**National Aeronautics and Space Administration**

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

## An Introduction To

# OceanWorks
## Ocean Science Platform

**Thomas Huang**
Jet Propulsion Laboratory
California Institute of Technology

*JPL Team*
Ed Armstrong, Joseph Jacob, Nga Quach, Vardis Tsontos, and Brian Wilson

*Florida State University Team*
Shawn Smith, and Mark A. Bourassa

*National Center for Atmospheric Research Team*
Steve J. Worley

*George Mason University Team*
Chaowei (Phil) Yang, Yongyao Jiang, and Yun Li

2017 ESTF

ESTO
Earth Science Technology Office

http://podaac.jpl.nasa.gov

- The **NASA Physical Oceanographic Distributed Active Archive Center (PO.DAAC)** at Jet Propulsion Laboratory is an element of the **Earth Observing System Data and Information System (EOSDIS)**. The EOSDIS provides science data to a wide communities of user for NASA's Science Mission Directorate.

- Archives and distributes data relevant to the physical state of the ocean

- **The mission of the PO.DAAC is to preserve NASA's ocean and climate data and make these universally accessible and meaningful.**

**Reality**

- With large amount of observational and modeling data, downloading to local machine is becoming inefficient
- Data centers are starting to provide additional services
  - Better searches – faceted, spatial, keyword, relevancy, etc.
  - Data subsetting – data reduction
  - Visualization – visual discovery

**2015 NASA ESTO/AIST Big Data Study Roadmap: Moving from Data Archiving to Data Analytics**

**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
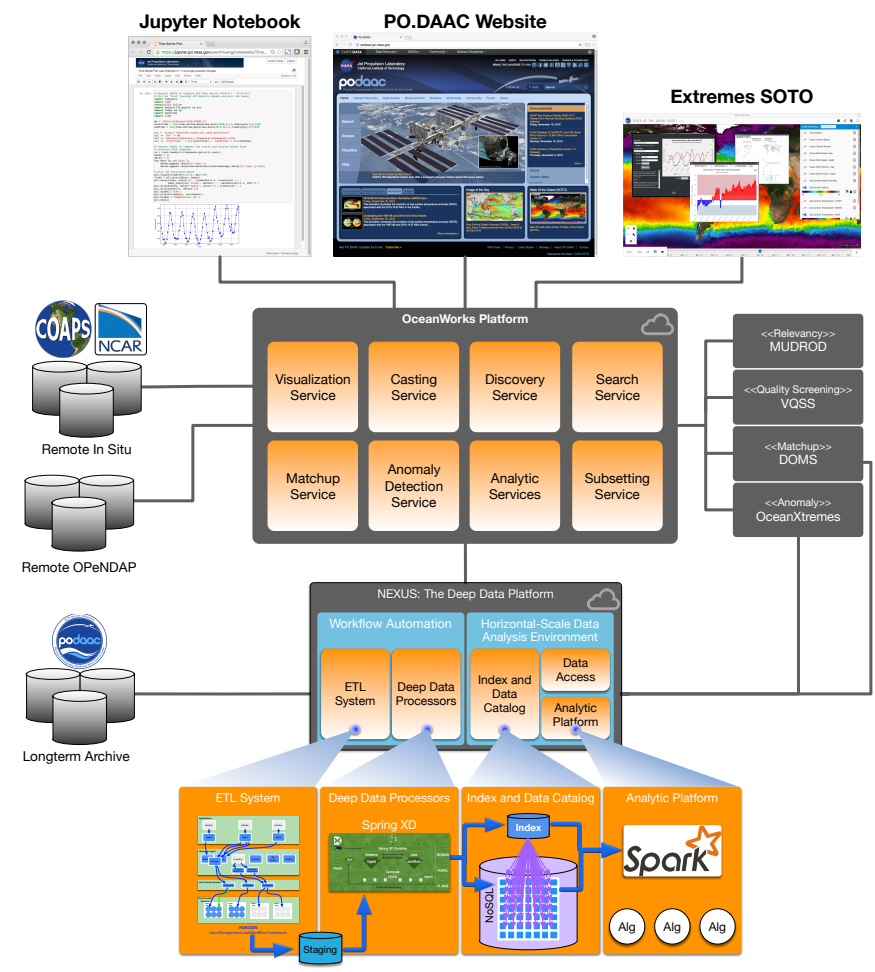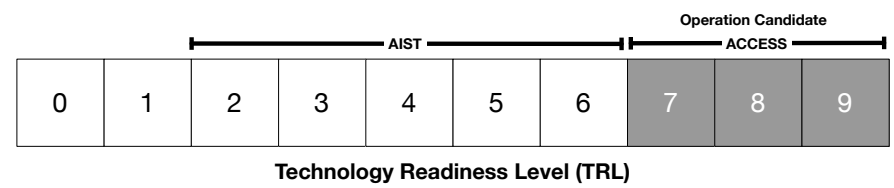- Apply computational and data science across the lifecycle

**Scalable Data Management**
- Capturing 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

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
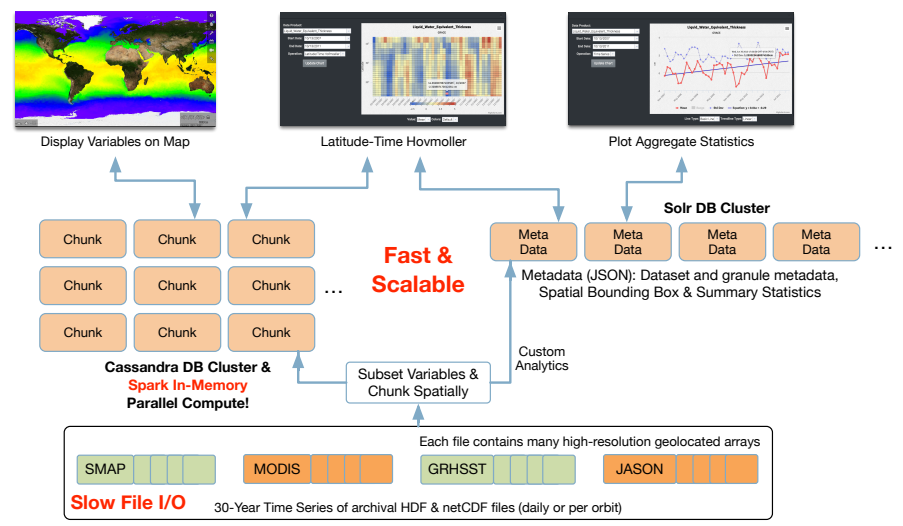California Institute of Technology
Pasadena, California

- **OceanWorks** is to establish an integrated data analytic center at the NASA PO.DAAC for Big Ocean Science. It focuses on technology integration, advancement and maturity

- **Collaboration between JPL, FSU, NCAR, and GMU**

- Bringing together PO.DAAC-related big data technologies

  - **AIST-14 OceanXtremes (PI: Huang/JPL) – TRL 4**
    Anomaly detection and ocean science
  - **NEXUS (PI: Huang/JPL) – TRL 6**
    Deep data analytic platform
  - **AIST-14 DOMS (PI: Smith/FSU) – TRL 4**
    Distributed in-situ to satellite matchup
  - **AIST-14 MUDROD (PI: Yang/GMU) – TRL 7**
    Search relevancy and discovery
  - **ACCESS-13 VQSS (PI: Armstrong/JPL) – TRL 7**
    Virtualized Quality Screening Service

**Operation Candidate**

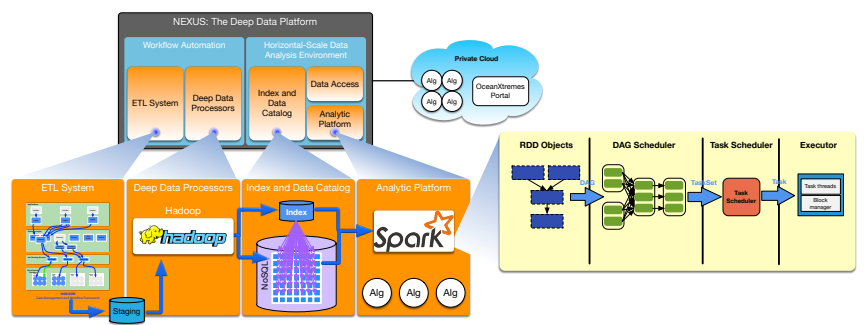| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|

**Technology Readiness Level (TRL)**

Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California Institute of Technology.

- **Improve Data Discovery**: help users to discover relevant data, anomalies and events, quality and uncertainty information, and provide multi-observing system matchup data. **PO.DAAC Search** webservice will be developed to demonstrate new data search and discover related matches and phenomenon

- **Subset and Distribute Data**: high performance data subsetting that supports quality screening

- **Identify and Catalog Ocean phenomenon**: registry of detected oceanographic phenomenon and published using datacasting technology

- **Matchup between Satellite and In-Situ Observations**: advanced matchup services that support in-situ to satellite colocation. It provides a mechanism for users to input a series of geospatial references for satellite observations and receive the in-situ observations that matched to the satellite data within a selectable temporal and spatial search domain.

- **Analyze Satellite Observations**: long time-series, correlation map, time averaged map, Hovmöller, climatological map, etc

- **Visualize and analyze Satellite Observation on the Web**: high performance data visualizations with linkages to actual source measurements, events and analytics. A **prototype integration** with PO.DAAC's new State of the Ocean (SOTO) web application will be developed, called **Extremes SOTO**, to demonstrate this objective.

- **API Integration** with Jupyter notebook to demonstrate working directly with OceanWorks' webserivce platform

2017 ESTF

- 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

**Open Source: Apache License 2**
https://github.com/dataplumber/nexus



Two-Database Architecture



Deep Data Computing Environment (DDCE)

2017 ESTF

National Aeronautics and Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

**SEA LEVEL CHANGE**
**Observations from Space**

# Data Analysis Tool   https://sealevel.nasa.gov
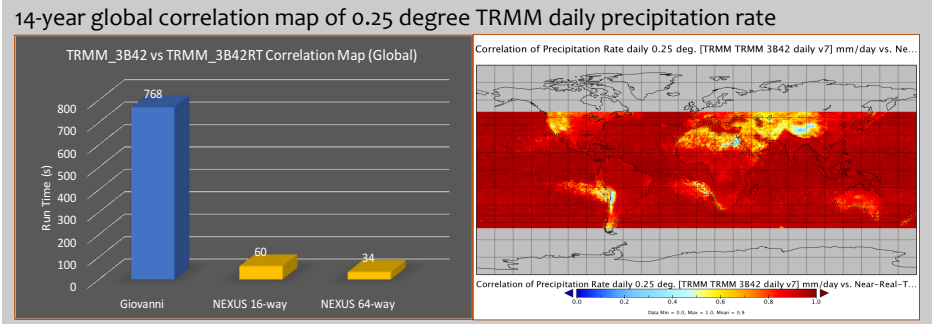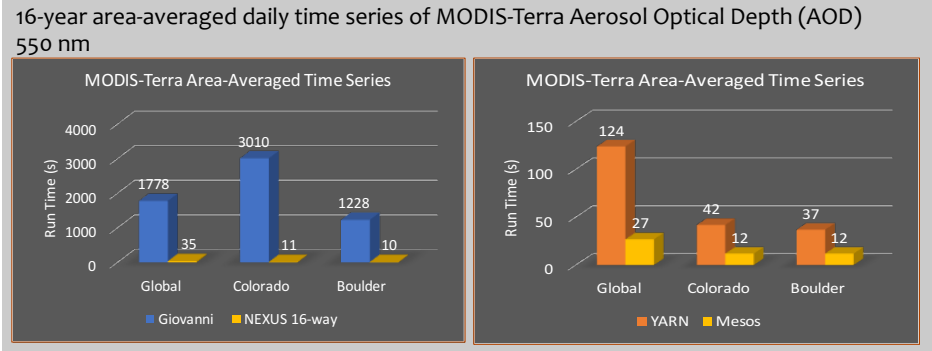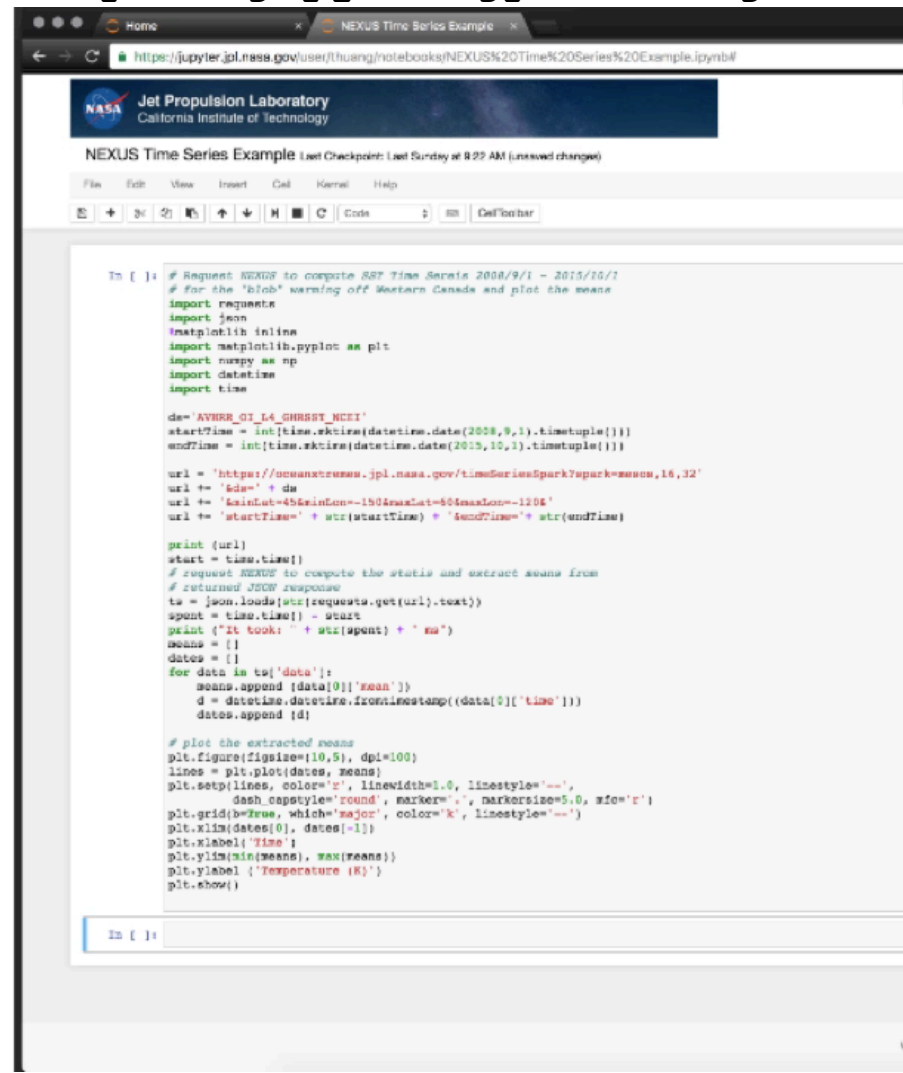


**Facebook: 28K followers**
**Twitter: 22K followers**

**"NASA Sea Level Change Website Offers Everything You Need To Know About Climate Change"**
http://www.techtimes.com/articles/147210/20160405/nasa-sea-level-change-website-offers-everything-need-know-climate.htm

**"NASA's new sea level site puts climate change papers, data, and tools online"**
http://techcrunch.com/2016/04/04/nasas-new-sea-level-site-puts-climate-change-papers-data-and-tools-online/

2017 ESTF

ESTO
Earth Science Technology Office

- **2016 Earthdata prototyping**
  - Twelve factor application
  - Benchmarking against Giovanni
  - Amazon cloud deployment
- **2016 ESTO Data Container Study**
  - Look at a variety of technologies for reorganizing and storing Earth science data to make them more tractable to full-scale science analysis.
  - To understand the strong points and tradeoffs of the different approaches to large-scale analysis
  - NEXUS' plug-an-play data storage
    - Apache Cassandra – Java implementation
    - ScyllaDB – C++ implementation
  - Task 1 – long time series
    - Point Time series for Boulder, CO
    - Area-averaged time series for the state of Colorado
    - Area Averaged time series for the globe
  - Task 2 – climatological map
  - Document ETL process, elapsed time, compute and storage
- **2017 Earthdata prototyping**
  - Auto ingestion
  - NASAcompilant Generation Application Platform (NGAP) Infusion
  - Cloud Analysis Toolkit to Enable Earth Science (CATEE)
  - 2017 ESIP Summer Workshops

16-year area-averaged daily time series of MODIS-Terra Aerosol Optical Depth (AOD) 550 nm



14-year global correlation map of 0.25 degree TRMM daily precipitation rate



17-year area-averaged time series over the continental United States

**National Aeronautics and Space Administration**

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

# `https://jupyter.jpl.nasa.gov`



```
# Request NEXUS to compute SST Time Series 2008/9/1 – 2015/10/1
# for the "blob" warming off Western Canada and plot the means
…
ds='AVHRR_OI_L4_GHRSST_NCEI'

url = … # construct the webservice URL request

# make request to NEXUS using URL request
# save JSON response in local variable
ts = json.loads(str(requests.get(url).text))

# extract dates and means from the response
means = []
dates = []
for data in ts['data']:
    means.append (data[0]['mean'])
    d = datetime.datetime.fromtimestamp((data[0]['time']))
    dates.append (d)

# plot the result
…
```
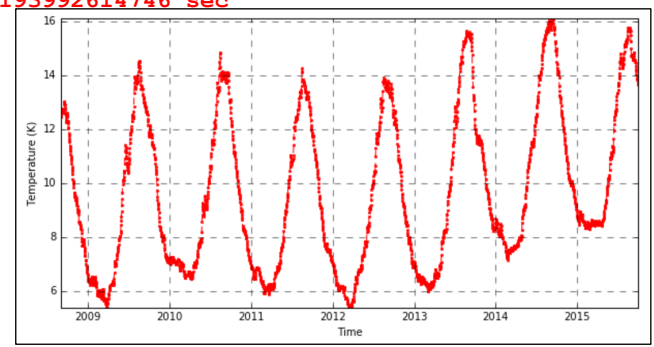
```
https://oceanxtremes.jpl.nasa.gov/timeSeriesSpark?
spark=mesos,
16,32&ds=AVHRR_OI_L4_GHRSST_NCEI&minLat=45&minLon=-150&ma
xLat=60&maxLon=-120&startTime=1220227200&endTime=14436576
00

It took: 6.909193992614746 sec
```

2017 ESTF

ESTO
Earth Science Technology Office

National Aeronautics and Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

- https://oceanxtremes.jpl.nasa.gov
- An oceanographic data-intensive anomaly detection and analysis portal
- Cloud-based big data analytic platform for
  - Climatology generation
  - On-the-fly daily difference computation
  - Anomaly registry and publication
  - On-the-fly data analytics
- Recent highlights
  - **Recreated identification of "The Blob"**
    - **The Blob** is the name given to a large mass of relatively warm water in the Pacific ocean off the coast of North America. It was first detected in late 2013 and continued to spread throughout 2014 and 2015.
    - SST anomaly = SST – SST Climatology at each location to compare with standard deviation  - Dr. Chelle Gentemann, Senior Scientist at Earth & Space Research
  - **Recreated the El Niño 3.4 regional signal**
    - Dec. 2010 – May 2016
    - **El Niño** is a phenomenon in the equatorial Pacific Ocean characterized by a five consecutive 3-month running mean of sea surface temperature (SST) anomalies in the Niño 3.4 region that is above (below) the threshold of +0.5°C (-0.5°C). This standard of measure is known as the Oceanic Niño Index (ONI).



The Blob



El Niño 3.4 regional signal

2017 ESTF

Earth Science Technology Office

Hurricane Katrina passed to the southwest of Florida on Aug 27, 2005. The ocean response in a 1 x 1 deg region is captured by a number of satellites. The initial ocean response was an immediate cooling of the surface waters by 2 deg C that lingers for several days. Following this was a short intense ocean chlorophyll bloom a few days later. The ocean may have been "preconditioned' by a cool core eddy and low sea surface height.

The SST drop is correlated to both wind and precipitation data. The Chl-A data is lagged by about 3 days to the other observations like SST, wind and precipitation.

*A study of a Hurricane Katrina–induced phytoplankton bloom using satellite observations and model simulations*
Xiaoming Liu, Menghua Wang, and Wei Shi

JOURNAL OF GEOPHYSICAL RESEARCH, VOL. 114, C03023, doi:10.1029/2008JC004934, 2009

http://shoni2.princeton.edu/ftp/lyo/journals/Ocean/
phybiogeochem/Liu-etal-KatrinaChlBloom-JGR2009.pdf



Hurricane Katrina
TRMM overlay SST Anomaly



Powered By NEXUS

2017 ESTF

# OceanXtremes: Identify . Analyze . Share

- Visualize parameter
- Compute daily differences against climatology
- Analyze time series area averaged differences
- Replay the anomaly and visualize with other measurements
- Document the anomaly
- Publish the anomaly

2017 ESTF

- https://doms.jpl.nasa.gov
- Distributed Oceanographic Matchup Service
- Typically data matching is done using one-off programs developed at multiple institutions
- A primary advantage of DOMS is the reduction in duplicate development and man hours required to match satellite/in situ data
  - Removes the need for satellite and in situ data to be collocated on a single server
  - Systematically recreate matchups if either in situ or satellite products are re-processed (new versions), i.e., matchup archives are always up-to-date.
- In situ data nodes at JPL, NCAR, and FSU operational.
- Provides data querying, subset creation, match-up services, and file delivery operational.
- Prototype graphical user interface (UI) and APIs accessible for external users.
- Plugin architecture for in situ data source using EDGE
  - Extensible Data Gateway Environment is an Apache License 2 open source technology
  - https://github.com/dataplumber/edge
- Defined specification for packaging matchup results. Working with Unidata and ESDSWG's data interoperability and standard groups

2017 ESTF

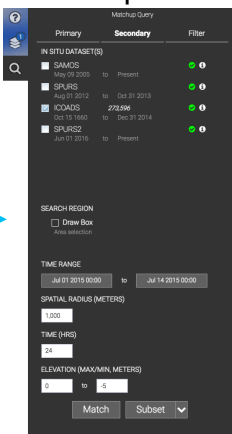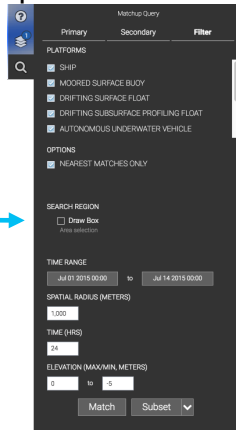National Aeronautics and Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

Primary Dataset

Match-up In-Situ
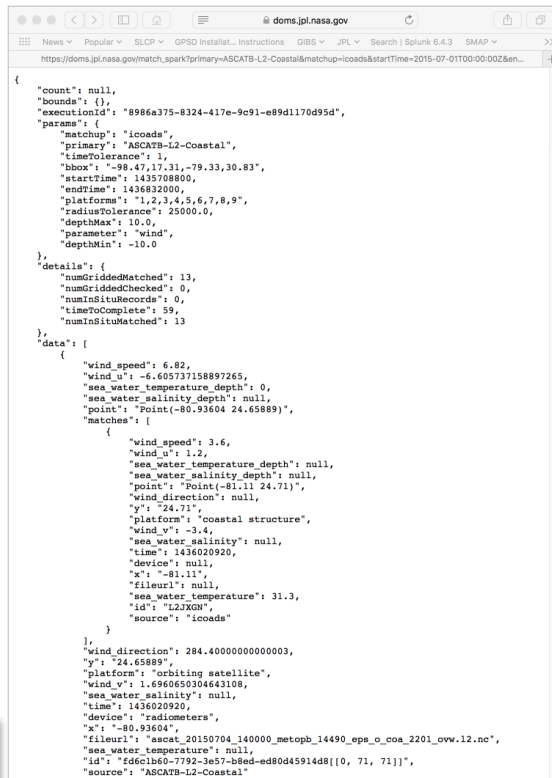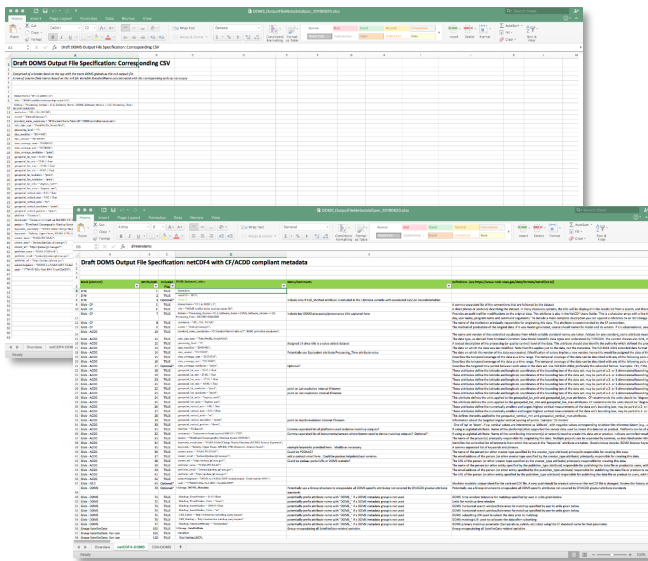
Optional Platform



- **Workflow driven**
- **Fast on-the-fly subseting satellite and in-situ**
- **Open webservice API**
- **ON-the-fly data analysis using NEXUS**
- **Cloud ready**

**PO.DAAC Satellite Data**
- AVHRR_OI_L4_GHRSST_NCEI
- JPL-L4-GHRSST-SSTfnd-MUR-GLOB-v02.0fv04.1
- SMAP_L2B_SSS
- ASCATB-L2-Costal

**In-situ Data**
- JPL SPURS 1 and 2
- FSU SAMOS
- NCAR ICOADS

**Drafted and implemented matchup specification**
- Promote CF and ACDD standards
- Self-contained
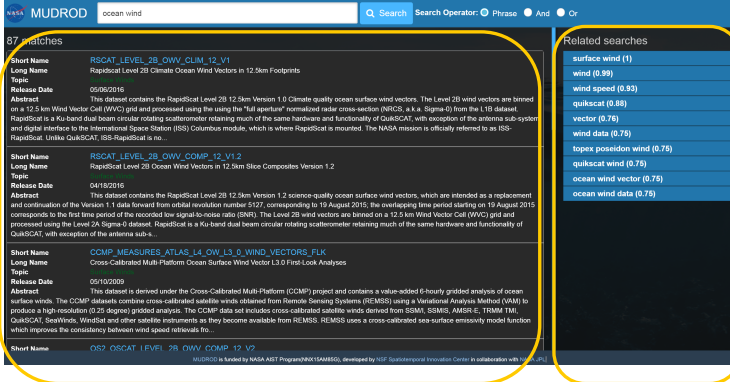- Simple to use by other tools and services

2017 ESTF

Earth Science Technology Office

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

**MUDROD:**
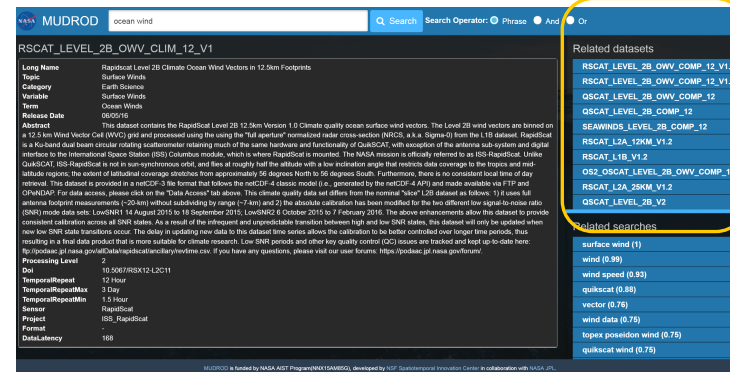**Mining and Utilizing Data Relevancy from Oceanographic Dataset**

- Mining and Utilizing Dataset Relevancy from Oceanographic Dataset

- **Search** – look for something you expect to exist
  - Information tagging
  - Indexed search technologies like Apache Solr or ElasticSearch
  - The solution is pretty straightforward

- **Discovery** – find something new, or in a new way
  - This is non-trivial
  - Traditional ontological method doesn't quite add up
  - The strength of semantic web is in inference
  - What happen when we have a lot of `subClassOf, equivalentClassOf, sameAs`?
  - How wide and deep should we go?

- **Relevancy**
  - It is domain-specific
  - It is personal
  - It is temporal
  - It is dynamic

- MUDROD analyzes web logs to discover user knowledge (the connections between datasets and keyword)

- Construct knowledge base by combining semantics and profile analyzer

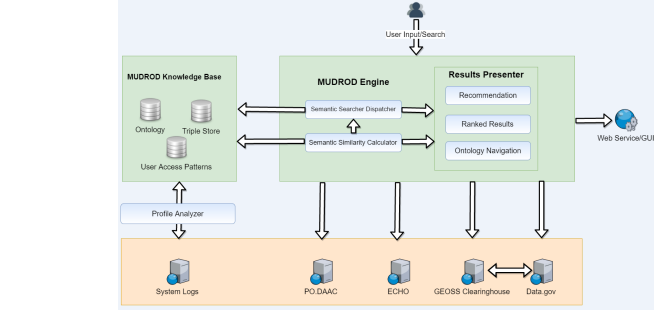- Improve data discovery by better ranked results
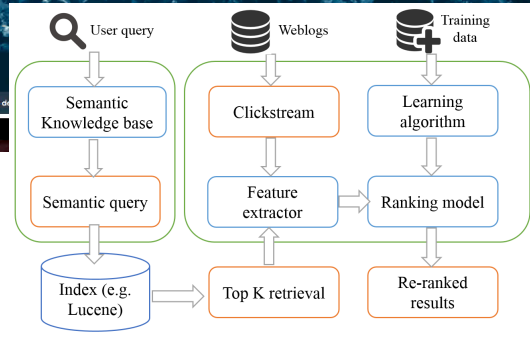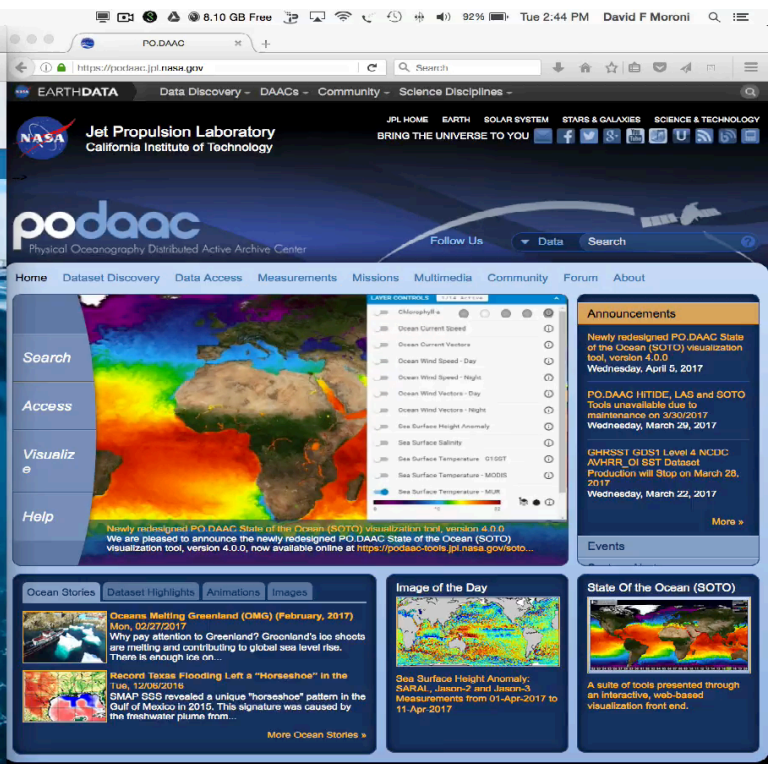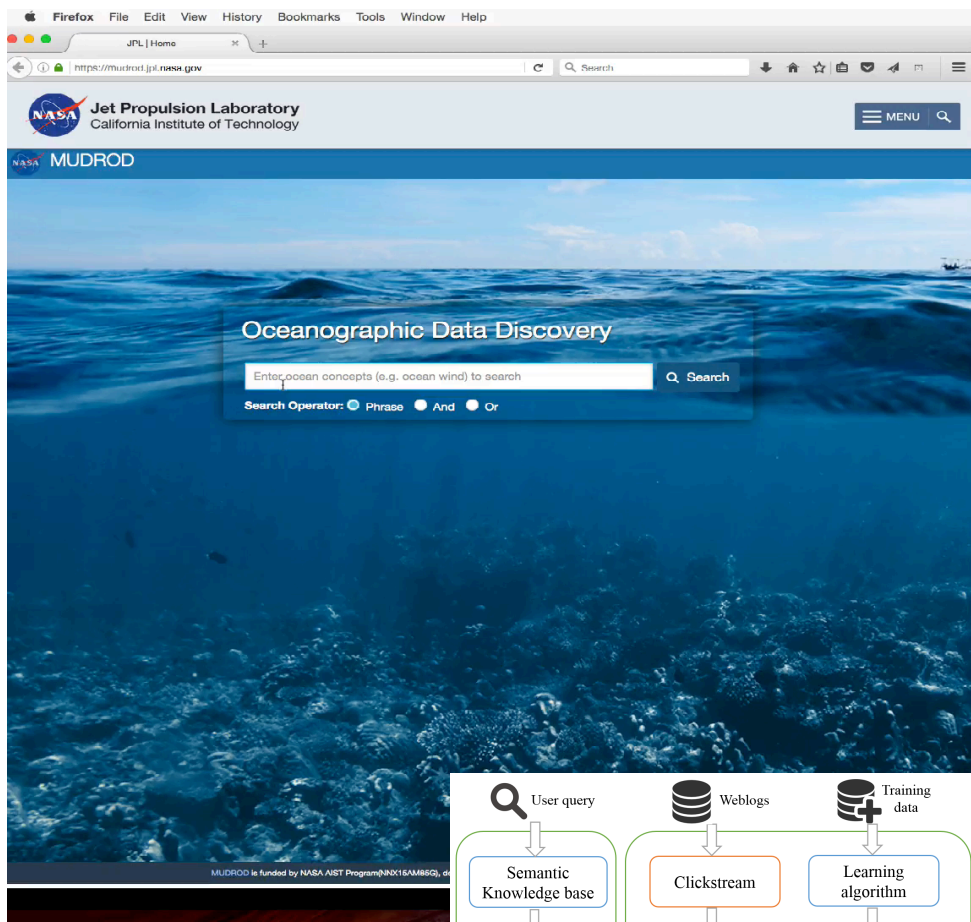
https://mudrod.jpl.nasa.gov

**Search Ranking**
Based on a machine learning model (RankSVM) which takes a number of features, such as vector space model, version, processing level, release date, all-time popularity, monthly-popularity, and user popularity.

**Search Recommendation**
Based on dataset metadata content and web session co-occurrence

2017 ESTF

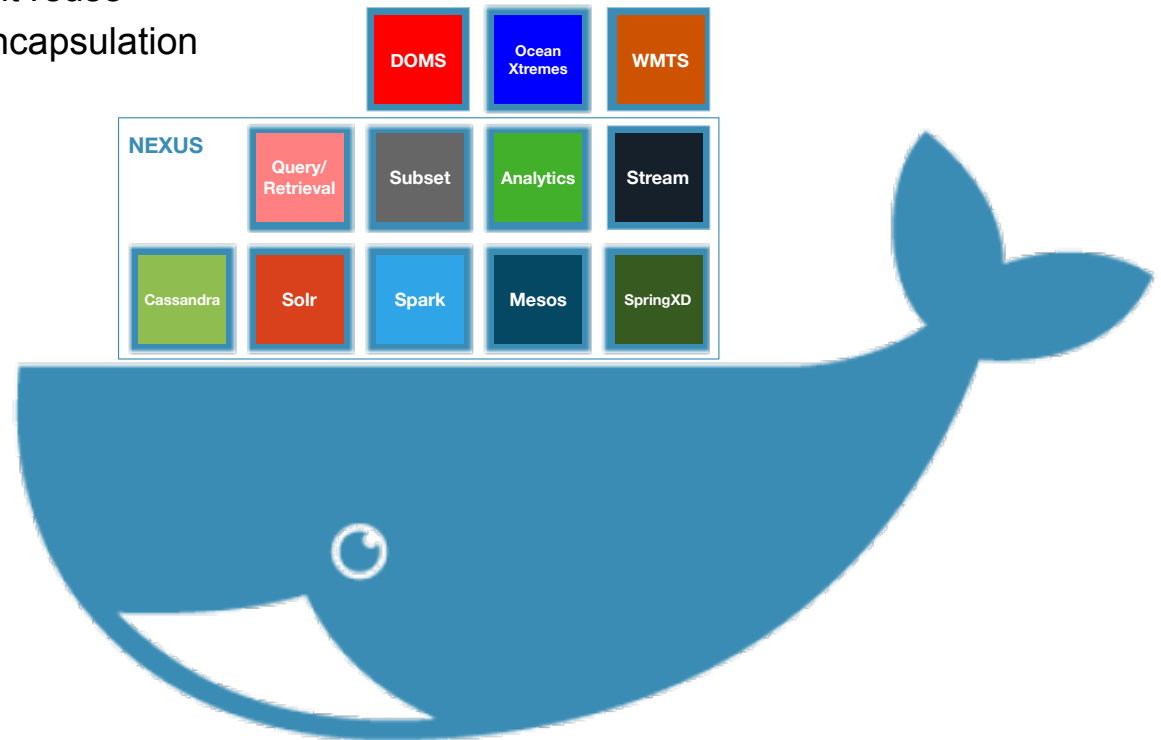Earth Science Technology Office

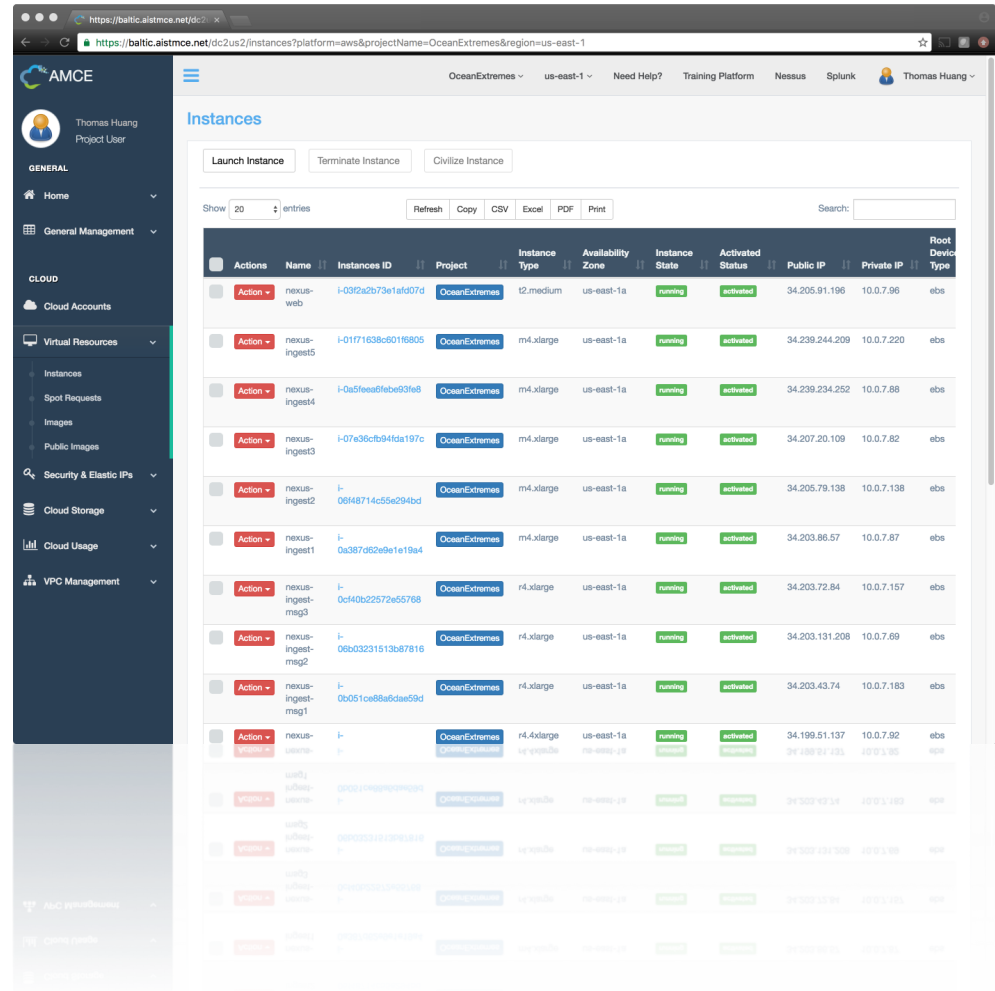# Speed . Relevant . Discovery

## Features

- Fast web log ingestion and processing using Apache Spark, in-memory MapReduce
- Session reconstruction
- Vocabulary semantic relationship extraction
- Machine Learning Search ranking
- Integration with SWEET Ontologies for semantic-driven search and recommendations
- Recommendation

2017 ESTF

- Dockerizing all services
- **Why?**
  - Rapid application deployment
  - Portability across machines
  - Application-centric vs machine/server-centric
  - Version control and component reuse
  - Secure due to isolation and encapsulation
  - Sharing
  - Lightweight footprint
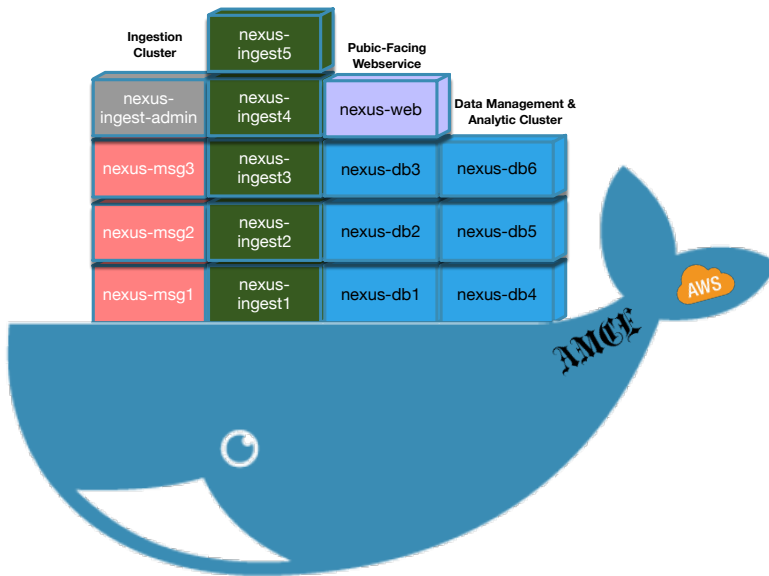  - Minimal overhead
  - Simplified maintenance

- OceanWorks will leverage the AIST Managed Cloud Environment (AMCE) for development – the AIST-provisioned Amazon Cloud environment
- OceanXtremes and NEXUS deployment
  - Multi-container, multi-cluster deployment
  - Leverage public DockerHub
  - Working toward dev-test pipeline automation
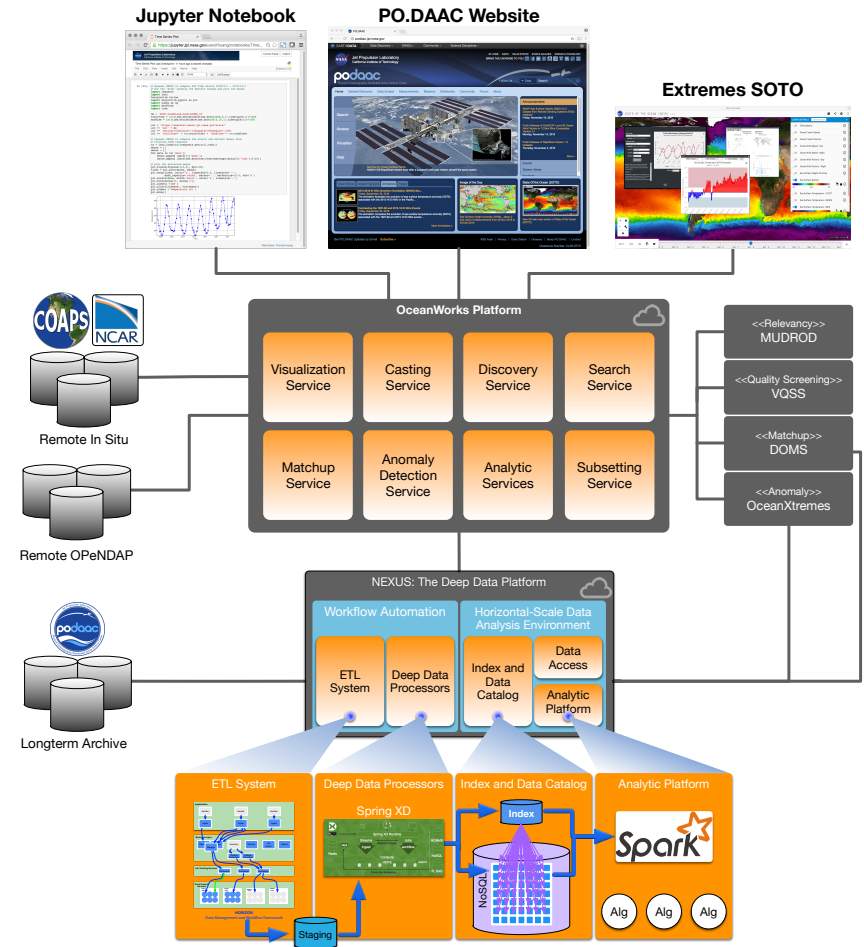  - Deployed under 16 Amazon instances. Needs testing

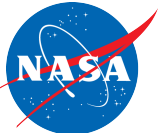OceanXtremes AMCE Deployment

- Lives on the Cloud
- Auto-Scaling
- One-The-Fly multi-parameter data analysis
- Access and matchup with in-situ measurements
- Smart subsetting
- Anomaly detection and registration
- Sharing of analytic results
- Lightning speed searches
- Discover relevant data, services, news and publications like never before
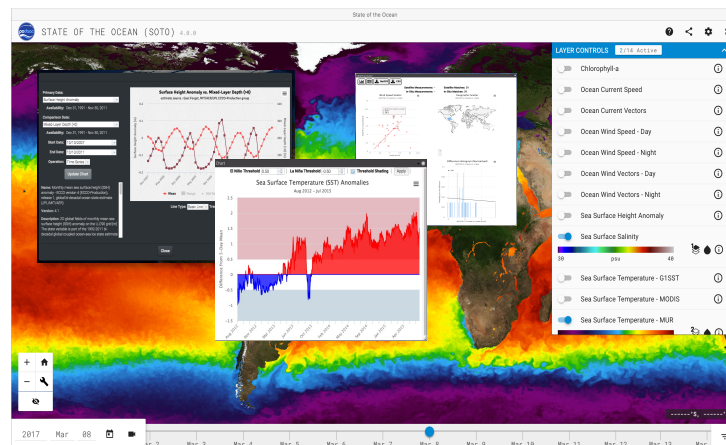- Fully Open Source

# OceanWorks: Ocean Science Data Platform

## PI: Thomas Huang, JPL

## Objective

Develop an integrated ocean science platform to provide a virtual environment for the PO.DAAC community and applications by leveraging NEXUS, OceanXtremes, DOMS, and MUDRODS to satisfy the following goals:

- Improve data service discovery
- Subset and distribute data
- Identify and catalog ocean phenomenon
- Coordination between satellite and in-situ observations
- Analyze satellite observations
- Visualize and analyze satellite observation on the web
- Enable webservice API integration



OceanWorks State of the Ocean (SOTO) Visualization and analysis

## Approach

- Define integrated system architecture and information model
- Prototype integration with PO.DAAC's existing visualization solution, SOTO, add develop web-based analytic capabilities
- Review by PO.DAAC UWG
- Improve data and service search and discovery
- Develop new data subsetting capability
- Demonstrate Jupyter notebook integration
- Facilitate validation by 4 organization and improve prformance
- PO.DAAC User Acceptance Testing environment (UAT) deployment

**Co-Is/Partners**: E. Armstrong, J. Jacob, N. Quach, V. Tsontos, B. Wilson, JPL; S. Smith, M. Bourassa, FSU; S. Worley, NCAR; C. Yang, Y. Jiang, Y. Li, GMU

## Key Milestones

| | |
|---|---|
| • Complete system design and dataset selection | 09/17 |
| • Perform OceanWorks system-level testing | 12/17 |
| • PO.DAAC User Working Group (UWG) CDR | 04/18 |
| • Complete discovery analysis services and performance optimization | 09/18 |
| • Data subsetting capability | 09/18 |
| • Integrate Jupyter notebook capability | 12/18 |
| • PO.DAAC UWG acceptance review and UAT deployment | 04/19 |

$TRL_{in} = 4$    $TRL_{current} = 4$

Earth Science Technology Office

**Credits**

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

**NASA**
Marjorie Cole
Brandi Quam
Mike Little
Chris Lynnes
Kevin Murphy

**FSU COAPS**
Mark Bourassa
Jocelyn Elya
Shawn Smith
Adam Stallard

**NCAR**
Steve Worley
Ji Zaihua

**GMU**
Yongyao Jiang
Yun Li
Chaowei (Phil) Yang

**JPL NEXUS Engineers**
Stewart (Parker) Abercrombie
Kevin Gill
Frank Greguska
Joseph Jacob
Nga Quach
Brian Wilson

**JPL Science Contributors**
Ed Armstrong
Andrew Bingham
Carmen Boening
Mike Chin
Michelle Gierach
Ben Holt
Tony Lee
David Moroni
Rob Toaz
Vardis Tsontos
Victor Zlotnicki

**Questions, and more information**

Thomas.Huang@jpl.nasa.gov