

# Scalable, Reliable Marshalling and Organization of Distributed Large Scale Data Onto Enterprise Storage Environments\*

Joseph JaJa  
joseph@  
umiacs.umd.edu

Mike Smorul  
toaster@  
umiacs.umd.edu

Fritz McCall  
fmccall@  
umiacs.umd.edu

Yang Wang  
wpwy@  
umiacs.umd.edu

*Institute for Advanced Computer Studies  
University of Maryland, College Park*

## Abstract

*Emerging technologies in high speed NAS, hierarchical storage management systems, and networked systems that virtualize interconnected storage over IP and fiber-channel networks, promise to consolidate distributed data stores onto large-scale professionally managed enterprise storage environments. We describe the software architecture of the PAWN (Producer – Archive Workflow Network) environment that enables scalable, reliable marshalling and organization of distributed data into such enterprise storage environments. PAWN was initially developed to capture the core elements required for long term preservation of digital objects as identified by researchers in the digital library and archiving communities. In this paper, we show how PAWN can be extended to enable multiple clients at a number of distributed sites to prepare, organize, and bulk transfer large scale data onto clusters of servers that securely verify the integrity of the data, register the metadata, and store the data into an enterprise storage environment. PAWN allows detailed description, auditing, and organization of the data, and hence will allow for efficient management, access, and disaster recovery. The basic software components are based on open standards and web technologies, and hence are platform independent.*

## 1. Introduction

We are pursuing, in collaboration with the San Diego Supercomputer Center (SDSC) and the National Archives and Records Administration (NARA), a broad research program to address major components required to build a computing infrastructure for enabling long term archiving and preservation of digital assets. It is well known that

digital preservation is substantially more challenging than the traditional problem of archiving and preserving physical objects. This collaboration has resulted in the establishment of a pilot persistent archive currently consisting of three node servers at SDSC, University of Maryland, and NARA, linked through the SDSC Storage Request Broker (SRB) data grid, and managing several terabytes of significant NARA selected collections. We have outlined in [1] the software architecture of the Producer – Archive Workflow Network (PAWN), which provides secure ingestion of digital objects into the persistent archive. PAWN captures the essential features of the producer-archive interface methodology articulated in [2], which covers the first stage of the Ingest Process as defined by the Open Archival Information System (OAIS) reference model [3]. We make use of METS (Metadata Encoding and Transmission Standard) schema [4] to encapsulate content, structural, descriptive, and preservation metadata, leading to the specification of the digital object model as described in the OAIS model.

In this paper, we show how the PAWN architecture can be extended to a general-purpose scalable software that can efficiently marshal and organize large scale distributed data into emerging mass storage systems. PAWN will allow efficient staging, arrangement, and assembly of the data at various sites, followed by parallel bulk transfers into receiving servers whose loads are managed by a scheduler. The receiving servers will validate data and metadata, organize the data into collections, and register the metadata into a unified metadata database before storing the data into the enterprise storage management environment. We provide a brief description of the overall architecture and some of the software components in the next two sections. We are currently conducting experimental results to illustrate the scalability of our software using the storage systems available through the Global Land Cover Facility (GLCF) at Maryland.

\* Sponsored by the National Archives and Records Administration and the National Science Foundation under the PACI Program.







