# From Data Curation to Data Analytics: A Big Data Experiment in NASA Earth Science

NASA

Science Mission Directorate

**Tsengdar J. Lee, Ph.D.**
**Weather Data Analysis Program Manager**
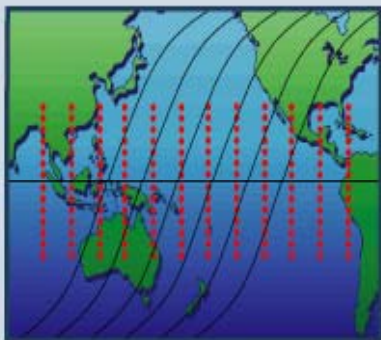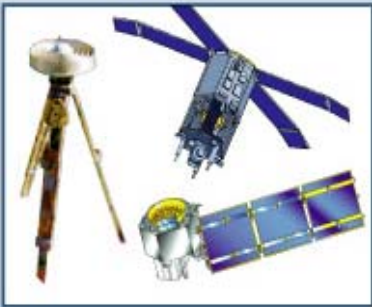**High-End Computing Manager**

# Turning Observations into Knowledge Products

**Downlink Speed**

## Petabytes $IO^{15}$

Multi-platform, multiparameter, high spatial and temporal resolution, remote & in-situ sensing
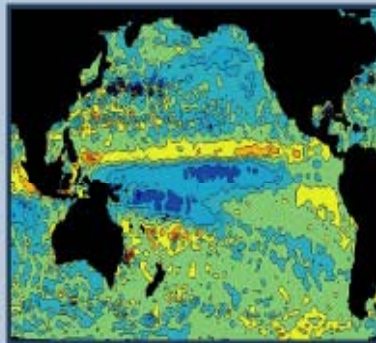
**Advanced Sensors**

## Terabytes $IO^{12}$

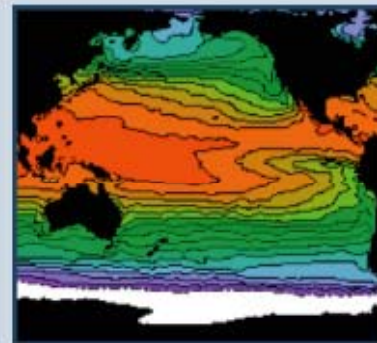Calibration, Transformation To Characterized Geophysical Parameters

**Data Processing & Analysis**

## Gigabytes $IO^{9}$

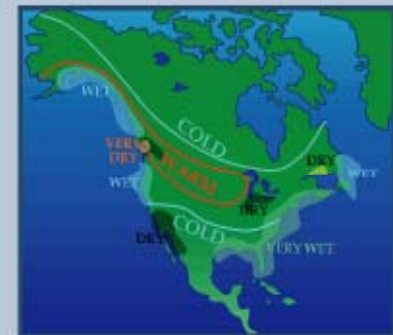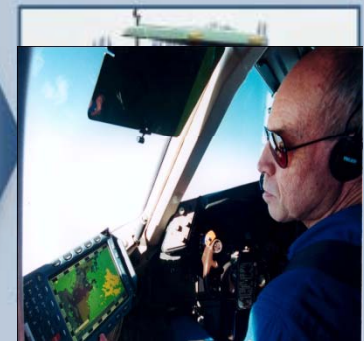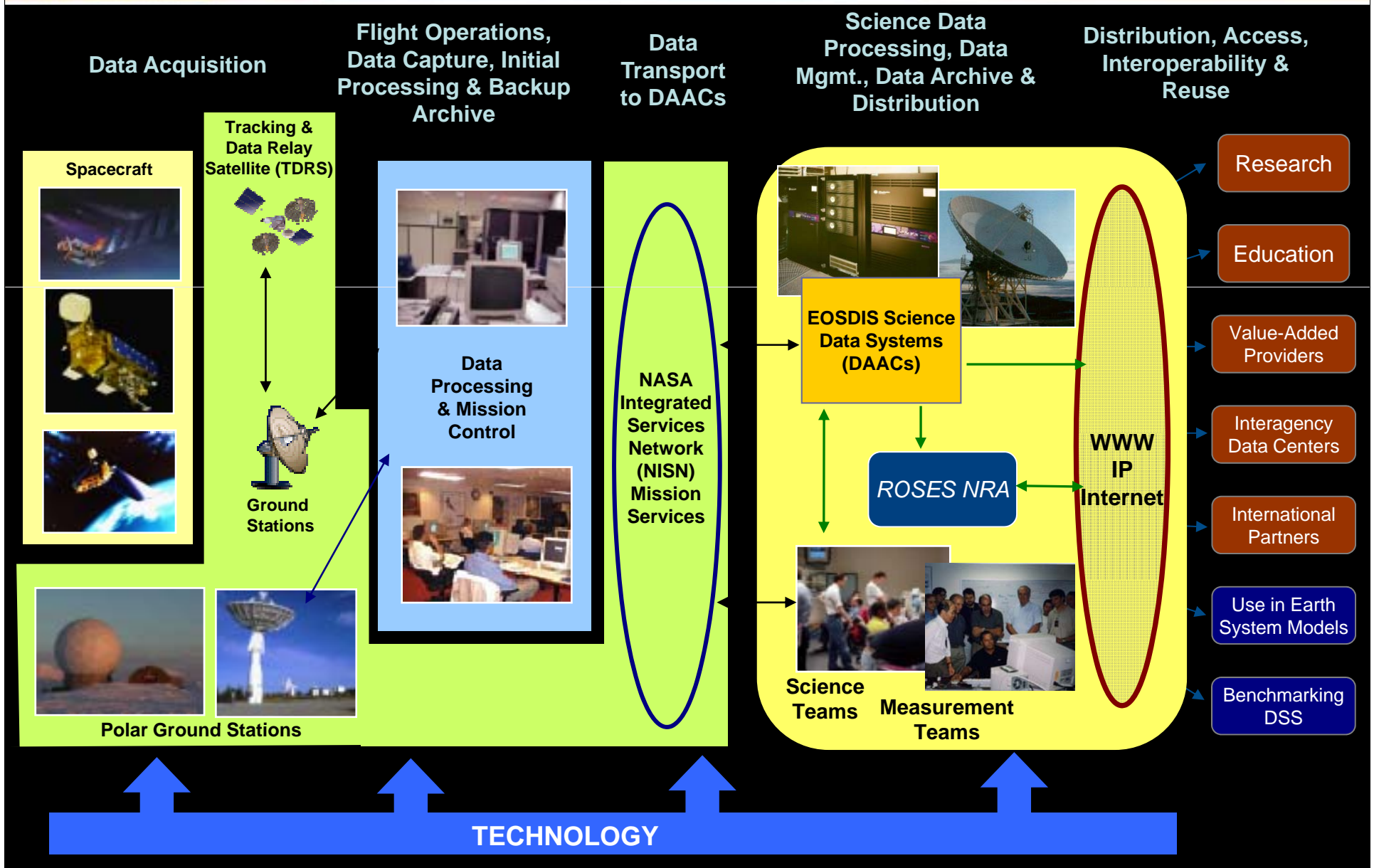Interaction Between Modeling/Forecasting and Observation Systems

**Information Synthesis**

## Megabytes $IO^{6}$

Interactive Dissemination and Predictions

**Access to Knowledge**

# Data Acquisition to Data Access

**Data Acquisition**

**Flight Operations, Data Capture, Initial Processing & Backup Archive**

**Data Transport to DAACs**

**Science Data Processing, Data Mgmt., Data Archive & Distribution**

**Distribution, Access, Interoperability & Reuse**

Spacecraft

Tracking & Data Relay Satellite (TDRS)

Ground Stations

Polar Ground Stations

Data Processing & Mission Control

NASA Integrated Services Network (NISN) Mission Services

EOSDIS Science Data Systems (DAACs)

ROSES NRA

Science Teams   Measurement Teams

WWW IP Internet

Research

Education

Value-Added Providers

Interagency Data Centers

International Partners

Use in Earth System Models

Benchmarking DSS

**TECHNOLOGY**

- Ian Foster talked about data processing workflow all the way to data distribution but why stop there?



Small science is struggling

More data, more complex data
Ad-hoc solutions
Inadequate software, hardware
Data plan mandates

- Michael Stonebraker talked about moving compute to data but there is a tremendous challenge. *(I can't even do it at one NASA center).*

# NEX Precursors and Development History

**1999** — TOPS

**2000-02** — IDU (AERO) *(Planning and Scheduling, SOA)*

**2003-07** — REASoN *(Automated data acquisition, first external user = NPS, NEED: Web Interface, Data summarization)*

**2008-09**
- R & A *(NEED: Flexible analysis capabilities)*
- Applied Sci *(NOAA, EPA, 20+ partners)*
- ACCESS *(OGC, Web-based maps)*

**2009-11**
- AIST *(Anomaly Detection and Analysis Framework)*
- ARRA/vTOPS/Water *(Social Network, Big Data, Technology prototyping)*

**2012-15**
- AIST (1) *(Workflow Capabilities)*
- NASA Earth Exchange NEX *(Community, Supercomputing, Collaboration)*
- R & A (11) CMS, NCA, LCLUC,CC,TE, IDS
- ACCESS (1) *(Data and tool management)*
- HECC (2) *(NEX Infrastructure)*
- Applied Sci. (9) *(ECOFOR, WATER)*

Investments

# NEX: the Big Data experiment and work completed to date

**Access to community/knowledge (240 members)**

**Access to ready-to-use data (24 products, 350TB)**



**Access to models/analysis tools**

Climate, weather, carbon, hydrology

**Access to workflows
to build upon
(VisTrails)**

**Ready for sophisticated users, needs proper integration for the rest…..**

**EARTH SCIENCE COMMUNITY USERS**

- **Access to Portal**

**NEX SCIENCE USERS**

- **Access to compute resources**

**NEX HPC USERS**

- **Access to NAS supercomputing resources**

## Portal

*Point of entry to NEX collaborative environment*

- **Project Information**
- **Collaboration and Social Networking**
- **Document Publication**
- **Resource Requests**
- **Data Discovery**

*Runs on NAS web servers*

## Sandbox

*Virtualized NEX compute environment*

### Domain Platform

- **Workflow Management**
- **Provenance**
- **Rich semantic search**
- **Data/Model/Tool access (API)**

### Infrastructure Platform

- **Virtualization Support**
- **Model and Analytic Tool Execution**

### Data Management

- **Data acquisition and pre-processing**
- **Data storage**

*Runs on dedicated NEX servers and storage*

## NAS HPC

*Environment for Computing at scale*

- **Execution of Jobs migrated from Sandbox**
- **Storage of results in NEX Data Management environment**

*Runs on NAS supercomputers and storage*

# Component Architecture

Search capabilities

Who is doing what where
Science network through
abstracts and papers
Who's the expert
Workflows
Archived seminars

Reporting capabilities

Annual reports
Highlights
Publications
Spatial distribution of funding

- Following the lead from Astrobiology and NASA Lunar Science Institutes to create a virtual institute, and offer
  - Summer short courses
  - Seminars
  - Conferences
  - Presentations
  - Have each funded project do one or two seminars that can be archived

# Benefits

## Science 2007

## GRL 2010

**Would have taken less than a week in NEX, Instead of one year it took to repeat the analysis**

*Efficient use of resources*
*Lowering barriers to entry*
*Interdisciplinary work*
*Transparency, repeatability, extensibility*
*Cost reduction*

Travel
Hardware
IT personnel
Data acquisition
Network costs

**Global Land Survey from Landsat (LCLUC)**

**30 years of change analysis (CC)**



1984-present

**Global monthly Landsat (WELD, MEASURES)**



Global WELD, Roy & Ju

**Biophysical products from Landsat (CMS)**



LAI

AGB, t/ha

**Crop water management with Landsat**



TOPS Satellite Irrigation Management Support

**WorldView-2, 50cm, 8 bands (NGA)**

*Space-time-resolution metric*
*(higher spatial resolution, larger extent,*
*longer time-series studies)*



Global WELD, Roy & Ju

# NEX Experiments Demonstrate Value of Collaborative Environment

- In a first application of NEX, a research team from around the U.S. used the environment to adjoin and atmospherically correct a mosaic of 9,000 Landsat Thematic Mapper scenes and retrieve global vegetation density at a 30-meter resolution.

- The entire processing of the nearly 340 billion pixels in the composite took just a few hours on the Pleiades supercomputer, allowing the team to experiment with new algorithms and products within just a few days.



- NEX's collaboration and knowledge-sharing platform for the Earth science community combines supercomputing, Earth system modeling, workflow management, and NASA remote sensing data feeds to deliver a complete work environment for users to explore/analyze large datasets, run modeling codes, collaborate, and share results.
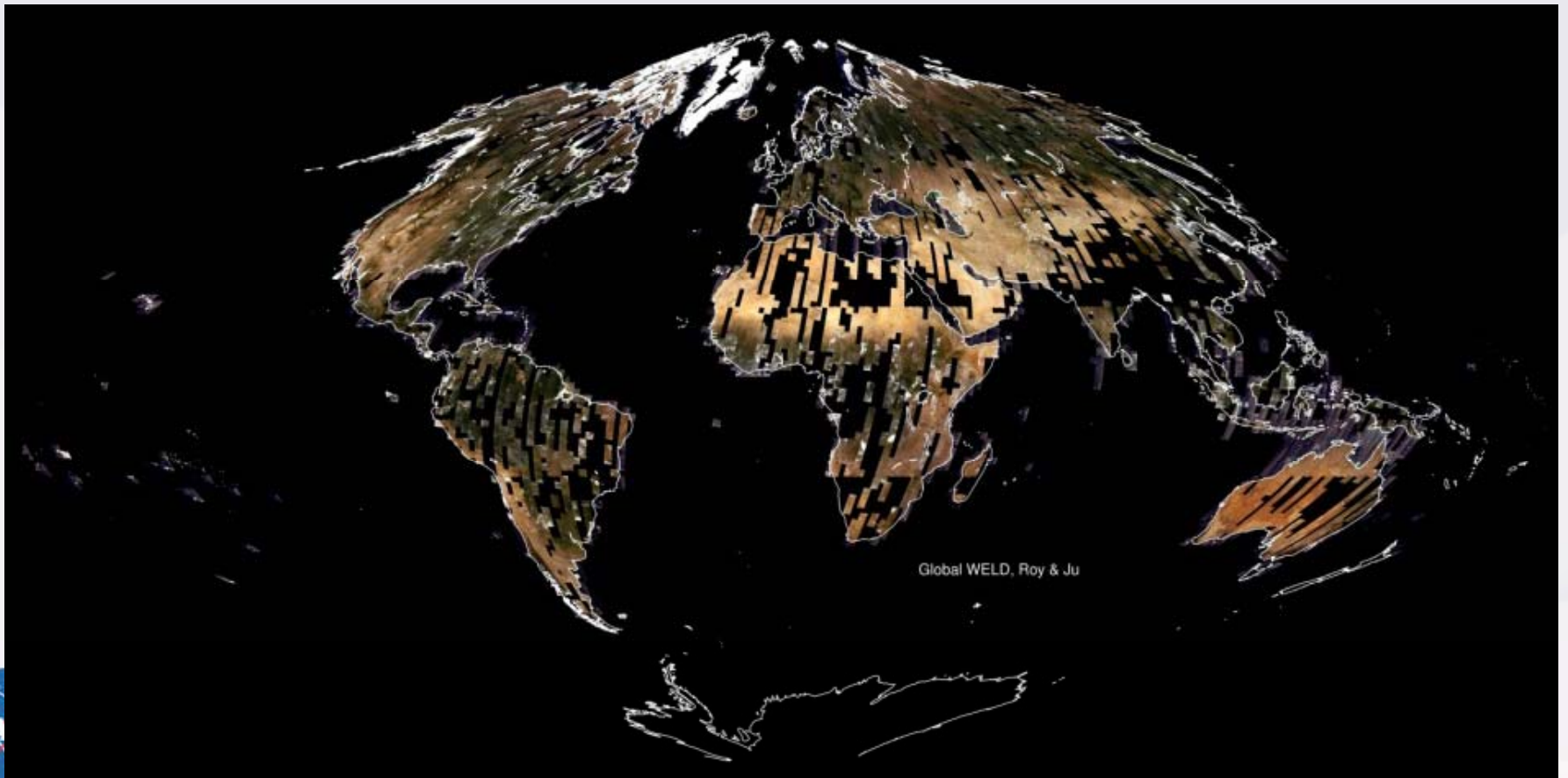
*POC: Petr Votava, petr.votava@nasa.gov, (650) 604-4675; Ramakrishna Nemani, ramakrishna.nemani@nasa.gov, (650) 604-6185, NASA Ames Research Center*

- Data Management
  - How to load and offload data in a multi-tier storage environment?
  - Distributed Multi-site Analytics?
- Workflow Reuse
  - Why create your own if a similar analytics was done before?
  - Use workflow to capture the data provenance?
- Workflow Discovery
  - Why not mine the literature to uncover the workflow?