# Long-term Information Preservation

CCSDS Contributions

# Introduction

- **We are the CCSDS Data Archive Ingest (DAI) WG**
  - Consolidates former DAI WG, Information Packaging and Registry (IPR) WG and Repository Audit and Certification (RAC) WG

- **What we've done**
  - Working within CCSDS, but is more broadly accepted than space industry
  - OAIS Reference Model
    - "now adopted as the de facto standard for building digital archives"
      Cyberinfrastructure Vision for 21st Century Discovery (http://www.nsf.gov/pubs/2007/nsf0728/nsf0728.pdf)
- **What we're doing**
  - Certification - ISO
  - Guidance on information to capture to aid preservation
- **What we plan to do**
  - Continued development of digital archiving standards
  - We'd like you to be an advocate for our work.
  - We'd appreciate inputs from you on suggested directions.

- **Details**      What we're doing and why we approach digital preservation the way we do

# Reference Model for an Open Archival Information System (OAIS)

- OAIS reference model provides: fundamental concepts for preservation
  - fundamental definitions so people can speak without confusion
    - Provides vocabulary – widely applicable
    - defines important **roles** in digital preservation
  - Conformance to OAIS Standard defined
  - OAIS approach to digital preservation:
    - covers all types of digitally encoded information
    - provides a way to test whether preservation is successful
    - does not require seeing into the future
    - does require transparency
    - but does not require "open access"
  - does not cover
    - social and organizational aspects
    - finance etc.
    - many detailed aspects noted in the "roadmap of follow-on standards"
- is not meant as a **design/blueprint** for a repository
  - it allows greater flexibility for implementations
- OAIS does provide a good basis for certification

# Challenges

- Wanted it to be applicable to all kinds of
  - Digital objects
  - Organizations
- Hardware
- Software
- Sources
- Users
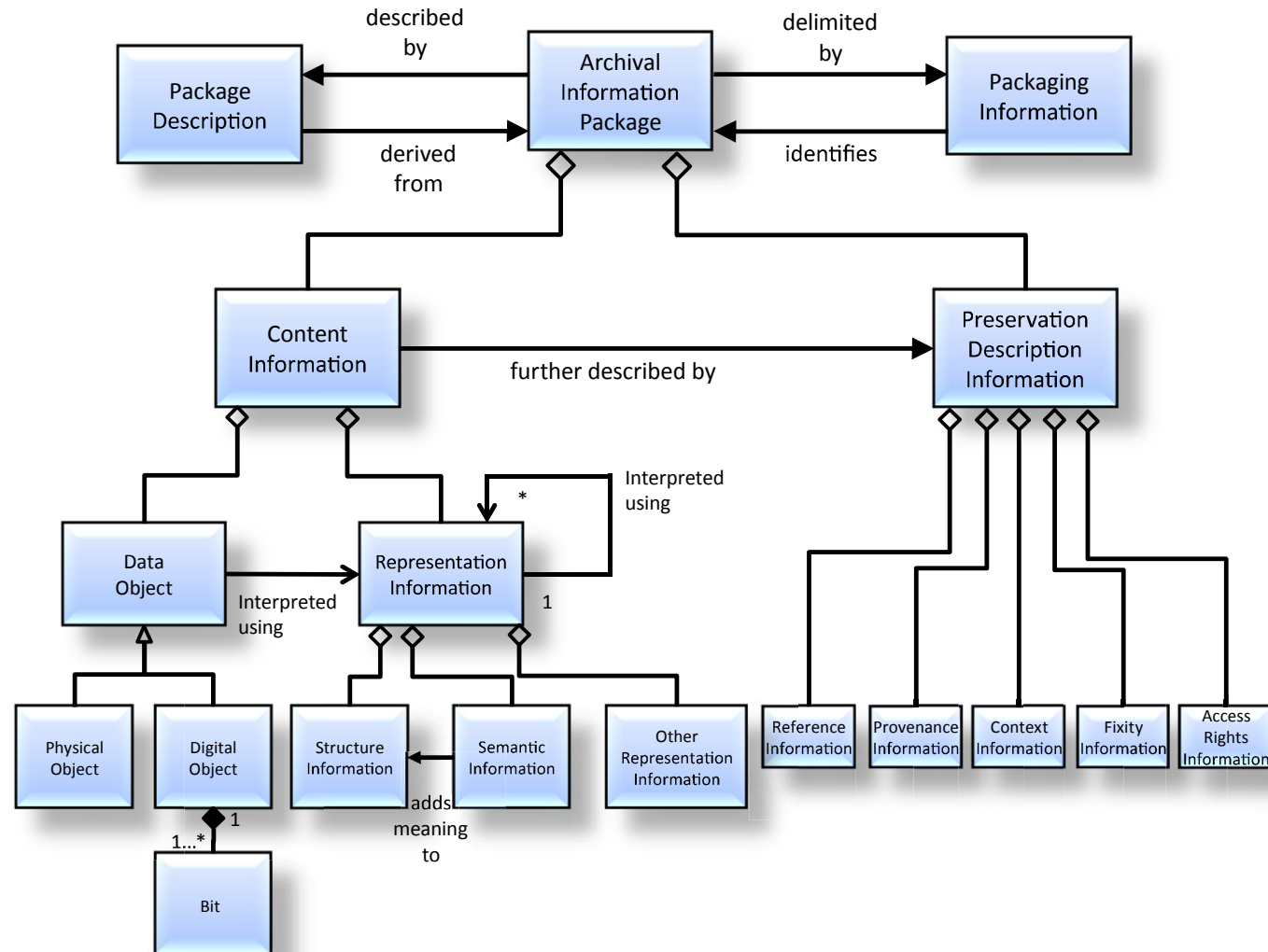- Legal systems
- All kinds of changes

# Key OAIS Concepts

- Claiming "This is being preserved" is untestable
  - Essentially meaningless
    - Except "BIT PRESERVATION"
- How can we make it testable?
  - Claim to be able to continue to "do something" with it
    - <span style="color:red">Understand/use</span>
      - Need extra information to help this (Representation Information )
- Still meaningless…
  - Things are too interrelated
    - Representation Information potentially unlimited
  - Need to define a target group (Designated Community) – those we guarantee can understand – so we can test
- Many other concepts identified
- Finer grained taxonomy than simply saying "metadata"
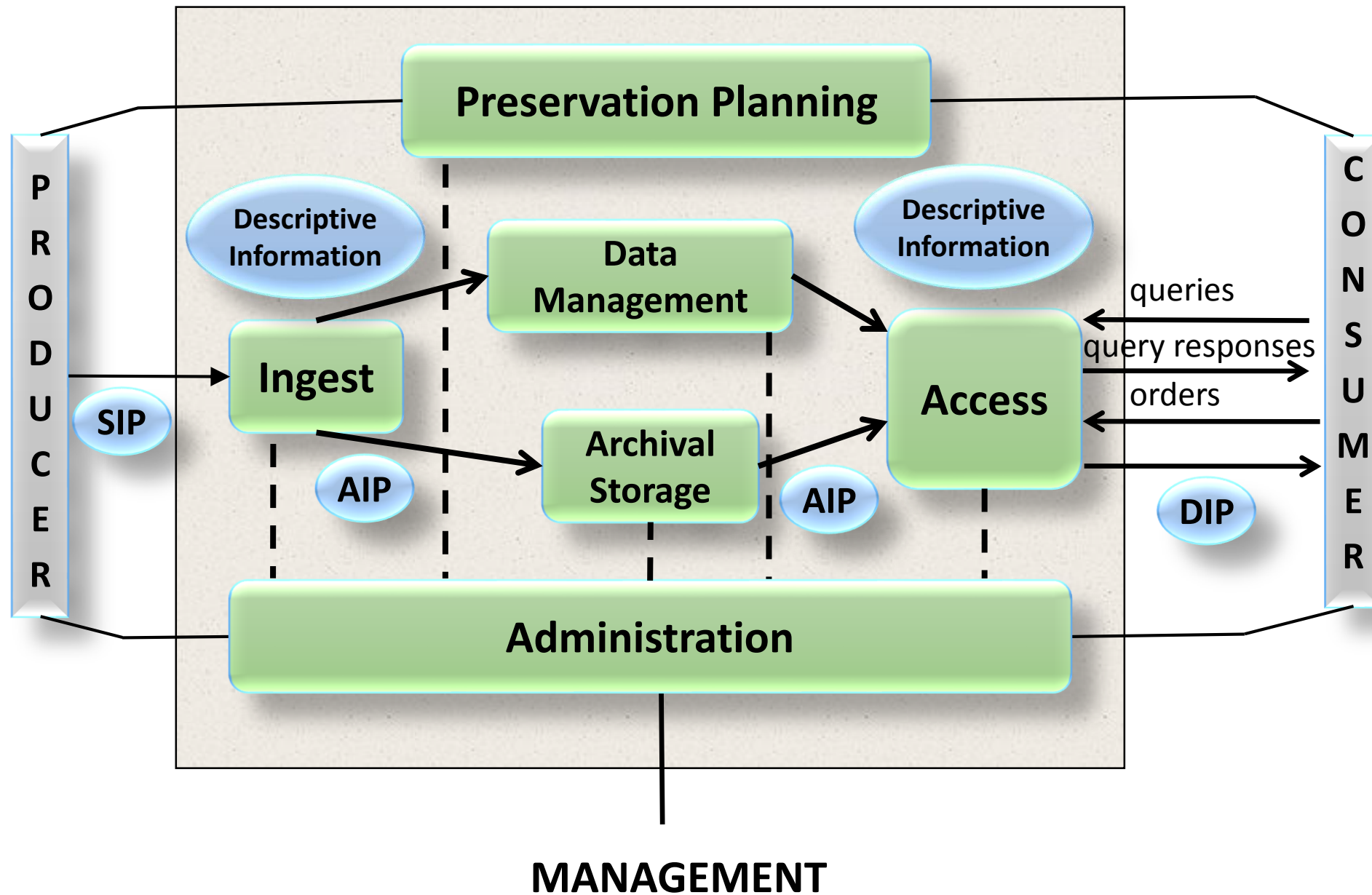  - Allows one to ask if one has all the required types

# OAIS Standard Overview

- Provides information model
- Provides mandatory requirements
- Provides framework for additional standards

- Follow-ons
  - Auditing  and Certification Standards
  - Detailed Processes
  - Some Protocols

- Builds on previous data description standards

# Archival Information Package



- Logical package e.g. using pointers
- Contains everything needed for long term preservation
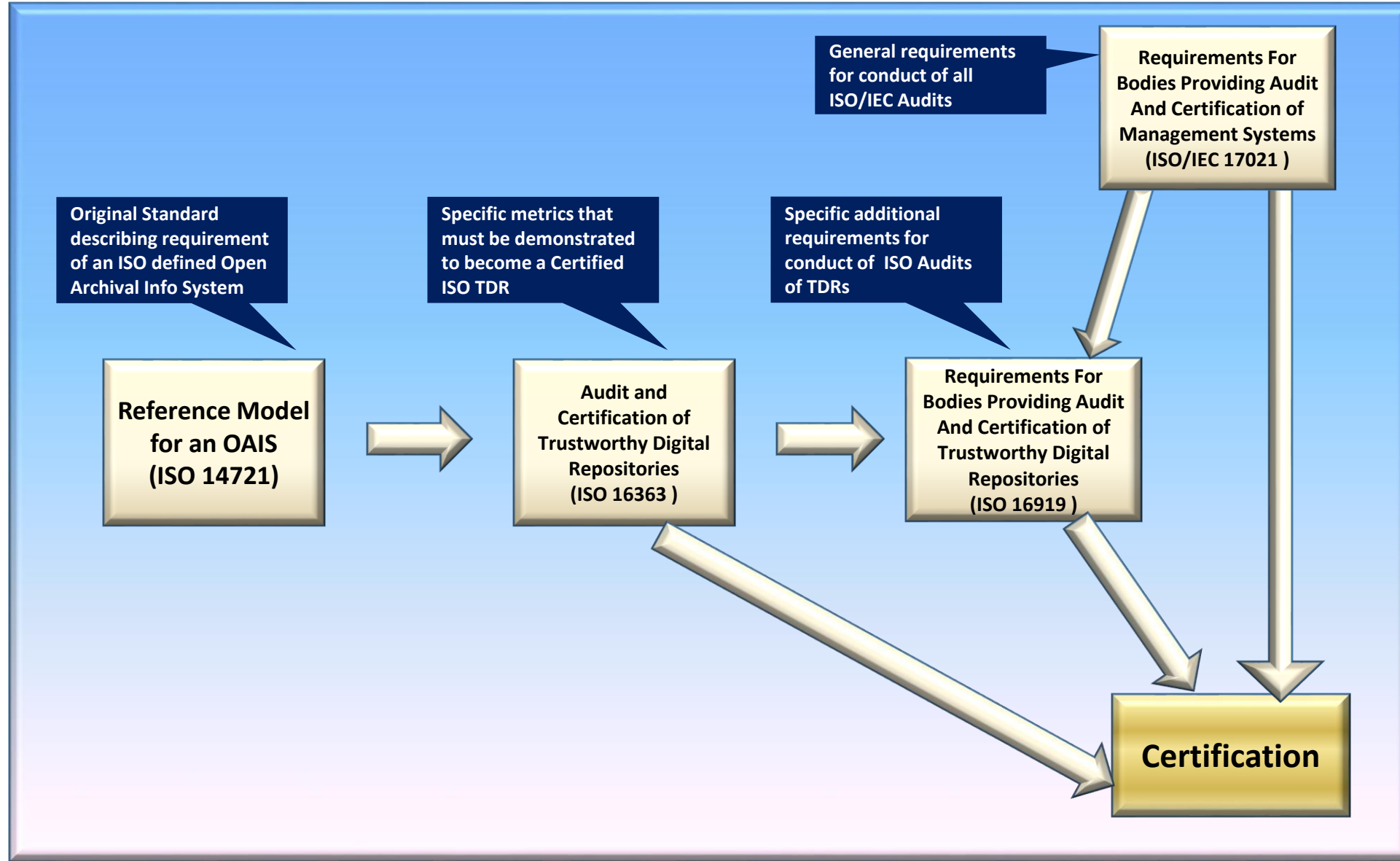  - Usability
  - Authenticity

# Menu

- What we've done – the standards
- Relationship to previous work – origins of digital preservation
- Examples of archives using OAIS
- Data – more details about the issues
- Preservation techniques
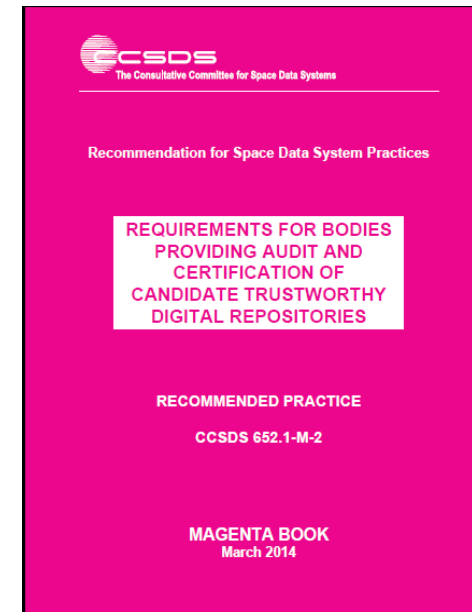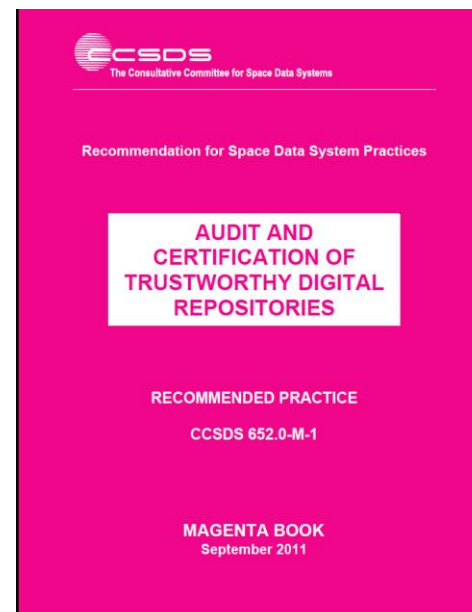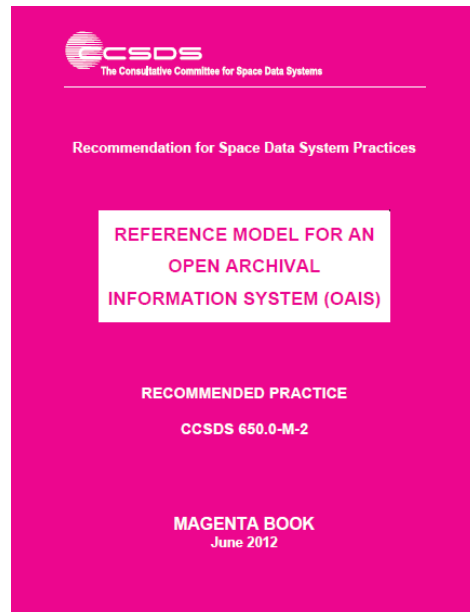- New standards in development
- Advocating

# What we've done – the standards

# Standards based Repository Audit and Certification (ISO 16363 and ISO 16919)



General requirements for conduct of all ISO/IEC Audits

Requirements For Bodies Providing Audit And Certification of Management Systems (ISO/IEC 17021 )

Original Standard describing requirement of an ISO defined Open Archival Info System

Specific metrics that must be demonstrated to become a Certified ISO TDR

Specific additional requirements for conduct of ISO Audits of TDRs

Reference Model for an OAIS (ISO 14721)

Audit and Certification of Trustworthy Digital Repositories (ISO 16363 )

Requirements For Bodies Providing Audit And Certification of Trustworthy Digital Repositories (ISO 16919 )
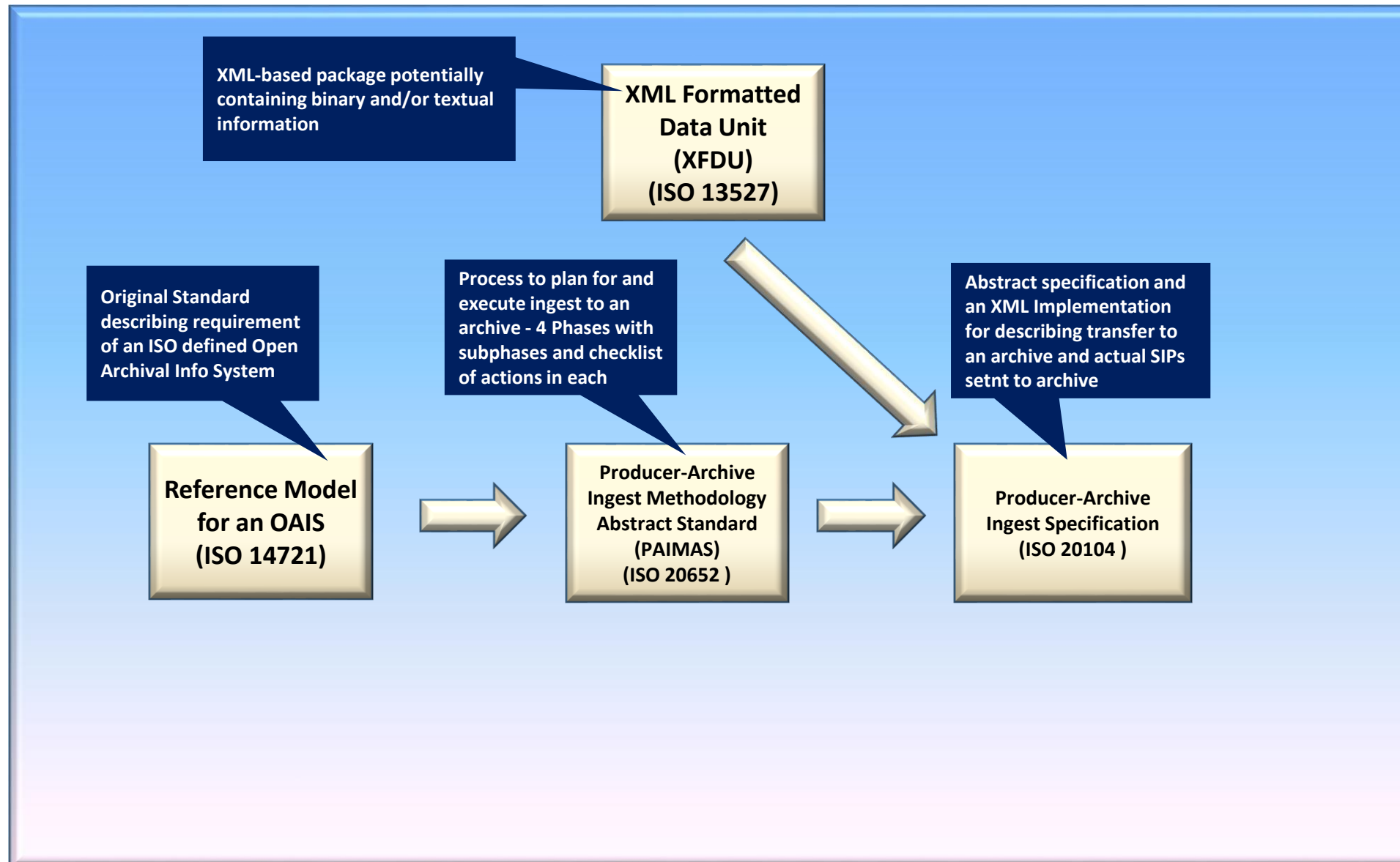
Certification

11

# Get Your Standards

- You can get all ISO Standards from ISO website at: http://www.iso.org/home/store/catalogue_tc/

- You can also download the CCSDS Magenta Book equivalents of the ISO docs for free
  - ISO 14721 (OAIS) equivalent here: http://public.ccsds.org/publications/archive/650x0m2.pdf
  - ISO 16363 (RAC metrics) equivalent here: http://public.ccsds.org/publications/archive/652x0m1.pdf
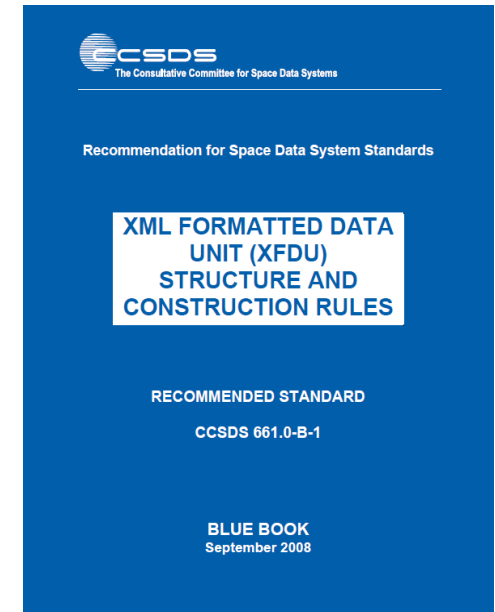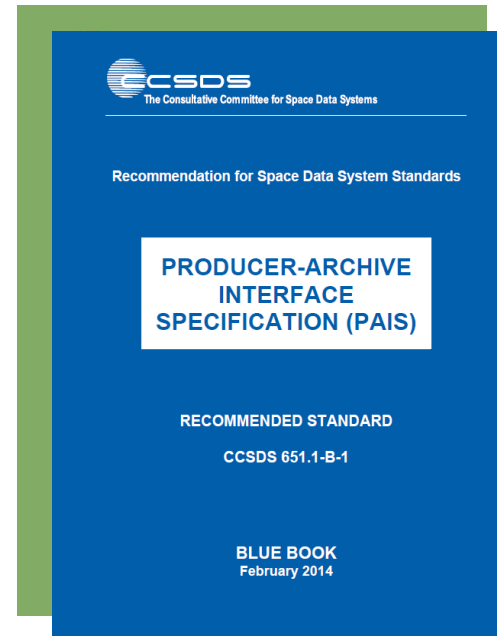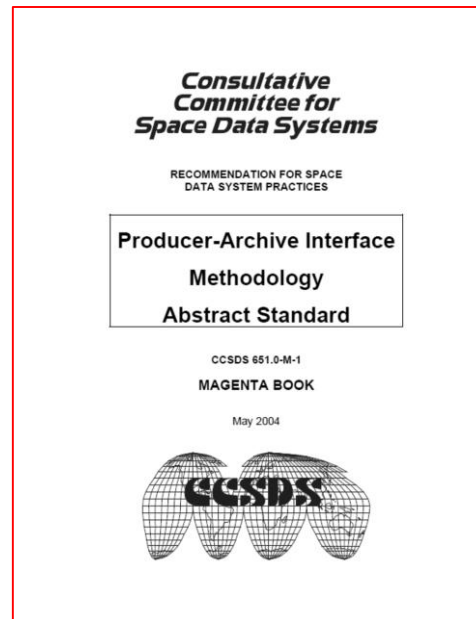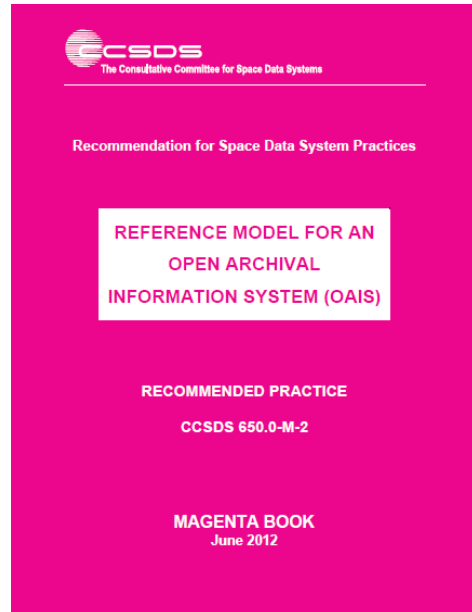  - ISO 16919 (RAC process) equivalent here: http://public.ccsds.org/publications/archive/652x1m2.pdf

# Archive Process and Protocol Standards (ISO 20652 and ISO 20104)



XML-based package potentially containing binary and/or textual information

**XML Formatted Data Unit (XFDU) (ISO 13527)**

Original Standard describing requirement of an ISO defined Open Archival Info System

Process to plan for and execute ingest to an archive - 4 Phases with subphases and checklist of actions in each

Abstract specification and an XML Implementation for describing transfer to an archive and actual SIPs setnt to archive

**Reference Model for an OAIS (ISO 14721)**

**Producer-Archive Ingest Methodology Abstract Standard (PAIMAS) (ISO 20652 )**

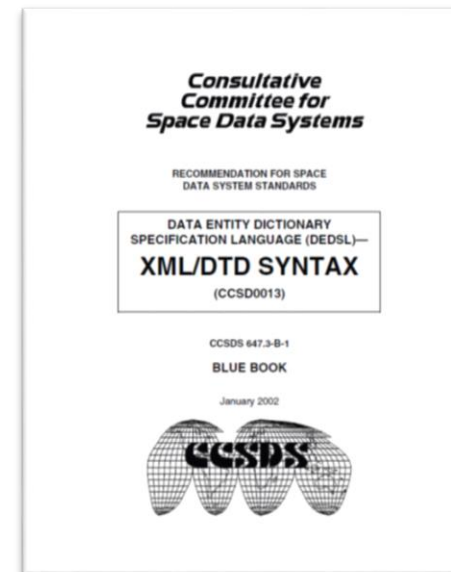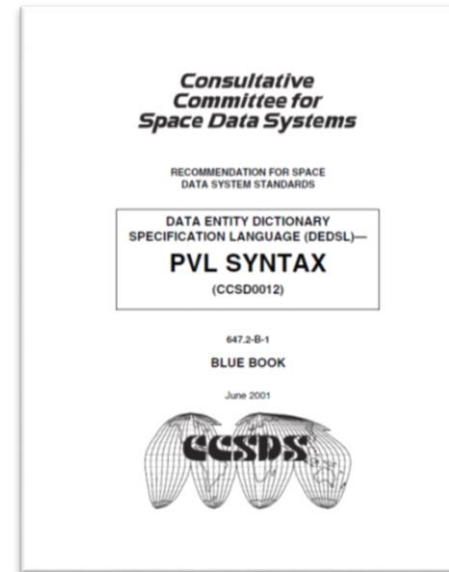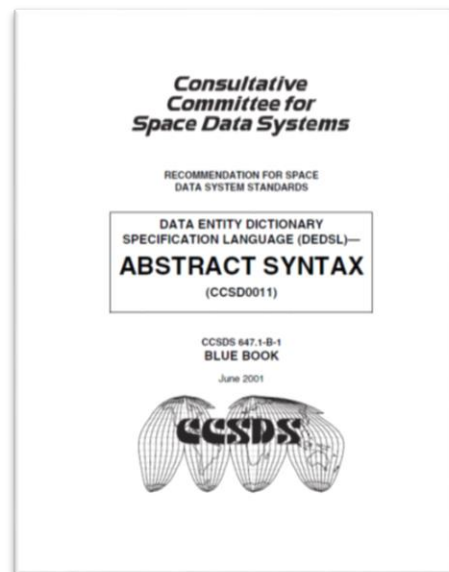**Producer-Archive Ingest Specification (ISO 20104 )**

# Get Your Standards

- You can get all ISO Standards from ISO website at: http://www.iso.org/home/store/catalogue_tc/

- You can also download the CCSDS Book equivalents of the ISO docs for free
  - ISO 14721 (OAIS) equivalent here:       http://public.ccsds.org/publications/archive/650x0m2.pdf
  - ISO 20652 (PAIMAS) equivalent here:   http://public.ccsds.org/publications/archive/651x0m1.pdf
  - ISO 20104 (PAIS) equivalent here:       http://public.ccsds.org/publications/archive/651x1b1.pdf
  - ISO 13527 (XFDU) equivalent here:       http://public.ccsds.org/publications/archive/661x1b1.pdf

# Other CCSDS DAI WG Standards

- EAST  (Data Format Structure Description Language)

- Data Entity Dictionary Specification Language (DEDSL)
  - Abstract definition
  - Parameter Value Language (PVL), XML/DTD, and XML Schema implementation

# Relationship to previous work – origins of digital preservation

# Who Begat Whom? Certification Standards

- Bruce Ambacher

# Examples of archives using OAIS

# Influence of OAIS-RM and Certification Audit on CIESIN Archive

- Bob Downs

# Influence of the OAIS-RM on PDS4

- The Planetary Data System (PDS4) is NASA's official archive for all Solar System Exploration science data.

- PDS4 is a major revision and transition to a modern system based on best practices for data system development while leveraging 20 years of lessons learned.
  - PDS4 has been operational since 2013.
  - PDS4 has been adopted by the international Planetary Science Community and considered to be a first-of-its-kind where a single system and set of standards has been adopted by an entire science community.

- The OAIS-RM heavily influenced the architecture and design of PDS4.
  - PDS4 complies with the *Mandatory Responsibilities* required to operate an OAIS archive.
  - The PDS4 system and services architecture maps to the *OAIS Functional Model*.
  - The PDS4 Information Model is based on OAIS-RM concepts.
    - The *Information Object*, consisting of a *Data Object* and its *Representation Information*, is a foundational concept.
    - *Preservation Description Information* (PDI) is required for all data.

- The PDS4 functional infrastructure is designed to respond to the OAIS-RM based PDS4 Information Model.
  - The model is independent from the implementation technology allowing both to evolve separately.
  - Model extracts are used to design and configure system software and services.

- The PDS is well positioned to be accredited against the ISO-16363 standards.

# Data – more details about the issues
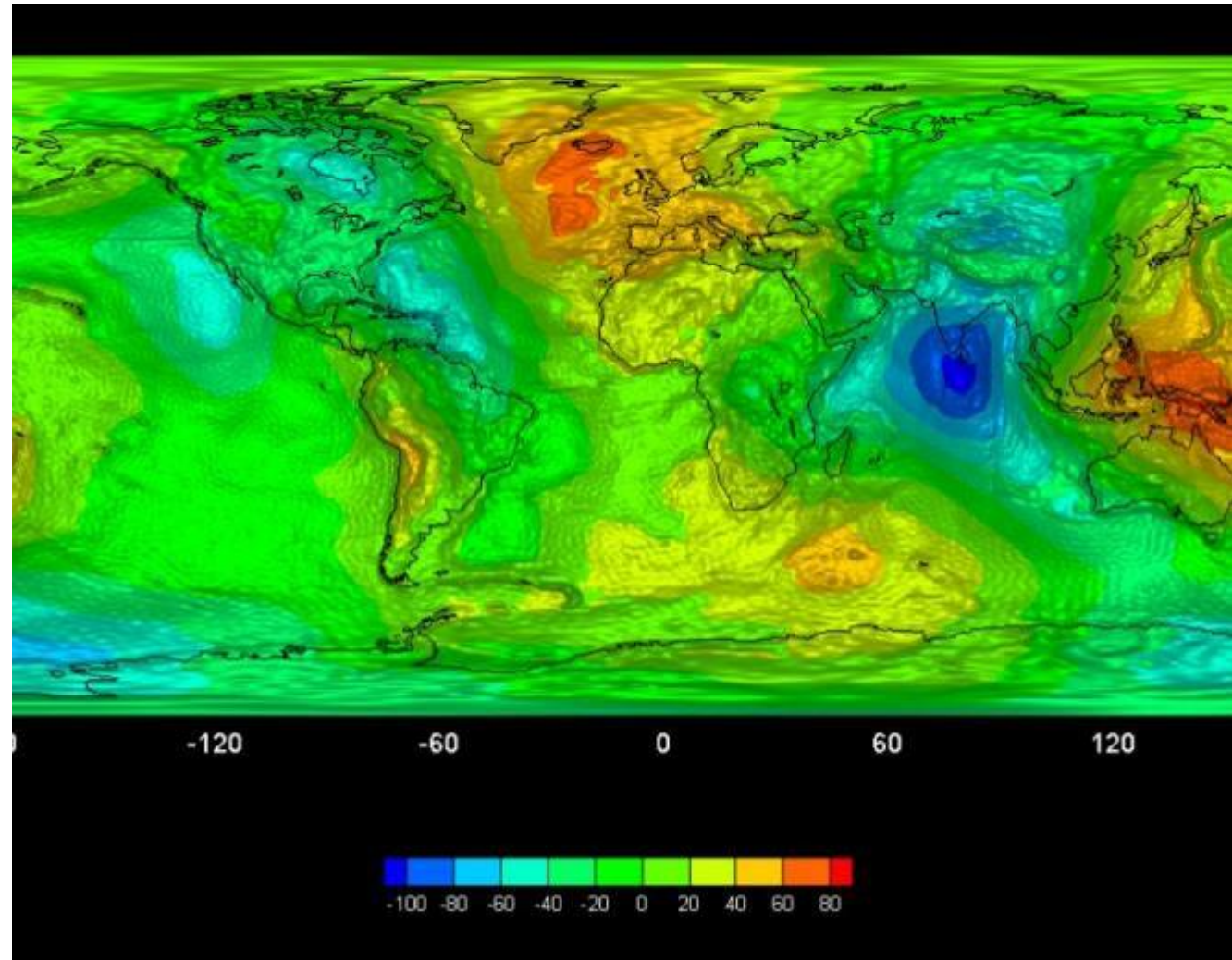
# Data contains numbers, etc.

May be digitally encoded
    But a lot of additional information is needed to understand it
        Layout, Bit Format, Structure
        Units, Field Names, Meaning,
        Semantics
        Etc.

To be combined and
processed to get this

# Preservation techniques

# Options

- **EITHER**
  - Hand on the information to some other custodian
- **OR**
  - Keep the information          – in which case one can
    - **EITHER**
      - Keep the bits unchanged
        - In which case one must add what OAIS calls "Representation Information" to ensure the information encoded in the bits can still be understood/used (by the Designated Community)
    - **OR**
      - Change the bits
        - In which case one must keep evidence about why this new object should be considered "authentic"
        - And enable the new bit sequences to be understandable/usable (add appropriate "Representation Information")

# Information



Data Object → (Interpreted using its) → Representation Information → (Yields) → Information Object

# Digital Preservation Methods



Figure 6.3 Digital Preservation Methods (From Thibodeau, 2002, p.19)
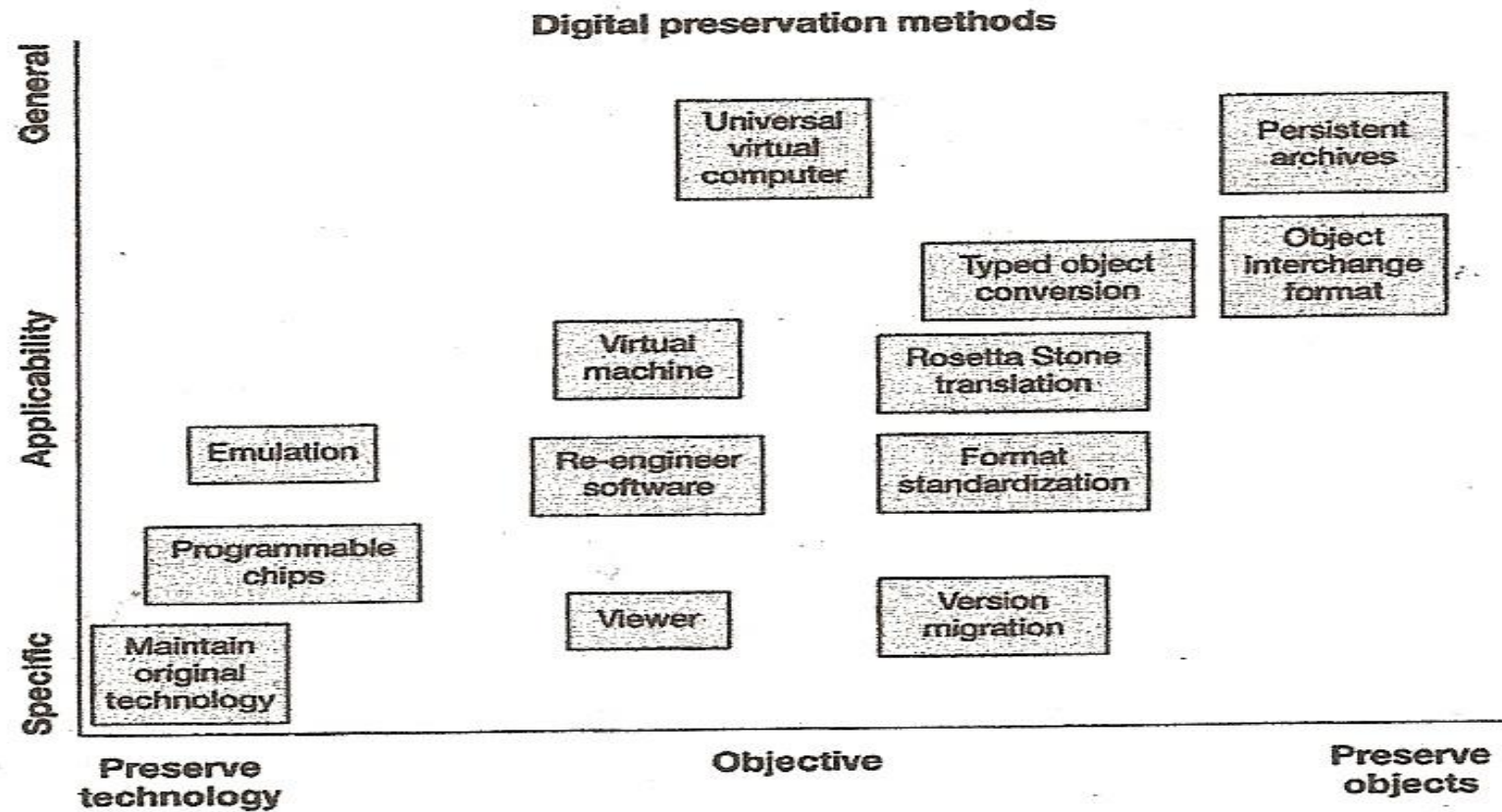
# Representation Information examples

- Emulation support

- Documentation

- Data Dictionaries

- Data descriptions

# Emulation – advantages and drawbacks

- Advantage
  - Can do what was do-able previously e.g.
    - display document in an format which is no longer supported by using emulation to enable use of  old software. Might need to emulate chip instruction sets, operating systems etc.
    - Analyse old data using old analysis software e.g. CERN/LHC analysis software to repeat analyses of data, perhaps tweaking parameters etc
- Disadvantage
  - Can ONLY do what was do-able previously
  - BUT
    - Difficult to use OLD data in **NEW** analysis software
    - Difficult to combine OLD data with NEW data

# Formal Descriptions of Structure and Semantics

- CNES EAST tools (http://east.cnes.f), OASIS, EAST C Library (reference implementation).

- Also DEBAT (BEST Tools) http://debat.c-s.fr/

- Data Request Broker (DRB) - http://www.gael.fr/drb/site/

- JNI Wrapper for EAST C Library (jnieast).

- DEDSL Abstract, PVL, and XML(DTD) syntax for defining some simple data semantics. RDF, RDFS and OWL.

- Interfaces for a more general data description language and semantics API (DSSIL).

- GUI Tools for capturing Object Oriented Semantics (RDF and RDFS) and Code Generation.

- DFDL  (Data Format Description Language from Open Grid Forum) https://www.ogf.org/ogf/doku.php/standards/dfdl/dfdl

# Advantages of Formal Descriptions of Structure

- Formal descriptions of structure gives a user the knowledge an ability to map data bits to data values in software.

- Gives a common language for syntactic information to describe many different data formats.

- One software API for many data formats:
  - may increase likelihood of data reuse

- Machine readable descriptions.

# Formal Descriptions of Structure (Examples)

# Virtualisation - building up data types...

- Building on simple virtualisation of tabular data into JTable based application
- Can combine data in different formats

# Advantages of Formal Descriptions of Semantics

- Semantic properties such as name, description and units add meaning to the data values and data objects.

- Object oriented view of the data combined with the structure and semantics virtualises the process of going from the bits to a usable data object in software.

- Combined with formal structure it encapsulates all the knowledge needed to go from the bits to a data object in software.

- Remove the need for data specific Access Software!

35

# DEDSL Example

# New standard in development

# Information Preparation to Enable Long Term Use



| Information Preparation to Enable Long Term Use | PAIMAS   PAIS | OAIS | AUDIT & CERTIFICATION |

Helps to ensure that the Additional Information needed is collected/created

Defines a mechanism to transfer the Primary Information and the Additional Information to the archive

Defines how the information should be preserved

Defines how to check that the information is being preserved

**Formulation**
- Propose activity.
- Obtain approval and resources for Data project

**Implementation**
- Preparation for activity
- Develop/ update systems (H/W, S/W, processes)

**Operation**
- Carry out activity to create/collect data

**Initial Exploitation**
- Extract initial value e.g. publications, commercial, social or economic
- Suggest other ways to exploit

Important for the development of Data Management Plans
- Working with RDA Active Data Management Plans Interest Group

# Brown Dog

- Presentation from Mark Conrad

# Advocating

# CCSDS

- Presentations in many fora
- Working with
  - National Accreditation Bodies e.g. ANAB
  - Certification Bodies
  - various industries e.g. aerospace
  - …..
- Need help

# Backup Slides

# Relationship between standards

**OAIS (ISO 14721)**

Trusted Digital Repositories: Attributes and Responsibilities

TRAC

**Requirements For Bodies Providing Audit And Certification (ISO 16919 )**

There is a hierarchy of ISO standards concerned with good auditing.
ISO 16919 is positioned within this hierarchy in order to ensure that these good practices can be applied to the evaluation of the trustworthiness of digital repositories using ISO 16363.
It covers principles needed to inspire confidence that third party certification of the management of the digital repository has been performed with impartiality, competence, responsibility, openness, confidentiality, and responsiveness to complaints

**Audit and Certification of Trustworthy Digital Repositories (ISO 16363 )**

Metrics concerning:
- **Organizational Infrastructure**
  - **e.g. The repository shall have a documented history of the changes to its operations, procedures, software, and hardware.**
- **Digital Object Management**
  - **e.g. The repository shall have access to necessary tools and resources to provide authoritative Representation Information for all of the digital objects it contains.**
- **Infrastructure and Security Risk Management**
  - **eg. The repository shall have procedures in place to evaluate when changes are needed to current software.**

Audit by external, **accredited**, auditors

**Certification**

IAF

TC    TC

WG

USA: ANSI    UK:    INDIA:
ANAB        UKAS   NABCB

Regional Group

Regional Group

National Accreditation Body (NAB)

ISO 17011(?)

NAB

NAB

National Accreditation Body appoints ASSESSORS to help decide whether or not to accredit

Certification Body

ISO 17021
(ISO) 16919

Repository

ISO 16363

45