# Preservation of Real-World Data: The Case for Preservation DataStores

Simona Cohen, **Michael Factor**, Dalit Naor,
Leeat Ramati, Petra Reshef, Shahar Ronen
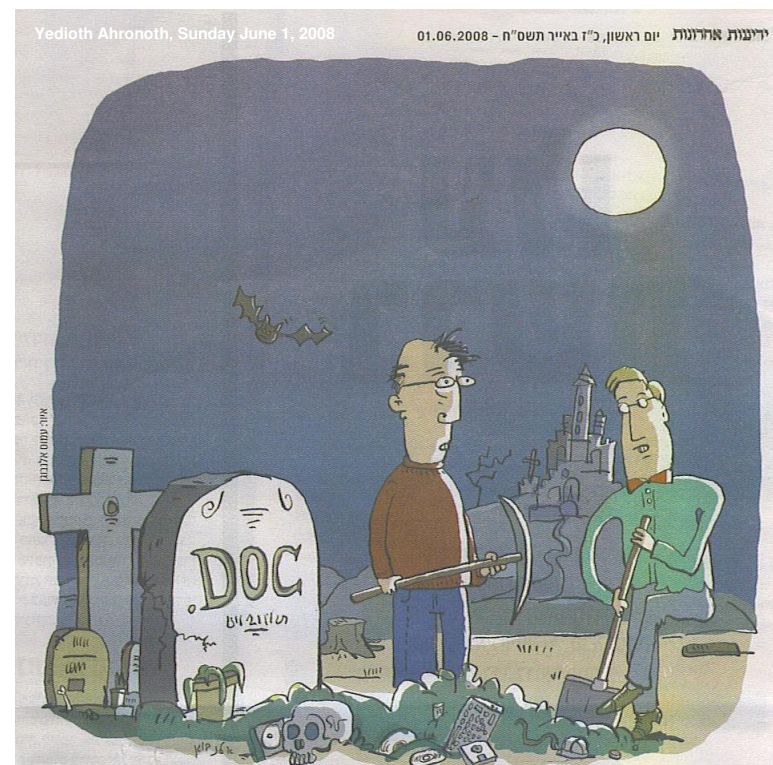
IBM Haifa Research Labs

http://www.haifa.il.ibm.com/projects/storage/ltdp/index.shtml

# What is Long Term Digital Preservation?

◈ *Long Term Digital Preservation (LTDP)* is a means of keeping digital information such that the same information can be used at some point in the future in spite of obsolescence of everything: hardware, software, processes, format, people, etc.

- ◈ *Bit Preservation* addresses obsolescence of hardware
- ◈ *Information or Logical Preservation* addresses obsolescence of everything else



Yedioth Ahronoth, Sunday June 1, 2008

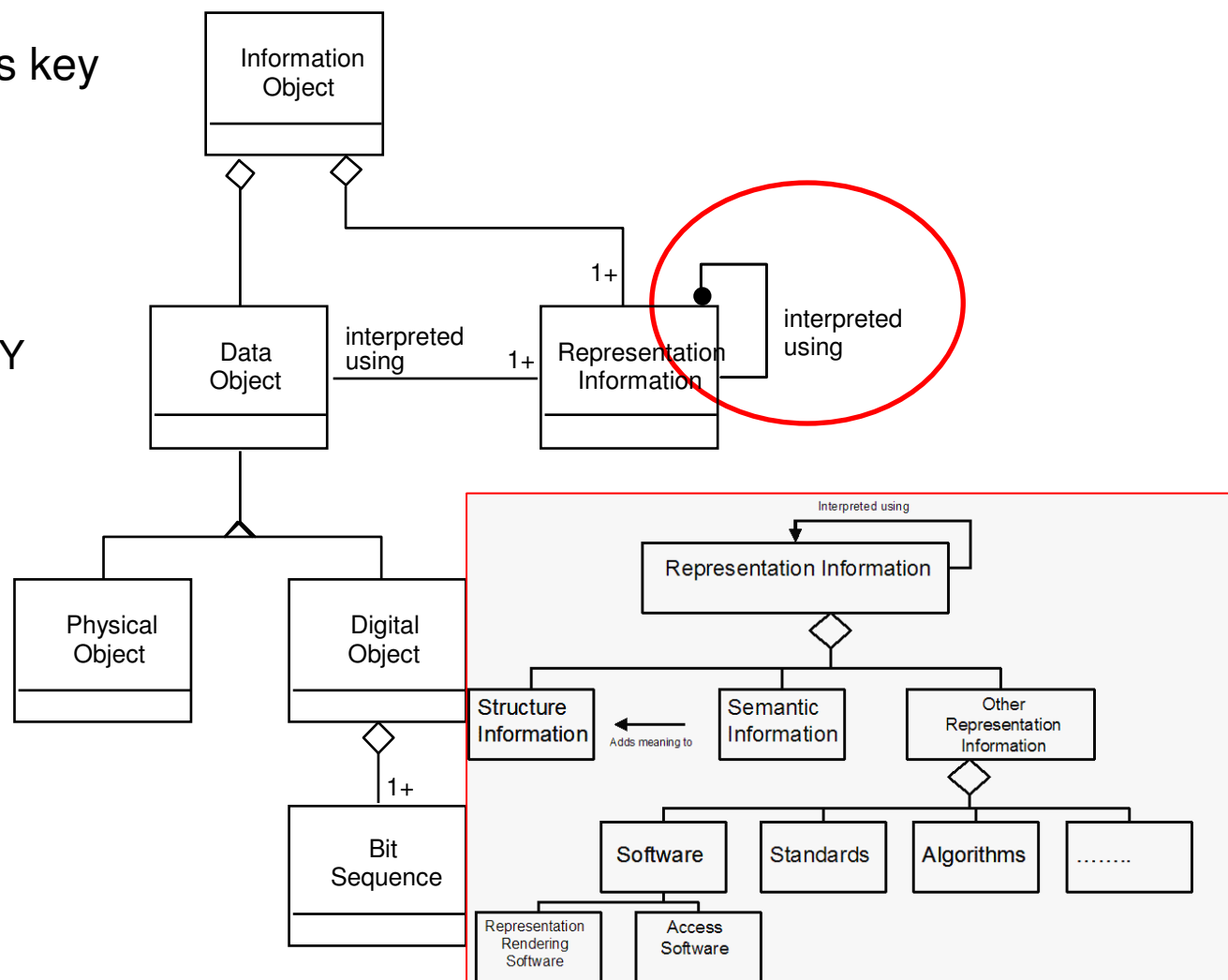01.06.2008 – ידיעות אחרונות יום ראשון, כ"ז באייר תשס"ח

# OAIS Information Model & Representation Information

The Information Model is key

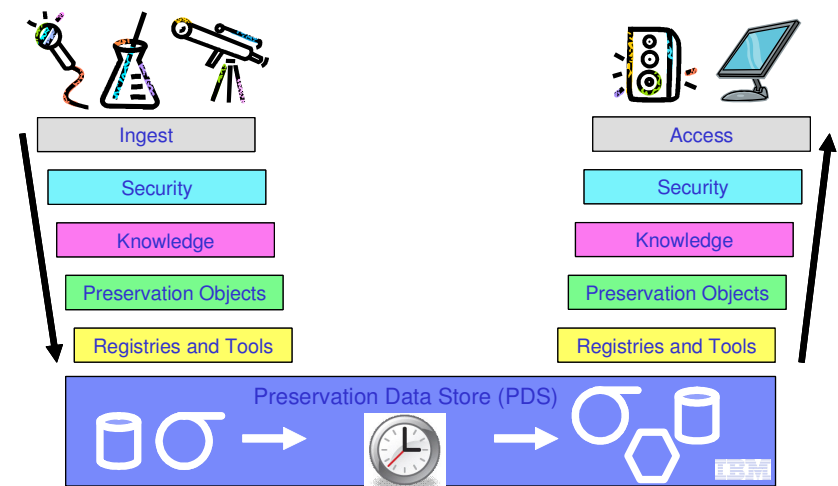Recursion ends at
KNOWLEDGEBASE of the
DESIGNATED COMMUNITY
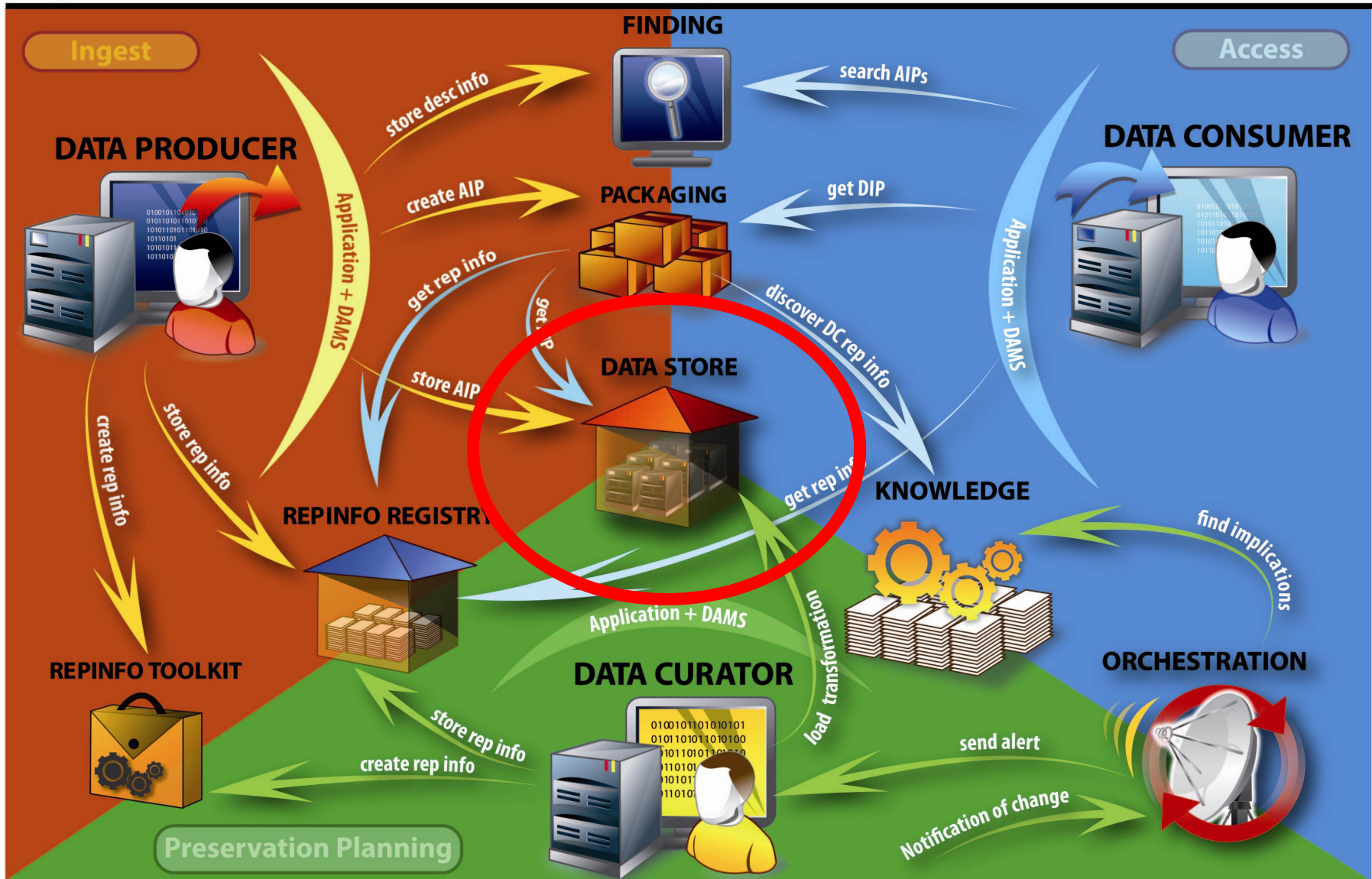
(this knowledge will change
over time and region)

**Information Object**

**Data Object** — interpreted using — 1+ — **Representation Information** — interpreted using

1+

1+

**Physical Object**

**Digital Object**

1+

**Bit Sequence**

Interpreted using

**Representation Information**

Structure Information ← Adds meaning to — Semantic Information — Other Representation Information

Software | Standards | Algorithms | ........

Representation Rendering Software | Access Software

Based upon presentation by David Giaretta at EVA/MINERVA 2006: The Third Annual Jerusalem Conference on the Digitisation of Cultural Heritage

3

© 2008 IBM Corporation

# CASPAR and Preservation DataStores

◈ Objective: Demonstrate validity of OAIS framework with heterogeneous data
  ◈ Open Archival Information System (OAIS) is an ISO Standard Reference Model for Long Term Preservation

◈ Preservation DataStores (PDS)
  ◈ IBM's responsibility
  ◈ Show how to build standards-based storage that is preservation aware:
    ◈ OAIS-required metadata
    ◈ Transform formats to avoid information obsolescence
    ◈ Manage media for bit obsolescence

◈ Other partners include data providers such as ESA, UNESCO, etc.

◈ **On June 24th, CASPAR (including PDS) successfully demonstrated preservation of heterogeneous data to the project officer of the European Union.**

http://www.casparpreserves.eu/ -- http://www.haifa.il.ibm.com/projects/storage/datastores/caspar.html

| Ingest | | Access |
| --- | --- | --- |
| Security | | Security |
| Knowledge | | Knowledge |
| Preservation Objects | | Preservation Objects |
| Registries and Tools | | Registries and Tools |

Preservation Data Store (PDS)

# CASPAR Testbed Data

◈ **Scientific data**

 ◈ European Space Agency (ESA) – IT; Science and Technology Facilities Council (STFC) – UK

 ◈ Complex digital objects, oriented towards processing, may be high-volume

◈ **Cultural heritage**

 ◈ UNESCO

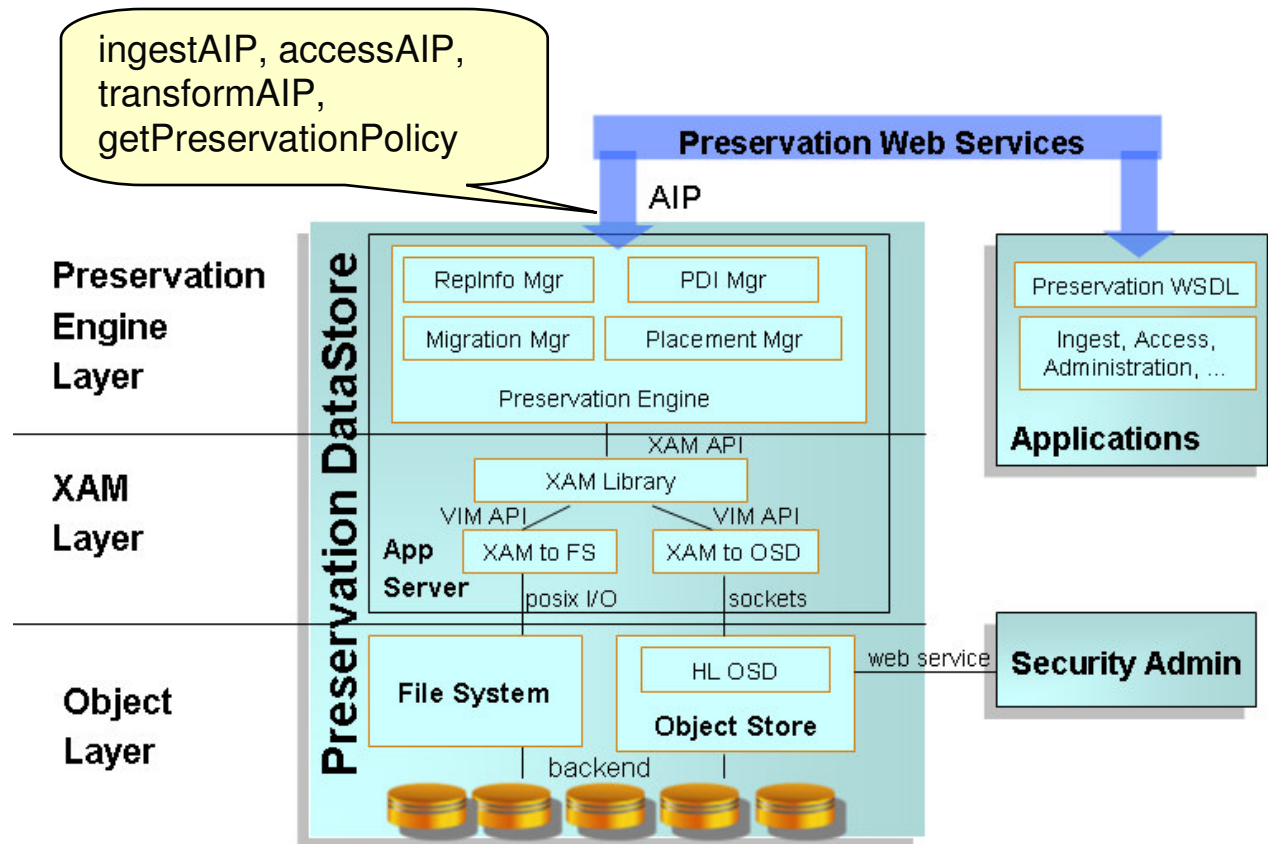 ◈ Dynamic interactive digital objects, oriented towards presentation and replay

◈ **Artistic data**

 ◈ INA, University of Leeds, CIANT, CNRS, Ircam

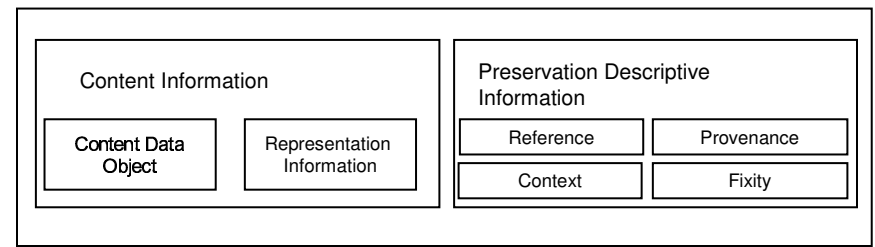 ◈ Virtual digital objects, spanning between processing and display

# PDS Architecture

◈ Layered approach based on open standards - OAIS, XAM, OSD

◈ In CASPAR, all layers are utilized. In other embodiments, only some layers can be used

◈ Utilize XAM to provide logical abstraction of containers (XSets)

◈ Offload preservation functionality to the storage



ingestAIP, accessAIP, transformAIP, getPreservationPolicy

**Preservation Web Services**

AIP

**Preservation DataStore**

**Preservation Engine Layer**
- RepInfo Mgr
- PDI Mgr
- Migration Mgr
- Placement Mgr
- Preservation Engine

XAM API

**XAM Layer**
- XAM Library
- VIM API
- App Server
- XAM to FS
- XAM to OSD
- posix I/O
- sockets

**Object Layer**
- File System
- HL OSD
- Object Store
- backend

**Applications**
- Preservation WSDL
- Ingest, Access, Administration, ...

web service **Security Admin**

# Preservation DataStores (PDS)

◈ OAIS–based preservation-aware storage, media-agnostic and generic storage to support logical preservation

◈ Manage preservation specific metadata
  ◆ Fixity computations
  ◆ Update technical provenance
  ◆ Manage the PDI RepInfo
  ◆ Ensuring referential integrity

| Content Information | Preservation Descriptive Information |
|---|---|
| Content Data Object / Representation Information | Reference / Provenance / Context / Fixity |

◈ Storlet container
  ◆ Module container that can execute restricted modules with predefined interfaces for data intensive functions, e.g., transformations, fixity calculation.
    ◈ Optimal scheduling
  ◆ Update PDS modules (e.g., fixity algorithm, packaging format)

◈ Managing availability/ data loss
  ◆ Physically co-locate data and metadata
  ◆ Cluster Related AIPs on the same media unit based upon their relative importance

◈ AIP identifier generation – Globally unique identifiers

# Preservation of an AIP

◈ Ingest AIP
- ◈ Storing an AIP in PDS

◈ Bit and logical preservation of an AIP
- ◈ Bit migrations
- ◈ Format migrations - data transformation
  - ◈ Transformation modules are packed as AIPs and preserved
  - ◈ Transformation result is a new version to the original AIP
- ◈ Migrations are documented as Provenance records
- ◈ During migrations PDS performs operations on AIP
  - ◈ Update PDI (e.g. Fixity calculation, additional Provenance events)
  - ◈ Execute previously loaded storlets

◈ Access AIP
- ◈ Retrieval of an AIP
  - ◈ By retrieval time the original AIP may have several versions and copies

# MST data and PDS



The Natural Environment Research Council (NERC) Mesosphere-Stratosphere-Troposphere (MST) Radar at Aberystwyth

◈ MST data ingested to PDS

  ◈ Highly complex atmospheric data preserved for long periods in self-describing NetCDF binary format

◈ PDS handles logical preservation

  ◈ Handle metadata

    ◈ Update provenance, compute fixity etc.

  ◈ Load and execute transformation

    ◈ The NASA-AMES format has become the preferred data format of the atmospheric scientists community

◈ Access

  ◈ By access time there are several versions to the original AIP

# Example: MST RepInfo Network

# Ingested MST AIP

NetCDF format

# Accessed MST data

NASA-AMES format

# Ingested MST AIP



# Access RepInfo of original AIP

NASA-AMES format

## Access new MST data

## Access new AIP provenance

# Some other things we are doing in IBM Research

◈ Storage Networking Industry Association (SNIA)

  ◈ LTACSI - Long Term Archive & Compliance Storage Initiative

  ◈ IBM co-chairs the Long Term Retention TWG

    ◈ A recently-formed working group focusing on a Self-contained Information Retention Format

◈ Long Term Digital Preservation Assessment

  ◈ Research tool to evaluate organization's ability to preserve its digital resources

  ◈ Based upon emerging standard audit checklists (ISO 14721)

# More Info

http://www.haifa.il.ibm.com/projects/storage/ltdp/index.shtml

Michael Factor: factor@il.ibm.com

IBM Haifa Research Lab