

Toward a Data Grid for Digital Archive

Eric Yen

Computing Centre, Academia Sinica

Sep. 22nd, 2002

Toward a DataGrid for Digital Archive

- Digital Archive demands for reliable storage systems for persistent digital objects, well-organized information structure for effective content management, efficient and accurate information retrieval mechanism, and flexible services for variant users needs. Hundreds of Petabytes digital information has been created and dispersed all over the internet since computers had been used for information processing, and the amount still grows in the rate of tens of petabyte per year. Transforming the information we collected into knowledge and keep on aggregating the knowledge we have, will lead to better and sustainable development of human life. Grid technologies enlightened a possible solution for processing of diversified and heterogeneous, petabyte scale digital archives. Metadata-based information representation makes specific and relative information retrieval more accurately, makes information resources interoperable, and paves the way for formal knowledge discovery. Taking advantage of advancing IT, semantic level information indexing, categorizing, analyzing, tracking, retrieving, and correlating could be implemented. DataGrid aims to set up a computational and data-intensive grid of resources for the analysis of data, and requires coordinated resource sharing, collaborative processing and analysis of huge amounts of data produced and stored by many institutions. And so that digital archives are not becoming information islands by themselves.
- In Taiwan, a National Digital Archive Project (NDAP) was initiated in 2002 with its pilot phase started in 2001. More than 40 Terabytes digital objects will be created and archived by 9 major content holders dispersed in Taiwan at the end of 2002. Not only delicate and gracious Chinese cultural assets could be accessed thru Internet, new paradigm of academic researches based on digital and integrated information resources will be devised and implemented. In this paper, idea for utilizing DataGrid infrastructure for NDAP will be depicted and discussed.

Outline

- Introduction
- Demands from Digital Archive
- Solutions – The Data Grid
- Archiving of Geospatial Contents in Academia Sinica
- Conclusions

Introduction to Digital Archive

- Digital Archive is a collection of digital objects.
- A digital object is defined as something (*e.g.*, an image, an audio recording, a text document, a movie, a map) that has been *digitally encoded and integrated with metadata* to support discovery, use, and storage of those objects.(ref. CDL)
- Goals for Digital Archive (functional point of view)
 - ✓ Protection of the original
 - ✓ Duplication for longevity
 - ✓ Flexible Search and Retrieval
 - ✓ Easy Access
 - ✓ Resource Sharing
 - ✓ Lower cost of maintenance and dissemination
 - ✓ Max. flexibility for integration of heterogeneous/homogeneous information resources
 - ✓ Providing abundant resources for knowledge discovery and knowledge construction

Issues for digital archives

- Architectures for persistent digital repositories (technology infrastructure, component layers of services, etc.);
- Preservation technologies and tools (Web harvesting methods, techniques for normalizing heterogeneous objects, algorithms for proving authenticity, audit mechanisms, interoperability, and large scale backup and recovery mechanism);
- Attributes of archived collections (methods for determining and validating completeness and closure of preserved digital objects);
- Policy and economic models (roles and responsibilities, access services, legal agreements, and balancing private interests with public good).

Important Issues of NDAP

- Intellectual property rights
- Time, Space and Language Coordination
- Multi-lingual issues
- Public information systems
- Technical Specifications and Standards
- Interoperability
- Meta-language and Documentation
 - ✓ *Metadata*
 - ✓ *Content Markup*
 - ✓ *References and Linking*
- Dissemination and Sharing
- Cooperation and collaboration
- Scalability, Adaptability and Durability

Demands of Digital Archive₁

- Persistent digital objects,
- Well-organized information structure for effective content management
- Efficient and accurate information retrieval mechanism
- Flexible services for variant users needs
- Consistency
- Integrate relationship management with information and data management
- High-performance remote data access
- Authentication and authorization
- Resource discovery and monitoring

Demands of Digital Archive₂

- Reliable and efficient storage system
 - ✓ Reliable replication system → replica locating mechanism
 - ✓ Reduced query latency → query routing scheme
 - ✓ Load sharing
 - ✓ Robust, high availability
 - ✓ Min. Access latency
 - ✓ Manageability
 - ✓ High Throughput
 - ✓ Adaptive
 - ✓ Transparency of location and protocol
 - ✓ Decentralization/centralization

The Unprecedented Ten Years

- Networking from 100Kbps to Gbps
- Computing from 100MFLOPS to TeraFLOPS
- Storage from 100GigaBytes to PetaBytes
- We are producing 3×10^{18} Bytes of data each year
- Most business processes, research, learning, commerce, socialising, etc. may be conducted on the Internet
- Internet and Digital Technology together bring in revolutionary ways to communicate, deal with information and collaborate

Effective Management System for Huge Volume of Data

- Remote sensing data: 2TB/day; And will accumulate to 5 Peta Byte in 2005 ◦
- According to the statistics of EU Space Center
 - ✓ Raw data from satellite : 100GB/day, 500GB/day (after Feb. 2002)
 - ✓ 800 TB data had been archived
- Big Challenge of IT for cataloging, searching, retrieval, management, identification, knowledge discovery, and integration ◦
- Trading off between decentralization and consolidation on cost,
 - ✓ Convergent to multi-centers of information resources in Internet
 - ✓ Think about how to facilitate the collaboration among those centers – Community and virtual organization
- Demands for complete architecture and services → Data Grid

What's the Solution

- Support sharing and coordinated use of diverse resources in dynamic “virtual organizations” – **Grid !**
- Good technical solutions for key problems, such as
 - ✓ Security enhancement like authentication and authorization
 - ✓ Resource discovery and monitoring
 - ✓ Reliable remote service invocation
 - ✓ High-performance remote data access
 - ✓ -- **Grid !**
- Good quality reference implementation, multi-lingual support, interfaces to many systems, large user base, industrial support, etc. – **Grid !**
- Persistent Web Services – **Grid !**

Evolution of Grid Technologies

- Initial exploration (1996-1999; Globus 1.0)
 - ✓ Extensive application in experiments; core protocols
- Data Grids (1999-?; Globus 2.0+)
 - ✓ Large-scale data management and analysis
- Open Grid Service Architecture (2001-?; Globus 3.0)
 - ✓ Integration with Web services, hosting environments, resource virtualization
 - ✓ Databases, higher-level services
- Radically scalable systems (2003-?)
 - ✓ Sensors, wireless, ubiquitous computing

Open Grid Services Architecture

- Service orientation to virtualize resources
- From Web services – standard interface definition mechanisms: multiple protocol bindings, multiple implementations, local/remote transparency
- Building on Globus Toolkit:
 - ✓ Grid service: semantics for service interactions
 - ✓ Management of transient instances (& state)
 - ✓ Factory, registry, discovery other services
 - ✓ Reliable and secure transport
- Web service with specified interfaces & behaviors
- Multiple hosting targets: J2EE, .net, C, ...

Data Grid

➤ Definition

- ✓ In broad sense, means effective sharing and integration of distributed resources
- ✓ Transparent remote access to heterogeneous data resources in grid environments

➤ Approaches

- ✓ Brand new design of system architecture
- ✓ Development of Middleware
 - ❖ Resource management and workload scheduling
 - ❖ Data management of distributed resources, e.g., caching, file replication, file migration, etc.
 - ❖ Automatic system management scheme
 - ❖ API for repository and specification for import/export
- ✓ Application Development
- ✓ Construction of testing platform

➤ Resource-On-demand

➤ Toward Virtual Data

- ✓ Consistent ways of access, with transparency of physical storage

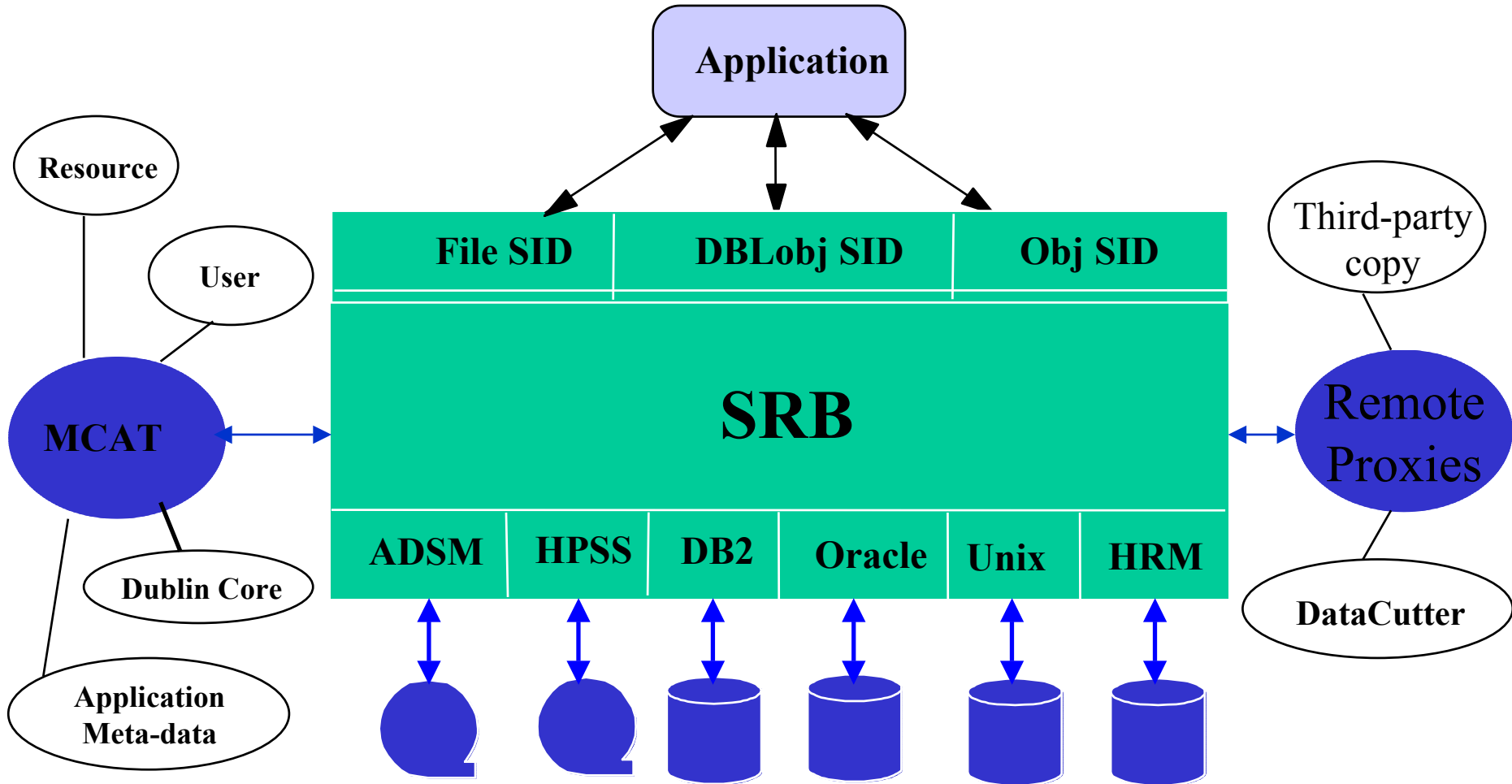
- ✓ Automatic routing for minimum latency

European Data Grid

- Characteristics for the requirements of scientific research in 21st century
 - ✓ Accumulation, storage and management of huge data
 - ✓ collaborative processing)
 - ✓ Integration and sharing of distributed resources
- Objectives : Building up an application environment for scientific research, with TeraFlops computing capability, and PetaByte scale storage management system ◦
- Leading by CERN, with 10 nations and 20 institutes , have invested EU\$ 9.8M in 2001~2003.
- Challenge : providing effective management of distributed resources based on current Internet infrastructure
- Application Fields : High Energy Physics, Processing of medical images, Survey & monitoring of Earth
- Core Services
 - ✓ job scheduling
 - ✓ data management
 - ✓ grid monitoring
 - ✓ Operation and Maintenance
- Delivery : Open Software Model

Storage Resource Broker & Metadata Catalog

提供屬性為基礎的遠端資料管理與擷取機制



Integration and Retrieval of Information Resources

- Models
 - ✓ Registry-based
 - ✓ Agent-based
- Criteria and Standard
 - ✓ UDDI
 - ✓ XML/WSDL
 - ✓ SOAP(HTTP/SMTP)
 - ✓ METS
- Domain specific-based, cross disciplinary collaboration
- Resource Broker
 - ✓ Resource Transparency
 - ✓ Location Transparency
 - ✓ Cross-Domain Authentication
 - ✓ Replicated Data Management
 - ✓ Data Redirection, Load Balancing, Fault Tolerance
 - ✓ Data Discovery
 - ✓ Uniform API, Protocol

Data Management in Digital Archive

- Persistent identifiers
 - ✓ Pertain consistent access method even physical data has been moved to other repository
- Backup of Digital Objects
 - ✓ Backup management: reliability, availability and consistency
- Backup during Work Flow Process
 - ✓ Version control
- Persistent archives
 - ✓ Completeness of an archived digital object
 - ✓ Scheduled migration and upon changes of technology/medium

WDC: CIESIN Virtual Data Center

➤ 動機與誘因

- ✓ 許多主要資料來源分散、定期更新，且願意分享(雖然初期以資料目錄為主)
- ✓ 使用者對數位化資料的接受度以及即時性需求已大幅提高
- ✓ 資料擁有單位希望大幅減低資料釋出與分享之服務成本
- ✓ 以空間為基礎之多重資料整合與地圖呈現需求大幅增加
- ✓ 資訊應用基礎架構日趨完善

➤ 執行情序

- ✓ 提昇各中心技術能力
 - ❖ 建立主題化資料目錄
 - ❖ 善用搜尋與檢索工具
 - ❖ 資料格式一致化與標準化
 - ❖ 促進國際交流與合作
- ✓ 建立鏡射站台(mirror site)
- ✓ 提昇資料擷取便利性與重複使用度
- ✓ 區分資料重要性，排定處理順序
- ✓ 整合資料發掘(discovery)與擷取

- ✓ 資料加值與整合
- ✓ 分散式資源整合

Replica Location Problem in Data Grid₁

- Given a unique logical identifier for desired data content, we need a mechanism to determine the physical locations of one or more copies of this content.
- In other words, a replica location mechanism for Data Grids might have to serve requests for many or all replicas corresponding to a give logical identifier.
- It's impossible to provide a complete consistent system view, and thus impossible to server “all replicas” requests reliably in a decentralized manner.
- One of the solutions
 - ✓ Flat overlay network of nodes: obtaining genuine decentralization and resilience when facing network and node failure
 - ✓ Probabilistic representations of replica location information: achieving space and bandwidth reduction
 - ✓ Soft-state protocols: decouple node state and achieve robustness

Replica Location Problem in Data Grid₂

➤ Functional Requirements

- ✓ Autonomy
- ✓ Best-effort consistency
- ✓ Adaptiveness

➤ Scale Requirements

- ✓ Total no. of files/replicas
- ✓ No. of storage sites and their geographical distribution
- ✓ Aggregated query and update rates

➤ Overall Benefits

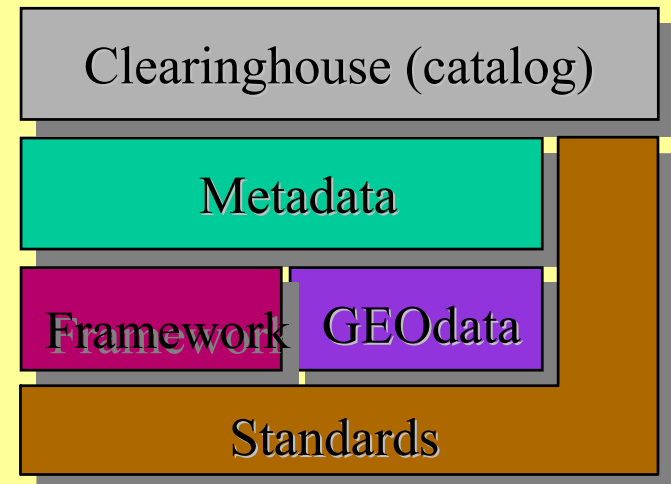
- ✓ Low query latency
- ✓ Adaptive
- ✓ Robust and high availability
- ✓ Manageability
- ✓ High throughput

Geolibrary

- Objective: Lower the barriers for applying GIScience technologies
- Approaches
 - ✓ Collecting and providing basic georeferenced spatial data/knowledge persistently
 - ✓ Building up application environment and tools for utilization of spatiotemporal knowledge and technologies
 - ✓ Development of spatiotemporal-based technologies for multi-disciplinary contents integration, aggregation, knowledge discovery in map-metaphor
- Focus & Approach
 - ✓ Construction of the System Infrastructure for Spatial and Temporal Information Technology
 - ✓ Development of Core Technology
 - ✓ Establishment of Effective Service Model for Research Support

Clearinghouse

- An instance of implementation of interoperability
- Functionality
 - ✓ Locating the required resources/services
 - ✓ Maintaining a persistent catalog of resources/services for sharing
 - ✓ Exchange of information content
 - ✓ Format transformation
 - ✓ Compilation
 - ✓ Integration



Partnerships

**Collaboration in the Digital Age, PNC 2002
Osaka City U., Japan, Sep. 20-22, 2002**

Applications in Academia Sinica

- High Energy Physics
- Analysis of Genetic Sequences
- Biodiversity and Long Term Ecological Research
- Data Center of Earthquake
- Digital Archives
- Astronomical Observation
- Geospatial Contents and Survey of the Earth

Conclusion and Future Works

- Networking Bandwidth and connectivity is one of the major reference for national strength in the 21st century
- Enhancing the capability of information infrastructure
 - ✓ Culture for collaboration and for sharing knowledge and expertise
 - ✓ Knowledge of the always-changing essence of IT, and deploy appropriate technology at the very right time
 - ✓ Collaboration, cooperation and
- Toward Knowledge Grid