

CASPAR Cultural, Artistic and Scientific knowledge for Preservation, Access and Retrieval

Preservation Information Conceptual Model (and Knowledge Management)

Presentation by: Yannis Marketakis

Yannis Tzitzikas, Vassilis Christophides, Grigoris Antoniou, Dimitris Kotzinos, Giorgos Flouris, Yannis Marketakis, Stamatis Karvounarakis, Dimitris Andreou, Yannis Theoharis
 FORTH-ICS
 CASPAR

CASPAR

Outline

- Motivation & Introduction
- Aspects of Preservation
- Preservation of Intelligibility (and OAIS)
 - Formalizing the Problem
 - Intelligibility Gap
- Modeling and Preserving Provenance
 - CIDOC CRM
- Knowledge Manager (one of the CASPAR key components)
 - Responsibilities
 - SWKM
 - GapManager
- Available Data and Current Deployments
- GapManager and Data Holders

CASPAR

Important Questions


- Important questions
 - *What digital information preservation is?*
 - *What kind (and how much) representation information do we need? How this depends on the Designated community?*
 - *What kind of automation could we offer?*
- CASPAR contribution
 - *A formalization of intelligibility and intelligibility gap through the notion of dependency. This approach can provide some answers to the above questions.*

CASPAR


The problem of Preserving the Intelligibility and Provenance

Phaistos disk (dated to 1700 BC)

We still cannot understand it (the meaning has not been preserved)



We still don't know how the pyramids were constructed. (the process has not been preserved)

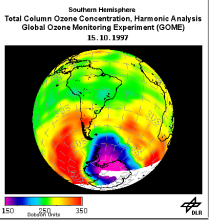


CASPAR

The problem of Preserving the Intelligibility and Provenance

How can we be sure that in the future one would be able to understand this byte stream?

100110110000110111011011101110010111100111



Southern Hemisphere
 Total Column Ozone Concentration, Harmonic Analysis
 Global Ozone Monitoring Experiment (GOME)
 15. 10. 1997

- *How this image has been derived?*
- *When and by whom it was taken?*
- *How the satellite image was processed (by what algorithms and with what parameters)?*

CASPAR

Aspects of Preservation

But what should we preserve?

- For sure we have to preserve the **bits** of the digital objects

We should also try to preserve the **information** carried by the digital objects

- Their accessibility
- Their integrity
- Their authenticity
- Their provenance
- Their intelligibility (by human or artificial actors)

OASIS and Intelligibility

- According to OASIS, metadata are distinguished to various categories.
- One very important is that of **Representation Information**
 - Aim at **enabling the conversion of bits to something useful**

7

Modules and Dependencies

- In order to abstract from the various domain-specific and time-varying details, we introduce the general notions of **Module** and **Dependency**.
- Module**
 - We adopt a very general definition. A module could be:
 - a piece of software/ hardware module.
 - a knowledge model expressed explicitly and formally (e.g. an Ontology)
 - a knowledge model not expressed explicitly (e.g. GreekLanguage)
 - (the only constraint is that modules need to have a unique identity)
- Dependency**
 - A module t depends on t' , written $t > t'$, if t requires t'
 - The meaning of a dependency $t > t'$
 - t cannot function/be understood/managed without the existence of t'
 - in general we can support several dependency types (they are goal-directed)
- Note:
 - We model the RI requirements of OASIS as **dependencies** between modules.

8

Modules and Dependencies: Examples

9

Modules and Dependencies: Examples (Semantic Web data)

10

Formalizing Actor/Community knowledge (in terms of modules and dependencies)

- Each actor or community u can be characterized by a **profile** T_u that contains those modules that are **assumed to be available/known to u** .
- Formalization: $T_u \subseteq T$ (where T is the set of all modules)

Examples

- u is an artificial agent
 - T_u may include the software/hardware modules available to it
- u is a human,
 - T_u may include modules that correspond to implicit knowledge

11

The notion of closure (of modules and profiles)

- Closure of a module t : $C(t)$ = all modules on which it depends
- Closure of a set of modules S : $C(S) = \cup \{ C(t) \mid t \in S \}$
- Required modules of t : $C^-(t) = C(t) - \{t\}$

$C(t_x) = C(t_x) - \{t_x\}$
 $C(t_y) = C(t_y) - \{t_y\}$
 Closure of T_u

Let user u be a user with profile T_u
 Let t be a module
Intelligibility Gap:
 The smallest set of extra modules that u needs to have in order to understand a module t .
 Notation: **Gap(t,u)**

12

Intelligibility and Intelligibility Gap (I)

- u can understand t iff: $C^*(t) \subseteq C(Tu)$
- The intelligibility gap: $Gap(t,u) = C^*(t) - C(Tu)$

This means that:

- if we want to preserve a digital object t for a community with profile Tu then we need to get and store only $Gap(t,u)$ plus an id that denotes Tu.
- if we want to deliver an object t to an actor with profile Tu, then the only extra modules we should deliver to him in order to return him something intelligible, is the set $Gap(t,u)$.

$Gap(ty,u) = \emptyset$ $Gap(tx,u) = \{t1, t2, t4, t5\}$

Exploiting DC Profiles for constructing the "right" AIPs (intelligible and redundancy free)

DC profiles could be exploited so that to be able to derive different AIPs and DIPs for different DCCommunities (if the dependencies are available this could be done automatically)

AIP of o1 wrt DC1: Object = o1, DCprofile = DC1, deps = {t1,t3}
 AIP of o1 wrt DC2: Object = o1, DCprofile = DC2, deps = {t1,t2,t4}
 AIP of o1 wrt DC3: Object = o1, DCprofile = DC3, deps = {t1,t2,t3,t4,t5,t6}

Scenario: Intelligibility-aware Packaging

DC Profiles

- P1 = {FITS} // for astronomers
- P2 = {PDF, XML} // for casual users
- P3 = {C3D, DirectX, MAX/MSP}

Objects

- o1 // a FITS file
- o2 // a pdf document
- o3 // a zip file containing multimedia performance data

Scenario: Intelligibility-aware Packaging

- $Gap(o2,P1) = \emptyset$
- $Gap(o2,P2) = \{FITS_STANDARD, FITS_DICTIONARY, DICTIONARY_SPECIFICATION\}$
- $Gap(o2,P3) = \{FITS_STANDARD, FITS_DICTIONARY, DICTIONARY_SPECIFICATION, PDF_STANDARD, XML_SPECIFICATION, UNICODE_SPECIFICATION\}$
- $Gap(o3,P3) = \{ZIP\}$
- $Gap(o3, P2) = \{ZIP, C3D, DirectX, MAX/MSP\}$

Example from cultural testbed

ASCII grid file

```

#NCSA 10
#NAME 10
#L1200000 213860 2167633
#C1200000 482500 1074054
#CELLSIZE 0.369176
#NCSA_VFILE -NONE
#
20.242000 20.493000 20.489000 20.492000 20.794000 20.798000 20.792000 20.796000 20.488000 20.492000 20.494000
20.492000 20.496000 20.494000 20.498000 20.487000 20.491000 20.493000 20.489000 20.495000 20.497000 20.493000
20.494000 20.498000 20.496000 20.500000 20.486000 20.490000 20.492000 20.488000 20.494000 20.496000 20.494000
20.496000 20.500000 20.498000 20.502000 20.485000 20.489000 20.491000 20.487000 20.493000 20.495000 20.493000
20.498000 20.502000 20.504000 20.506000 20.484000 20.488000 20.490000 20.486000 20.492000 20.494000 20.492000
20.500000 20.504000 20.506000 20.508000 20.483000 20.487000 20.489000 20.485000 20.491000 20.493000 20.491000
20.502000 20.506000 20.508000 20.510000 20.482000 20.486000 20.488000 20.484000 20.490000 20.492000 20.490000
20.504000 20.508000 20.510000 20.512000 20.481000 20.485000 20.487000 20.483000 20.489000 20.491000 20.489000
20.506000 20.510000 20.512000 20.514000 20.480000 20.484000 20.486000 20.482000 20.488000 20.490000 20.488000
20.508000 20.512000 20.514000 20.516000 20.479000 20.483000 20.485000 20.481000 20.487000 20.489000 20.487000
20.510000 20.514000 20.516000 20.518000 20.478000 20.482000 20.484000 20.480000 20.486000 20.488000 20.486000
20.512000 20.516000 20.518000 20.520000 20.477000 20.481000 20.483000 20.479000 20.485000 20.487000 20.485000
20.514000 20.518000 20.520000 20.522000 20.476000 20.480000 20.482000 20.478000 20.484000 20.486000 20.484000
  
```

Modeling Modules and Dependencies (usage of SW languages for extensibility / inheritance)

Extensible typology of modules and dependency types

Example: Extending the typology of dependencies

The **goal** determines the dependencies. For this reason we allow specializing the dependency relation

```

    graph LR
      ajava[a.java] -- _compile --> javac[javac]
      ajava -- _edit --> ASCII[ASCII]
      aclass[a.class] -- _run --> JVM[JVM]
      subgraph dependsOn
        _run --> dependsOn
        _compile --> dependsOn
        _edit --> dependsOn
      end
  
```

Assume a DC profile $pr1 = \{ASCII\}$
 $gap(a.java, pr1) = \{javac\}$ // not type of dependency is specified (so all count)
 $gap(a.java, pr1, _edit) = \emptyset$ // here the type of dependency is specified

19

Aspects of Preservation

- For sure we have to preserve the **bits** of the digital objects
- We should also try to preserve the **information** carried by the digital objects
 - Their accessibility
 - Their integrity
 - Their authenticity
 - Their **provenance**
 - Their intelligibility (by human or artificial actors)

20

Modeling Descriptive Metadata (Provenance, Context, etc)

- Contributions in extending and specializing CIDOC CRM
 - CIDOC Conceptual Reference Model ISO/FDIS 21127
- Can be used as the conceptual backbone for descriptive metadata

21

Example of modeling provenance

Change of custody chain

22

Example of modeling provenance

Conversion - case 1

Figure: JPEG2000 Converter

23

Cultural, Artistic and Scientific knowledge for Preservation, Access and Retrieval

(Knowledge Manager) Component Description

24

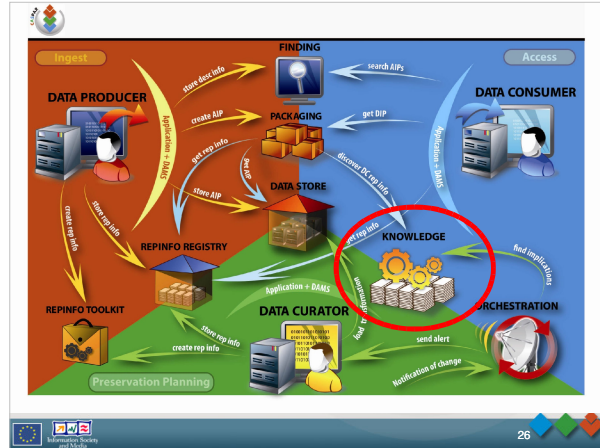
Component responsibility

Responsibilities

- Capture Higher level Semantics
- Manage Designated Community Knowledge profile
- Identify ReplInfo Gaps
- Manage Ontologies and Metadata

It comprises two layers

- A) **SWKM** (Semantic Web Knowledge Middleware) which is the lower layer providing a set of core services for managing Semantic Web data.
- B) **GapManager** which is the upper layer providing high level services based on the OAIS model aiming at offering an abstraction useful for preservation information systems.



Available Data and Current Deployments

Available data

ontologies and data

- Intelligibility-related
 - Core Ontology for the Gap Manager
 - Instantiations of the ontology
 - exported from the PRONOM registry
 - from various testbeds
 - from the Registry
- Descriptive metadata
 - CIDOC CRM ontology
 - Specializations of CIDOC CRM (e.g. FRBRoo, testbed-related)
 - Instantiations of the above
 - contemporary arts testbed (Distance Liquide, Avis ...), cultural testbed

Available Data and Current Deployments

<http://mariateresa.isti.cnr.it:8093/CasparGui/>
http://139.91.183.8:3027/GapManagerGWT_SWKM/

GapManager

In brief, Gap Manager can aid the following tasks

- the decision of what metadata need to be captured and stored
- the identification of the data objects that are in danger in case a module (e.g. a software component or a format) is becoming obsolete (or has been vanished),
- the reduction of the metadata that have to be archived, or delivered to the users (as response to queries) on the basis of the designated community.

(for more details see D2102, the slides <http://wiki.casparpreserves.eu/pub/Main/PreparationForMonth24EuRevi> and the video tutorial)

GWT-based prototype

- GapManager is implemented as software component offering a plethora of methods for managing dependencies and profiles
- The implementation of an indicate prototype application using GWT is ongoing

GapManager and Data Holders

**GWT-based prototype
Ingestion of Modules and specification of their Dependencies**

Module Profile: Intelligibility Services

Insert Module

Identifier:

Name:

Version:

Type:

Module Software Structural Binary

Dependency Types: Modules

Search Remove

Dependencies

Typology of dependency types

Typology of modules

31

**GWT-based prototype
List All Profiles**

Module Profile: Intelligibility Services

Search Results

Composer/Performer Profile <http://www.icsim.org.uk/profile#composerPerformer> PDF

Multimedia User Profile <http://athena.ics.forth.gr/profile#multimediaUser> CSD Direct MAX_MSP

Developer/Operator Profile <http://www.icsim.org.uk/profile#developerOperator> CFG MAF TRC

Multimedia User Profile <http://www.icsim.org.uk/profile#multimediaUser> MOV AIF

Astronomer Profile <http://athena.ics.forth.gr/profile#astronomer> FITS_Module

Casual User Profile <http://athena.ics.forth.gr/profile#casualUser> PDF_Standard_HTML_Specification

Known modules by a DC profile

32

**GWT-based prototype
Intelligibility-aware Packaging (DIPs)**

Module Profile: Intelligibility Services

Modules

Profiles

Selection of module

Selection of DC profile

Computation of gap

33

GWT-based prototype

Module Profile: Intelligibility Services

MODULE MAX_MSP

Edit

Identifier: <http://athena.ics.forth.gr/>

Name: MAX_MSP

Version: 1.0

Types: Modules Software/Binary Software

Direct Dependencies: 23

Get Closure

Direct Dependents: [Multimedia.zip](#)

Allows the curator to see what objects will be in danger if a module is no longer available

34

GapManager and Data Holders

Methodologies that could be (directly) adopted

[SIMPLE] Very simple and easy

- 1/ Each information object is recorded to gap manager
 - only its id and its name is required
- 2/ The dependencies of these objects are specified
 - dependencies may point to modules that are already recorded (more than 2200 modules already exist)
- 3/ Appropriate DC profiles are defined
 - Just selecting the modules that are assumed to be known

[Advanced]

- 0/ Identify Tasks => Dep types => Specialization of the Core Ontology
- Then continue steps 1/ .. 3/ of the [Simple] methodology

35

**Cultural, Artistic and Scientific knowledge
for Preservation, Access and Retrieval**

thanks for your attention

36

CASPAR Cultural, Artistic and Scientific knowledge for Preservation, Access and Retrieval

backup slides

37

CASPAR Cultural, Artistic and Scientific knowledge for Preservation, Access and Retrieval

CASPAR Research on Knowledge Management

Outline

Research Results of the 1st year (in brief)

Research Results of the 2nd year: on Knowledge Evolution

- On Comparing RDF Knowledge Bases
- On Ontology Evolution
- On Archiving RDF Knowledge Bases

38

CASPAR Component vs State of the Art FORTH-ICS

- Gap Manager
 - Enables Intelligibility-aware Services. Support of DC profiles, extensible typologies (of modules and dependency types), computation of closures and gaps, import/export in various formats
 - More rich and expressive functionality comparing to other registries
- SWKM
 - Complete and scalable suite of basic services for validating, storing, querying, updating and exporting descriptive metadata expressed in RDF/S.
 - All services are based on a common knowledge repository enabling the consistency of its contents.
 - Advanced evolution services (declarative update languages, comparison services, versioning, ontology evolution)

39

CASPAR 1st year Results FORTH-ICS

Contributions

- A formalization of intelligibility and intelligibility gap through the notion of dependency. This view could be used for providing answers to some basic questions (who much RI we need, how this depends on DC, etc)

- Y. Tzitzikas, "Dependency Management for the Preservation of Digital Information", 18th International Conference on Database and Expert Systems Applications, DEXA'2007, Regensburg, Germany, September 2007
- Y. Tzitzikas and G. Flouris, "Mind the (Intelligibility) Gap", 11th European Conference on Research and Advanced Technology for Digital Libraries, ECDL'2007, Budapest, Hungary, September 2007
- Yannis Tzitzikas, Dependency Management for the Preservation of Digital Information, International Conference PV'2007 (Ensuring the Long-Term Preservation and Value Adding to Scientific and Technical Data), Oberpfaffenhofen/Munich, Germany, October 2007

The implementation activities of the 2nd year have exploited these results

40

CASPAR Motivation for Research on Knowledge Evolution and 2nd year Research Results FORTH-ICS

Motivation

- Evolution is a key concept because preservation is a dynamic process (world evolves, software/hardware evolves, metadata schemas evolve, digital objects evolve, community knowledge evolves).
- Digital objects have to be preserved not only against hardware and software technology changes, but also against changes in its Designated Community

Requirements

- Bulk metadata updates, Versioning metadata and ontologies, Ontology Evolution, Change log management, Comparison operators, etc

Research Results

- On Comparing Knowledge Bases
- On Ontology Evolution
- On Archiving

41

CASPAR CASPAR RESEARCH RESULTS Comparing RDF Knowledge Bases FORTH-ICS

- D. Zeginis, Y. Tzitzikas, V. Christophides, "On The Foundations of Computing Deltas between RDF models", ISWC'2007 (Nov 2007) (BEST PAPER AWARD!)

Topic: What is the delta between two knowledge bases?
Can we reduce the number of change operations needed to be applied on K so that to reach K'?

K

K'

Delta(K->K') = ?

42

CASPAR RESEARCH RESULTS
Comparing RDF Knowledge Bases

- Why it is useful in preservation
 - Consider the case where PDI is represented as an RDF/XML file (that instantiates one or more ontologies e.g. CIDOC CRM)
 - Suppose a new version of that file arrives
 - Issue: How the SWKM repository (that is used by Finding Aids) is updated?
 - Solution: Diff(file1, file2) gives a set of change operations that can be forwarded to SWKM to update it
- Research/Development Status
 - A prototype implementation over SWKM (as a Web Service) is ongoing.
 - Also a main memory implementation is already available

CASPAR RESEARCH RESULTS
Ontology Evolution

G. Konstantinides and G. Flouris, G. Antoniou and V. Christophides, "A Formal Approach for RDF/S Ontology Evolution", ECAI'2008 (accepted for publication)

Topic: How to reflect ontology changes to the underlying data?

CASPAR RESEARCH RESULTS
Archiving Versions> Synopsis

Y Tzitzikas, Y. Theoharis, D. Andreou, "On Storage Policies for Semantic Web Repositories that Support Versioning", European Semantic Web Conference 2008, ESWC'2008 (June, 2008)

- This is the first work that focuses on the storage aspect of SW repositories that support versioning.
- We proposed an index called **POI (Partial Order Index)**, we verified the space gains of this index experimentally and we provided an efficient version insertion algorithm with acceptable main memory space requirements.
- From our experiments, POI can be 180 times more space economical compared to **IC (Individual Copies-approach)** and 18 times compared to **CB (Change-based approach)** for parallel version tracks. Moreover, POI allows performing efficiently various cross-version operations.
- POI is an advantageous approach for archiving set-based data (e.g. an RDF KB is a set of triples), especially good for inclusion-related and "oscillating" data

Cultural, Artistic and Scientific knowledge
 for Preservation, Access and Retrieval

thanks for your attention

Cultural, Artistic and Scientific knowledge
 for Preservation, Access and Retrieval

more backup slides

CASPAR RESEARCH RESULTS
Exploiting DC profiles for having intelligibility aware AIPs and DIPs

According to OAIS:

- AIP (Archival Information Package): It is actually a format which consist of
 - Data Object
 - Rep Info
 - PDI (Preservation Description Information): e.g. provenance, context, fixity
- DIP (Dissemination Information Package): Is the version of the information package delivered to the Consumer in response to an access request. May differ in form (e.g. TIFF to JPEG) or content (e.g. amount of metadata supplied) to that which resides in the archival store.

DC profiles could be exploited so that

- to be able to derive different AIPs and DIPs for different DCCommunities
 - If the dependencies are available this could be done automatically

Exploiting DC Profiles for constructing the "right" DIPs (intelligible and redundancy free)

DIP of o1 wrt DC1

Object = o1
deps = {t1,t3}

DIP of o1 wrt DC2

Object = o1
deps = {t1,t2,t4}

DIP of o1 wrt DC3

Object = o1
deps = {t1,t2,t3,t4,t5,t6}

49

Repackaging

KMGr (Gap Manager)

Preservation Data Store

Object = o5
Object = o4
Object = o3
DC = DC1
Object = o2
Object = o1
DCprofile = DC1
deps = {t1,t3}

And of course, the availability of dependencies and DC profiles could aid in transforming the stored AIPs to DIPs of different profiles (than those of the AIPs)

50

When DC profiles evolve

AIP of o1 wrt DC2

Object = o1
DCprofile = DC2
deps = {t1,t2,t4}

AIP of o1 wrt DC2'

Object = o1
DCprofile = DC2
deps = {t1}

51

Intelligibility-related dependencies and provenance

Naive Modeling

CIDOC CRM modeling

52