**National Aeronautics and
Space Administration**

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

# Analysis Optimized Data Storage in Apache Science Data Analytics Platform

**Thomas Huang**

thomas.huang@jpl.nasa.gov

Group Supervisor - Computer Science for Data-Intensive Applications

Strategic Lead - Interactive Data Analytics
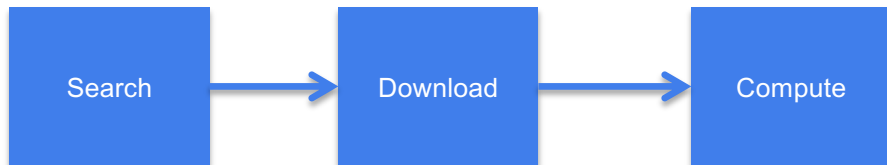
Jet Propulsion Laboratory

California Institute of Technology

4800 Oak Grove Drive, Pasadena, CA 91109-8099, U.S.A.

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

NASA

# Background

- **NASA has historically focused on systematic capture and stewardship of data for observational Systems**
- **With large amount of observational and modeling data, finding and downloading is becoming inefficient**
- **Reality with large amount of observational and modeling data**
  - Downloading to local machine is becoming inefficient
  - Search has gotten a lot faster.  Too many matches.
  - Finding the relevant measurement has becoming a very time consuming process "*Which SST dataset I should use?*"
  - Analyze decades of regional measurement is labor-intensive and costly
- **Increasing "big data" era is driving needs to**
  - Scale computational and data infrastructures
  - Support new methods for deriving scientific inferences
  - Shift towards integrated data analytics
  - Apply computational and data science across the lifecycle
- **Scalable Data Management**
  - Capture well-architected and curated data repositories based on well-defined data/information architectures
  - Architecting automated pipelines for data capture
- **Scalable Data Analytics**
  - Access and integration of highly distributed, heterogeneous data
  - Novel statistical approaches for data integration and fusion
  - Computation applied at the data sources
  - Algorithms for identifying and extracting interesting features and patterns
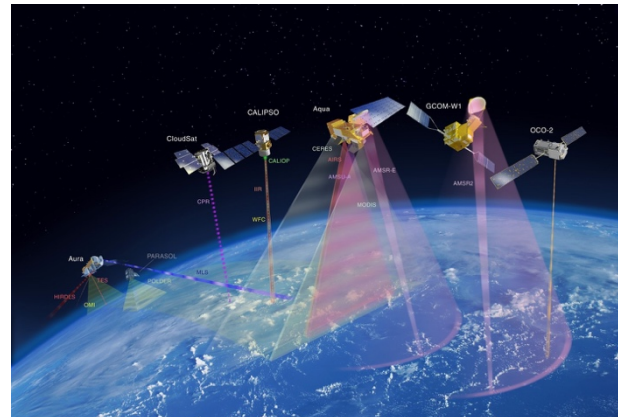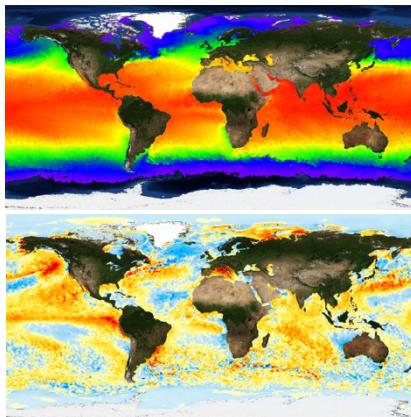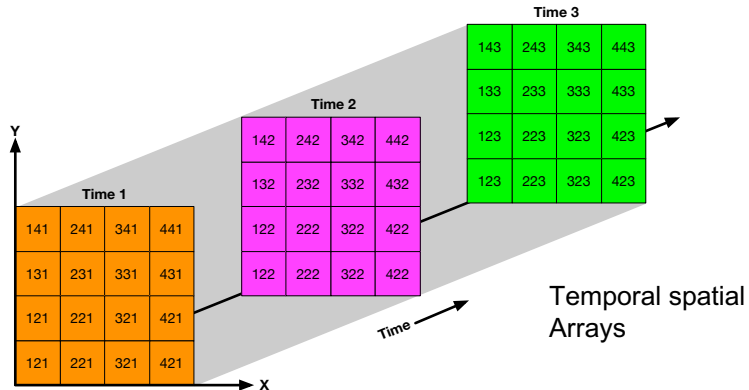
# Data Centers

- **Mainly focus on archives and distributions**
- **With additional services**
  - Better searches – faceted, spatial, keyword, ranking, etc.
  - Data subsetting – home grown, OPeNDAP, etc.
  - Visualization – visual discovery, PO.DAAC's SOTO, NASA Worldview, etc.
- **Limitations**
  - Little to no interoperability between tools and services: metadata standard, keyword, spatial coverage (0-360 or -180..180), temporal representation, etc.
  - Making sure the most relevant measurements return first
  - Visualization is nice, but it doesn't provide enough information about the event/phenomenon captured in the image.
  - With large amount of observational data, data centers need to do more than just storing bits
    - "Is the red blob in the middle of Pacific normal this time of the year?"
    - "Any relevant news and publications relate to what I am looking at?"
    - "What other measurements, phenomena, news, publications relate to the period and location I am looking at?"
    - "I can see the observation from satellite, are there any relevant in situ data I can look at?"

# Traditional Method for Analyze Satellite Measurements

Search → Download → Compute

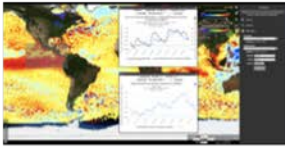

Temporal spatial Arrays

## Observation

- Traditional methods for data analysis (time-series, distribution, climatology generation) can't scale to handle large volume, high-resolution data. They perform poorly

- Performance suffers when involve large files and/or large collection of files

- A high-performance data analysis solution must be free from file I/O bottleneck

- Depending on the data volume (size and number of files)

- It could take many hours of download – (e.g. 10yr of observational data could yield thousands of files)

- It could take many hours of computation

- It requires expensive local computing resource (CPU + RAM + Storage)

- After result is produced, purge downloaded files
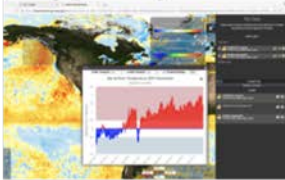
# Enabling Next Generation of Earth Science Tools and Services
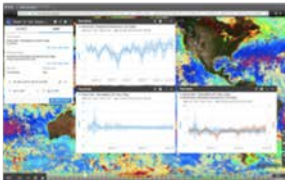
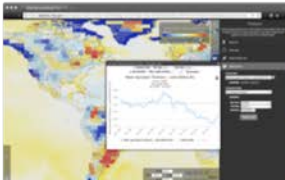

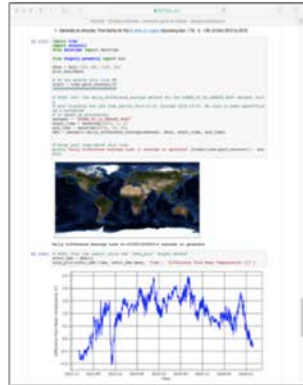NASA Sea Level Change Portal

Oceanographic Anomaly Detection
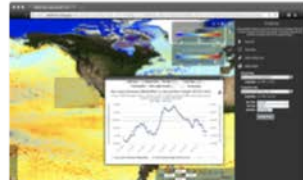
PO.DAAC State Of The Ocean

Hydrological Basin Analysis

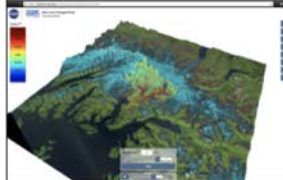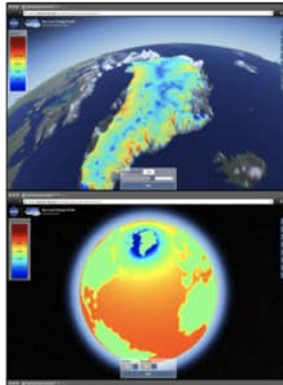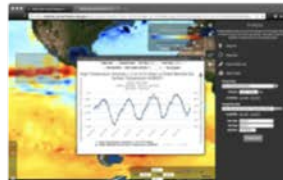Jupyter Notebook - Interactive Workbench
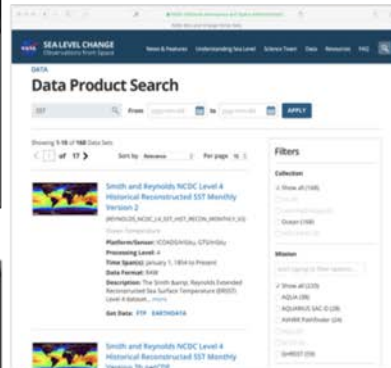
Mobile Analysis

In Situ Data Analysis

Model Simulations
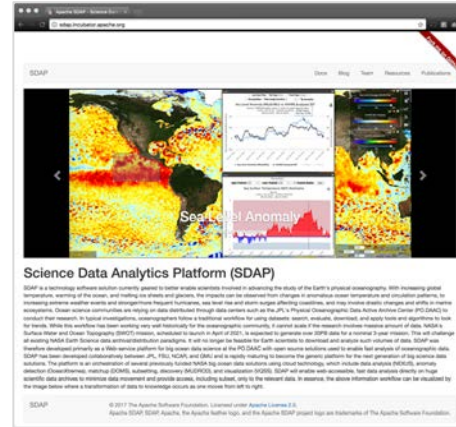
Model - Observation Comparison
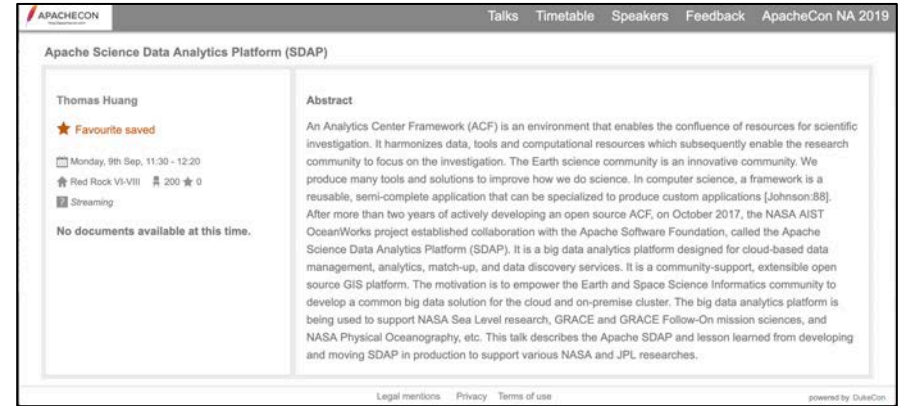
Integrated Search and Discovery

# Managing Open Source

- After more than two years of active development, on October 2017 the **NASA ESOT/AIST OceanWorks** team established Apache Software Foundation and established the **Science Data Analytics Platform (SDAP)** in the **Apache Incubator**
- Technology sharing through Free and Open Source Software (FOSS)
- Why? Further technology evolution that is restricted by projects / missions
- It is more than GitHub
    - Quarterly reporting
    - Reports are open for community review by over 6000 committers
    - SDAP has a group of appointed international mentors
- **SDAP and many of its affiliated projects are now being developed in the open**
    - Support local cluster and cloud computing platform support
    - Fully containerized using Docker and Kubernetes
    - Infrastructure orchestration using Amazon CloudFormation
    - Satellite and model data analysis: time series, correlation map,
    - In situ data analysis and collocation with satellite measurements
    - Fast data subsetting
    - Upload and execute custom parallel analytic algorithms
    - Data services integration architecture
    - OpenSearch and dynamic metadata translation
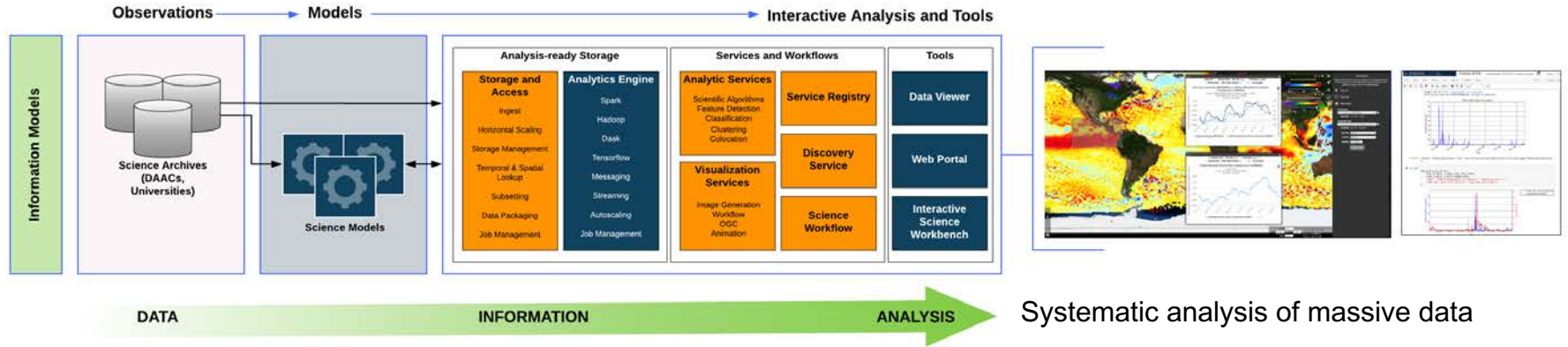    - Mining of user interaction and data to enable discovery and recommendations



http://sdap.apache.org

# Integrated Science Data Analytics Platform
## SaaS and PaaS for Science Tools and Services



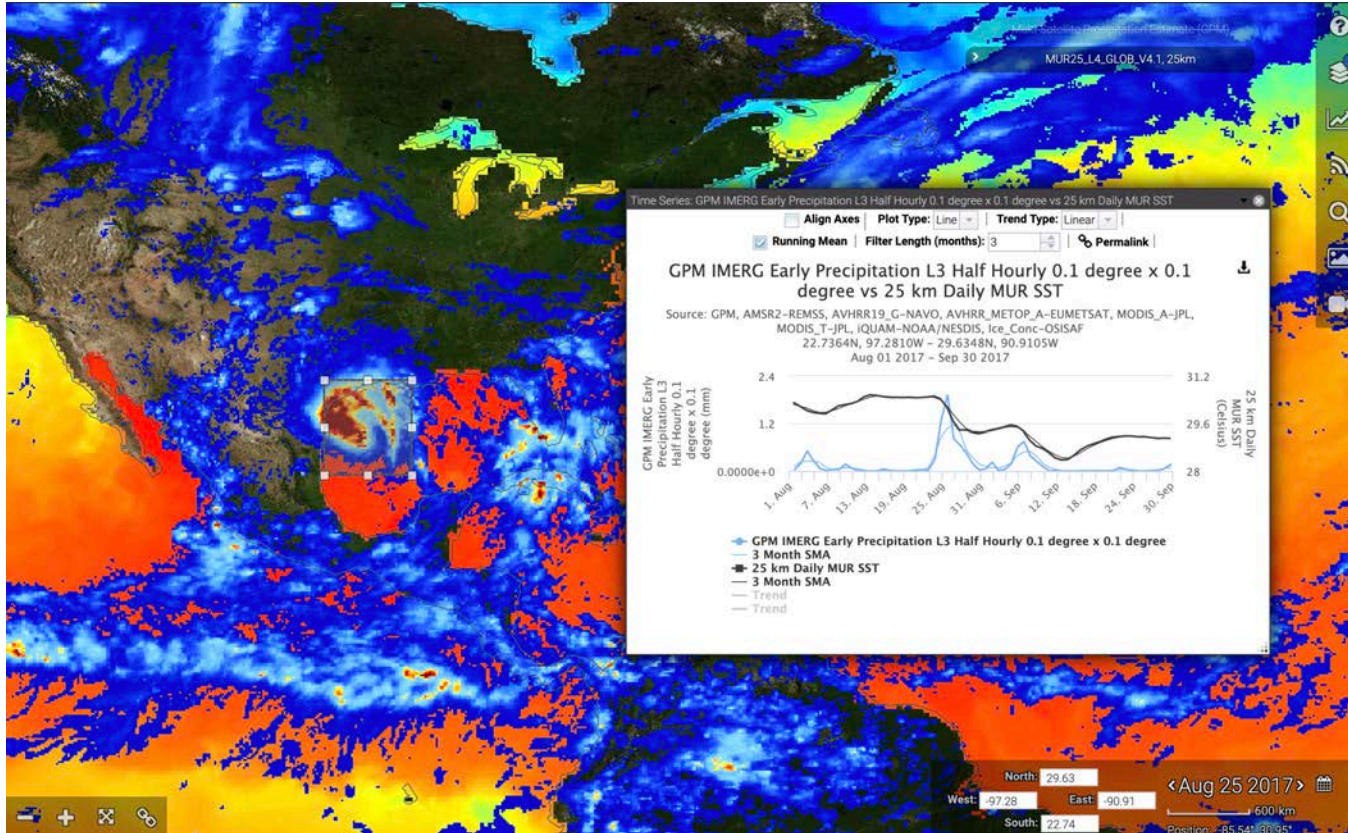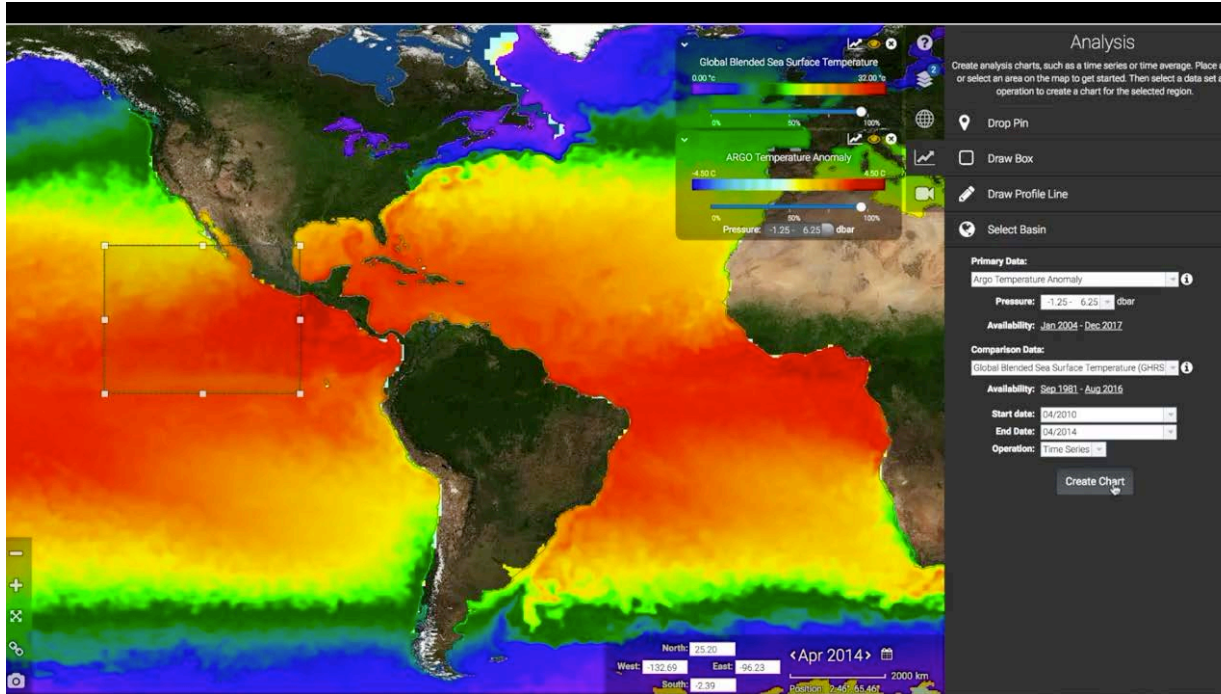Systematic analysis of massive data

- **An effort funded by the NASA's Advanced Information Systems Technology (AIST) program**
- **Integrated Science Data Analytics Platform**: an analytic center framework to provide an environment for conducting a science investigation
  - Enables the confluence of resources for that investigation
  - Tailored to the individual study area (physical ocean, sea level, etc.)
- Harmonizes data, tools and computational resources to permit the research community to focus on the investigation
- Scale computational and data infrastructures
- Shift towards integrated data analytics
- Algorithms for identifying and extracting interesting features and patterns

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
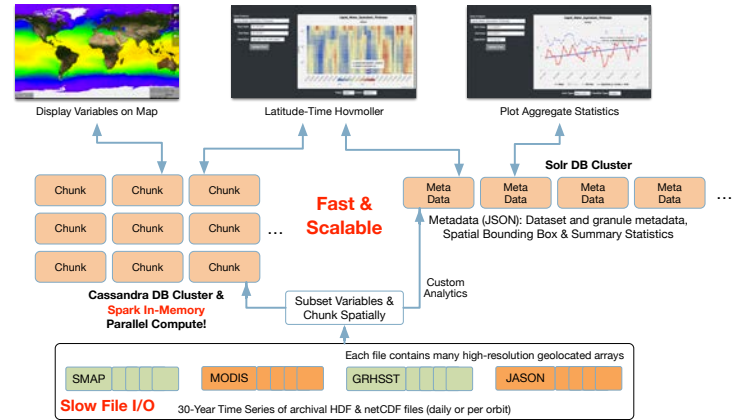California Institute of Technology
Pasadena, California



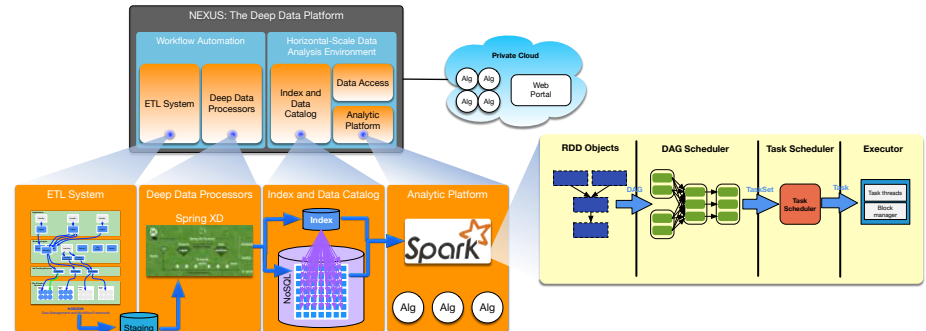Analyze *in situ* and satellite observations

Analyze Sea Level
on mobiles

# Scalable Data Analytic Solution

- SDAP's analytics engine (a.k.a NEXUS) is a data-intensive analysis solution using a new approach for handling science data to enable large-scale data analysis
- Streaming architecture for horizontal scale data ingestion
- Scales horizontally to handle massive amount of data in parallel
- Provides high-performance geospatial and indexed search solution
- Provides tiled data storage architecture to eliminate file I/O overhead
- A growing collection of science analysis webservices using Apache Spark: parallel compute, in-memory map-reduce framework
- Pre-Chunk and Summarize Key Variables
  - Easy statistics instantly (milliseconds)
  - Harder statistics on-demand using Spark (in seconds)
  - Visualize original data (layers) on a map quickly (Cassandra store)
- **Algorithms** – Time Series | Latitude/Time Hovmöller| Longitude/Time Hovmöller| Latitude/Longitude Time Average | Area Averaged Time Series | Time Averaged Map | Climatological Map | Correlation Map | Daily Difference Average



Two-Database Architecture

**National Aeronautics and Space Administration**
**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

**Dataset**: MODIS AQUA Daily
**Name**: Aerosol Optical Depth 550 nm (Dark Target) (MYD08_D3v6)
**File Count**: 5106
**Volume**: 2.6GB
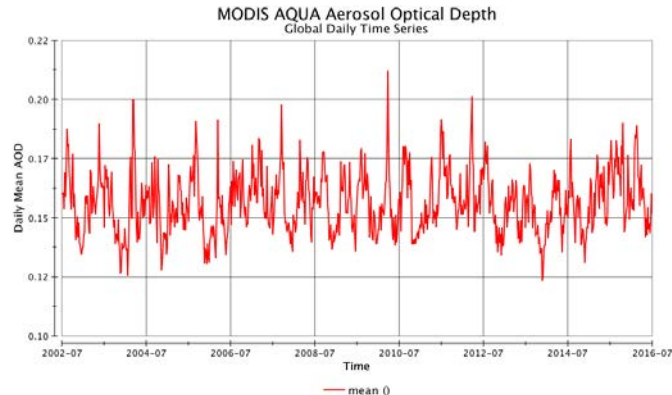**Time Coverage**: July 4, 2002 – July 3, 2016

**File-based**: A web-based application for visualize, analyze, and access vast amounts of Earth science remote sensing data without having to download the data.

- Represents current state of data analysis technology, by processing one file at a time
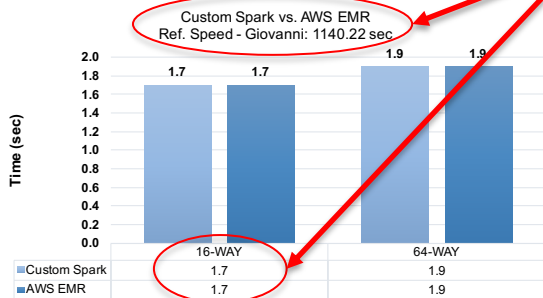- Backed by the popular NCO library. Highly optimized C/C++ library

**AWS EMR**: Amazon's provisioned MapReduce cluster



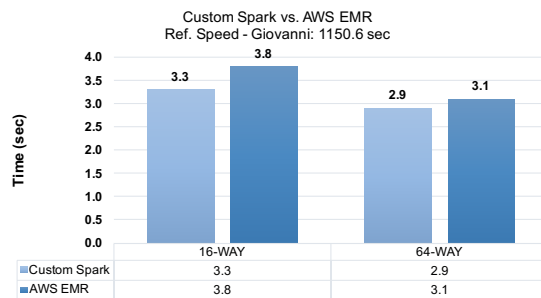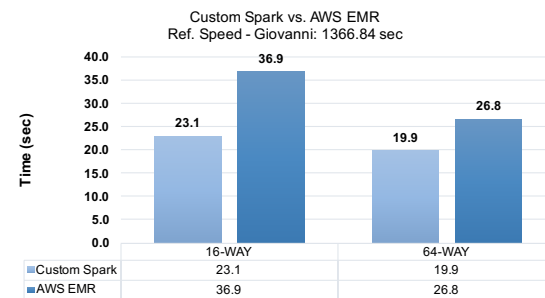MODIS AQUA Aerosol Optical Depth
Global Daily Time Series

**File-based: 20 min**
**NEXUS: 1.7 sec**

**Area Averaged Time Series on AWS - Boulder**
July 4, 2002 - July 3, 2016
NEXUS Performance

Custom Spark vs. AWS EMR
Ref. Speed - Giovanni: 1140.22 sec

| | 16-WAY | 64-WAY |
|---|---|---|
| Custom Spark | 1.7 | 1.9 |
| AWS EMR | 1.7 | 1.9 |

**Area Averaged Time Series on AWS - Colorado**
July 4, 2002 - July 3, 2016
NEXUS Performance

Custom Spark vs. AWS EMR
Ref. Speed - Giovanni: 1150.6 sec

| | 16-WAY | 64-WAY |
|---|---|---|
| Custom Spark | 3.3 | 2.9 |
| AWS EMR | 3.8 | 3.1 |

**Area Averaged Time Series on AWS - Global**
July 4, 2002 - July 3, 2016
NEXUS Performance

Custom Spark vs. AWS EMR
Ref. Speed - Giovanni: 1366.84 sec

| | 16-WAY | 64-WAY |
|---|---|---|
| Custom Spark | 23.1 | 19.9 |
| AWS EMR | 36.9 | 26.8 |

Algorithm execution time

# Analyze Large Collection of Observational Data Directly … across the ocean



**Retrieval of a single river time series**

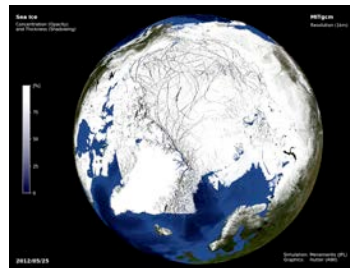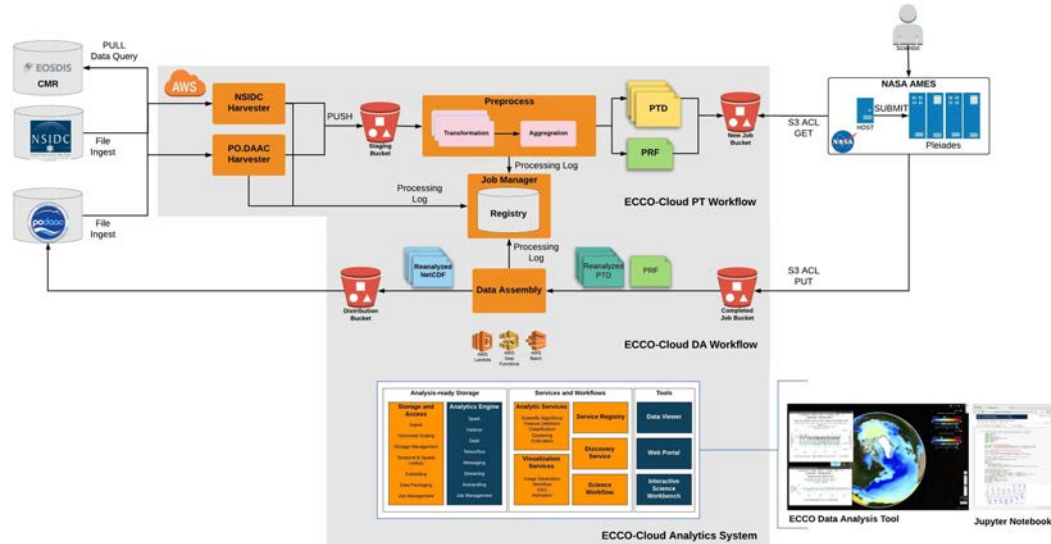**Retrieval of time series from 9 rivers**

**Time series coordination between TRMM and river**

- Running Jupyter from Germany and interacts with analytics services hosted on Amazon and at JPL
- Simulated hydrology data in preparation for SWOT hydrology
- **River data**: ~3.6 billion data points. 3-hour sample rate. Consists of measurements from ~600,000 rivers
- **TRMM data**: 17 years, .25deg, 1.5 billion data points
- Sub-second retrieval of river measurements
- On-the-fly computation of time series and generate coordination plot

# Data Access and the ECCO Ocean and Ice State Estimate
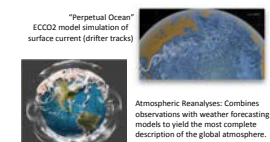## NASA ACCESS Program | PI: Patrick Heimbach; Co-Is: Ian Fenty and Thomas Huang

- **Estimating the Circulation and Climate of the Ocean** (ECCO) is a consortium endeavors to produce the best possible estimates of ocean circulation and its role in climate

- Combining state-of-the-art ocean circulation models with global ocean and sea-ice data in a physically and statistically consistent manner

- ECCO products are being used in studies on ocean variability, biological cycles, coastal physics, water cycle, ocean-cryosphere interactions, and geodesy

- **Goals**
  - Expand and accelerate in a sustainable and scalable manner the integration of NASA Earth system data into ECCO through automated preprocessing and transformation
  - Automate generation of ECCO reanalysis products into CF-compliant NetCDF products
  - Radically streamline the integration of updated ECCO products into NASA's Earth Observing System Data and Information System (EOSDIS)
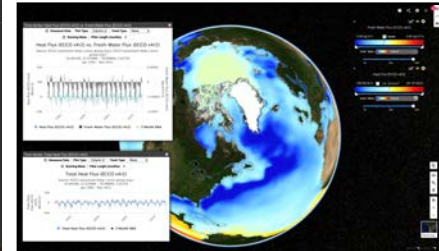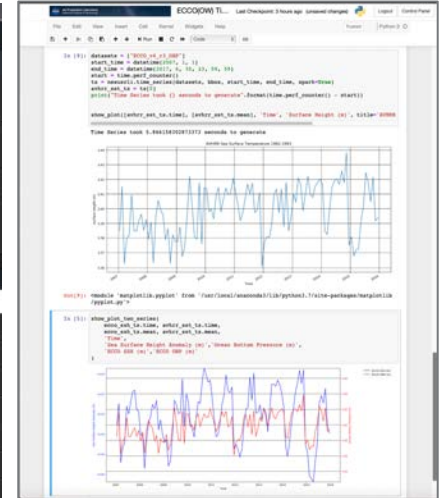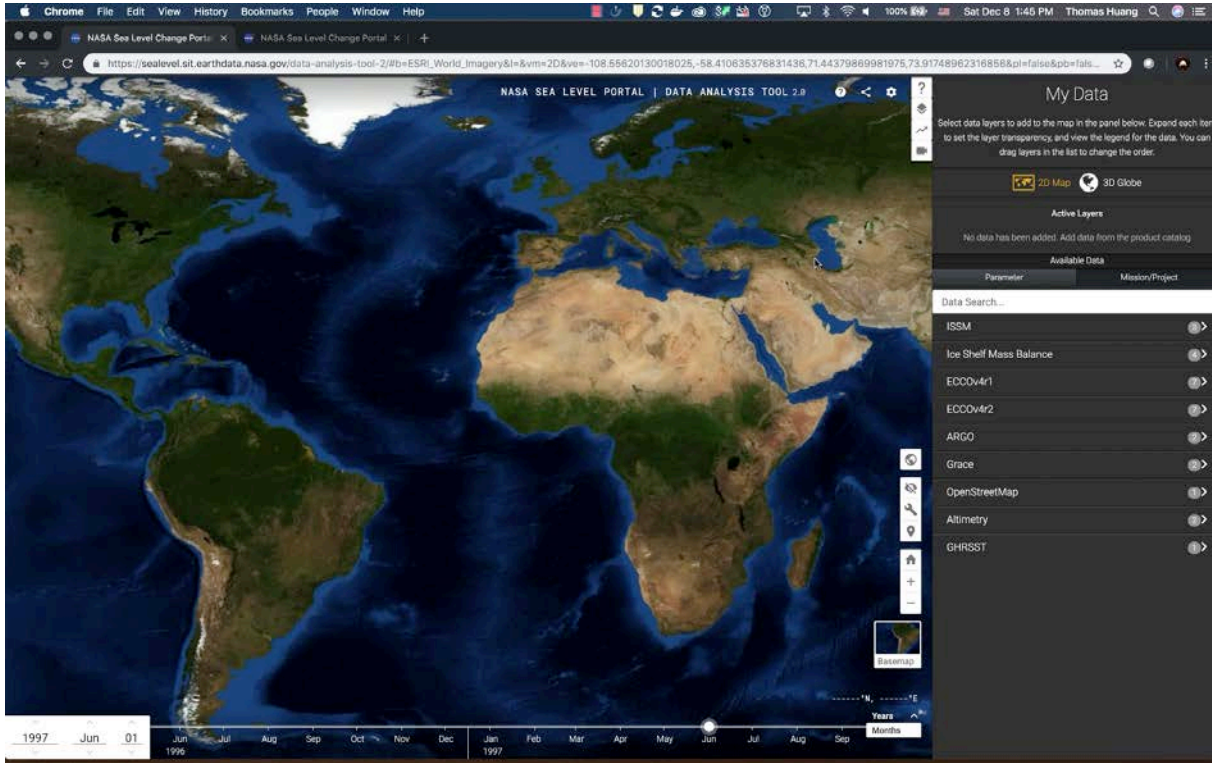




Models

General circulation models provide complete descriptions of the ocean, motivating their use as a "curve" to fit the observations.

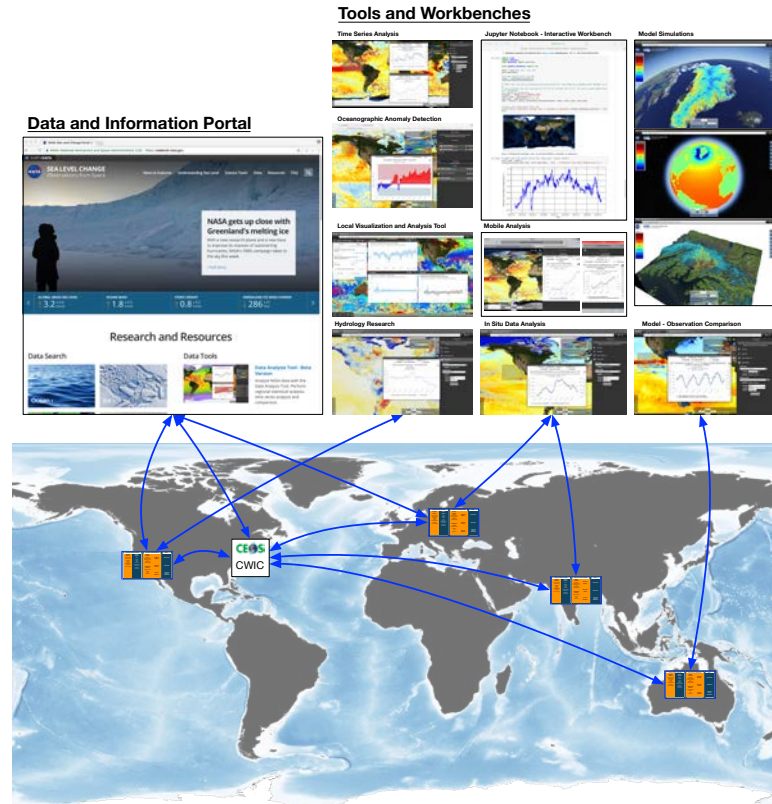"Perpetual Ocean" ECCO2 model simulation of surface current (drifter tracks)

Atmospheric Reanalyses: Combines observations with weather forecasting models to yield the most complete description of the global atmosphere. e.g., ERA-5 relative vorticity (FZ Juelich)

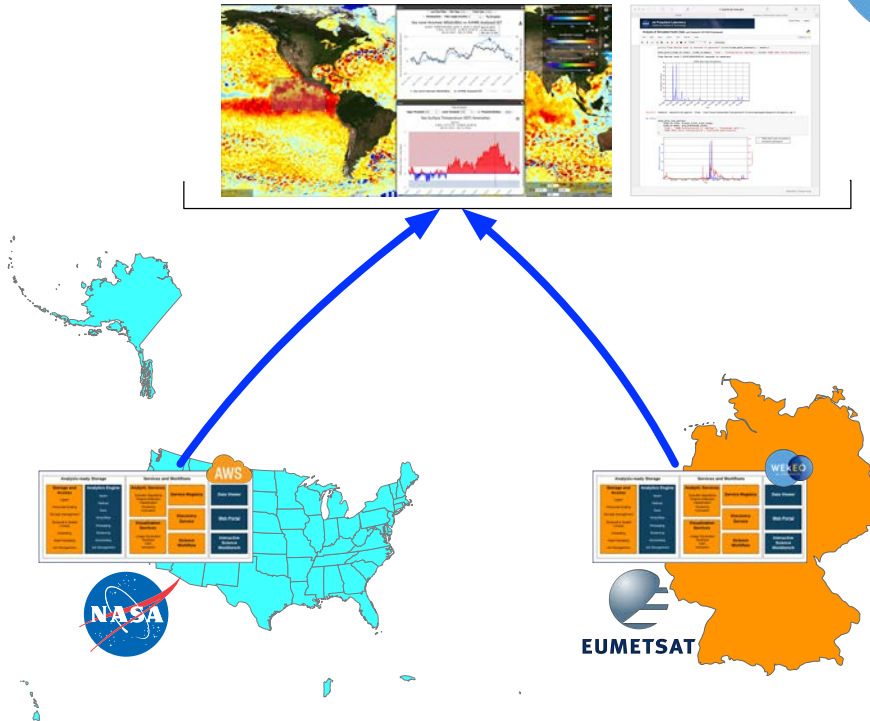# Interactive Analysis of ECCO Products

# Distributed Analytics Center Architecture

- **Committee of Earth Observation Satellites (CEOS) Ocean Variables Enabling Research and Applications for GEO (COVERAGE) Initiative**

- Seeks to provide **improved access** to **multi-agency ocean remote sensing data** that are **better integrated with in-situ and biological observations**, in support of **oceanographic and decision support applications** for societal benefit.

- A community-support open specification with common taxonomies, information model, and API (maybe security)

- Putting value-added services next to the data to eliminate unnecessary data movement

- Avoid data replication. Reduce unnecessary data movement and egress charges

- Public accessible RESTful analytic APIs where computation is next to the data

- Analytic engine infused and managed by the data centers perhaps on the Cloud

- Researchers can perform multi-variable analysis using any web-enabled devices without having to download files



Tools and Workbenches

Data and Information Portal

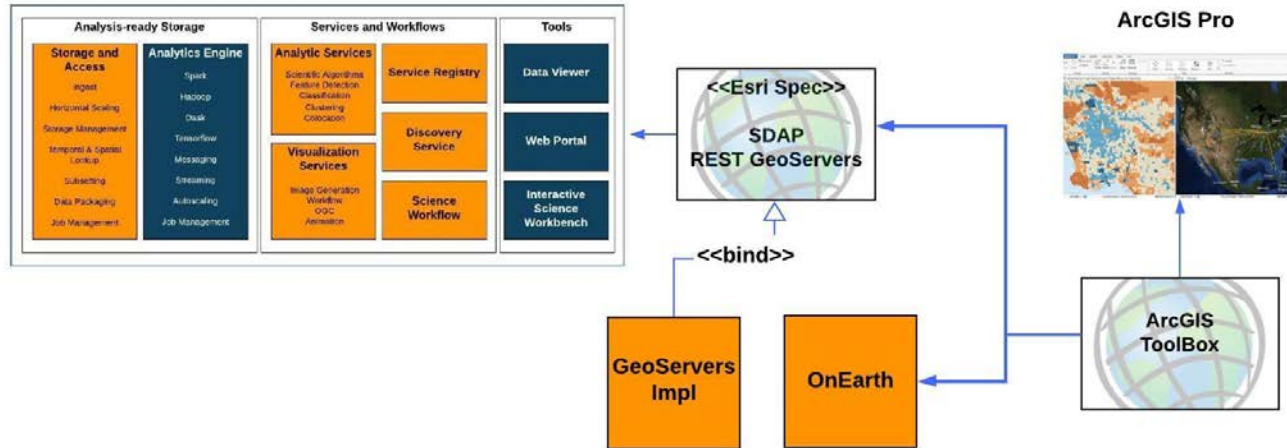# COVERAGE Phase B

- WEkEO
  - Copernicus Data and Information Access Services (DIAS)
    1. Copernicus Data
    2. Virtual Environment and Tools
    3. User Support
  - Harmonized Data Access for Satellite data and Services
  - Virtualized infrastructure for personal sandboxes
  - Pre-configured tools
- COVERAGE Phase B
  - Establish US Node on Amazon Cloud
  - Establish EU Node on WEkEO at EUMETSAT
  - Establish COVERAGE data portal and analysis tool powered by the COVERAGE Nodes at US and EU
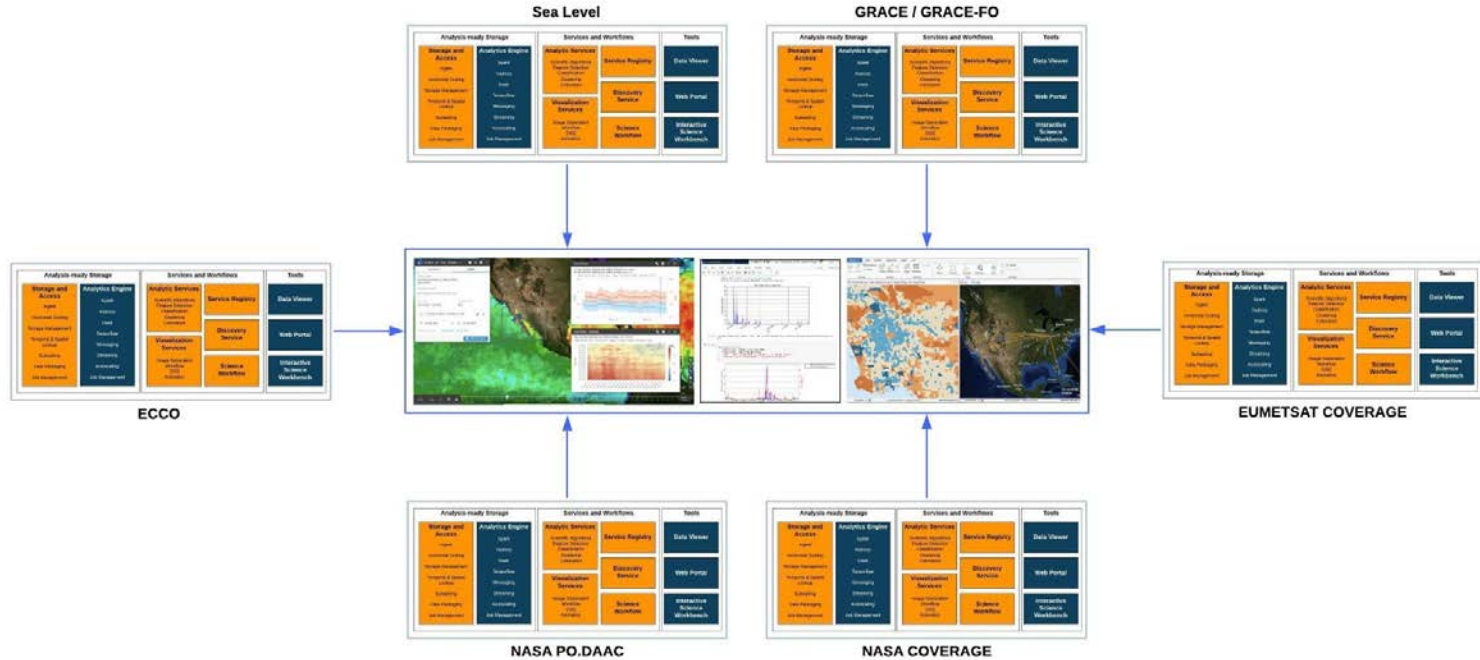
# Integration with Esri ArcGIS

- NASA's Advanced Information System Technology (AIST) effort for their big ocean science analytics solution using Apache SDAP (a.k.a OceanWorks)
- Enable scientists to use OceanWorks data and analytics within ArcGIS.
- Many scientists already use ArcGIS for their day-to-day work. Enabling them to use OceanWorks data and cloud analytics from within their familiar ArcGIS tools will enable them to perform analyses that cut across several NASA Ocean science datasets, and will expand the reach and impact of the OceanWorks data analytic system within the Ocean Science community

# Connecting to a Federation of SDAPs



Instances of Apache SDAP already in production or will be in production within the next 12 months

- Deliver solutions to establish coherent platform solutions
- Embrace open source software
- Community validation
- Evolve the technology through community contributions
- Share recipes and lessons learned
- Technology demonstrations
- Host webinars, hands-on cloud analytics workshops and hackathons


2019 EGU – NASA Hyperwall


Big Data Analytics and Cloud Computing Workshop, 2017 ESIP Summer Meeting, Bloomington, IN


2019 JPL Data Science Showcase

- **The gap between visionary to pragmatists is significant**. – Geoffrey Moore
- Become an expert in the production environment and devote resources in automations
- Give project engineering team early access to the PaaS
- Deliver all technical documents and work with project system engineering
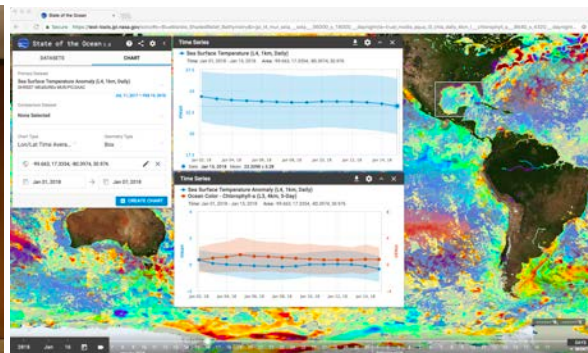- Provide project-focused trainings



NASA's Sea Level Change Team



CEOS SIT Technical Workshop





NASA's Physical Oceanography Distributed Active Archive Center (PO.DAAC)

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

Presentation | Evaluation | Collaboration

# In Summary

- **You've got to think about big things while you're doing small things, so that all the small things go in the right direction** – Alvin Toffler

- Climate research requires Autonomously Sustainable Solutions

- Open source and a web of analytics centers should be the architecture for climate science

- Focus on delivering professional quality open source solutions that enables end-to-end data and computation architecture, and the total cost of ownership

- Open source should not be a destination, it should be in place from the beginning

- How a technology is being managed will determine how far it can go

If you want to go fast, go Alone. If you want to go far, go Together.

African Proverb

**Thomas Huang**

thomas.huang@jpl.nasa.gov

Jet Propulsion Laboratory

California Institute of Technology