
| 27 | dbo.DatabaseLog | U | 1595 | 6.555 | - | 0.102 | == | 6.453 | = | 0.047 | + |
| 28 | dbo.ErrorLog | U | 0 | 0.000 | - | 0.000 | == | 0.000 | = | 0.000 | + |
| 29 | HumanResources.Department | U | 16 | 0.031 | - | 0.000 | == | 0.031 | = | 0.023 | + |
| 30 | HumanResources.Employee | U | 290 | 0.195 | - | 0.008 | == | 0.188 | = | 0.133 | + |
| 31 | HumanResources.EmployeeDepartmentHistory | U | 296 | 0.063 | - | 0.000 | == | 0.063 | = | 0.047 | + |
| 32 | HumanResources.EmployeePayHistory | U | 316 | 0.031 | - | 0.000 | == | 0.031 | = | 0.016 | + |
| 33 | HumanResources.JobCandidate | U | 13 | 0.227 | - | 0.070 | == | 0.156 | = | 0.031 | + |
| 34 | HumanResources.Shift | U | 3 | 0.047 | - | 0.000 | == | 0.047 | = | 0.039 | + |

Customers : Table

| CustomerID | CompanyName | ContactName | ContactTitle | Address | City |
|---|---|---|---|---|---|
| ALFKI | Alfreds Futterkis | Maria Anders | Sales Represen | Obere Str. 57 | Berlin |
| ANATR | Ana Trujillo Emp | Ana Trujillo | Owner | Avda. de la Con | México D.F. |
| ANTON | Antonio Moreno | Antonio Moreno | Owner | Mataderos 231 | México D.F. |
| AROUT | Around the Horn | Thomas Hardy | Sales Represen | 120 Hanover Sq | London |
| BERGS | Berglunds snab | Christina Berglu | Order Administr | Berguvsvägen 8 | Luleå |
| BLAUS | Blauer See Deli | Hanna Moos | Sales Represen | Forsterstr. 57 | Mannheim |
| BLONP | Blondesddsl pèr | Frédérique Citea | Marketing Mana | 24, place Klébe | Strasbourg |
| BOLID | Bólido Comidas | Martín Sommer | Owner | C/ Araquil, 67 | Madrid |
| BONAP | Bon app' | Laurence Lebiha | Owner | 12, rue des Bou | Marseille |
| BOTTM | Bottom-Dollar M | Elizabeth Lincol | Accounting Mar | 23 Tsawassen E | Tsawassen |
| BSBEV | B's Beverages | Victoria Ashwor | Sales Represen | Fauntleroy Circu | London |
| CACTU | Cactus Comidas | Patricio Simpso | Sales Agent | Cerrito 333 | Buenos Aires |
| CENTC | Centro comercia | Francisco Chan | Marketing Mana | Sierras de Gran | México D.F. |
| CHOPS | Chop-suey Chin | Yang Wang | Owner | Hauptstr. 29 | Bern |
| COMMI | Comércio Minei | Pedro Afonso | Sales Associate | Av. dos Lusíada | Sao Paulo |
| CONSH | Consolidated Ho | Elizabeth Browr | Sales Represen | Berkeley Garde | London |
| DRACD | Drachenblut Del | Sven Ottlieb | Order Administr | Walserweg 21 | Aachen |
| DUMON | Du monde entie | Janine Labrune | Owner | 67, rue des Cinc | Nantes |
| EASTC | Eastern Connec | Ann Devon | Sales Agent | 35 King George | London |
| ERNSH | Ernst Handel | Roland Mendel | Sales Manager | Kirchgasse 6 | Graz |
| FAMIA | Familia Arquiba | Aria Cruz | Marketing Assis | Rua Orós, 92 | Sao Paulo |
| FISSA | FISSA Fabrica I | Diego Roel | Accounting Mar | C/ Moralzarzal, | Madrid |
| FOLIG | Folies gourmand | Martine Rancé | Assistant Sales | 184, chaussée | Lille |

Record: 1 of 91

BObjectType-N'UVITSQ', BShowInternalTable-N'Yes', BOrderBy-N'T', BUpdateUsage-1

Results | Messages

| | Object Name | Ty... | Rows | Total(MB) | - | Unused(MB) | == | Used(MB) | = | Index(MB) | + |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 23 | sys.xml_index_nodes_162099618_32000 (Production.ProductModel) | IT | 350 | 0.133 | - | 0.078 | == | 0.055 | = | 0.055 | + |
| 24 | sys.xml_index_nodes_162099618_32001 (Production.ProductModel) | IT | 650 | 0.195 | - | 0.086 | == | 0.109 | = | 0.109 | + |
| 25 | sys.xml_index_nodes_270624007_32000 (Sales.Store) | IT | 9113 | 0.570 | - | 0.055 | == | 0.516 | = | 0.516 | + |
| 26 | dbo.AWBuildVersion | U | 1 | 0.016 | - | 0.000 | == | 0.016 | = | 0.008 | + |
| 27 | dbo.DatabaseLog | U | 1595 | 6.555 | - | 0.102 | == | 6.453 | = | 0.047 | + |
| 28 | dbo.ErrorLog | U | 0 | 0.000 | - | 0.000 | == | 0.000 | = | 0.000 | + |
| 29 | HumanResources.Department | U | 16 | 0.031 | - | 0.000 | == | 0.031 | = | 0.023 | + |
| 30 | HumanResources.Employee | U | 290 | 0.195 | - | 0.008 | == | 0.188 | = | 0.133 | + |
| 31 | HumanResources.EmployeeDepartmentHistory | U | 296 | 0.063 | - | 0.000 | == | 0.063 | = | 0.047 | + |
| 32 | HumanResources.EmployeePayHistory | U | 316 | 0.031 | - | 0.000 | == | 0.031 | = | 0.016 | + |
| 33 | HumanResources.JobCandidate | U | 13 | 0.227 | - | 0.070 | == | 0.156 | = | 0.031 | + |
| 34 | HumanResources.Shift | U | 3 | 0.047 | - | 0.000 | == | 0.047 | = | 0.039 | + |
| 35 | Person.Address | | | | | | | | | | |

**Record #10**
**Record #11**

**Customers : Table**

| CustomerID | CompanyName | ContactName | ContactTitle | Address | City |
|---|---|---|---|---|---|
| ALFKI | Alfreds Futterkis | Maria Anders | Sales Represen | Obere Str. 57 | Berlin |
| ANATR | Ana Trujillo Emp | Ana Trujillo | Owner | Avda. de la Con | México D.F. |
| ANTON | Antonio Moreno | Antonio Moreno | Owner | Mataderos 231 | México D.F. |
| AROUT | Around the Horn | Thomas Hardy | Sales Represen | 120 Hanover Sq | London |
| BERGS | Berglunds snab | Christina Berglu | Order Administr | Berguvsvägen 8 | Luleå |
| BLAUS | Blauer See Deli | Hanna Moos | Sales Represen | Forsterstr. 57 | Mannheim |
| BLONP | Blondesddsl pèr | Frédérique Citea | Marketing Mana | 24, place Klébe | Strasbourg |
| BOLID | Bólido Comidas | Martín Sommer | Owner | C/ Araquil, 67 | Madrid |
| BONAP | Bon app' | Laurence Lebiha | Owner | 12, rue des Bou | Marseille |
| BOTTM | Bottom-Dollar M | Elizabeth Lincol | Accounting Mar | 23 Tsawassen E | Tsawassen |
| BSBEV | B's Beverages | Victoria Ashwor | Sales Represen | Fauntleroy Circu | London |
| CACTU | Cactus Comidas | Patricio Simpso | Sales Agent | Cerrito 333 | Buenos Aires |
| CENTC | Centro comercia | Francisco Chan | Marketing Mana | Sierras de Gran | México D.F. |
| CHOPS | Chop-suey Chin | Yang Wang | Owner | Hauptstr. 29 | Bern |
| COMMI | Comércio Minei | Pedro Afonso | Sales Associate | Av. dos Lusíada | Sao Paulo |
| CONSH | Consolidated Ho | Elizabeth Browr | Sales Represen | Berkeley Garde | London |
| DRACD | Drachenblut Del | Sven Ottlieb | Order Administr | Walserweg 21 | Aachen |
| DUMON | Du monde entie | Janine Labrune | Owner | 67, rue des Cinc | Nantes |
| EASTC | Eastern Connec | Ann Devon | Sales Agent | 35 King George | London |
| ERNSH | Ernst Handel | Roland Mendel | Sales Manager | Kirchgasse 6 | Graz |
| FAMIA | Familia Arquiba | Aria Cruz | Marketing Assis | Rua Orós, 92 | Sao Paulo |
| FISSA | FISSA Fabrica I | Diego Roel | Accounting Mar | C/ Moralzarzal, | Madrid |
| FOLIG | Folies gourmand | Martine Rancé | Assistant Sales | 184, chaussée | Lille |

Record: 1 of 91

**person**

| **id** | **name** | **birth** | **city_id** |
|:------:|:--------:|:---------:|:-----------:|
| 1 | Mary | 1986-03-28 | 5 |
| 2 | Phillip | *NULL* | 6 |

**city**

| **id** | **name** | **mayor** | **country_id** |
|:------:|:--------:|:---------:|:--------------:|
| 5 | Payne Springs | 1 | 16 |
| 6 | Rosenhayn | *NULL* | 16 |

**country**

| **id** | **name** |
|:------:|:--------:|
| 16 | United States |

- Tables
- Column data types
- Relations
- Constraints
- Projections (views)
- Behaviour (triggers)
- Other (users, permissions, etc.)

# SIARD 2.0

## Specification for archiving relational databases

It's a vendor independent database archiving format

Stands for Software Independent Archiving of Relational Databases

Version 2 is a joint effort of the Swiss Federal Archives, the E-ARK project and the eCH (eGovernment standards in Switzerland)


Based on international standards

Unicode, SQL:1999, XML, XML Schema, URI


SIARD 2 is E-Government Standard (eCH-0165)

Replaced version 1.0 in 2016

What it does…

# Preserves information — not layout or interaction

The application and the business logic are not preserved

# Preserves data — not code

Stored procedures, functions and other code-like properties are stored but not preserved

# Preserves tables and relations —not dynamic data

Views are not be preserved in SIARD, but they can be materialised

Some technical details…

# Data is stored in a folder-like structure

Optionally, the folder can be compressed as ZIP for storage saving purposes

# A header folder stores database metadata

General information about the archived database and database structure stored as XML

# A content folder stores database content

Tabular data stored as multiple XML files

# Binary objects are stored in 3 different ways

Inline, inside, outside

# file.siard

📁 header
- 📄 metadata.xml
- 📄 metadata.xsd
- 📄 metadata.xsl

📁 content
- 📁 schema0
  - 📁 table0
    - 📄 table0.xml
    - 📄 table0.xsd
  - 📁 table1
    - 📄 table1.xml
    - 📄 table1.xsd
    - 📁 lob0
      - 📄 record0.bin
      - 📄 record1.bin
      - 📄 …
  - 📁 table2
    - …
  - 📁 …

```xml
<table>
  <name>address</name>
  <folder>table2</folder>
  <description>This table contains addresses</description>
  <columns>
    <column>
      <name>address_id</name>
      <type>SMALLINT</type>
      <typeOriginal>SMALLINT UNSIGNED</typeOriginal>
      <nullable>false</nullable>
      <description>The address unique id.</description>
    </column>
    <column>
      <name>address</name>
      <type>CHARACTER VARYING(50)</type>
      <typeOriginal>VARCHAR</typeOriginal>
      <nullable>false</nullable>
      <description>First address line</description>
    </column>
    <column>
      <name>address2</name>
      <type>CHARACTER VARYING(50)</type>
      <typeOriginal>VARCHAR</typeOriginal>
      <nullable>true</nullable>
      <description>Second address line</description>
    </column>
    <column>
      <name>district</name>
      <type>CHARACTER VARYING(20)</type>
      <typeOriginal>VARCHAR</typeOriginal>
      <nullable>false</nullable>
      <description>Address district</description>
    </column>
    <column>
      <name>city_id</name>
      <type>SMALLINT</type>
      <typeOriginal>SMALLINT UNSIGNED</typeOriginal>
      <nullable>false</nullable>
      <description>Address city (id)</description>
```

Information about the database and its structure

```xml
<?xml version="1.0" encoding="UTF-8"?>
<table
  xsi:schemaLocation="http://www.admin.ch/xmlns/siard/2.0/schema1/table2.xsd table2.xsd"
  xmlns="http://www.admin.ch/xmlns/siard/2.0/schema1/table2.xsd"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" >
  <row>
    <c1>1</c1>
    <c2>47 MySakila Drive</c2>
    <c4>Alberta</c4>
    <c5>300</c5>
    (other columns omitted)
  </row>
  <row>
    <c1>2</c1>
    <c2>28 MySQL Boulevard</c2>
    <c4>QLD</c4>
    <c5>576</c5>
    (other columns omitted)
  </row>
  <row>
    <c1>3</c1>
    <c2>23 Workhaven Lane</c2>
    <c4>Alberta</c4>
    <c5>300</c5>
    (other columns omitted)
  </row>
  (remaining rows omitted)
</table>
```

Table content

# Toolset

database preservation toolkit — Data extraction

RODA — Data archiving

database visualization toolkit

**database preservation** toolkit — Data extraction

**○ RODA** — Data archiving

**database visualization** toolkit — Data visualisation

Use case scenario

An information system has been in use for 15 years

It is about to be decommissioned as it has been replaced by a new, more advanced system

There is an interest in maintaining the data produced by the legacy information system for legal and historical reasons

This means that the system's database has been selected for long-term archival

Let's see how it works…

Database

Database

Database

database preservation toolkit

Archival format

SIARD 2.0

Database

database preservation toolkit

Archival format

SIARD 2.0

Transfer & ingest

EARK SIP

RODA in
SIP creation tool

Database

database preservation toolkit

Archival format

SIARD 2.0

Transfer & ingest

RODA
Long-term digital repository

EARK SIP

RODA in
SIP creation tool

Years later, a user wants to access the data…

Discovery services

Discovery services

RODA
Long-term digital repository

SIARD 2.0

Discovery services

RODA
Long-term digital repository

SIARD 2.0

database
visualization toolkit

Web access
to data

Discovery
services

RODA
Long-term digital repository

database
visualization toolkit

SIARD 2.0

Using the
Database
Preservation
Toolkit

Current live
database
system

In summary…

# The Database Preservation Toolkit

Extracts data from the original RDBMS and stores it in a long-term archival format, i.e. SIARD 2

# The SIARD 2 archival format

Can be submitted to a long-term preservation repository for continuous preservation, monitoring and access

Or, alternatively, can be written to a storage infrastructure

If years later, a user wants to access the archived data…

# Lookup the database

Search the repository catalogue to find the right database

# View database on a light-weight viewer

Using the Database Visualization Toolkit

# Export parts of the database

For a shallow analysis or simple analytics using, e.g. Microsoft Excel

# Load the database into a new RDBMS

Perform more advanced data analysis, e.g. OLAP

# IT department is happy!
as it no longer needs to maintain the legacy database system

# Management is happy!
because costs have been greatly reduced

# Demo

# Summary

database preservation toolkit

# Extracts/loads data from/to relational databases

Stores data in preservation formats: SIARD 1, SIARD DK, SIARD 2

# Command-line tool

Enables using the tool over SSH

# Multi-platform

Runs on Windows, Linux, MacOSX, and any Java compatible operating system

# Lightweight and highly performant

Tests reveal a throughput of 88.000 records/second on a large text database

# Viewing and navigating on relational data

Lists data, follows relationships, shows structure, users, data, triggers, stored procedures, functions, metadata, etc.

# Searching, filtering, column hiding, etc.

AND, OR, NOT and "exact phrase"

# Export data

After filtering, the user can export data to CSV

# Copes with millions of records

Supported by NoSQL horizontally scalable technologies

In conclusion…

Standards and tools are available right now

Documentation and source code published on GitHub

Examples, video tutorials are also available

Tools have been piloted in several institutions

Pilots report will be available on the E-ARK web site

http://www.eark-project.com/resources/project-deliverables/97-d25-1/file

www.database-preservation.com

# Questions?