

# Building the Archives of the Future: Self-Describing Records



Kenneth Thibodeau

Director, Electronic Records Archives Program

*National Archives and Records Administration*

July 18, 2001

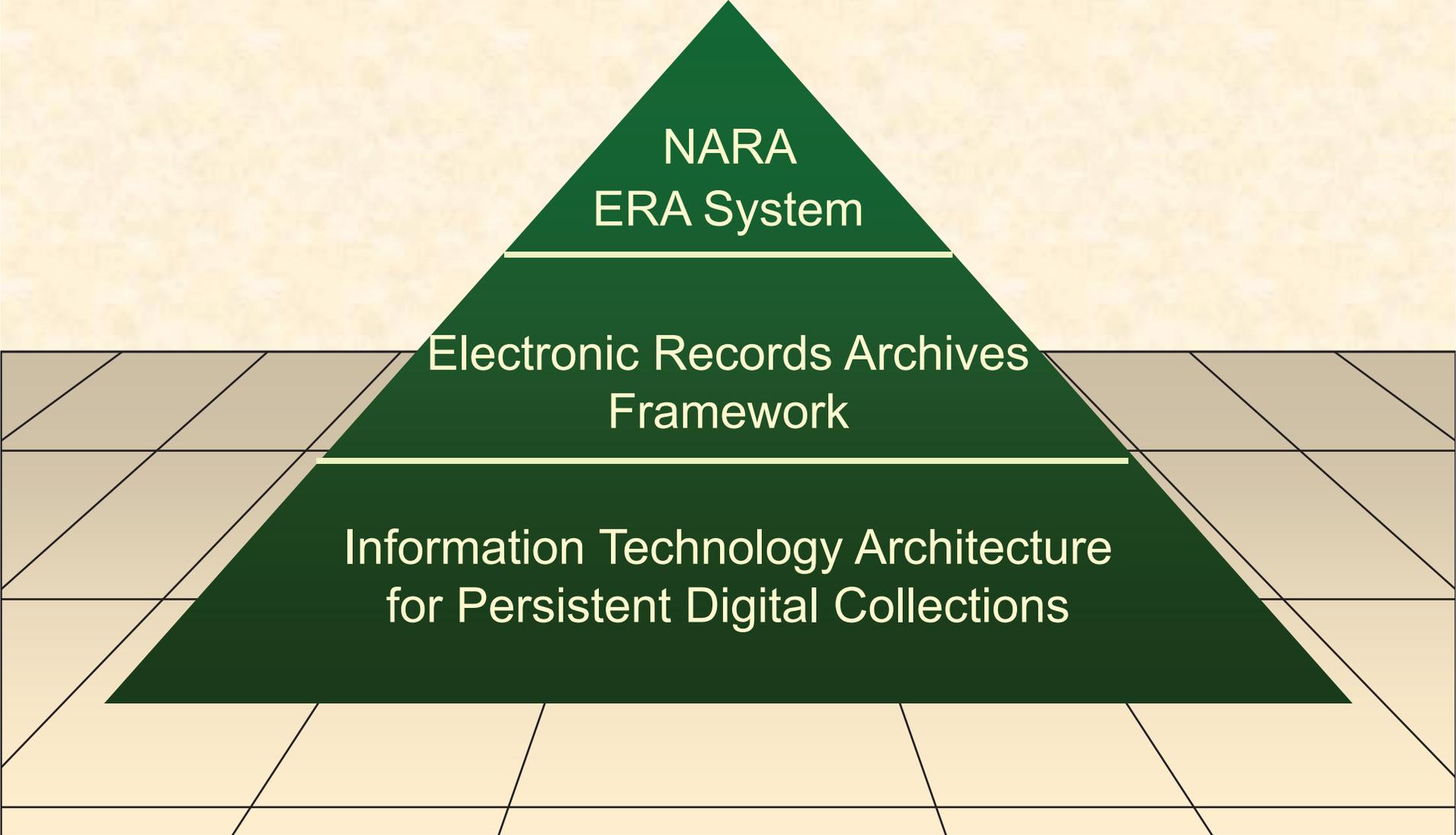
# The Electronic Records Archives Vision

- Overcome technological obsolescence in a way that preserves demonstrably authentic records.
- Build a dynamic solution that incorporates the expectation of continuing change in information technology and in the records it produces.
- Find ways to take advantage of continuing progress in information technology in order to maintain and improve both performance and customer service

# Critical Challenge

- Proven methods for preserving digital information across generations of technology are limited to the simplest formats
- Available methods are increasingly inadequate
- The market has not delivered solutions.

# How will we develop the Electronic Records Archives?



NARA  
ERA System

Electronic Records Archives  
Framework

Information Technology Architecture  
for Persistent Digital Collections

# ERA Infrastructure Concept

SCALABLE

Gb/sec Internet

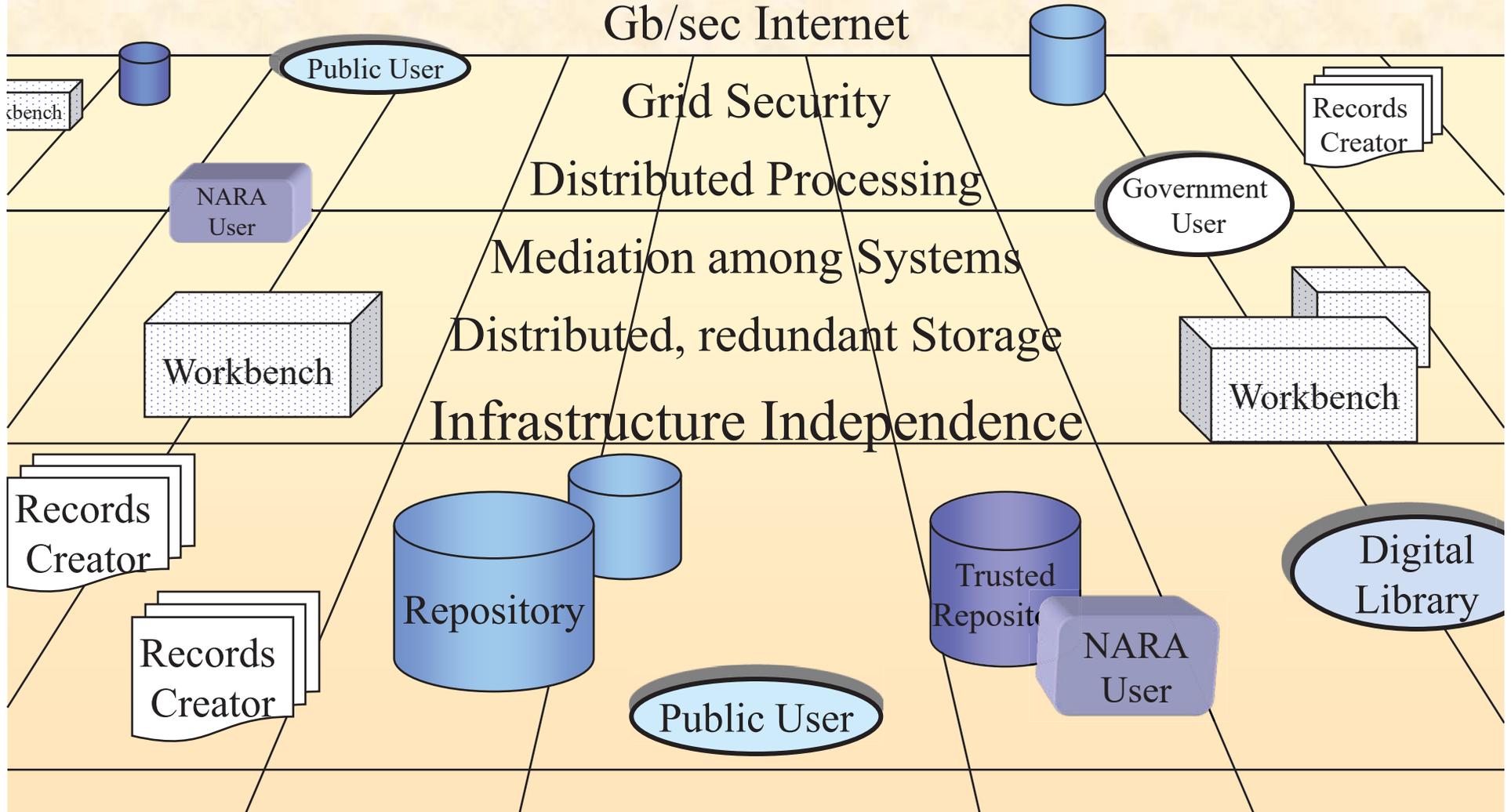
Grid Security

Distributed Processing

Mediation among Systems

Distributed, redundant Storage

Infrastructure Independence



# ERA Infrastructure

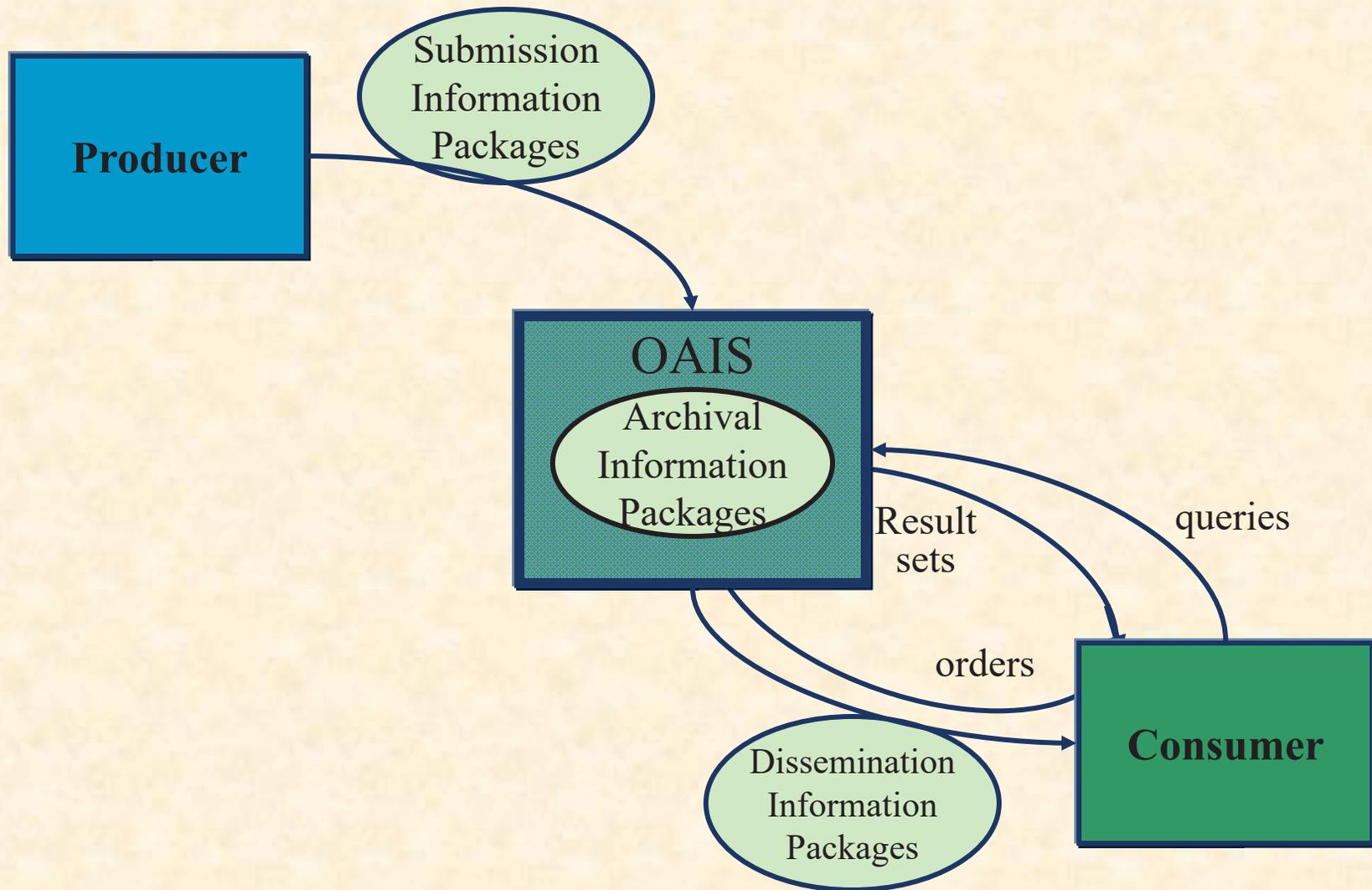
- In NARA (using NARANET)
  - **Archival workstations for staff**
  - **Reference workstations for researchers**
- On the National Information Infrastructure, under NARA's control
  - **ERA Ingest & Distribution portal (Internet & Media)**
  - **POP repositories (Normal, Trusted, Special)**
  - **Affiliated Archives**
- On the National Information Infrastructure
  - **Agency systems with access to NARA portal**
  - **Digital Libraries**
  - **Public Users**

# NARA Partnerships

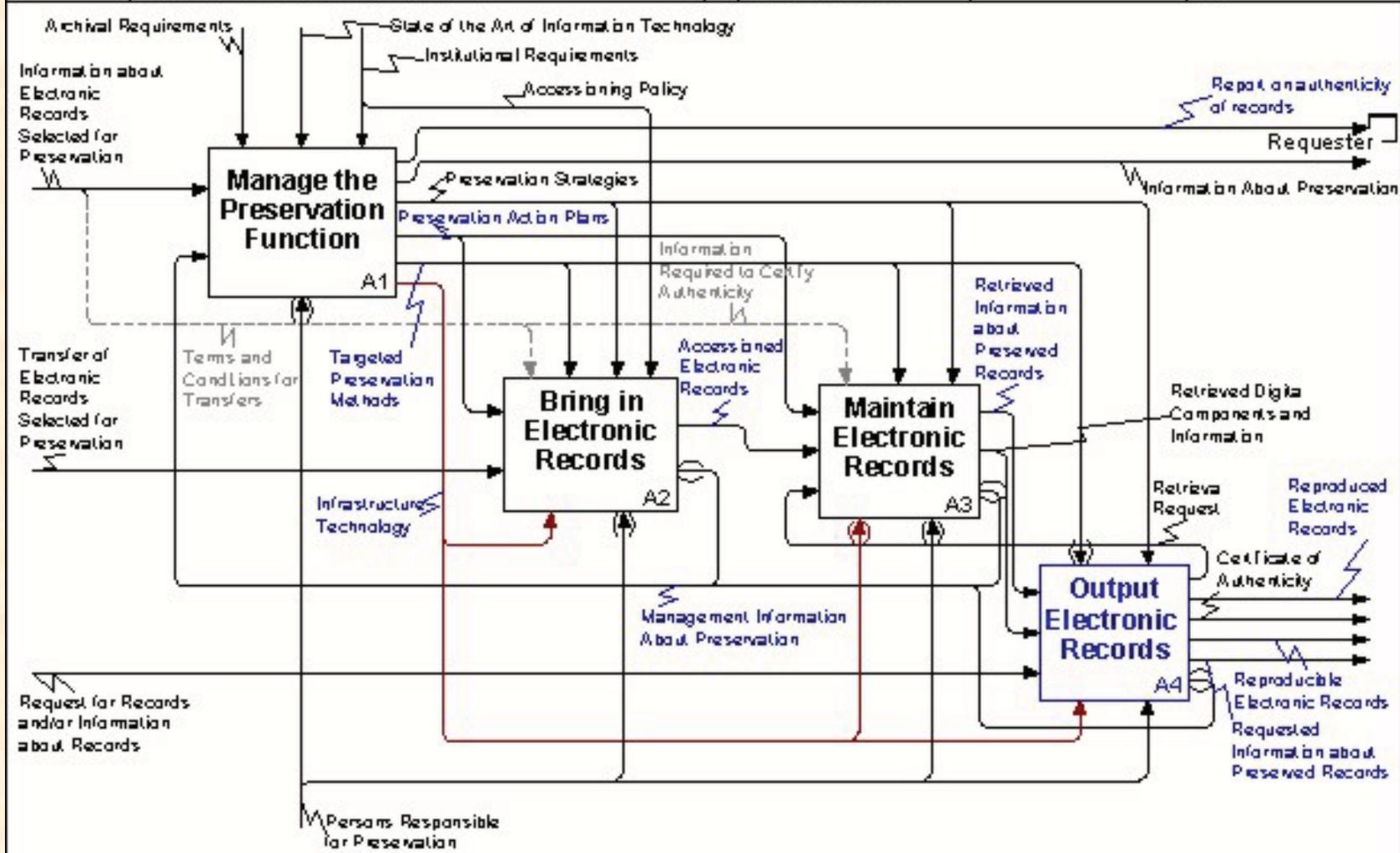
- **Open Archival Information System (OAIS) Reference Model**
  - NASA, Consultative Committee on Space Data Systems
- **Distributed Object Computation Testbed (DOCT)**
  - Defense Advanced Research Projects Agency, U.S. Patent and Trademark Office
- **National Partnership for Advanced Computational Infrastructure (NPACI)**
  - National Science Foundation
- **Presidential Electronic Records Processing Operational System (PERPOS)**
  - Army Research Laboratory, Georgia Tech Research Institute
- **Archivist's Workbench**
  - NHPRC Grant to San Diego Supercomputer Center
- **International research on Permanent Authentic Records in Electronic Systems (InterPARES)**
  - 7 international, multidisciplinary research teams, 10 national archives

# ERA Functional Model

## An Open Archival Information System Implementation

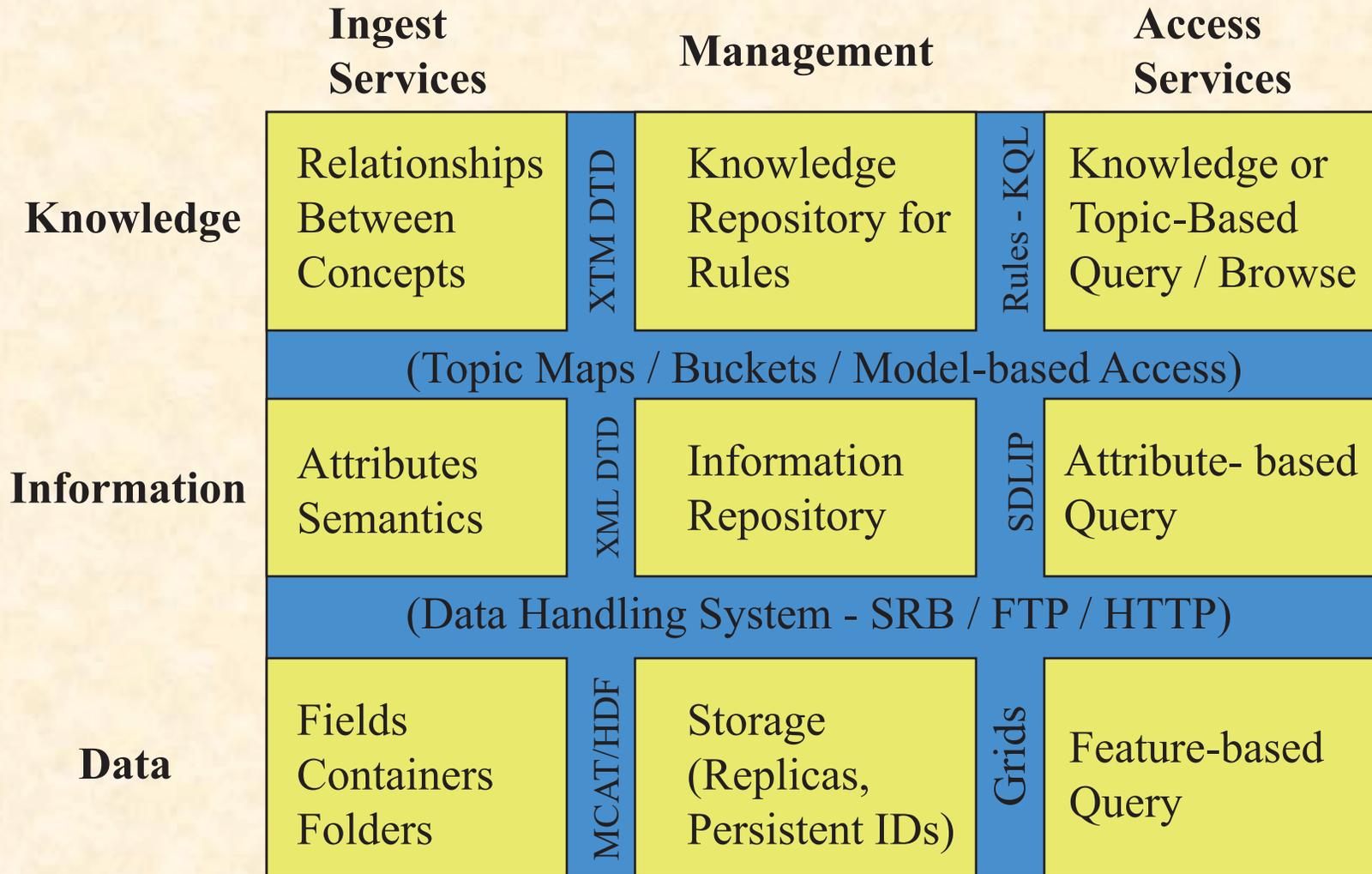


USED AT: Workshop 8	AUTHOR: Preservation Task Force	DATE: 2/17/2000	WORKING	READER	DATE	CONTEXT:
	PROJECT: InterPARES Project	REV: 7/6/2001	DRAFT			
			RECOMMENDED			
			PUBLICATION			A-0
NOTES: 1 2 3 4 5 6 7 8 9 10						

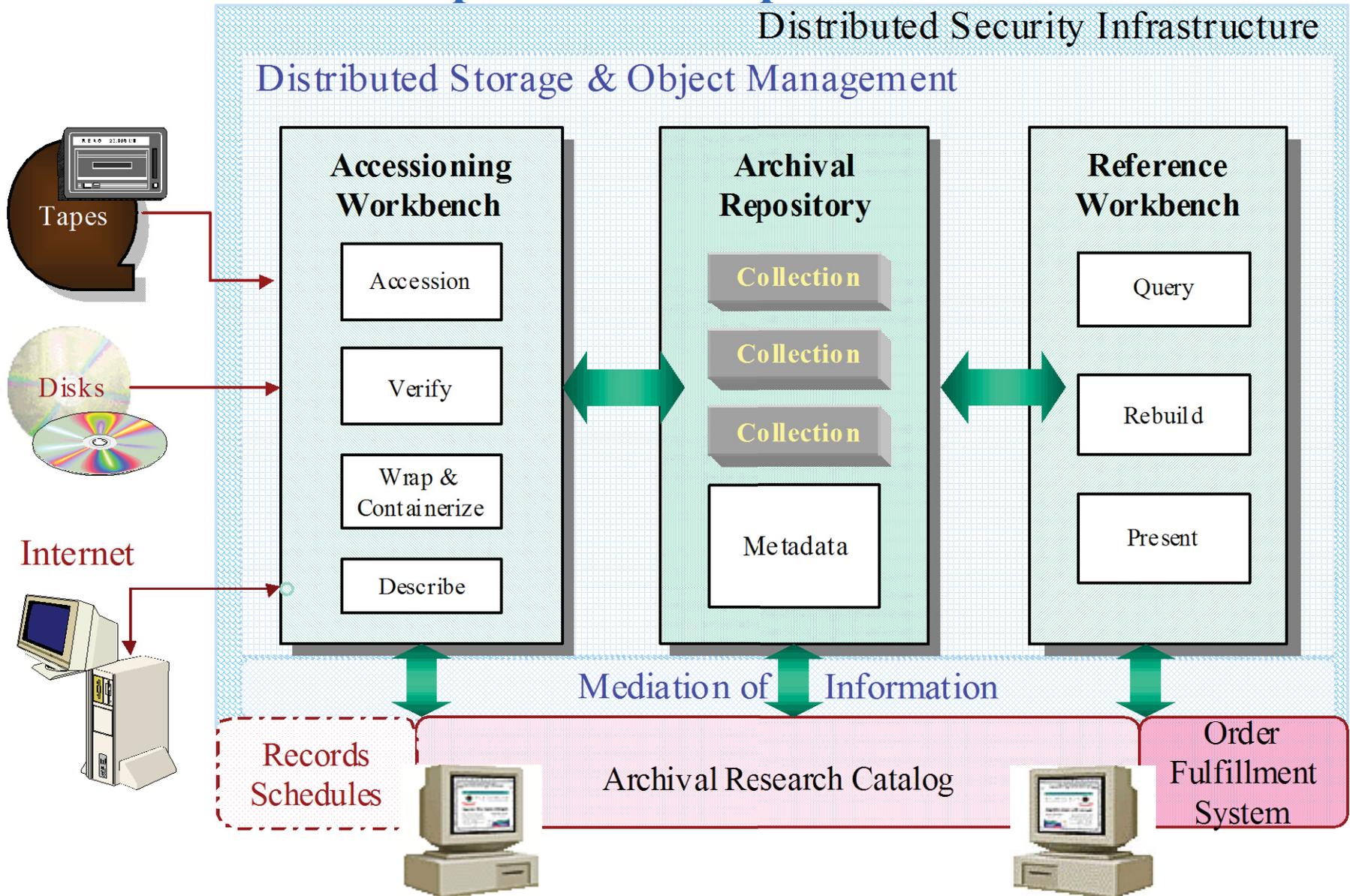


NODE: A0	TITLE: Preserve Electronic Records	NUMBER: v 5b draft
-------------	---------------------------------------	-----------------------

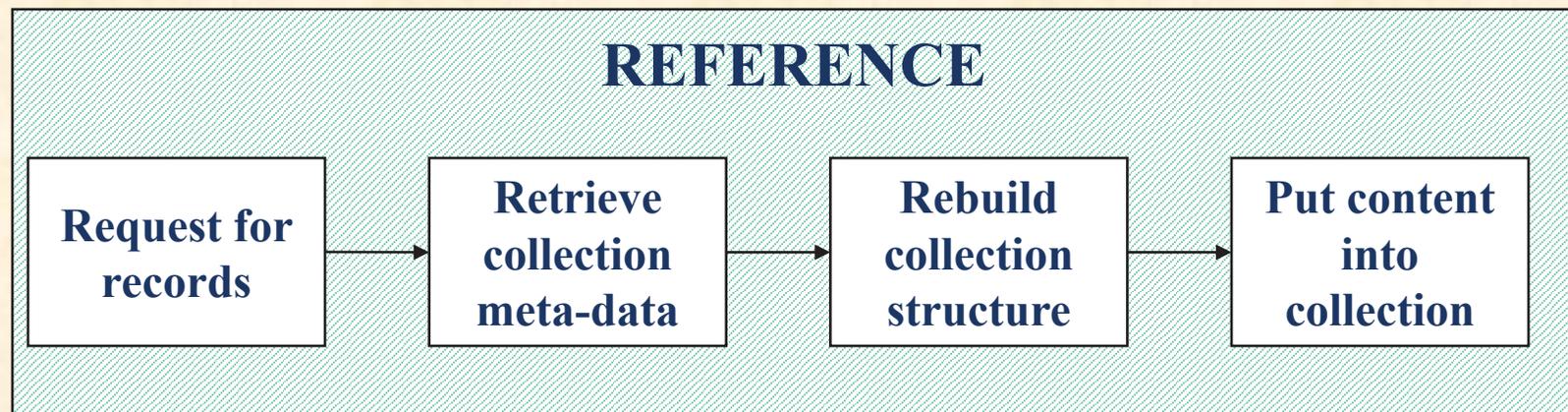
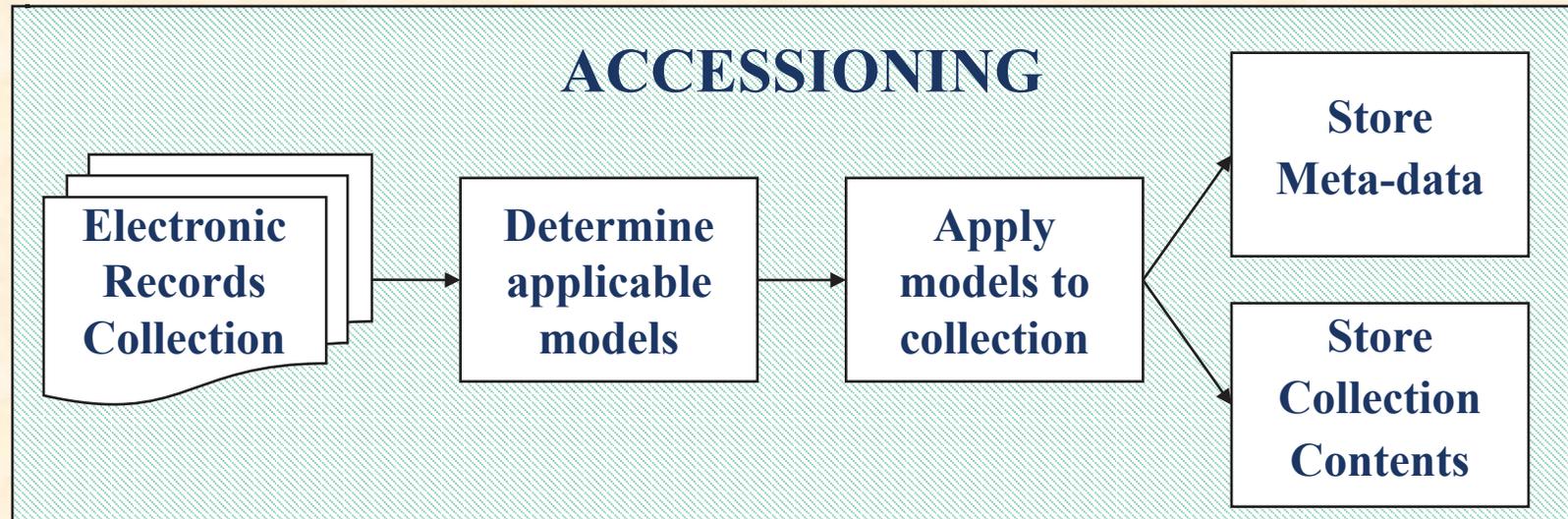
# Information Management Architecture for Persistent Object Preservation



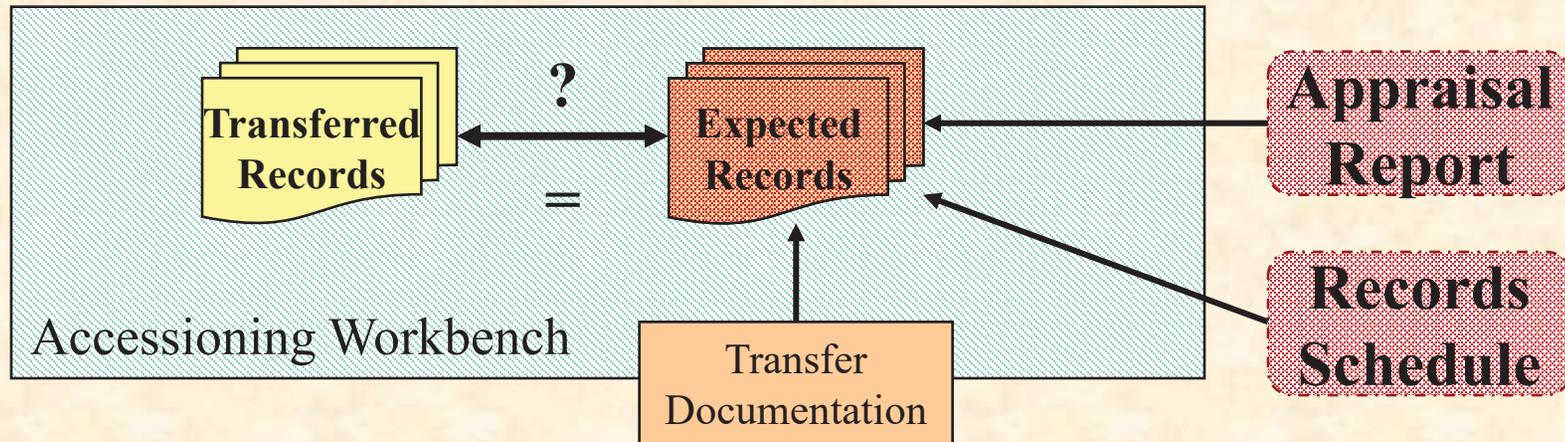
# ERA: Archival Components Concept



# ERA Processes



# Accept an Accession?



- What should the agency have transferred?
- What did the agency say it transferred?
- What was transferred?

# SELF-DESCRIBING

- Records
- Files
- Series
- Record  
Systems

# Persistent Object Method

- Characterize significant properties of the things that are to be preserved.
- Express these properties in formal models
- Encapsulate objects in metadata defined in the models.
- Use software “mediators” to enable future technologies to interpret the models and metadata
  - to rebuild and repopulate collections
  - to re-present the records
  - support information discovery and delivery.

# E-mail: Groupwise view



From: Yigal Arens <arens@ISI.EDU>

CC:

To:

BC: Ken Thibodeau

Subject: Announcing DG Online, the magazine of digital governmentresearch

Message:

DG Online <<http://www.dgrc.org/dg-online/>>.

DG Online: The Magazine for Digital Government Research is the new online quarterly of dg.o (DigitalGovernment.Org), a national consortium of government agencies, computer science researchers, the IT industry, and civic organizations concerned with improving online government operations and services.

DG Online presents the latest developments in advanced computer and IT research for Digital Government along with news and viewpoints on the most important DG issues:

- \* Cooperation among federal, state, and local government agencies
- \* Privacy and security
- \* Universal access--bridging the digital divide
- \* Streamlined delivery and tracking of public services
- \* Electronic voting, taxes, the Census, and other sensitive online data collection
- \* IT user friendly interfaces

Attach:



Close



Reply



Forward



Info



Delete



## E-mail: MIME-aware view

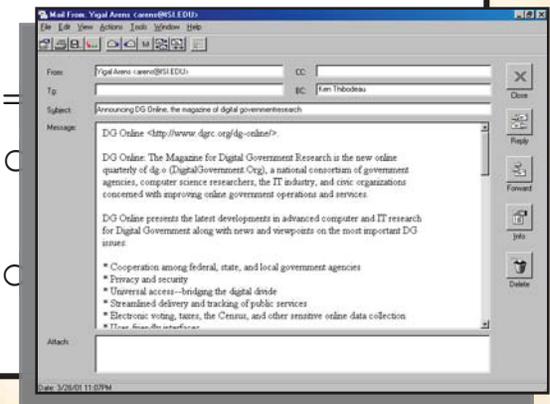
```
Message-Id: <p05010429b6e7e938e0b4@[10.2.68.205]>
X-Organization: USC/Information Sciences Institute
X-Phone: (310) 822-1511 ext. 766
X-Fax: (310) 822-0751
Date: Wed, 28 Mar 2001 22:22:30 -0500
To: Readers:;
From: Yigal Arens <arens@ISI.EDU>
Subject: Announcing DG Online, the magazine of
digital government
research
Content-Type: multipart/alternative;
boundary="===== _ -
1226283781==_ma======"
|
--===== _ -1226283781==_ma======"
Content-Type: text/plain; charset="us-ascii" ;
format="flowed"

DG Online <http://www.dgrc.org/dg-online/>.

DG Online: The Magazine for Digital Government
Research is the new online
```

# Tagged MIME E-mail Message

```
<Message-Id> p05010429b6e7e938e0b4@[10.2.68.205]
</Message-Id>
<X-Organization>USC/Information Sciences Institute</X-
Organization>
<X-Phone>(310) 822-1511 ext. 766</X-Phone>
<X-Fax>(310) 822-0751</X-Fax>
<Date>Wed, 28 Mar 2001 22:22:30 -0500</Date>
<To>Readers: ;</To>
<From>Yigal Arens {arens@ISI.EDU}</From>
<Subject>Announcing DG Online, the magazine of digital
government research</Subject>
<Content-Type: multipart/alternative;
boundary="=====-1226283781==_ma====="
-----=-1226283781==_ma-----
<Content-Type: text/plain; charset="us-asc
format="flowed">
<Message_Body>DG Online {http://www.dgrc.c
....
```



# Structure of E-mail Message aka: Document Type Definition

```
<!ELEMENT Email_Message (Header, Message Body, Attachment*)>
  <!ELEMENT Header (Internal Header, External Header)>
    <!ELEMENT Internal_Header (Message_Id, X-
      Organization, X-Phone, X-Fax)>
      <!ELEMENT Message-Id>
      <!ELEMENT X-Organization>
      <!ELEMENT X-Phone>
      <!ELEMENT X-Fax>
    <!ELEMENT External_Header (Date, To, From,
      Subject)>
      <!ELEMENT Date (Weekday, Day_of_Month, Month,
        Time)>
      <!ELEMENT To (#PCDATA)+>
      <!ELEMENT From (#PCDATA)>
      <!ELEMENT Subject (#PCDATA)*>
  <!ELEMENT Message_Body (#PCDATA)*>
```



# eXtensible Business Reporting Language (XBRL) Example

## REPORT OF INDEPENDENT ACCOUNTANTS

To the Board of Directors and Stockholders of Great Plains Software, Inc.

In our opinion, the consolidated financial statements listed in the accompanying index present fairly, in all material respects, the financial position of Great Plains Software, Inc. and its subsidiaries at May 31, 1999 and 1998, and the results of their operations and their cash flows for each of the three years in the period ended May 31, 1999, in conformity with generally accepted accounting principles. In addition, in our opinion, the financial statement schedules listed in the accompanying index present fairly, in all material respects, the information set forth therein when read in conjunction with the related consolidated financial statements. These financial statements and financial statement schedules are the responsibility of the Company's management; our responsibility is to express an opinion on these financial statements and financial statement schedules based on our audits. We conducted our audits of these statements in accordance with generally accepted auditing standards, which require that we plan and perform the audit to obtain reasonable assurance about whether the financial statements are free of material misstatement. An audit includes examining, on a test basis, evidence supporting the amounts and disclosures in the financial statements, assessing the accounting principles used and significant estimates made by management, and evaluating the overall financial statement presentation. We believe that our audits provide a reasonable basis for the opinion expressed above.

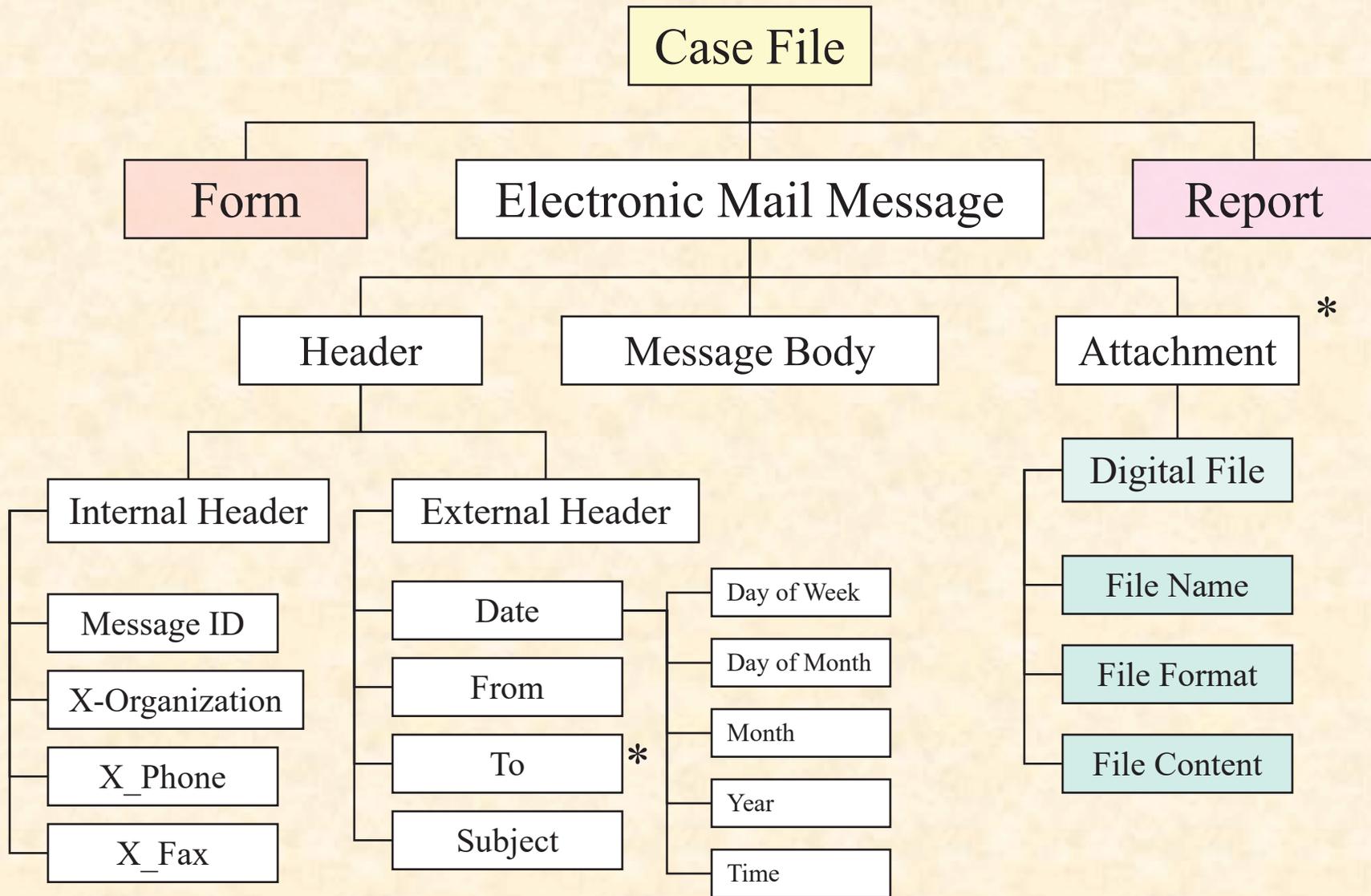
/s/ PricewaterhouseCoopers LLP



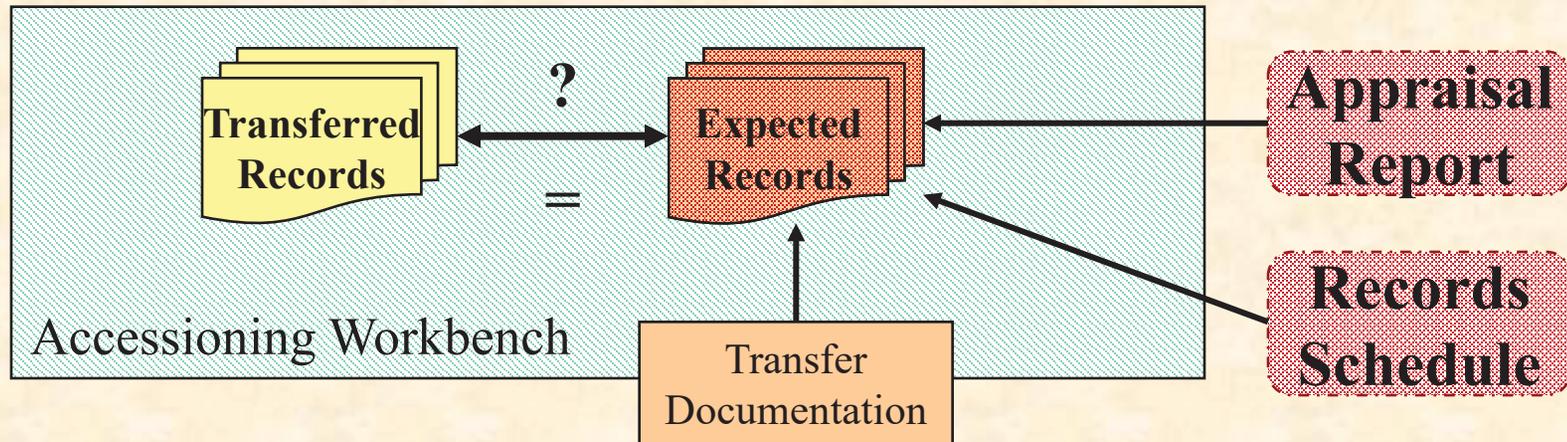
# XBRL DTD

Structure	Values
[COMMENT]	SECTION: AccountantReport
ITEM	PricewaterhouseCoopers LLP
TYPE	accountant_sReport.independent
ITEM	/s/ PricewaterhouseCoopers LLP
TYPE	accountant_sReport.accountantSignature
ITEM	Minneapolis
TYPE	accountantSignature.city
ITEM	Minnesota
TYPE	accountantSignature.state
ITEM	REPORT OF INDEPENDENT ACCOUNTANTS
TYPE	accountant_sReport.titleOfAccountantsReport
ITEM	To the Board of Directors and Stockholders of Great Plains Software, Inc.
TYPE	accountant_sReport.addressee
ITEM	June 25, 1999
TYPE	reportDate.date
ITEM	Unqualified
TYPE	accountant_sReport.typeOfOpinion
ITEM	US GAAP
TYPE	reportingMethod.generallyAcceptedAccountingPrinciples
ITEM	In our opinion, the consolidated financial statements listed in the accompanying index present fairly, in all material respects, the f...
TYPE	scopeOfWorkPerformed.auditedFinancialStatements
[COMMENT]	SECTION: BalanceSheet
ITEM	26983
ID	BS-01
TYPE	cashAndCashEquivalents.cashEquivalents
PERIOD	1999-05-31
ITEM	18197
ITEM	96700
ITEM	48721
ITEM	12593
ITEM	8790

# Structure expressed as Tree



# Accept an Accession?

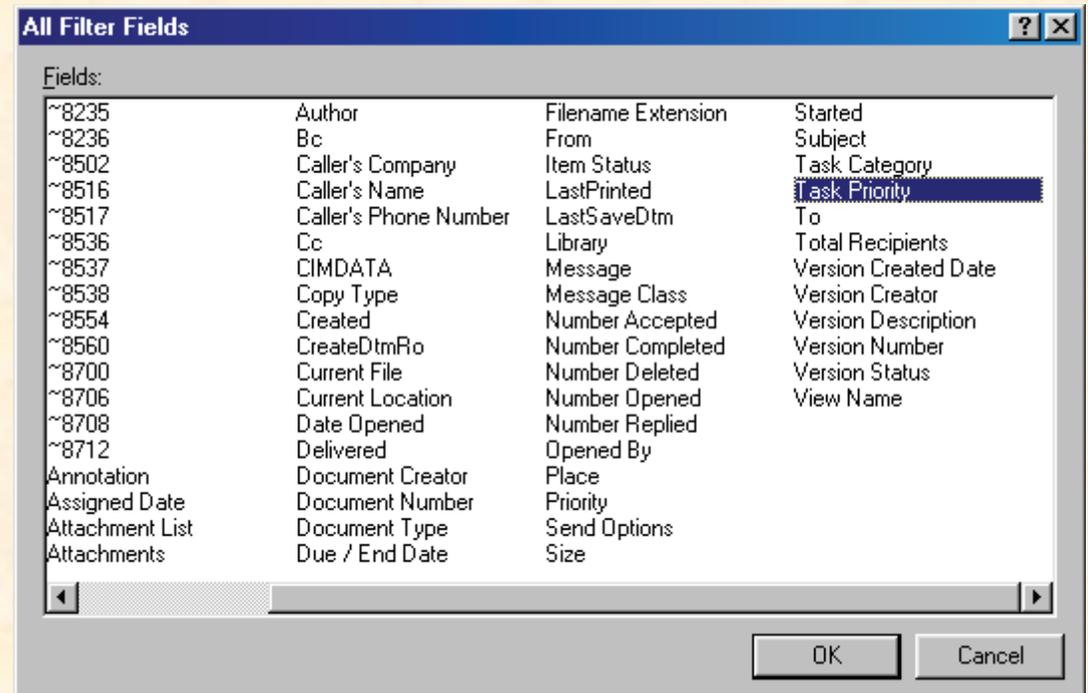


## How does ERA determine the dates of records?

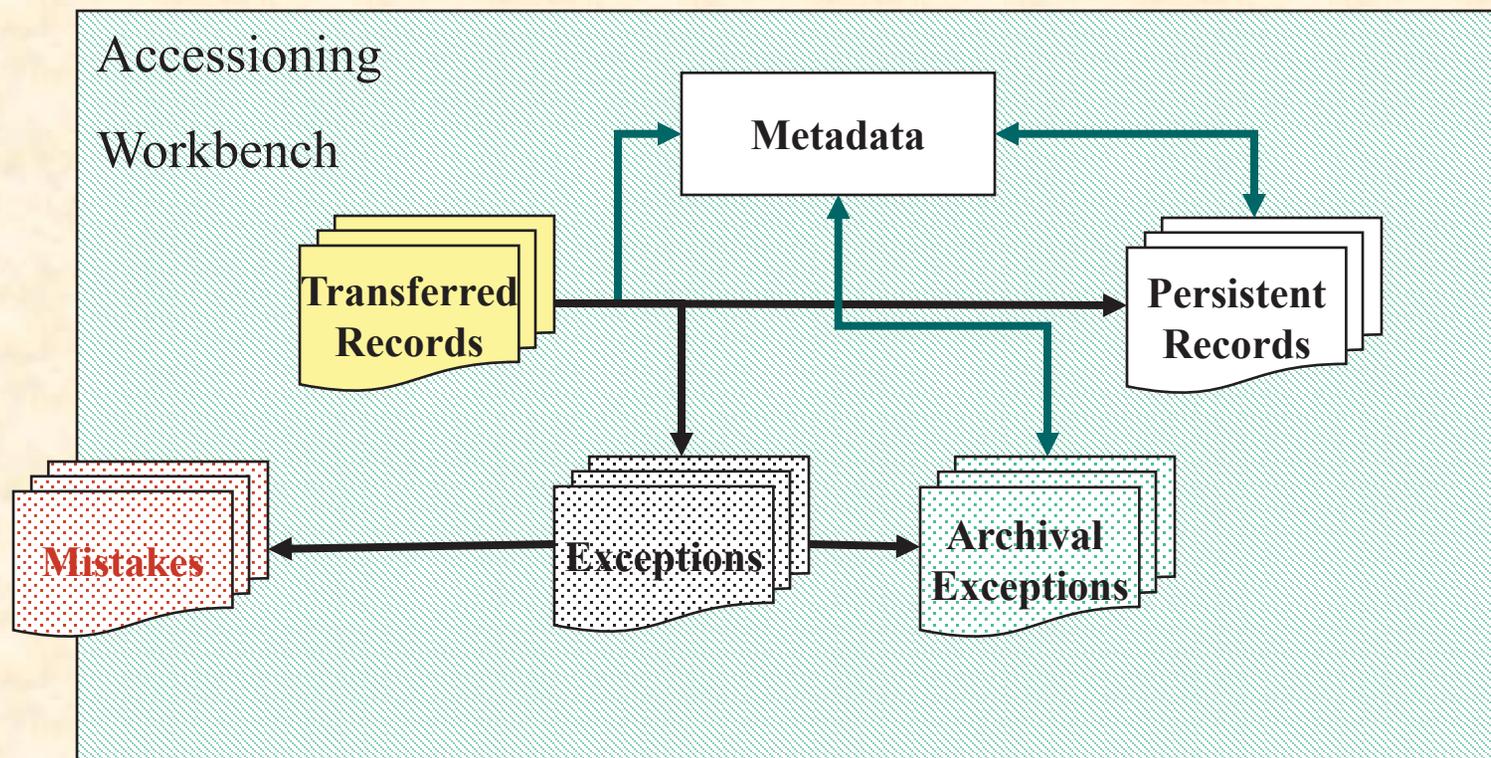
- E-mail
  - All e-mail contains a field indicating the date it was sent. For the sender, that is the date of the record. ERA needs to search the date-sent fields.
    - *(Technology solution)*
- Attachments to e-mail messages
  - Attachments to a record are parts of that record. The date of the message is the date of the record.
    - *(Archival principle)*
- Records forwarded, via e-mail, for filing in a recordkeeping system
  - E-mail is used only to transmit a record to the system. The date of the attached record depends on the record
    - *(Archival principle)*

# Defining models for electronic records e.g. E-mail

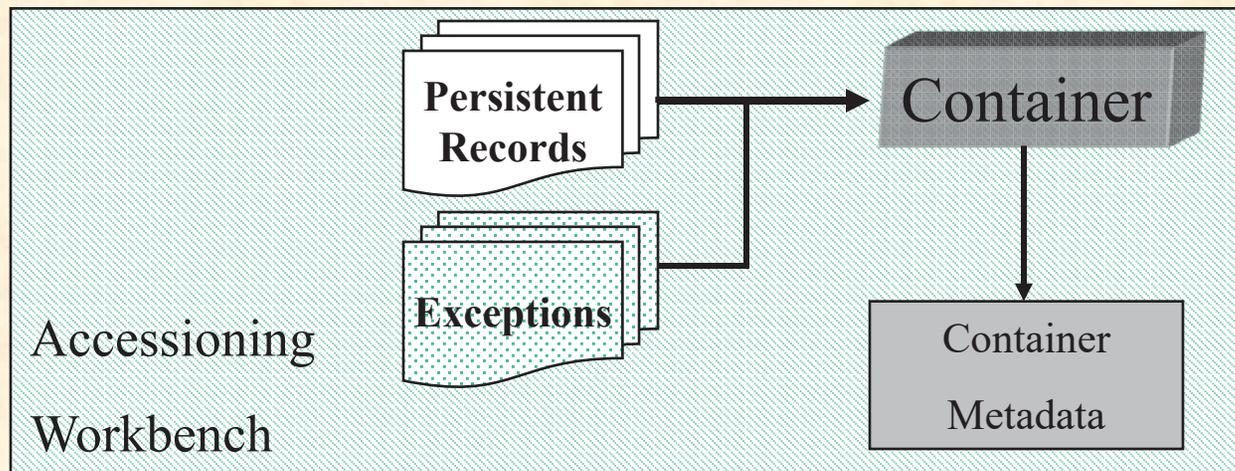
- All E-mail
  - Groupwise mail
  - cc:mail
  - USENET mail
    - User defined fields
  - .....



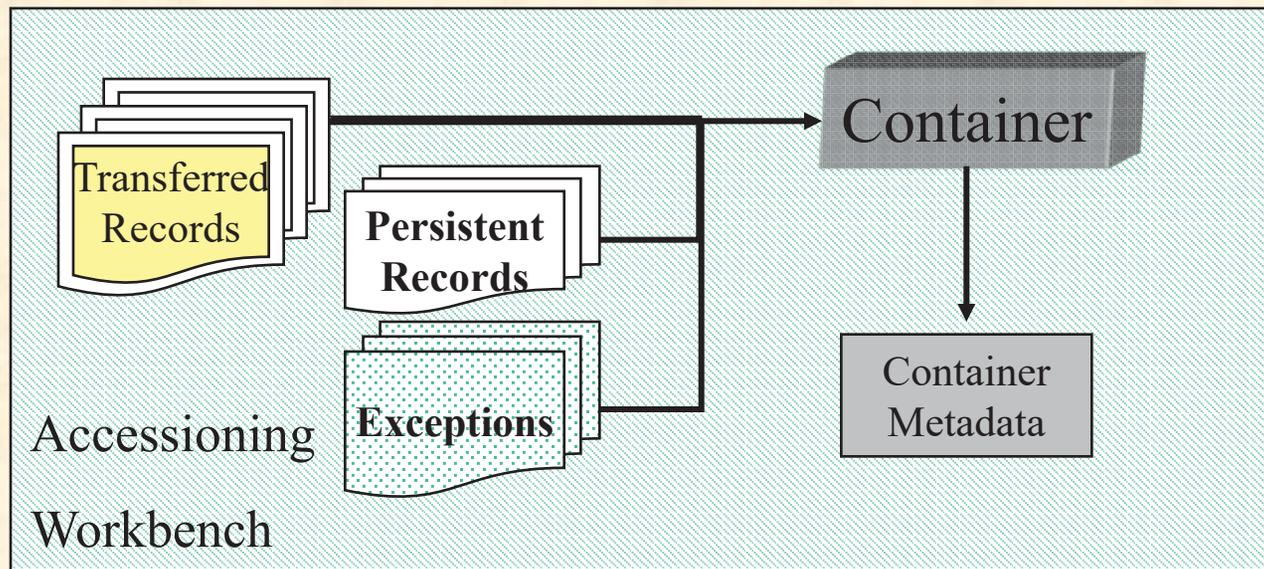
# Transformation



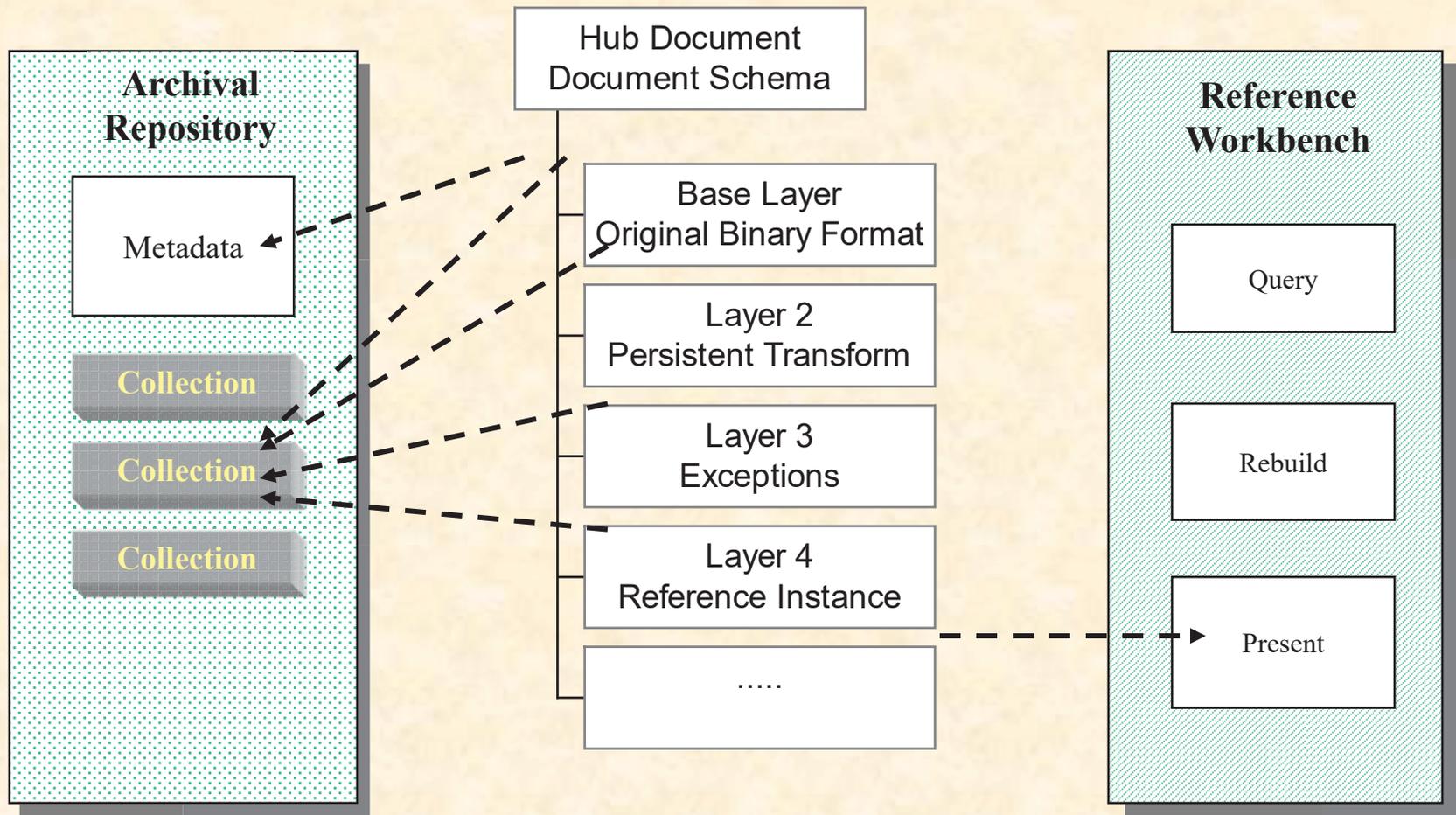
# Aggregation



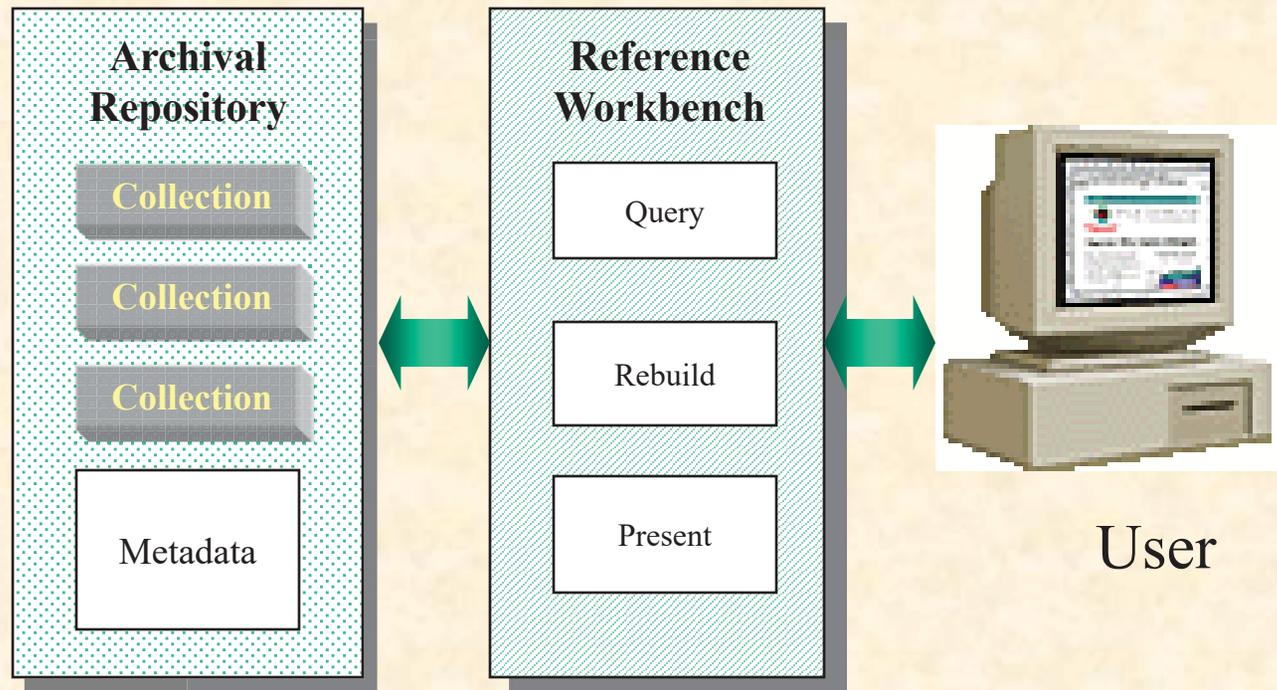
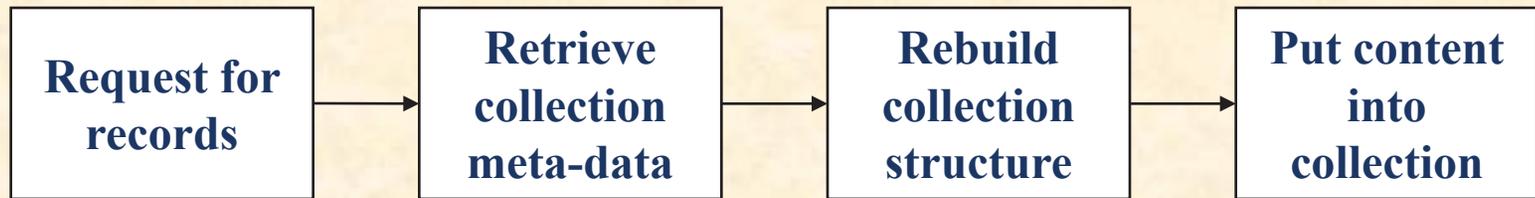
# Aggregation: risk management



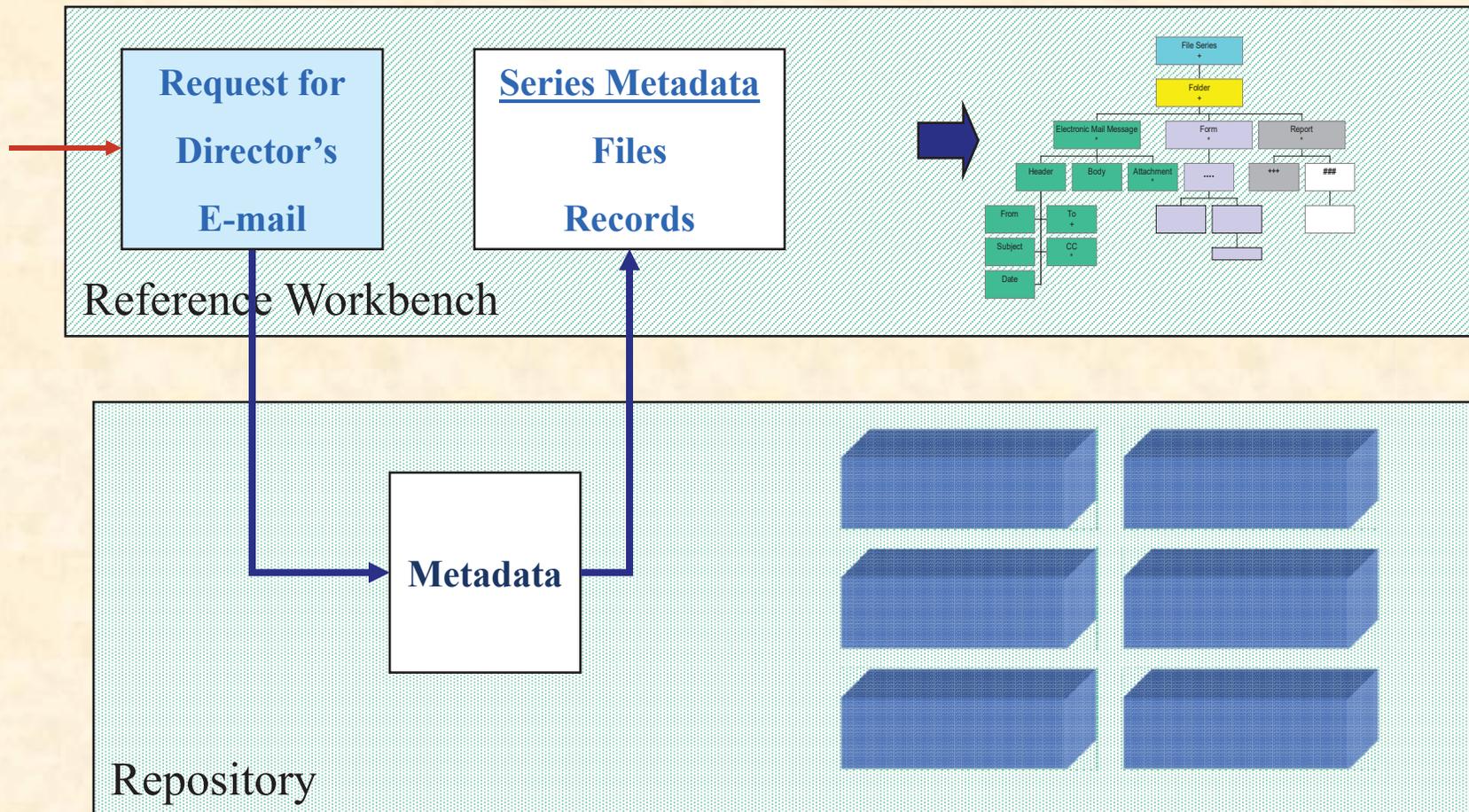
# Risk Management: Multi-Valent Documents



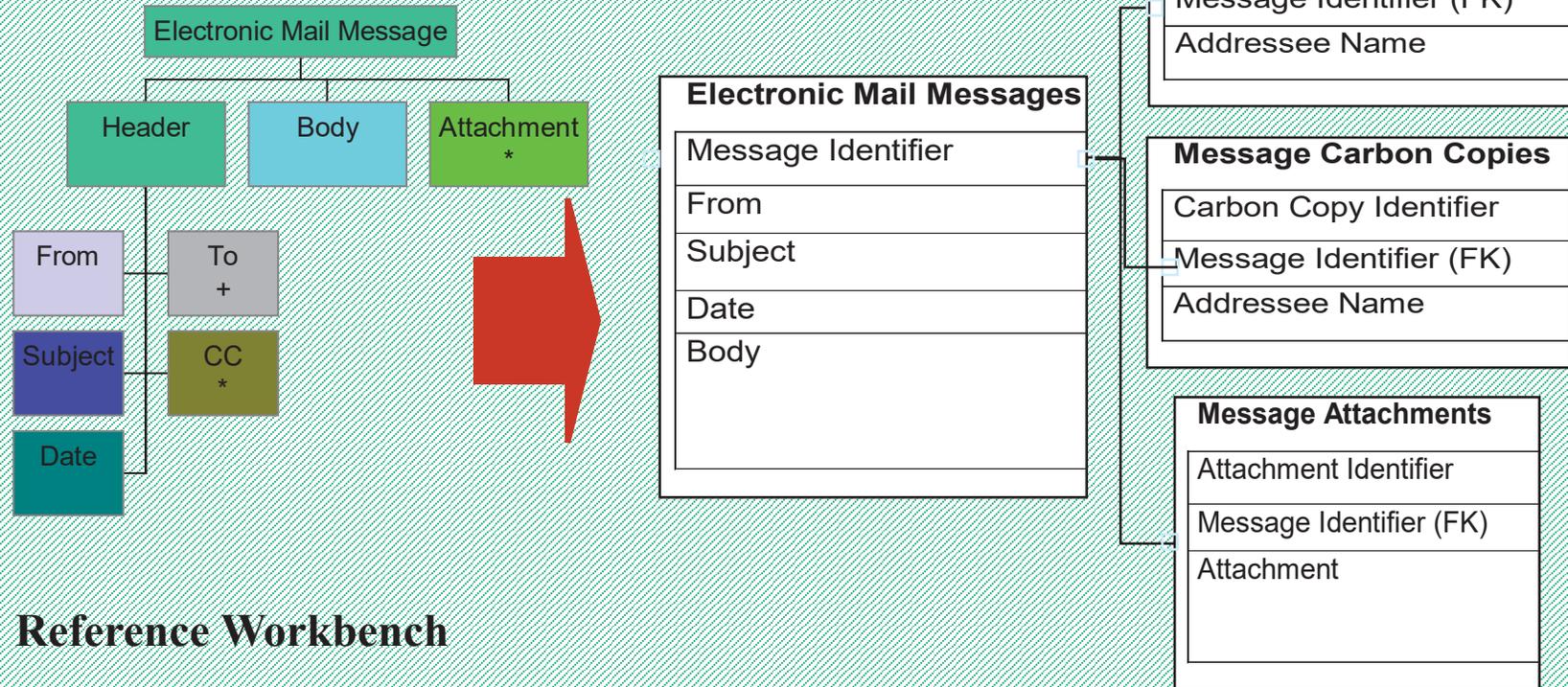
# ERA: Reference Process



# Process: Check metadata for the series to identify relevant DTDs

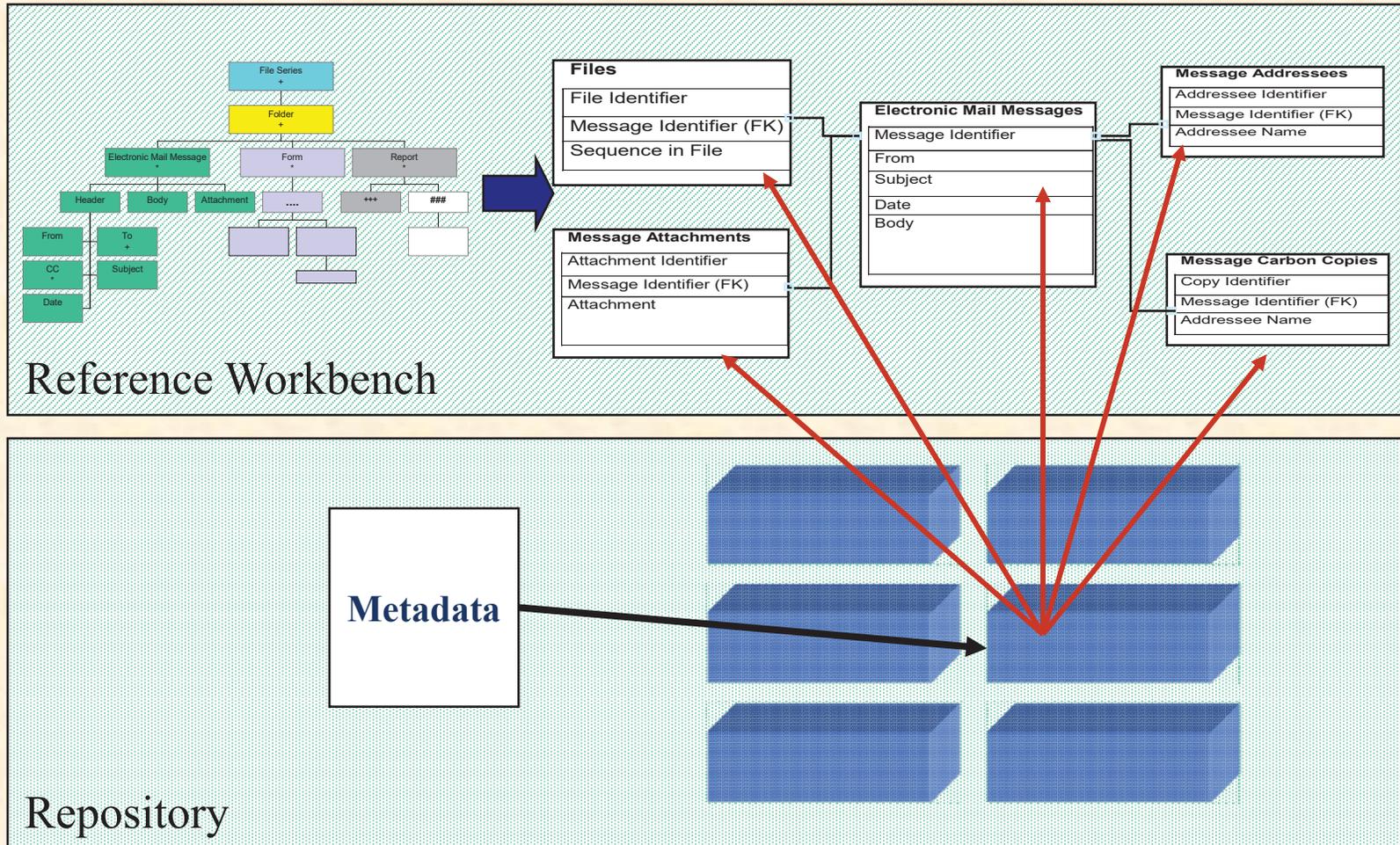


# Translate E-mail DTD to Relational Database Structure



Reference Workbench

# Process: Retrieve the records and place in the target structure



# Persistent Object Preservation

- + Aims at independence of technological infrastructure
  - + Reduce threats to integrity and authenticity by minimizing changes over time.
- + Embeds changes in a comprehensive information management architecture designed for preservation
- + Inherently extensible
- + Facilitates use of future, advanced technologies, without requiring change in what is preserved.
- Currently beyond state of the art of information technology.

# Self-describing Objects for Records Management

- Facilitate management, exchange, and disposition of records
  - explicitly identify the content of records, files, series,...
  - express how content is organized
  - allow the content to be stored once and used in different documents
  - separate, but link, management of content and presentation
  - capture the relationships among documents and collections of documents
  - and support multiple views of a collection of documents
  - all in plain language

# Thank you.

For more information:

[www.nara.gov/era](http://www.nara.gov/era)

