



# InterPARES 2 Project

International Research on Permanent Authentic Records in Electronic Systems

**Title:** General Study 01 Final Report:  
Building Preservation Environments  
with Data Grid Technology<sup>†</sup>

**Status:** Final (public)

**Version:** 1.0

**Submission Date:** October 2004

**Release Date:** September 2007

**Author:** The InterPARES 2 Project

**Writer(s):** Reagan W. Moore  
San Diego Supercomputer Center

**Project Unit:** Focus 2

**URL:** [http://www.interpares.org/display\\_file.cfm?doc=ip2\\_gs01\\_final\\_report.pdf](http://www.interpares.org/display_file.cfm?doc=ip2_gs01_final_report.pdf)

---

<sup>†</sup> A complimentary report with the same title was subsequently published by this author in 2006 in *American Archivist* 69(1): 139-158.

## Abstract

The preservation of digital records requires many of the same types of processes as used for preservation of paper records. The InterPARES Project has identified standard processing steps and characterized them in a functional description. The InterPARES functions can be mapped onto software infrastructure based on modern data management technology. The software infrastructure can be assembled into a preservation environment that supports the technical stabilization and physical and intellectual protection of data, documents, and records through time. I present an analysis of the software infrastructure needed to preserve digital records, and identify the new processes and concepts required for permanent preservation, the stable uninterrupted management of data without a foreseeable end.

## Persistent Archive Description

The preservation of digital records is of interest to all communities managing digital records, from archivists, to computer scientists, to librarians.<sup>2</sup> The fundamental issue is the preservation of authenticity of digital records while the technology in the supporting infrastructure evolves.<sup>3</sup> Whenever digital records are preserved for time periods greater than the lifetime of technology, mechanisms are required to support migration to new technology or emulation of old technology. The challenge of technology evolution must be understood and managed for long-term preservation to be achievable.

Preservation environments can be characterized by functional descriptions, such as that produced by the project on International Research on Permanent Authentic Records in Electronic Systems (InterPARES), by dataflow descriptions, which track the processing steps that are applied to electronic records, and by infrastructure descriptions, which identify the software systems needed to manage electronic records.

The InterPARES 1 Preservation Model<sup>4</sup> specifies the activities associated with ensuring the preservation of authentic electronic records.<sup>5</sup> A *preservation action plan* describes the preservation actions to be taken for the transfer of records to the archives, for accessioning the records and for maintaining the records. Preservation actions are implemented using preservation methods. The methods rely on software for generic preservation methods such as integrity checks, methods for packaging or archiving many files as one, for refreshing media, for data base management and for archival storage. They also include specific preservation methods, for example, for reproducing records, for converting proprietary formats to standard formats and for converting digital objects in proprietary formats to persistent objects. The actions in a preservation action plan trigger methods associated with an archival process.

In this paper, I examine the mapping from the operations required by the InterPARES 1 Preservation Model to the capabilities provided by Data Grids.<sup>6</sup> Data Grids provide interoperability mechanisms across

---

<sup>2</sup> Thibodeau, K., "Building the Archives of the Future: Advances in Preserving Electronic Records at the National Archives and Records Administration," U.S. National Archives and Records Administration, <http://www.dlib.org/dlib/february01/thibodeau/02thibodeau.html>.

<sup>3</sup> Moore, R. (2000), "Knowledge-based Persistent Archives," in *Proceedings of La Conservazione Dei Documenti Informatici Aspetti Organizzativi E Tecnici*, in Rome, Italy, October, 2000.

<sup>4</sup> Preservation Task Force, "Appendix 5 – A Model of the Preservation Function," in *The Long-term Preservation of Authentic Electronic Records: Findings of the InterPARES Project*, Luciana Duranti, ed. (San Miniato, Italy: Archilab, 2005), 253-292. PDF version available at <http://www.interpares.org/book/index.cfm>. [Note: At the time this report was written, the InterPARES 2 Project's *Chain of Preservation Model* (see [http://www.interpares.org/ip2/ip2\\_models.cfm](http://www.interpares.org/ip2/ip2_models.cfm)) was still under development and thus unavailable for use in this report.]

<sup>5</sup> Underwood, W. E., "As-Is IDEF0 Activity Model of the Archival Processing of Presidential Textual Records," TR CSITD 98-1, Information Technology and Telecommunications Laboratory, Georgia Tech Research Institute, December 1, 1998.

<sup>6</sup> Moore, R. and C. Baru (2003), "Virtualization Services for Data Grids," chapter 16 in *Grid Computing: Making the Global Infrastructure a Reality*, F. Berman, A. J. G. Hey and G. C. Fox, eds. (John Wiley & Sons), 409-436.

storage repositories, information repositories and authentication environments. Data Grids also provide persistent naming conventions that permit the migration and replication of digital entities onto new storage systems. The lowest-level processes in the InterPARES model specify explicit capabilities that can be supported with Data Grid technology. I present an overview of archival processes, the capabilities provided by Data Grids, an assessment of alternate preservation approaches and a characterization of each of the principal Data Grid capabilities.

A preservation environment that is based upon Data Grid technology is called a Persistent Archive.<sup>7</sup> I apply the term “persistent” to the concept of an archive to represent the management of the evolution of the software and hardware infrastructure over time. In this report, I describe the capabilities that are needed by preservation environments to automate both the management of technology evolution and the application of archival processes by archivists. I note that the terminology used by the preservation community conflicts with the terminology from the Data Grid community. Throughout the paper, references are made to the management of digital entities in persistent archives. Archives, within a Data Grid, refer to the storage systems used for long-term storage. Archives, as used by archivists, are the organized non-current records of an institution. I use the term digital entity, to represent any sequence of bits that constitute an entity that will be preserved. Records, as the term is used by archivists, are a class of digital entities with unique attributes and constraints. Records, in addition to their data bits or content, also require metadata to describe their provenance (source, submitting agency, original form) and their authenticity (digital signatures to verify they have not been changed, audit trails to track operations by archivists). The associated metadata defines the preservation context. A record is the combination of the content (digital entity) with the context (preservation metadata). The preservation of records requires the ability to manage the consistency of the archival context relative to the original content. The archival context in turn is created by the application of archival processes to the content.

A persistent archive maintains not only the data bits comprising digital entities, but also the context that defines the provenance, authenticity and structure of the digital entities. The context, from the perspective of the data grid, is managed as attributes that are organized into an authoritative catalog. I use the term archival form of a digital entity to represent the data bits, a definition of the structure of the digital entity and the associated preservation attributes. For example, the Open Archival Information System (OAIS) specifies an Archival Information Package (AIP) for defining the context of a digital entity.<sup>8</sup> The archival collection is the aggregation of the archival forms of the digital entities.

## InterPARES 1 Preservation Model

The InterPARES 1 preservation model contains three main processing steps: A2 – Bring in Electronic Records, A3 – Maintain Electronic Records and A4 – Output Electronic Records. Each main process is further decomposed into subsidiary steps, based upon version 6.0 of the InterPARES 1 IDEF0 model for the preservation of electronic records.<sup>9</sup> For the main subsidiary steps, the functionality that is needed

---

<sup>7</sup> Moore, R., C. Baru, A. Rajasekar, B. Ludascher, R. Marciano, M. Wan, W. Schroeder and A. Gupta (2000), “Collection-Based Persistent Digital Archives – Part 1,” *D-Lib Magazine* (March), <http://www.dlib.org/dlib/march00/moore/03moore-pt1.html>; and Moore, R., C. Baru, A. Rajasekar, B. Ludascher, R. Marciano, M. Wan, W. Schroeder and A. Gupta (2000), “Collection-Based Persistent Digital Archives – Part 2,” *D-Lib Magazine* (April), <http://www.dlib.org/dlib/april00/moore/04moore-pt2.html>.

<sup>8</sup> OAIS – “Preservation Metadata and the OAIS Information Model,” the OCLC working group on Preservation Metadata, June 2002, <http://www.oclc.org/research/pmwg/>.

<sup>9</sup> Underwood, W. E., “The InterPARES Preservation Model: A Framework for the Long-Term Preservation of Authentic Electronic Records,” paper presented at the *Choices and Strategies for Preservation of the Collective Memory*, International Conference, 25-29 June 2002, Toblach/Dobbiaco, Bolzano Province, Italy. Published in *Archivi per la Storia*.

within data management software is described in Appendix A. The InterPARES definition is given for each processing step, followed by the mechanisms provided by Data Grids.

A mapping of each of the Data Grid capabilities to the InterPARES 1 Preservation Model is shown in Table 1. All mechanisms that are needed for each activity are listed under the heading “Core Capabilities.” Five major categories are listed: 1) Logical name space for persistent identifiers, 2) Storage repository abstraction for migrating to new storage technologies, 3) Information repository abstraction for migrating to new databases, 4) Distributed resilient scalable architecture for performance and 5) Virtual data grid for managing archival processes.

Note that a given Data Grid capability is used multiple times across the InterPARES 1 preservation activities. Processing steps bring electronic records into a processing workbench, create the archival context for each digital content component and then manage consistency of the content and context. The content is written to storage repositories and the context is written to an authoritative catalog. Both content and context may be replicated to ensure the ability to recover from disasters. Archival processes are executed under authentication and authorization controls, and all operations are tracked using audit trails.

## Persistent Archive Functionality Requirements

The requirements for a persistent archive can be expressed in general as “transparencies” that hide virtual data grid implementation details.<sup>10</sup> Examples include digital entity name transparency, data location transparency, platform implementation transparency, encoding standard transparency and authentication transparency for single sign-on systems. The capabilities of a persistent archive can be characterized as the set of “transparencies” needed to manage technology evolution. Implementations exist in data grids for at least four key functionalities or transparencies that simplify the complexity of accessing distributed heterogeneous systems:

- Name transparency – The ability to identify a desired digital entity without knowing its name can be accomplished by queries on descriptive attributes, organized as a collection. Persistent archives are inherently archives of collections of digital entities that map from unique attribute values to a global, persistent identifier.
- Location transparency – The ability to retrieve a digital entity without knowing where it is stored can be accomplished through use of a logical name space that maps from the global, persistent, identifier to a physical storage location and physical file name. If the data grid owns the digital entities (stored under the custodian user ID), the administrative attributes for storage location and file name can be self-consistently updated every time the digital entity is moved.
- Platform implementation transparency – The ability to retrieve a digital entity from arbitrary types of storage systems can be accomplished through use of a data grid that provides a storage repository abstraction. The data grid maps from the protocols needed to talk to the storage systems to the operations defined by the storage repository abstraction. Every time a new type of storage system is added to the persistent archive, a new driver is added to the data grid to map from the new storage access protocol to the data grid data transport protocol. Similar platform transparency is needed for the information repository in which the persistent archive stores the

---

<sup>10</sup> Moore, R., C. Baru, A. Rajasekar, R. Marciano and M. Wan (1999), “Data Intensive Computing,” chapter 5 in *The Grid: Blueprint for a New Computing Infrastructure*, I. Foster and C. Kesselman, eds. (San Francisco: Morgan Kaufmann), 105-130.

Table 1. Mapping of the InterPARES 1 Preservation Model to functionality provided by Data Grids

<b>Core Capabilities</b>	A2.1	A2.2	A2.3	A2.4	A3.1	A3.2	A3.3	A4.1	A4.2	A4.3	A4.4	A4.5
Storage repository abstraction	x		x			x	x	x	x			
Storage interface to at least one repository	x		x			x	x					
Standard data access mechanism	x		x			x	x	x				
Standard data movement protocol support	x		x			x	x	x			x	
Containers for data	x					x	x					x
Logical name space	x	x	x	x	x	x	x	x	x	x	x	
Registration of files in logical name space	x		x		x	x	x					
Retrieval by logical name			x			x	x					
Logical name space structural independence from physical name space	x				x	x	x					
Persistent handle					x	x						
Information repository abstraction	x	x	x	x	x	x	x	x	x	x		x
Custodian owned data	x	x	x	x	x	x	x	x				
Collection hierarchy for organizing logical name space	x		x		x			x	x			
Standard metadata attributes (controlled vocabulary)	x	x	x	x	x	x	x	x	x	x		x
Attribute creation and deletion			x									
Scalable metadata insertion	x		x	x	x	x	x					
Access control lists for logical name space to control who can see, add and change metadata	x	x	x	x	x	x	x	x			x	
Attributes for mapping from logical file name to physical file names	x		x		x	x	x	x	x			
Encoding format specification attributes		x	x				x		x	x		
Data referenced by catalog query								x				
Containers for metadata	x				x	x	x		x			x
Distributed resilient scalable architecture	x				x	x	x	x				
Specification of system availability	x				x	x		x	x			
Standard error messages	x				x	x	x	x	x	x	x	x
Status checking	x		x		x	x		x	x			
Authentication mechanism	x	x	x	x	x	x	x	x			x	
Specification of reliability against permanent data loss						x	x					
Specification of mechanism to validate integrity of data and metadata			x						x		x	x
Specification of mechanism to assure integrity of data and metadata					x	x	x	x				
Virtual Data Grid			x				x			x		x
Knowledge repositories for managing collection properties		x	x									x
Application of transformative migration for encoding format			x				x			x		
Application of archival processes	x	x	x		x	x	x			x		x

- collection context. An information repository abstraction is defined for the set of operations needed to manipulate a catalog in an information repository, or database.
- Encoding standard transparency – The ability to display a digital entity requires understanding the associated data model and encoding standard for information. If infrastructure independent standards are used for the data model and encoding standard (non-proprietary, published formats), a persistent archive can use transformative migrations to maintain the ability to display the digital entities. The transformative migrations will need to be defined between the original encoding standard and the contemporary infrastructure independent data model standard.

Possibly the unique capability that must be present in a persistent archive is the ability to preserve authenticity. This implies an environment in which only authorized actions can take place. Every operation within the persistent archive should be tracked and the corresponding metadata updated to guarantee consistency of the metadata to preserve authenticity. This can be most easily implemented by having the digital entities stored under the control of the Data Grid. This forces access to be done through the data grid, making it possible to track all operations that are done on the digital entities, from transformative migrations, to media migrations, to replication, to accesses. Data Grids implement restricted access through the use of collection-based ownership of the registered digital entities.

One can think of a Data Grid as the set of abstractions that manage differences across storage repositories, information repositories, knowledge repositories and execution systems. Data Grids also provide abstraction mechanisms for interacting with the objects that are manipulated within the grid, including digital entities (logical namespace), processes (service characterizations or application specifications) and interaction environments (portals). The Data Grid approach can be defined as a set of services and the associated application programming interfaces (APIs) and protocols used to implement the services. The Data Grid is augmented with portals that are used to assemble integrated work environments to support specific applications or disciplines. An example is an archivist workbench, which provides separate functions for each of the archival processes.

## **Persistent Archive versus Persistent Storage**

There is a distinction between Persistent Archives and Persistent Storage. Persistent storage systems provide archival media that have a very long shelf life, such as heavy-ion beam encoded disk, film, etc. A standard encoding is chosen (such as ASCII) that is assumed readable at an arbitrary date in the future. The technology to extract meaning from the archived material is based on the ability to parse ASCII. Persistent archives recognize that there is a cost benefit that can be obtained by migrating to new technology, including minimization of floor space through higher density media, lower cost storage media, elimination of obsolete equipment and improved access.

Both systems need universal identifiers, the ability to manage descriptive, authenticity and preservation metadata for the archived material, policy management systems to control the archival workflow, and audit trails of transactional activity and user history. Grid technology provides the mechanisms that make it possible to migrate to new versions of technology and to new archival services. The distributed state information that is managed by grid technology can be employed to support more extended applications for in-depth re-purposing beyond general catalog functions.

## Prior Conceptual Models

Conceptual models can be used from prior research efforts to evaluate the completeness of the proposed approach to persistent archives based on data grid technology. The models are selected from the archives and computer science domains, and include: traditional archive procedures, a reference model for BAC (business acceptable communications) developed by the University of Pittsburgh, the records continuum model proposed by the Monash University in Australia, the reference model for an open archival information system (OAIS) designed by the CCSDS of NASA, and the ISO/IEC 11179 standard for data element composition for inclusion in metadata registries. A comparison with these models illustrates the multiple characterizations of archival processes that have been used. The comparison also demonstrates the importance of policy management issues.

## Traditional Archival Procedures

The capabilities of virtual Data Grids can be used to implement the traditional archival processes of appraisal, accession, arrangement, description, preservation and access. These can be mapped to the InterPARES 1 Preservation Model.

## Records Continuum Model

The records continuum model uses four processes for preservation, namely: create, capture, organize and pluralize.<sup>11</sup> Frank Upward states: "The four continua I chose to represent as sets within a spacetime continuum model were identity, transactionality, evidentiality and recordkeeping containers [which I more normally refer to these days as recordkeeping objects]."<sup>12</sup> The records continuum model was originally developed as a teaching tool to communicate evidence-based approaches to archives and records management. Upward states that: "It [the records continuum model] can never provide complete or satisfying views of detailed practice, but that is not what a worldview does. It provides an overview for re-organising our detailed knowledge and applying our skills in contexts framed by the task at hand...As a view it presents a multi-layered and multi-faceted approach which can be used to re-organise knowledge and deploy skills. It is more in tune with electronic communications and technological change than a life cycle view."<sup>13</sup> Data management systems that provide mechanisms to manage technological change are consistent with the records continuum model.

## BAC Reference Model

The Reference Model for BAC is also based on a distributed environment.<sup>14</sup> Thus the Data Grid approach is consistent with the BAC model, except for three layers: terms and conditions layer, contextual layer and user history layer. These layers comprise policy management (for terms and conditions), a knowledge layer (for defining the context) and an access layer (for describing user interaction history). These layers are expected to become grid mechanisms that in the future will be part of virtual data grids. The implication is that the proposed Data Grid model must continue to evolve to include future grid services that are appropriate for preservation.

---

<sup>11</sup> Records Continuum Model, Records Continuum Research Group, <http://rcrg.dstc.edu.au/>.

<sup>12</sup> Upward, F. (2000), "Modeling the records continuum as a paradigm shift in record keeping and archiving processes, and beyond - a personal reflection," *Records Management Journal* 10(3): 115-139.

<sup>13</sup> *Ibid.*, 128.

<sup>14</sup> BAC, Business Acceptable Communications, <http://www.phila.gov/records/divisions/rm/units/perp/presentations/nagara/nagara96/sld005.html>.

## OAIS Model

The OAIS system specifies a reference model for describing the processes associated with preservation from the viewpoint of submission information packages (SIPs), archival information packages (AIPs) and dissemination information packages (DIPs).<sup>15</sup> The OAIS reference model can be implemented on top of data grid technology through the specification of the interaction and information packaging mechanisms. The OAIS reference model specifies the information that should be associated with each procedure. The information can be stored with the digital entities, or stored in a metadata repository that can be queried, or stored in both places. This implementation choice is left to the persistent archive.

From the OAIS point of view, "The OAIS model provides a theoretical framework for an archival system, and integrates its conceptual approach with a hierarchical structure for organizing information." The model does not specify an implementation strategy; instead it provides guidelines to address digital archiving concepts both from functional and information model. The OAIS model describes an Archival Information Package (AIP) as an aggregation of four types of Information Objects:

1. Content information object: includes the data object as well as representation information (structural and semantic information about the object).
2. Preservation description information object: comprised of reference information, provenance information, context information and fixity information.
3. Packaging information object: information that is stated as being used to bind and identify the components of an Information Package. An example is the ISO-9660 volume and directory information used on a CD-ROM to define the content of several files containing Content Information and Preservation Description Information.
4. Descriptive information object.

Besides the OAIS metadata, technical metadata (refers to the administrative, structural and preservation metadata related to digital objects) is needed to facilitate management and access to archival objects. The technical metadata needs to be independent of the object itself to support interoperability and preservation. The OAIS reference model provides a very good reference not just for the design of a long-term digital archive, but also provides several real use cases in its appendix for verification.

## ISO/IEC 11179

ISO/IEC 11179 specifies basic aspects of data element composition for inclusion in metadata registries.<sup>16</sup> The standard applies to the formulation of data element representations and meaning as shared among people and machines. Metadata registries are authoritative semantic maps with associated procedures for storing and registering detailed metadata from multiple sources and diverse organizations in a common structured form. Extensions to the formats are recorded, as are agreed-upon mappings between diverse formats. Use of the ISO /IEC standard and participation in metadata registries promotes access, understanding and sharing of data across time and space, and use of this structure makes it easier to check the metadata for consistent application.

## Other models

There are many possible levels of granularity and different ways to categorize information. Working within the OAIS framework and ISO/IEC 11179 is a sound strategy because it makes possible improved

---

<sup>15</sup> OAIS - Reference Model for an Open Archival Information System (OAIS). Submitted as ISO draft, <http://www.ccsds.org/documents/pdf/CCSDS-650.0-R-1.pdf>, 1999.

<sup>16</sup> ISO/IEC 11179, Specification and Standardization of Data Elements, <http://www.diffuse.org/oii/en/meta.html#ISO11179>.



communication among divergent digital applications. Other projects are also addressing the characterization of preservation systems.

- The National Library of the Netherlands Long Term Preservation Study distinguishes between Intellectual Preservation, Media Preservation and Technology Preservation.<sup>17</sup>
- The Making of America project distinguishes between descriptive, administrative and structural metadata.<sup>18</sup>
- The METS schema classifies administrative metadata into four types: technical metadata, intellectual property rights metadata, source metadata and digital provenance metadata.<sup>19</sup>
- The IEEE Learning Object Metadata draft standard includes metadata categories for general; lifecycle; meta-metadata; technical; educational; rights; relation; annotation; classification.<sup>20</sup>
- The NDAP National Digital Archive Project proposes core capabilities for preservation, including linkage to the original object to keep complete information about how the digital entity was created (ways of digitization, equipment used, workflow, accuracy, data quality and specifications for the digitization work); metrics to specify the quality and completeness (in terms of information lost ratio) of the digitized entity; support for knowledge level information discovery; content analysis based on complete metadata analysis to provide structure and organization for the contents (metadata schema design); descriptions of each archived collection in the content space, which is composed of space, time and linguistics perspectives; flexible presentation with the content (representation) separated from presentation through a content management framework (CMF) constructed to manage the workflow from authoring to publishing; and authentication for data and owner.<sup>21</sup>

## Managing Technology Evolution

A major design issue for the creation of persistent archives is the development of an approach in which the digital entities can be preserved in an unchanged format, while still making it possible for future presentation applications to display the digital entity. The challenge is that the encoding format interpreted by future applications will not be the same as the encoding format used to create the original digital entity. Three approaches are being considered within the archival community to resolve this challenge.

1. Migration of digital entities applies an archival process to create an infrastructure independent representation of the digital entity by changing the encoding format to a non-proprietary standard. In the process, the bits of the digital entity must be changed to the standard encoding format. The expectation is that the transformative migration will need to be done at an infrequent interval.
2. Emulation preserves the original digital entity by migrating the presentation application onto new technology. Instead of migrating the digital entity to new encoding formats, the presentation application is migrated to new operating systems. This requires migrating onto new technology the applications that were used to create or view each digital entity. The result is a system that preserves the look and feel of the original software, but at the same time makes it very difficult to apply any new techniques to the interpretation of the digital entities. An emulator can be characterized as the set of

<sup>17</sup> National Library of the Netherlands, Koninklijke Bibliotheek, <http://www.konbib.nl/>.

<sup>18</sup> Making of America II, <http://sunsite.berkeley.edu/MOA2/>.

<sup>19</sup> METS – “Metadata Encoding and Transmission Standard,” <http://www.loc.gov/standards/mets/>.

<sup>20</sup> IEEE Learning Object Metadata, [http://ltsc.ieee.org/doc/wg12/LOM\\_1484\\_12\\_1\\_v1\\_Final\\_Draft.pdf](http://ltsc.ieee.org/doc/wg12/LOM_1484_12_1_v1_Final_Draft.pdf).

<sup>21</sup> NDAP, National Digital Archives Project, Taiwan.

operations that the original application must be able to perform through an operating system. This characterization is typically specified as a set of operating system calls. An emulator maps from the system calls used by the original application to the system calls provided by current operating systems. The Dyninst system<sup>22</sup> is an example of software that supports the dynamic insertion of new system calls into existing code, and can be viewed as enabling infrastructure for the development of emulators.

3. Migration of digital ontologies, characterizations of the data structure and data model that specify how to manipulate a digital entity.<sup>23</sup> Emulation and migration capabilities can be combined by creating a digital ontology that organizes the relationships present within a digital entity. A digital entity can be viewed as a sequence of bits onto which structural, procedural and semantic relationships are applied. These relationships include the structural relationships that define how to turn the bits into binary arrays, or words, or tables. Logical relationships are used to apply semantic tags to the structures. Spatial relationships are used to map binary arrays to coordinate systems. Temporal relationships are used to apply time stamps to structures. The digital ontology specifies the order in which the relationships need to be applied to correctly interpret the information and knowledge content.

The digital entity is kept in its original encoding format. Instead of changing the encoding format of the digital entity to a non-proprietary standard, a digital ontology is created that defines the relationships present within the digital entity. The digital ontology is migrated onto new encoding standards for relationships over time. For instance, a digital ontology can be represented using the Resource Description Framework syntax. In the future, when a new syntax is used to specify relationships, the digital ontology can be migrated from the old syntax to the new syntax, without modifying the original digital entity.

The presentation application is emulated as the set of operations that can be performed on the defined relationships. The set of operations can be kept fixed on the original set, or they can be expanded over time as new capabilities are created (such as causal queries on time stamps). In effect, the presentation application is emulated as operations on a digital ontology, and the digital ontology is migrated forward in time onto new encoding formats.

All references to migration in this report can be interpreted as either migration of digital entities onto new encoding formats for display by future applications, or migration of digital ontologies onto new encoding formats for display through a standard set of operations.

## Example Persistent Archive

An example persistent archive has been constructed using the San Diego Supercomputer Center Storage Resource Broker data grid.<sup>24</sup> The persistent archive components based upon the SRB include:

- Logical name space implemented in the Metadata Catalog (MCAT).<sup>25</sup> The logical names are chosen by the archivist. The archivist will use the logical names as they are defined in the record

---

<sup>22</sup> Dyninst – a machine-independent interface to permit the creation of tools and applications that use runtime code patching, <http://www.paradyn.org/release3.3/>.

<sup>23</sup> Moore, R. (2002), "The San Diego Project: Persistent Objects," in *Proceedings of the Workshop on XML as a Preservation Language*, Urbino, Italy, October, 2002.

<sup>24</sup> SRB - "The Storage Resource Broker Web Page," <http://www.npaci.edu/DICE/SRB/>.

<sup>25</sup> MCAT - "The Metadata Catalog," <http://www.npaci.edu/DICE/SRB/mcat.html>.

collection whenever possible. A mapping is maintained from the logical name to the physical file location. The logical names are infrastructure independent and are organized in a collection hierarchy, allowing the specification of different descriptive metadata for each sub-collection. Soft links and shadow links are supported for the logical organization and registration of digital entities. Digital entities may include files, URLs, SQL command strings, directories and database tables. Distributed state information is mapped onto the logical name space as attributes.

- Storage repository abstraction implemented in the SRB. The set of operations that are supported include Unix file system operations (create, open, close, unlink, read, write, seek, sync, stat, fstat, mkdir, rmdir, chmod, opendir, closedir and readdir), latency management operations (aggregation of data, I/O commands and metadata) and metadata manipulation (extraction, registration) through use of remote proxies. Containers are used to physically aggregate digital entities before storage into archives. Both digital entities and containers can be replicated. The storage repository abstraction is used to manage data within Unix file systems, archives, object-relational databases, object ring buffers, storage resource managers, FTP sites, GridFTP sites and Windows file systems.
- Information repository abstraction implemented in the MCAT. Mechanisms are supported for schema extension through addition of new attributes, table restructuring and metadata import and export through XML files. Soft links are supported for logical reorganization of digital entities within a collection hierarchy. Metadata attributes are maintained for provenance attributes (Dublin core), administrative metadata (file location), descriptive metadata (user-defined attributes) and authenticity metadata (audit trails, digital signatures). The information repository abstraction is used to manage metadata in both proprietary and non-proprietary databases including DB2, Oracle, Sybase, Informix, SQLServer and Postgresql.
- Distributed resilient architecture implemented through a federated client server architecture. Servers are installed in front of each storage repository and in front of the information repository. Access to the system results in the creation of a service instance that manages further interactions for the request. The service instance retrieves all required distributed state information from the MCAT catalog that is needed to complete the request, and interacts with remote servers as needed to access the data. The system has been designed to minimize the number of message sent over wide area networks to improve performance and increase reliability. Data retrieval requests are automatically retried on a replica when a storage repository does not respond. All error messages generated by the network, storage repository and information repository are returned to the user. Consistency constraints on distributed state information are explicitly integrated into the software through use of write locks and synchronization flags. This makes it possible to update a file that has been aggregated into a container and replicated into an archive, lock out competing activities to avoid over-writes and then synchronize all replicas to the new state. When additional records for a record series are received, they can be appended to the container holding the records that have already been accessioned. Changes to digital entities within a container are made by marking the original digital entity as deleted and appending the new form of the digital entity to the end of the container. The addition of digital entities to an archival collection can also be done through soft links within the logical name space, making it possible to link digital entities into an existing collection, while simultaneously organizing the new digital entities in a separate sub-collection. All system metadata is automatically generated and updated by the SRB on each request.

- Virtual data grid implemented through use of remote proxies and external process management systems. The SRB provides a mechanism to process data remotely, before it is sent over a network. The Ohio State University DataCutter technology is used to filter data.<sup>26</sup> External process management systems can control the generation of derived data products through application of remote proxies or the DataCutter filters. Interactions with databases can be expressed through SQL command strings that are registered into the logical name space. The SRB is able to apply simple transformative migrations such as unit conversion and reformatting of query results into HTML or XML. More complex transformations require the use of a process management system.

## Summary

A proposed set of core capabilities can be defined for minimizing the labor required to implement, manage and evolve a persistent archive. The capabilities are present within implementations of current data grids. Many of the capabilities are general properties that have been implemented across almost all existing data grids. A characterization of each capability has been defined. This characterization can be used as the set of requirements for defining a persistent archive architecture that supports the InterPARES 1 Preservation Model.

## Acknowledgements

The results presented here were supported by the NSF NPACI ACI-9619020 (NARA supplement), the NSF NSDL/UCAR Subaward S02-36645, the DOE SciDAC/SDM DE-FC02-01ER25486 and DOE Particle Physics Data Grid, the NSF National Virtual Observatory, the NSF Grid Physics Network and the NASA Information Power Grid. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the National Science Foundation, the National Archives and Records Administration, or the U.S. government.

The definition of terms used by archivists was provided by Mark Conrad of the National Historical Publications and Records Commission, and the comparison with OAIS technology was provided by Han-wei Yen, Simon C. Lin, Ya-ning Chen and Shu-jiun Chen of the Computing Centre of Academia Sinica, Taiwan. Margaret Hedstrom provided valuable comments on preservation terminology. William Underwood provided the example assessments for mapping archival processes onto data grid capabilities, an analysis of the records continuum model and the comparison of data grid capabilities with the InterPARES 1 Preservation Model.

More information about the example persistent archive based on the San Diego Supercomputer Center Storage Resource Broker can be found at [http://www.sdsc.edu/srb/index.php/Main\\_Page](http://www.sdsc.edu/srb/index.php/Main_Page) and <http://www.sdsc.edu/NARA/>.

---

<sup>26</sup> Beynon, M.D., T. Kurc, U. Catalyurek, C. Chang, A. Sussman and J. Saltz (2001), "Distributed Processing of Very Large Datasets with DataCutter," *Parallel Computing* 27(11): 1457-1478.

## Glossary

Two sets of terminology are used in the report; one from the preservation community and one from the grid community. In some cases, the same word is used in two different contexts.

### Terms used within the preservation community

#### ACCESS

The right, opportunity, or means of finding, using, or approaching documents or information (SAA).  
<http://rpm.lib.az.us/alert/thesaurus/terms.asp?letter=a>

Access to archival documents is provided through an archival reference service. Key steps in an archival reference service are:

- Querying the researcher to draw out the specific nature of the subject as well as secondary aspects of the subject that can serve as leads to documentation sources.
- Translating the terms and concepts of the inquiry into the terms and concepts of the archives' reference apparatus.
- Explaining finding aids, archival methodology, and the nature of manuscripts and records documentation
- Guiding the researcher to the appropriate finding aids and/or records.
- Retrieving the records that appear to be relevant to the researcher's inquiry.
- Informing the researcher of policies and practices for making copies and handling documents to ensure that the records are not damaged or disarranged.
- Consulting with the researcher during and after the visit to determine how well the records answered the question or led to new questions.

<http://web.library.uiuc.edu/ahx/define.htm>

#### ACCESSION

(v.) To transfer physical and legal custody of documentary materials to an archival institution.

(n.) Materials transferred to an archival institution in a single accessioning action.

[http://www.archives.gov/research\\_room/alic/reference\\_desk/archives\\_resources/archival\\_terminology.html](http://www.archives.gov/research_room/alic/reference_desk/archives_resources/archival_terminology.html)

#### APPRAISAL

(n.) the process of determining the value and thus the disposition of records based upon their current administrative, legal, and fiscal use; their evidential and informational value; their arrangement and condition; their intrinsic value; and their relationships to other records (SAA).

<http://rpm.lib.az.us/alert/thesaurus/terms.asp?letter=a>

#### ARCHIVES

The organized non-current records of an institution or organization retained for their continuing value in providing a) evidence of the existence, functions, and operations of the institution or organization that generated them, or b) other information on activities or persons affected by the organization.

Derived from the Greek word for "government house," the term "archives" also refers to the agency responsible for selecting, preserving, and making available non-current records with long-term value and to the building or part of the building housing them.

<http://web.library.uiuc.edu/ahx/define.htm>

#### ARRANGEMENT

The body of principles and practices which archivists follow to group records in such a way as to reflect the manner in which they were held and used by the office or person creating the records. It involves the fundamental principles of respect des fonds, provenance, and sanctity of original order. The key units in archival arrangement are: record groups, sub-groups, and record series.

<http://web.library.uiuc.edu/ahx/define.htm>

#### AUTHENTIC RECORD

A record that is what it purports to be and that is free from tampering or corruption.

[http://www.interpares.org/book/interpares\\_book\\_q\\_gloss.pdf](http://www.interpares.org/book/interpares_book_q_gloss.pdf)

#### AUTHENTICATION

A declaration of a record's authenticity at a specific point in time by a juridical person entrusted with the authority to make such a declaration.

[http://www.interpares.org/book/interpares\\_book\\_q\\_gloss.pdf](http://www.interpares.org/book/interpares_book_q_gloss.pdf)

#### AUTHENTICATION CERTIFICATE OF TRUSTED THIRD PARTY

An attestation issued by a trusted third party for the purpose of authenticating the ownership and characteristics of a public key. It appears in conjunction with the digital signature of the author of a record, and is itself digitally signed by the trusted third party.

[http://www.interpares.org/book/interpares\\_book\\_q\\_gloss.pdf](http://www.interpares.org/book/interpares_book_q_gloss.pdf)

#### AUTHENTICITY

The quality of being authentic, or entitled to acceptance. As being authoritative or duly authorized, as being what it professes in origin or authorship, as being genuine.

[http://www.interpares.org/book/interpares\\_book\\_q\\_gloss.pdf](http://www.interpares.org/book/interpares_book_q_gloss.pdf)

#### COLLECTION

The hierarchy of archival collections generally goes from the fonds to the Record Group to the Record Sub-group to the Record Series to the Record sub-series to the file or folder to the individual record. Collection is often used when talking about a non-archival group of documents that were artificially put together by a collector and were not created by an institutional process.

#### CONTEXT

The circumstances of creation and history of ownership and usage of an archival collection, as well as the collection's original arrangement or filing structure. A clear context gives a collection enhanced legal and research value as it indicates that the collection's integrity was respected during a continuous chain of custody (ownership). The evidence in the collection remains intact. The collection was not rearranged or inappropriately added to or weeded. Historians may depend upon the inferences they draw from the collection's authentic filing structure. See also original order and provenance

<http://crm.cr.nps.gov/archive/22-2/22-02-19.pdf>

## DESCRIPTION

The process of recording information about the nature and content of the records in archival custody. The description identifies such features as provenance, extent, arrangement, format, and contents, and presents them in a standardized form.

<http://www.sfu.ca/archives/glossary.html>

## FONDS

The whole of the records, regardless of form or medium, automatically and organically created and/or accumulated and used by a particular individual, family, or corporate body in the course of that creator's activities or functions.

<http://www.sfu.ca/archives/glossary.html>

## PRESERVATION

Preservation encompasses the activities that prolong the usable life of archival records. Preservation activities are designed to minimize the physical and chemical deterioration of records and to prevent the loss of informational content.

[http://www.archives.gov/preservation/about\\_preservation.html](http://www.archives.gov/preservation/about_preservation.html)

## PROVENANCE

The principle of archival arrangement according to which each deposit of records should be placed within an overall arrangement or classification scheme that reflects its origin and relation to other deposits from the same administrative body.

<http://web.library.uiuc.edu/ahx/define.htm>

## RECORDS

Documents, regardless of form, produced or received by any agency, officer, or employee of an institution or organization in the conduct of its business. Documents include all forms of recorded information, such as: correspondence, computer data, files, financial statements, manuscripts, moving images, publications, photographs, sound recordings, drawings, or other material bearing upon the activities and functions of the institution or organization, its officers, and employees. A document becomes a record when it is placed in an organized filing system for use as evidence or information. It becomes archival when transferred to a repository for preservation and research use.

<http://web.library.uiuc.edu/ahx/define.htm>

## RECORD GROUP

A body of organizationally related records, normally large in size and established on the basis of provenance to accommodate the records of major organizational units and functions of an institution.

<http://web.library.uiuc.edu/ahx/define.htm>

## RECORD SUB-GROUPS

Smaller (than record groups) bodies of organizationally related records placed within a record group to correspond to the subordinate administrative units that collectively form the record group.

<http://web.library.uiuc.edu/ahx/define.htm>

## RECORD SERIES

A systematic gathering of documents that have a common arrangement and common relationship to the functions of the office that created them. Record series are the filing units created by offices at all levels in an institutional hierarchy. Each series will be arranged internally according to a system established and modified by its creators. Boundaries between one record series and the next are

sometimes razor-sharp and sometimes fuzzy. Typical record series include subject files, project files, chronological correspondence files, client files, applicant files, financial records files, voucher files, and minutes and agenda files.

<http://web.library.uiuc.edu/ahx/define.htm>

#### RESPECT DES FONDS

The principle of archival arrangement according to which each deposit (fonds) should be maintained as a separate entity, even if other fonds cover the same or similar subjects. It requires archivists to respect the integrity of the body of records at the time it is deposited in the archives.

<http://web.library.uiuc.edu/ahx/define.htm>

#### SANCTITY OF THE ORIGINAL ORDER

The principle of archival arrangement according to which the creator's arrangement of files and documents within a deposit should be maintained.

<http://web.library.uiuc.edu/ahx/define.htm>

### **Terms used within the Data Grid community**

#### BULK METADATA LOAD

The ability to import attribute values for multiple objects registered within the logical name space from a single input file.

#### COLLECTION-OWNED DATA

The storage of digital entities under a Unix user ID that corresponds to the collection. Access to the data is then restricted to a server running under the collection ID.

#### COLLECTIONS

The organization of the metadata attributes that are managed for the digital entities registered into the logical name space

#### CONTAINER

Aggregation of multiple digital entities into a single file, while retaining the ability to access and manipulate each digital entity within the file. A container is the digital equivalent of a cardboard box.

#### CURATION CONTROL

The administration tasks associated with creating and managing a logical collection

#### DERIVED DATA PRODUCTS

The result of execution of processes under the control of a virtual data grid. For persistent archives, derived data products can be transformative migrations of digital entities to new encoding formats. A data collection can be thought of as a derived data product that results from the application of archival processes to a group of constituent documents.

#### DIGITAL ONTOLOGIES

The organization of the set of semantic, structural, spatial, temporal, procedural and functional relationships that are present within a digital entity. The digital ontology specifies the order in which the relationships need to be applied to correctly display or manipulate the digital entity.



## INFORMATION REPOSITORY

A software system that is used to manage combinations of semantic tags (attribute names) and the associated attribute data values. Examples are relational databases, XML databases, object-relational databases, etc.

## INFORMATION REPOSITORY ABSTRACTION

The set of operations performed on an information repository for the manipulation of a catalog or collection.

## KNOWLEDGE

Relationships between attributes, or relationships that characterize properties of a collection as a whole. Relationships can be cast as inference rules that can be applied to digital entities. An example is the set of structural relationships used to parse metadata from a digital entity in metadata extraction.

## LOGICAL FOLDERS

Sub-collections within a collection hierarchy that are equivalent to directories in a file system, but are used to manage different sets of metadata attributes.

## LOGICAL NAME SPACE

A naming convention for labeling digital entities. The logical name space is used to create global, persistent identifiers that are independent of the storage location. Within the logical name space, information consists of semantic tags that are applied to digital entities. The logical name space can be organized as a collection hierarchy, making it possible to associate different metadata attributes with different sets of digital entities within the collection. This is particularly useful for accession, arrangement and description.

## METADATA

The combination of semantic tags and the associated tagged data, typically managed as attributes in a database. Metadata is called data about data.

## REGISTRATION

Addition of an entry to the logical name space, creation of a logical name and storage of a pointer to the file name used on the storage system.

## REPLICAS

Copies of a file registered into the logical name space that may be stored on either the same storage system or on different storage systems.

## SHADOW LINKS

Pointers to objects owned by individuals, used to register individually owned data into the logical name space, without requiring creation of a copy of the object on storage systems managed by the logical name space.

## SOFT LINKS

The cross registration of a single physical data object into multiple folders or sub-collections within the logical name space

### STORAGE REPOSITORY

A storage system that holds digital entities. Examples are file systems, archives, object-relational databases, object-oriented databases, object ring buffers, FTP sites, etc.

### STORAGE REPOSITORY ABSTRACTION

The set of operations performed on a storage repository for the manipulation of data.

### TRANSFORMATIVE MIGRATIONS

The processing of a digital entity to change its encoding format. The processing steps required to implement the transformative migration can themselves be characterized and archived, and then applied later.

### USER ACCESS

Consists of authentication to the data grid, checking of access controls for authorization, and then retrieval of the digital entity by the data grid from storage through the collection ID for transmission to the user.

### VIRTUAL DATA GRID

The automation of the execution of processes. References to the output of a process can result in the application of the process, or direct access to the output.

## Appendix 1: InterPARES 1 Preservation Model

Each task in the preservation model is listed, along with the mechanisms that are used in the data grid to support the functionality.

### **A2 Bring in Electronic Records**

Process each transfer of electronic records into accessioned electronic records, producing information about each transfer of electronic records

#### **A 2.1 Register Transfer**

Capture information about the transfer, such as submitter's name, record creator's name, and the date of receipt of the transfer in a Record of the Transfer, and establish basic control over the materials transferred by identifying what has been transferred and where it is located. Inspect what was received in order to ensure that the physical transfer has been accomplished correctly.

Within the Data Grid, content is registered from a remote site into the logical name space used by the archival processing environment (workbench), and the content is then copied onto a storage repository managed by the workbench. The content is stored under the control of the process custodian. The attributes that describe the source of the transfer are added to an authoritative catalog as part of the context for each content component, using standard archival metadata. Administrative metadata is updated to record the location of the content on the workbench. Status information is used to track completion of the transfer. Support for bulk registration and loading of content is provided through use of containers for both data and metadata.

#### **A2.2 Verify that the Transfer is Authorized**

Determine if the transfer is authorized; that is, it comprises the records that have been selected for preservation, and those records have been submitted either by the records creator or an agent acting for the creator. Verify information about the records, their digital components, and the basis for asserting the authenticity of the records as received; and that the materials transferred are of the correct types and in the specified formats.

Within the Data Grid, authentication systems are used to identify the submitter, and control access to the workbench. The authoritative catalog is accessed for information that describes the expected encoding formats, structural relationships, and provenance metadata.

#### **A2.3 Examine Electronic Records**

Determine if the transfer actually includes all records and aggregates of records specified in the terms and conditions of transfer and that these records and aggregates are adequately and accurately described in the accompanying information to enable their preservation, reproduction in authentic form, and interpretation. Identify any actions required to preserve both the individual records transferred and the archival sets in which these records belong. Initiate technical or other preservation actions that should be taken immediately and schedule preservation actions that should be taken at a later date. Sub-processes include “A2.3.1 – Map Records and Digital Components within Transferred Materials,” “A2.3.2 – Verify that the Records in the Transfer Can Be Preserved and Reproduced,” and “A2.3.3 – Take Action Needed to Preserve the Record.”

Within the Data Grid, the records are read from storage repository using a standard access mechanism, the records are analyzed to identify the digital components, the logical name space is accessed for each

record, and the associated structural metadata is written into the authoritative catalog. To verify the records, the submitted content is parsed and checks are done to verify that the desired archival form can be created. This requires reading the content, assigning metadata attributes to record the status of the analysis, updating the authoritative catalog, updating audit trails, and managing the parsing process. To preserve the record, the archival form that associates preservation context with the material content is created. The context is updated in the authoritative catalog, along with the audit trails.

#### **A2.4 Accession Electronic Records**

Formally accept responsibility for preserving a transferred body of records. Create a Record of the Accession using Retrieved Information about the Presumption of Authenticity.

Within the Data Grid, the authoritative context is updated to reflect the status of the processing, and attributes are set to assert accession of the content.

#### **A3 Maintain Electronic Records**

Apply preservation method(s) by maintaining the digital components of accessioned electronic records, along with related information necessary to reproduce the records, certify their authenticity, and enable correct interpretation of the records.

##### **A3.1 Manage Information About Records**

Collect and maintain information necessary to carry out the Preservation Strategy, including information about the digital components, the archival aggregates they comprise, their authenticity, their interpretation, and the preservation activities performed on them. Include Storage Information identifying the files, locations, and other relevant data about the digital components of the Accessioned Electronic Records when they are placed in storage and subsequently when storage parameters are changed. Sub-processes include “A3.1.1 – Maintain Information About Records,” “A3.1.2 – Retrieve Information About Records,” and “A3.1.3 – Retrieve Information About Digital Components.”

Within the Data Grid, administrative metadata is managed that tracks the location of the content in storage repositories, audit trails are updated when content is moved, and information about relationships between digital components is updated on transformative migrations. If context is replicated between multiple catalogs, a persistent handle is used to assert equivalence between authoritative versions. Scalable mechanisms are used for bulk insertion of metadata. To retrieve preservation metadata and administrative metadata, the authoritative catalog is accessed using the logical name space, or records of interest may be identified through queries on descriptive metadata. The person issuing the request is authenticated and checked for access authorization.

##### **A3.2 Manage Storage of Digital Components of Records**

Place the digital components of Accessioned Electronic Records into storage. In response to a Request for Digital Components, retrieve the requested components and output them. Provide Updated Storage Information about the identities, locations and other relevant parameters of stored digital components whenever components are updated or other changes, such as media refreshment, are made in storage. Sub-processes include “A3.2.1 – Place Record Components in Storage,” “A3.2.2 – Refresh Storage,” “A3.2.3 – Monitor Storage,” “A3.2.4 – Correct Storage Problems,” and “A3.2.5 – Retrieve Components from Storage.”

Within the Data Grid, content is stored using the storage repository abstraction for managing heterogeneous storage systems. Administrative metadata is updated in the authoritative catalog, along

with the audit trails. The data integrity mechanism replicates both content and the authoritative context through use of Data Grid Federation mechanisms. Data Grids store the content under custodian control, and use containers for aggregating small files, to avoid overloading the archive name space. Scalable metadata insertion mechanisms are used to manage large numbers of digital entities. To refresh storage, content is migrated to new technology by accessing the original storage repository, reading the content, replicating the content to the new storage repository, and updating the administrative metadata for the location of the content. Audit trails track the operations on the content. Monitoring consists of tracking the status of the storage repositories, the specification of reliability against data loss, and analyzing which content is at risk. The results are recorded as administrative attributes in the authoritative catalog. Authenticity attributes identify partial completion of tasks such as synchronization across replicas. A standard data movement protocol is used to transport the data to the user after authentication and authorization against access controls maintained for each digital record. The system availability is accessed to decide which replica to retrieve.

### **A3.3 Update Digital Components**

Update Digital Components of a Record that cannot be preserved because of technological obsolescence, changes in Preservation Strategy, or similar factors. Examples of update processes include migration, standardization, and transformation to persistent form. Return the Updated Digital Components to Storage, providing Information about the Updated Digital Components to the 'Manage Information' process.

Within the Data Grid, content is transformed to new encoding formats, the updated content is stored, and the authoritative catalog is updated. Scalable management mechanisms are used to enable the update of entire collections.

### **A4 Output Electronic Record**

Produce an authentic copy of a record in response to a request, and a certificate attesting to the authenticity of the copy. Alternatively, if requested, produce a reproducible electronic record; that is, the digital component(s) of the record along with instructions for producing an authentic copy of the record and information necessary to interpret the record.

#### **A4.1 Manage the Request**

Register an incoming Request for a Record and/or Information about a Record. Translate the request into terms that can be executed in the preservation system, define Request Controls to ensure that the request is fulfilled, and report Problems with Retrieval.

Within the Data Grid, error returns are tracked when delivering output, alternate replicas are accessed when a storage system is unavailable, authentication and authorization of the requestor is performed for each record, and standard error messages are reported to requestors.

#### **A4.2 Review Retrieved Components and Information**

Determine whether all components and information necessary to satisfy a request for records have been received and can be processed for output.

Within the Data Grid, the authoritative catalog is accessed to associate components with records. The availability of the components is determined by tracking the status of each storage repository.

**A4.3 Reconstitute Record**

Apply the appropriate Targeted Preservation Method to Retrieved Digital Components to link or assemble the components as necessary to reproduce the record and output the Requested Reconstituted Record.

Within the Data Grid, the record is assembled by applying the structural metadata stored in the authoritative catalog.

**A4.4 Present Record**

Present the record with the appropriate extrinsic form, and if requested, produce a Certificate of Authenticity for the Reproduced Electronic Record.

Within the Data Grid, the record is transmitted to the requestor, along with audit trail information.

**A4.5 Package Output**

Combine Requested Digital Components with Information, including instructions on how to reproduce the record, into a package suitable for reproducing or presenting the record on an external system designated by the Requester.

Within the Data Grid, a record is reconstituted from component parts through specification of the processes that should be applied in the virtual data grid, and the result is packaged into a container for delivery.