# *Object Oriented Data Technology (OODT)*

April 23, 2003

Dan  Crichton
Sean Kelly
Jet Propulsion Laboratory
California Institute of Technology
National Aeronautics and Space Administration

# *Jet Propulsion Laboratory*

- NASA's lead center for robotic exploration of the solar system

- Has a dual character:

  - A unit of Caltech, staffed with Caltech employees;

  - A Federally-Funded Research and Development Center (FFRDC) under NASA sponsorship;

- Is a major national research and development (R&D) center supporting:

  - NASA programs;

  - Defense programs;

  - Civil programs of national importance compatible with JPL capabilities.

- Currently 5500 employees located in Pasadena, CA on 177 Acres

# Key Data Management Challenges of NASA Scientists and Engineers

- Search and retrieval of data sets across projects, missions and data centers

- Long term preservation of data

- Distribution of data to scientists

- Data sharing

    - Different formats, systems, access methods, etc

    - Data Policies for Data Release

- Data storage

- Automated data understanding

- Collaboration across multi-agencies

- Data Analysis and Correlation

# *Example: Difficulty Sharing Space Science Data*

- ☞ Space scientists cannot easily locate or use data across the hundreds if not thousands of autonomous, heterogeneous, and distributed data systems currently in the Space Science community.

- ☞ Heterogeneous Systems
  - ↗ Data Management - RDBMS, ODBMS, HomeGrownDBMS, BinaryFiles
  - ↗ Platforms - UNIX, LINUX, WIN3.x/9x/NT, Mac, VMS, …
  - ↗ Interfaces - Web, Windows, Command Line
  - ↗ Data Formats - HDF, CDF, NetCDF, PDS, FITS, VICR, ASCII, ...
  - ↗ Data Volume - KiloBytes to TeraBytes

- ☞ Heterogeneous Disciplines
  - ↗ Moving targets and stationary targets
  - ↗ Multiple coordinate systems
  - ↗ Multiple data object types (images, cubes,  time series,  spectrum, tables, binary, document)
  - ↗ Multiple interpretations of single object types
  - ↗ Multiple software solutions to same problem
  - ↗ Incompatible and/or missing metadata

# Evolution of Data Systems
### *(Trying to make order out of entropy)*

**JPL**

*Locally Centralized Data*

*Interoperable & Distributed Databases*

Data System Evolution

**Local Database**
- Local Tools
- No Data Sharing between Centers
- No Common Data Elements

**Limited Data Sharing**
- Manual Data Sharing
- Manual Correlation
- Export/Import Data
- Limited CDEs

**Full Data Sharing**
- Location Independence
- Data Interchange
- Data Sharing
- Common CDEs between centers
- Heterogeneous Systems

**Single Mission**

NASA Data Architecture

**Multi-Center, Multi Mission Environments**

# *Object Oriented Data Technology*

**JPL**

- Started in 1998 as a research and development task funded at JPL by the Office of Space Science to address
  - Application of Information Technology to Space Science
  - Provide an infrastructure for distributed data management
  - Research methods for interoperability, knowledge management and knowledge discovery
  - Develop software frameworks for data management to reuse software, manage risk, reduce cost and leverage IT experience
- OODT Initial focus
  - Data archiving – Manage heterogeneous data products and resources in a distributed, metadata-driven environment
  - Data location – Locate data products across multiple archives, catalogs and data systems
  - Data retrieval – Retrieve diverse data products from distributed data sources and integrate

# JPL/NIH Interagency Agreements

☞ September 2000, JPL/NIH signed an interagency agreement to explore infusion of space science data architectures and technologies into NIH research networks

  ↗ Agreement between JPL and Office of Science Policy, Office of the Director

☞ April 2002, JPL/NCI signed an interagency agreement

  ↗ Agreement to transfer technology and build a knowledge environment for data sharing across the Early Detection Research Network

# *OODT Projects*

- ***Technology Infrastructure for the Planetary Data System***

- Technology Infrastructure for the SeaWinds Earth Science Data System

- Basis for JPL Institutional Information Architecture

- Candidate framework driving standards for the International Consultative Committee of Space Science Data Systems (CCSDS)

- ***Technology Infrastructure for the NCI Early Detection Research Network (EDRN)***

- Technology Infrastructure for the Alaska State Government Denali Commission

- Future infrastructure for the Cassini Mission to Saturn

- Candidate Technology Infrastructure for a proposed Space Physics Archive System (SPASE)

- Proposed Technology Infrastructure for NASA Earth Science Data Systems

# OODT System Design Goals

- **Separate the technology and the data architecture**
- Encapsulate individual data systems to hide uniqueness
- Provide data system location independence
- Require that communication between distributed systems use metadata
- Define a standard data dictionary structure and approach for describing systems and resources
- Provide a scalable and extensible solution
- Provide a mechanism for data product exchange
- Allow systems using different data dictionaries and metadata implementations to be integrated
- Define an architecture that can leverage off of open standard approaches

# *Technology Architecture*

**JPL**

- ☞ Create intelligent middleware to capture and share data

- ☞ Implemented in Java

- ☞ Data layer implemented with the Extensible Markup Language (XML)

- ☞ Uses Java Remote Method Invocation (RMI)

- ☞ Secure Socket Layer (SSL) for data encryption

- ☞ Uses a standardized XML DTD messaging and querying language for communication

- ☞ Support a variety of client access methods

  - ↗ Java API

  - ↗ HTTP

# Middleware Data Encapsulation



Middleware can tie application, data, and user interfaces together and hide the unique interfaces

# Data Architecture

- ☞ **Use Extensible Markup Language (XML) for the data architecture**

    - ↗ Use XML metadata tags to describe data products

        - ☞ Metadata provides labels for describing data products

        - ☞ Metadata provides location information about products which can be stored remotely

    - ↗ Use XML for messaging between distributed computers

        - ☞ Standard for the exchange of information

        - ☞ A query language for locating and retrieving disparate data products

# *Metadata Development*

**JPL**

- ☞ Metadata has been identified as a critical component of capturing and sharing data
    - ↗ http://www.cio.gov/docs/metadata.htm
- ☞ Develop methods for managing the semantics of data that are shared within and between domains
    - ↗ Data Dictionary – Inventory of domain terms with definitions and other distinguishing attributes.
        - ☞ Common set of data elements used to describe information
    - ↗ XML for metadata registry and communication
- ☞ Use standards where appropriate
    - ↗ ISO/IEC 11179 – A framework for the Specification and Standardization of Data Elements
    - ↗ Dublin Core – A metadata element set intended to facilitate discovery of electronic resources.

# *OODT Component Framework*

☞ Java based software middleware component architecture that provides a software framework for archiving, search and retrieval, and data product exchange

    ↗ Archive Component – Archive Service
- ☞ Provides centralized data archiving and cataloging of data products
- ☞ Distributed

    ↗ Profile Metadata Component – Profile Service
- ☞ Manage metadata associated with resources (i.e. pointers to data products)
- ☞ Locate resources across geographically distributed data systems
- ☞ Distributed

    ↗ Data Product Exchange Component – Product Service
- ☞ Support interchange (data sharing) of data products
- ☞ Support heterogeneous implementations and systems
- ☞ Distributed

    ↗ Query Service Component – Query Service
- ☞ Ties search and product exchange services together
- ☞ Distributed

# Component Framework for OODT

**JPL**

OODT/Science Web Tools

Archive Client

**OBJECT ORIENTED DATA TECHNOLOGY FRAMEWORK**

| Archive Service | Profile Service | Product Service | Query Service | Bridge to External Services |

Navigation Service

Other Service 1

Other Service 2

Profile XML Data

Data System 1

Data System 2

# Solutions to Data Search

**JPL**

- Build metadata "profiles" that describe data system resources
  - Define using "XML"
  - Encapsulate individual data systems resources (Hide uniqueness)
  - Enable interoperability based on metadata compatibility
  - Refocus problem on metadata development
    - Communicate using metadata (Provide metadata with data)
  - Provide a core framework of software components to interconnect distributed data systems
- Define profiles using standard industry approaches
  - Use XML to describe profiles
  - ISO/IEC 11179 – A framework for the Specification and Standardization of Data Elements
  - Dublin Core – A metadata element set intended to facilitate discovery of electronic resources.

# *Profile DTD*

```
<!ELEMENT profiles
 (profile+)>

<!ELEMENT profile
 (profAttributes,
  resAttributes,
  profElement*)>

  <!ELEMENT profAttributes
   (profId, profVersion*, profTitle*, profDesc*, profType*,
    profStatusId*, profSecurityType*, profParentId*, profChildId*,
    profRegAuthority*, profRevisionNote*, profDataDictId*)>

  <!ELEMENT resAttributes
   (Identifier, Title*, Format*, Description*, Creator*, Subject*,
    Publisher*, Contributor*, Date*, Type*, Source*,
    Language*, Relation*, Coverage*, Rights*,
    resContext*, resAggregation*, resClass*, resLocation*)>

  <!ELEMENT profElement
   (elemId*, elemName, elemDesc*, elemType*, elemUnit*,
    elemEnumFlag*, (elemValue | (elemMinValue, elemMaxValue))*,
    elemSynonym*,
    elemObligation*, elemMaxOccurrence*, elemComment*)>
```

```
<profile>
 <profAttributes>
  <profId>OODT_PDS_DATA_SET_INV_82</profId>
<profDataDictId>OODT_PDS_DATA_SET_DD_V1.0</profDataDictId>
 </profAttributes>
 <resAttributes>
  <Identifier>VO1/VO2-M-VIS-5-DIM-V1.0</Identifier>
  <Title>VO1/VO2 MARS VISUAL IMAGING SUBSYSTEM DIGITAL …</Title>
  <Format>text/html</Format>
  <Language>en</Language>
  <resContext>PDS</resContext>
  <resAggregation>dataSet</resAggregation>
  <resClass>data.dataSet</resClass>
 <resLocation>http://pds.jpl.nasa.gov/cgi-bin/pdsserv.pl?…</resLocation>
 </resAttributes>
```

**JPL**

```xml
<profElement>
    <elemId>ARCHIVE_STATUS</elemId>
    <elemName>ARCHIVE_STATUS</elemName>
    <elemType>ENUMERATION</elemType>
    <elemEnumFlag>T</elemEnumFlag>
    <elemValue>ARCHIVED</elemValue>
  </profElement>
  <profElement>
   <elemId>TARGET_NAME</elemId>
   <elemName>TARGET_NAME</elemName>
   <elemType>ENUMERATION</elemType>
   <elemEnumFlag>T</elemEnumFlag>
   <elemValue>MARS</elemValue>
  </profElement>
</profile>
```

# Solutions to Data Product Exchange

- Extend framework to support common access to distributed data systems by creating a "Product Service Component"

  - Product Servers - Middleware that negotiates the interfaces between the data system implementations despite the heterogeneity

- Design the component to leverage off of consistent data architecture

- Provide data and location abstraction

- Provide a standard language for communication

# *Planetary Data System (PDS)*

- Official NASA "Active" Archive for all Planetary Data
  - Data ingestion required as part of Announcement of Opportunity (AO) for a mission
- 9 Nodes with data located at discipline sites
- Common Data Architecture
- Different data systems located at the sites
- Prior to October 2002, no ability to find and share data between PDS nodes
  - Data distribution via CD ROM
  - Limited electronic distribution

# *PDS for Mars Odyssey*



- ☞ Provide unified view across distributed science data archives

- ☞ Support online distribution of science data to scientists (up to 250 MB products)

  - ↗ Enable interoperability to distributed PDS data nodes

  - ↗ Internet as the primary means of distribution of data products

  - ↗ A unified web interface for accessing all PDS data products

  - ↗ Support real-time access to data products

- ☞ Provide a common messaging technology architecture allowing scientists and developers to link in their own tools

- ☞ Uses existing PDS databases and repositories

# Deployed PDS System



PDS-D D01 Architecture: Mars Odyssey

Science

Education

General Public

USER COMMUNITIES

Planetary Atlas | Default Browser Set Browser | Data Set View | IDL, WIPE

DISTRIBUTED CLIENTS

In:Query
Out::Data and Metadata

DITDOS    Standard Interfaces (OODT Middleware)

Data Products and Metadata

MARIE, PDS PPI

THEMIS ASU

Radio Science PDS GEO

Documents and Ancillary Files PDS CN

GRS PDS GEO

ACCEL PDS ATMOS

SPICE PDS NAIF

DISTRIBUTED DATA REPOSITORIES and CATALOGS

# PDS Nodes

Planetary Data System
Distributed Planetary Science Archive

Rings Node
Ames Research Center
Moffett Field, CA

Imaging Node
JPL and USGS
Pasadena, CA and Flagstaff, AZ

Central Node
Jet Propulsion Laboratory
Pasadena, CA

Planetary Plasma Interactions Node
University of California Los Angeles
Los Angeles, CA

Navigation Ancillary Information Node
Jet Propulsion Laboratory
Pasadena, CA

Atmospheres Node
New Mexico State University
Las Cruces, NM

Geosciences Node
Washington University
St. Louis, MO

Small Bodies Node
University of Maryland
College Park, MD

# NCI Early Detection Research Network - EDRN

- Funded by the National Cancer Institute

- Network consists of 18 Labs

  - DMCC (Fred Hutchinson)

  - Clinical Epidemiological Centers

  - Biomarker Development Labs

  - Biomarker Validation Labs

- Specimen data located at labs

- Data in validation studies

  - Captured and archived centrally

# EDRN Informatics Goals

- Develop a collaborative *knowledge* environment that
  - Provides *seamless access* to science data resources captured in EDRN studies
  - Allows investigators to *share* data using informatics tools
  - *Increases* the sample size of data resources by combining and correlating data from multiple EDRN sites
  - Provides data *standards* in the capture and exchange of critical data sets
  - *Use* existing IT infrastructures and tools located at EDRN PI sites
  - *Minimize* impact on IT systems already in place
  - Allows the IT environment to *evolve* as new data sets are available

# EDRN Informatics Key Challenges

- ☞ Data are *geographically distributed across heterogeneous* data systems making the location, retrieval and use of this data difficult
  - ↗ Data at each site is captured *differently* in
    - ☞ database systems
    - ☞ data formats
    - ☞ data definitions
  - ↗ Access to data at each site is *limited* to local tools and users

- ☞ *Different* levels of IT support and capabilities at each institution

- ☞ Data *sharing* and *privacy* issues

# EDRN Informatics Approach

- ☞ Develop a *cross-disciplinary* team of biomedical and computer science researchers

- ☞ Develop *Common Data Elements* to standardize data definitions for databases, forms, and communication

- ☞ Develop an *Informatics infrastructure* that allows for data located in disparate databases to exchange information

  - ↗ Leverage JPL/NASA's *experience* and software in developing IT infrastructures to support planetary science

  - ↗ Use *existing EDRN databases* without requiring changes (i.e. software handles translation between local database and EDRN)

  - ↗ Deploy common software at EDRN sites

  - ↗ Develop a *common IRB protocol*

- ☞ Develop a common science *portal* to provide a single point of entry to EDRN data resources

# *Benefits of Informatics Infrastructure*

☞ *Seamless* search and retrieval of data products
  - ↗ Users can access EDRN resources without knowing their location ("*one stop shopping"*)
  - ↗ *Integration* of EDRN Sites (one integrated system!)
  - ↗ Support *heterogeneous* data repositories
  - ↗ Support *geographically distributed* data repositories

☞ *Standard interfaces* for software developers to develop new bioinformatics tools

☞ Provide a *translation layer* between EDRN and the local institution's database

☞ *Plug-ins* for preferred tools (i.e. SAS)

☞ EDRN can *evolve* as basic information technology changes

# EDRN Informatics Tools

) **EDRN Secure Website**

- ↗ A *unified portal* allowing PIs to access shared information
- ↗ Restricted to EDRN *registered* users
- ↗ Uses the Internet as the primary means of access to the data
- ↗ A *collaborative website* for sharing of information among PIs

) **EDRN Resource Network Exchange (ERNE)**

- ↗ An *infrastructure* for sharing data resources across EDRN
- ↗ Supports *real time* (on demand) *distribution* of data to users
- ↗ First release - *Specimen sharing tool*

) **EDRN CDE Mapping Tool**

- ↗ Allow EDRN sites to *map local data definitions* to Common Data Elements (CDEs)

# EDRN Resource Network Exchange Tool



National Data Sharing Infrastructure
Supporting Collaboration In Biomedical Research For EDRN

Fred Hutchinson
Cancer Research Center, Seattle
(DMCC)

University of
Pittsburgh
(CEC)

Creighton
University
(CEC)

University
of Michigan
(CEC)

University of
Colorado
(CEC)

UT Health Science
Center, San Antonio
(CEC)

Moffitt Cancer
Center, Tampa
(BDL)

**As of September 2002**

# Rollout of EDRN Informatics Infrastructure



EDRN Resource Network Exchange

# EDRN Query Scenario

☞ Find *DNA blood specimens* for participants younger than 70 years old that have cancer

☞ Possible constraints

  ↗ Cancer Site

  ↗ Storage Mechanism

  ↗ Smoking

  ↗ Age

  ↗ Ethnicity

# ERNE Search Tool

**JPL**

- ☞ Connects to distributed databases

- ☞ Reports all available sites

- ☞ Allows user to select specific or all sites

# EDRN Query Example

- Bio-specimen search

- Based on CDEs

- Real time access to EDRN data

- Search performed locally at each institution

# EDRN Query Results

- Results from all applicable sites based on query

- Summary information of samples from each site

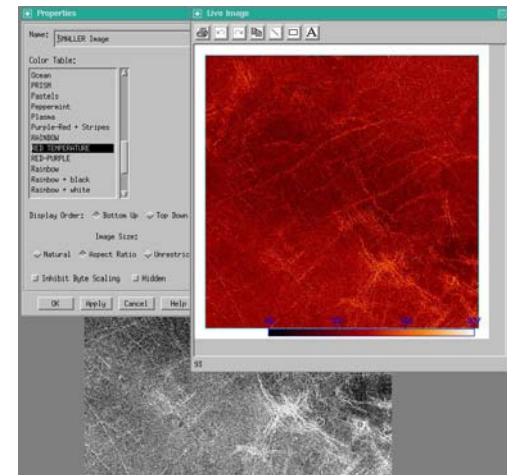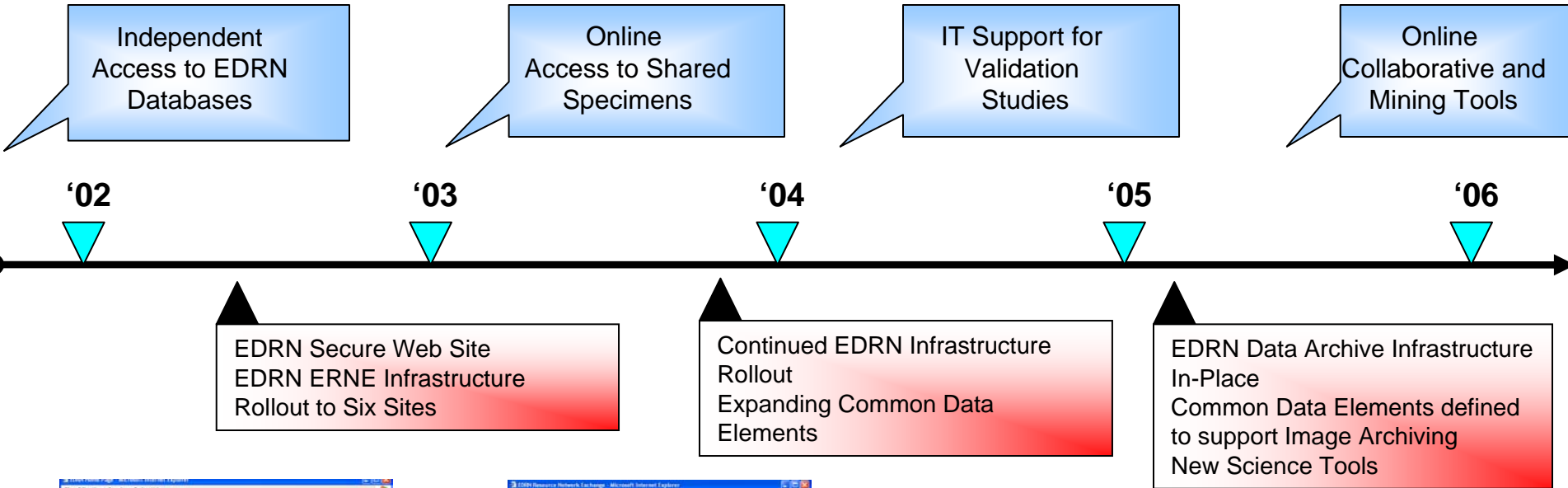- Ability to drill down through results

# Available Specimens

| Site | Specimen | Cancer Type |
|------|----------|-------------|
| Moffitt | Blood, Bone marrow, Sputum, Tissue | Various/Lung |
| San Antonio | Blood | Prostate/various |
| Creighton | Blood, Tissue | Various |
| GLNE | Blood, Tissue, Urine | Colon/various |
| Colorado | Blood, Sputum, Tissue, Urine | Various |
| Pittsburgh | Blood | Various |

# *EDRN Informatics Timeline*

**JPL**

| Independent Access to EDRN Databases | Online Access to Shared Specimens | IT Support for Validation Studies | Online Collaborative and Mining Tools |

**'02**     **'03**     **'04**     **'05**     **'06**

EDRN Secure Web Site
EDRN ERNE Infrastructure
Rollout to Six Sites

Continued EDRN Infrastructure Rollout
Expanding Common Data Elements

EDRN Data Archive Infrastructure In-Place
Common Data Elements defined to support Image Archiving
New Science Tools

# Key Accomplishments

- ☞ Deployed *science tools*

- ☞ *Multi-agency, multi-discipline* working groups and collaborations

- ☞ *National and International* Presentations and Publications

- ☞ Science-driven solutions benefiting both *cancer* and *planetary science* research

- ☞ *Seamless access* between seven EDRN research sites (including the DMCC)

# *Informatics Working Group Members*

**JPL**

- Data Management and Coordinating Center, Fred Hutchinson Cancer Research Center

- H. Lee Moffitt Cancer Center

- University of Texas, San Antonio

- Creighton University

- University of Colorado

- University of Pittsburgh

- University of Michigan/Dartmouth University (Great Lakes New England Consortium)

- Brigham and Womens Hospital

- New York University

- MD Anderson, University of Texas

- Cancer Biomarkers Group, NCI

- NASA Jet Propulsion Laboratory

# *More Information and References*

- Information about the JPL OODT Project (http://oodt.jpl.nasa.gov)
- Information about the Planetary Data System (http://pds.jpl.nasa.gov)
- Information about the Early Detection Research Network (http://edrn.nci.nih.gov)
- Dublin Core (http://purl.oclc.org/dc)
- Extensible Markup Language (http://www.w3c.org/XML)