**National Aeronautics and
Space Administration**

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

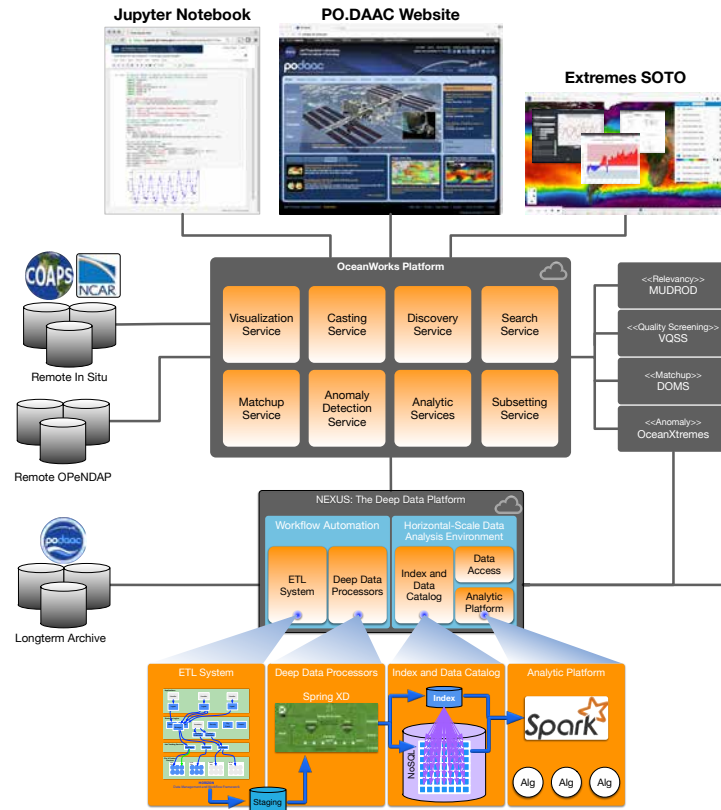# Apache Science Data Analytic Platform (SDAP)

**Thomas Huang**

Data Scientist | Principal Investigator | Technologist | Architect

thomas.huang@jpl.nasa.gov

Jet Propulsion Laboratory
California Institute of Technology
4800 Oak Grove Drive, Pasadena, CA 91109-8099, U.S.A.
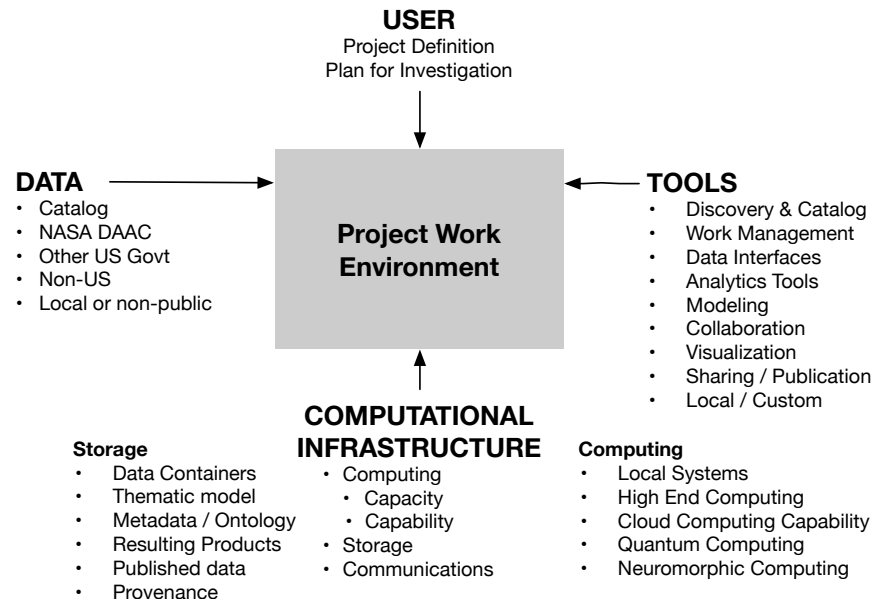
# Big Data and Data Centers

- **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 computation and data science across the lifecycle
- **For NASA Data Centers, with large amount of observational and modeling data, downloading to local machine 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
- **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

# AIST OceanWorks

- **OceanWorks** is to establish an **Integrated Data Analytic Center** at the NASA Physical Oceanography Distributed Active Archive Center (PO.DAAC) for Big Ocean Science
- Focuses on technology integration, advancement and maturity
- Collaboration between JPL, FSU, NCAR, and GMU
- Bringing together PO.DAAC-related big data technologies
  - **OceanXtremes –** Anomaly detection and ocean science
  - **NEXUS –** Big data analytic platform
  - **Data Container Studies**
  - **DOMS –** Distributed in-situ to satellite matchup
  - **MUDROD –** Search relevancy and discovery – linking datasets, services, and anomalies through recommendations
  - **VQSS –** Virtualized Quality Screening Service

# Integrated Data Analytic Center

- An environment for conducting a Science investigation
  - Enables the confluence of resources for that investigation
  - Tailored to the individual study area (ocean, atmospheric, sea level, etc.)
- Harmonizes data, tools and computational resources to permit the research community to focus on the investigation
  - Reduce the data preparation time to something tolerable
  - Catalog of optional resources
  - Semantic-enabled catalog of resources
  - Relevant publications
  - Provide established training data sets of varying resolution
  - Provide effective project confidentiality, integrity and availability
  - Single sign-on and unified financial tracking

**USER**
Project Definition
Plan for Investigation

**DATA**
- Catalog
- NASA DAAC
- Other US Govt
- Non-US
- Local or non-public

**Project Work Environment**

**TOOLS**
- Discovery & Catalog
- Work Management
- Data Interfaces
- Analytics Tools
- Modeling
- Collaboration
- Visualization
- Sharing / Publication
- Local / Custom

**COMPUTATIONAL INFRASTRUCTURE**

**Storage**
- Data Containers
- Thematic model
- Metadata / Ontology
- Resulting Products
- Published data
- Provenance

- Computing
  - Capacity
  - Capability
- Storage
- Communications

**Computing**
- Local Systems
- High End Computing
- Cloud Computing Capability
- Quantum Computing
- Neuromorphic Computing

Credit: Mike Little, NASA

# OceanWorks as an Analytic Center

**DATA**
- Earthdata CMR
- nonCMR DAAC
- PI Generated
  - ECCO
  - Altimetry
- In Situ
  - ICOADS
  - SAMOS
  - SPURS I & 2
- Satellite
  - Chlorophyll
  - Gravity
  - Salinity
  - SST
  - Winds

**PHYSICAL OCEANOGRAPHERS**
Project Definition
Plan for Investigation
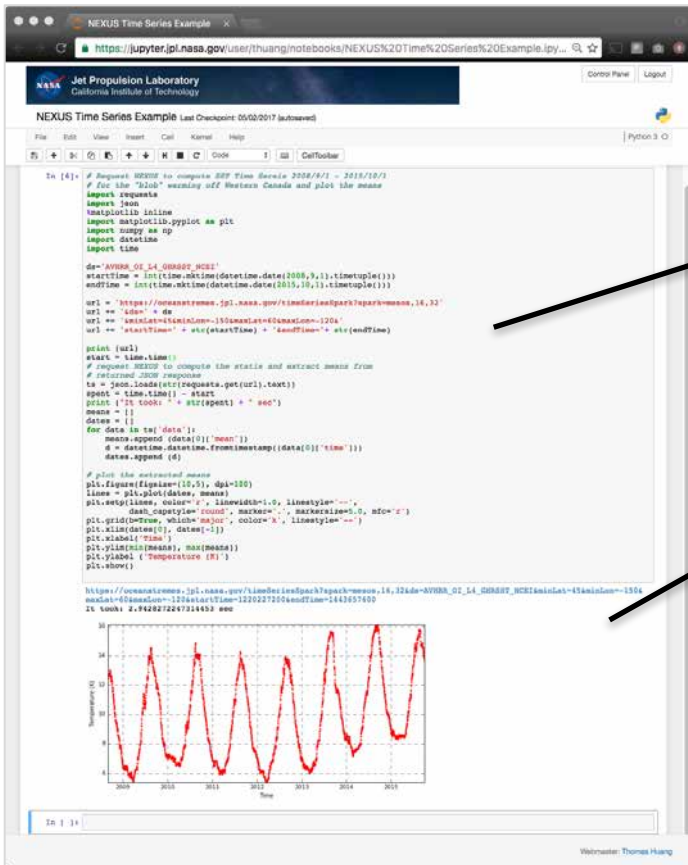
**Project Work Environment**

**TOOLS**
- EDGE and MUDROD: Metadata, Search & Discovery
- Services
  - Area Averaged Time Series
  - Time Averaged Map
  - Correlation Map
  - Anomaly: Daily Differences
  - Matchup (single satellite - multiple in situ)
- Workflow
  - AWS Lambda, Step Functions, Batch
  - SpringXD
  - Jupyter Notebook
- Visualization
  - CMC (GIS)
  - OnEarth
- Deployment
  - Bamboo
  - Jenkins
  - Docker
  - AWS CloudFormation
- Collaboration
  - Confluence, JIRA, GIT
  - Apache wiki
  - Smartsheet and Google Office
  - Slack

**COMPUTATIONAL INFRASTRUCTURE**

**Storage**
- NEXUS
- Apache Solr
- Amazon S3

**Computing**
- Local Systems
- Amazon
- AMCE Cloud Computing
- NGAP
- JPL on Premises Cloud

# End User Applications

# Enable Science without File Download



```python
# 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
…
```
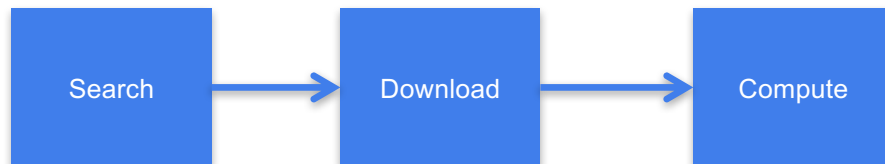
```
https://oceanxtremes.jpl.nasa.gov/timeSeriesSpark?spark=me
sos,16,32&ds=AVHRR_OI_L4_GHRSST_NCEI&minLat=45&minLon=-
150&maxLat=60&maxLon=-
120&startTime=1220227200&endTime=1443657600

It took: 2.9428272247314453 sec
```
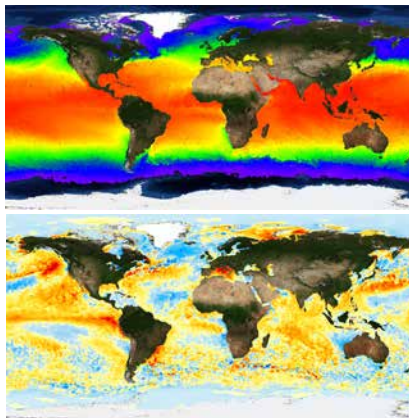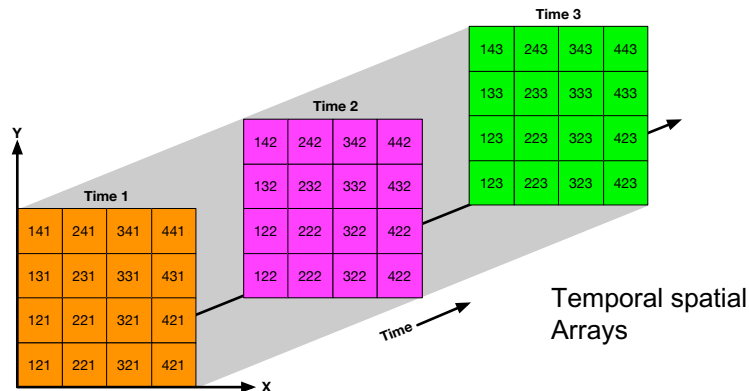
# Traditional Method for Analyze Satellite Measurements

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

Search → Download → Compute

- 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

**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



Temporal spatial Arrays

# NEXUS Performance: Custom Spark vs. AWS EMR

**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)
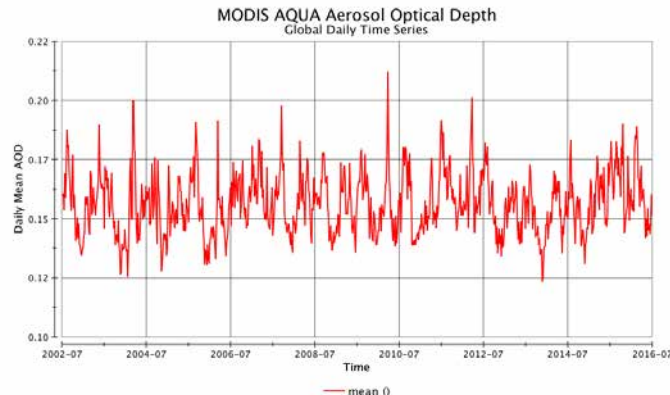**File Count**: 5106
**Volume**: 2.6GB
**Time Coverage**: July 4, 2002 – July 3, 2016

**Giovanni**: 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



**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

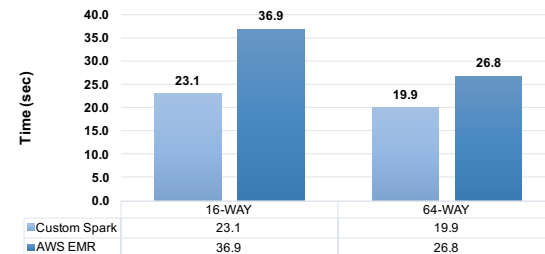|  | 16-WAY | 64-WAY |
|---|---|---|
| Custom Spark | 3.3 | 2.9 |
| AWS EMR | 3.8 | 3.1 |



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

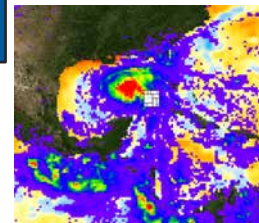Custom Spark vs. AWS EMR
Ref. Speed - Giovanni: 1366.84 sec

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

# Hurricane Katrina Study



Powered by NEXUS

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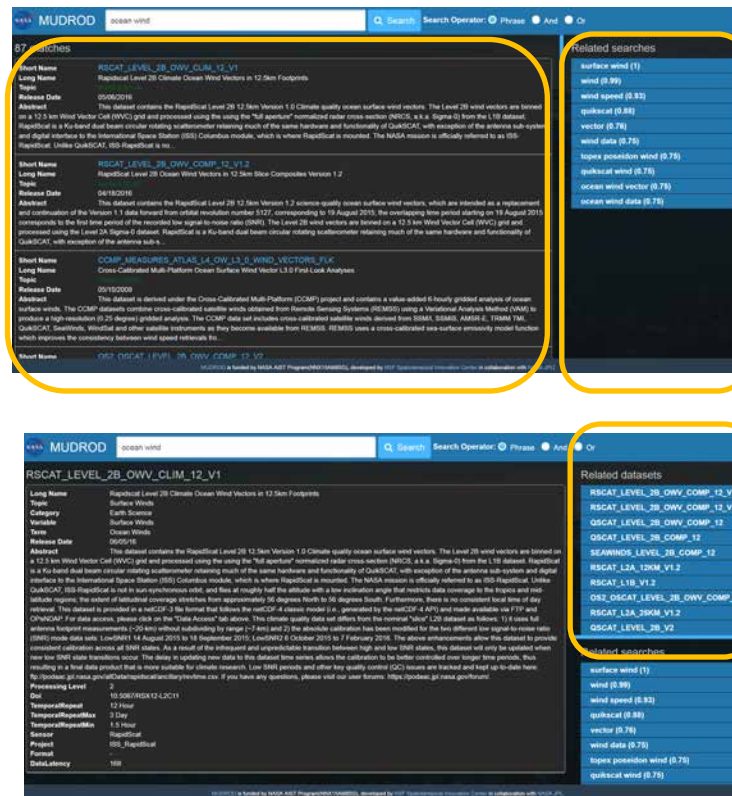
The SST drop is correlated to both wind and precipitation data. The Chl-A data is lagged by about 3 days to the other observations like SST, wind and precipitation.



Hurricane Katrina TRMM overlay SST Anomaly

*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

- **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

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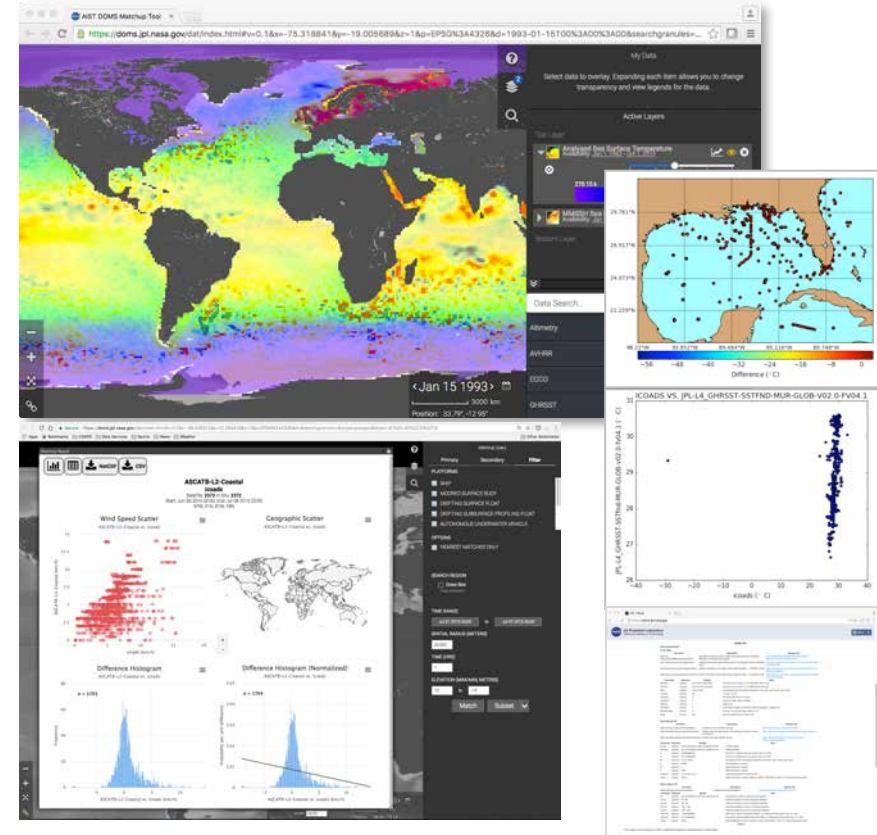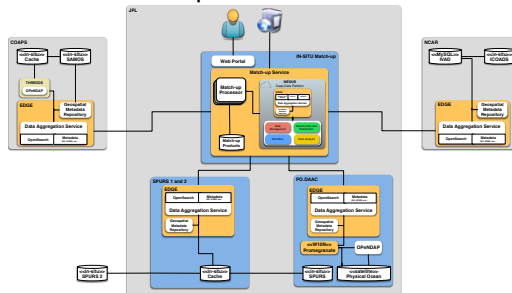
**Search Recommendation**
Based on dataset metadata content and web session co-occurrence

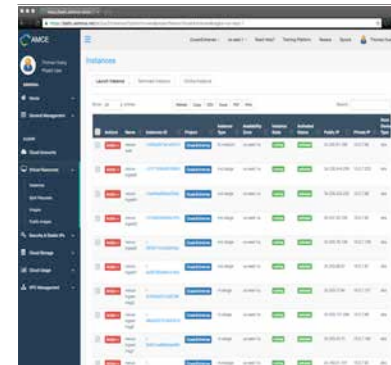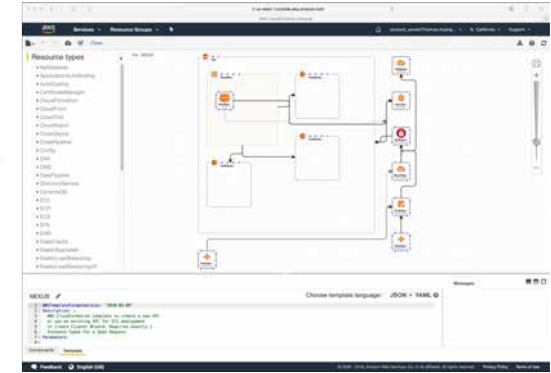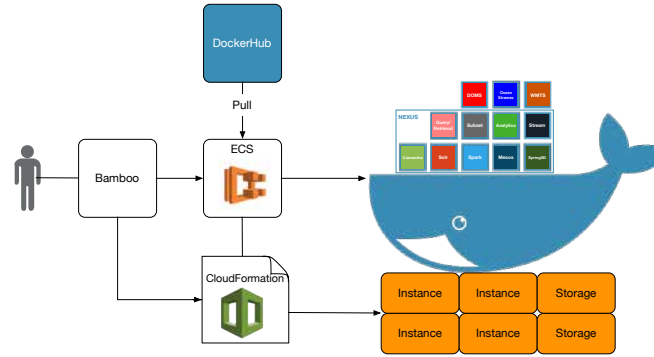# Developing Information Discovery Solutions

# In Situ to Satellite Matchup

- 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.
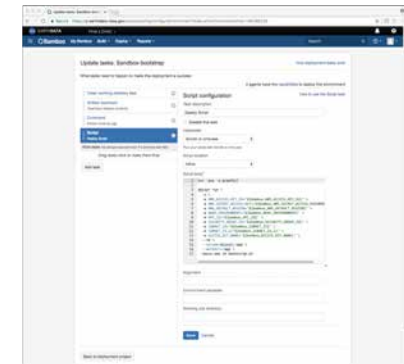- Plugin architecture for in situ data providers

# Deployment Automation

- Cloud Deployment is nontrivial
- Infrastructure Definition
  - Various machine instances
  - Storage and buckets
- Software Deployment.. manually
  - Build
  - Package
  - Install
  - Configure
  - Shell login (security issues)
- Best Practice: Deployment Automation
  - Script Infrastructure Definition (e.g. Amazon CloudFormation)
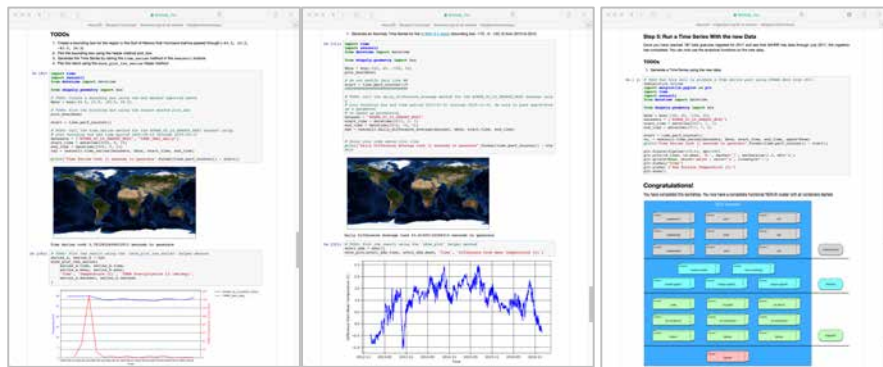  - Container-based Deployment (e.g. Amazon ECS and DockerHub)



AMCE Deployment      NGAP Deployment

# Working with both Science and Informatics Communities

- Established Apache Incubator project
- OceanWorks is developed in the open
- Target Apache top-level project by 2019.
- Public hands-on workshops
- Organize technical sessions at conferences
- Invited speaker and panelist
- Lead Editor: 2018 Wiley Book on **Big Earth Data Analytics in Earth, Atmospheric and Ocean Sciences**
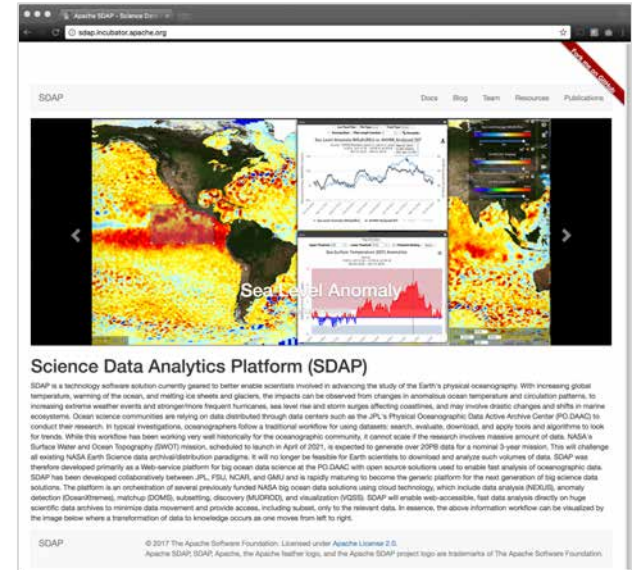


Analyze Hurricane Katrina by comparing SST and TRMM time series

Generate daily difference average
"The Blob" is an oceanographic anomaly

Each participant deployed 3 computing clusters, a total of 24 containers on EC2

# Open Source

- Technology sharing through Free and Open Source Software (FOSS)
- Further technology evolution that is restricted by projects / missions
- **Science Data Analytic Platform (SDAP)**, the implementation of **OceanWorks**, in **Apache Incubator**
  - Cloud platform
  - Analyzing satellite and model data
  - In situ data analysis and coordination with satellite measurements
  - Fast data subsetting
  - Mining of user interactions and data to enable discovery and recommendations
  - Streamline deployment through container technology
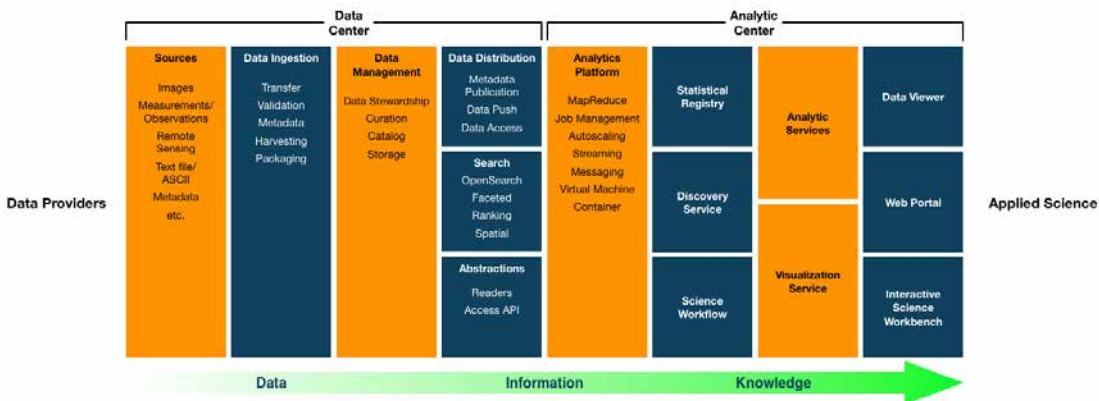


http://sdap.incubator.apache.org

# In Summary

- Traditional method for scientific research (search, download, local number crunching) is unable to keep up
- Think beyond the archive
- Connected information enables discovery
- Community developed solution through open sourcing
- Thanks to the NASA ESTO/AIST and Sea Level Rise programs, and the NASA ESDIS project
- Investment in data and computational sciences
- Data Centers might want to be in the business of Enabling Science!
- OceanWorks infusion 2018 – 2019
  - Watch for changes to the Sea Level Change Portal
    - Even faster analysis capabilities
    - More variety of measurements – satellites, in situ, and models
    - Event more relevant recommendations
  - NASA's Physical Oceanography Distributed Active Archive Center (PO.DAAC)

**Transforming Data to Knowledge**

**National Aeronautics and Space Administration**

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



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

**JPL Team**
Ed Armstrong, Frank Greguska, Joseph Jacob, Lewis McGibbney, Nga Quach, Vardis Tsontos, and Brian Wilson

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

**National Center for Atmospheric Research Team**
Steve J. Worley, Tom Cram, Zaihua Ji

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