## AIST-14: OceanXtremes
## Oceanographic Data-Intensive Anomaly Detection and Analysis Portal
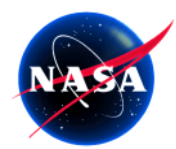
### Earth Science Technology Forum

PI: Thomas Huang

Co-Is: Ed Armstrong, George Chang, (Mike) Toshio Chin, and Brian Wilson

Engineers: Kevin Gill, Frank Greguska, Joseph Jacob, and Nga Quach

June 13, 2016

## Objective

Develop an anomaly detection system which identifies items, events or observations which do not conform to an expected pattern

- Mature and test domain-specific, multi-scale anomaly and feature detection algorithms.
- Identify unexpected correlations between key measured variables.

Demonstrate value of technologies in this service:

- Adapted Map-Reduce data mining.
- Algorithm profiling service.
- Shared discovery and exploration search tools.
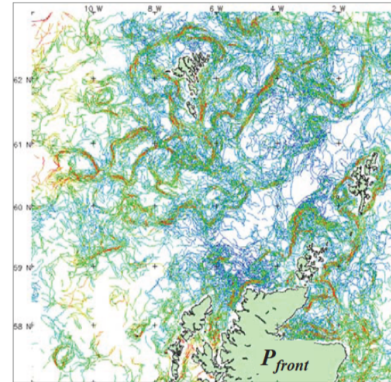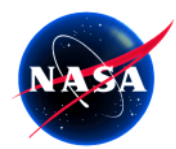- Automatic notification of events of interest.



Illustration of future OceanXtremes analysis capability showing sea surface temperature (SST) gradients from AVHRR imagery (warner colors indicate higher gradient persistence)

## Approach

- Setup on-premise Cloud environment.
- Select dataset and algorithm for anomaly detection.
- Design and develop OceanXtremes backend.
- Validate OceanXtremes using selected datasets and algorithms.
- Design, develop and integrate web portal to backend system.
- Integrate datacasting and visualization capability.
- Expand the number of datasets and algorithms supported within OceanXtremes
- Conduct end-to-end demonstration.

**Co-Is**: E. Armstrong, G. Chang, T. Chin, B. Wilson, JPL

## Key Milestones

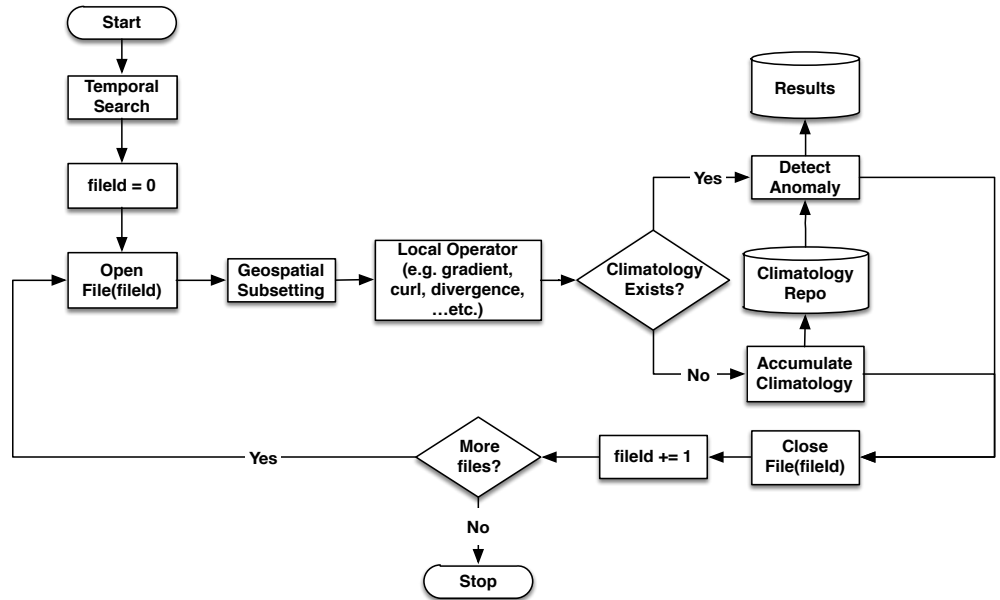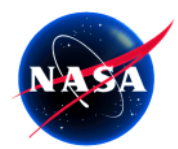| Milestone | Date |
|---|---|
| Complete backend system design | 12/15 |
| Complete testing of backend system | 05/16 |
| Complete web portal design | 08/16 |
| Integrate web portal and backend system | 11/16 |
| Integrate datacasting and visualization capability | 02/17 |
| Collect benchmarking data | 04/17 |
| Conduct end-to-end demonstration | 05/17 |

$TRL_{in} = 2$       $TRL_{current} = 3$

Earth Science Technology Office

# Motivation

- Anomaly detection is a process of identifying items, events or observations outside the "norm" or expected patterns

- Current and future oceanographic missions and our research communities present us with challenges to rapidly identify features and anomalies in increasingly complex and voluminous observations

- Typically this is a two-stage procedure
  1. Determine a long-term/periodic mean ("climatology")
  2. Deviations from the mean are searched. Step 1 could be omitted in cases where a climatology data set already exists.

AIST-14: OceanXtremes

# OceanXtremes Architecture

**Xtremes Ingester**
   Real-time ingestion system
**Xtremes Climatology**
   Batch-oriented climatology
   computation service
**Xtremes Processor**
   Horizontal-scale system for
   anomaly computation and
   detection
**Xtremes Analyzer**
   Webservice to access data
   and anomalies
**Xtremes Visualizer**
   Web service for data
   visualization
**Xtremes Speaker**
   Feed generation and
   management system
**Xtremes Explorer**
   Web-based data visualization
   and analysis



Xtremes Explorer

# Deep Data Computing Environment (DDCE)



**DDCE is OceanXtremes' development environment. It consists of**

- **CloudWork**: A Mirantis OpenStack private cloud computing environment
- **DeepData**: A high-performance data cluster with locally attached storages
- High speed switches

Earth Science Technology Office

# NEXUS Deep Data Analytics: One-Minute Summary

**NEXUS** is an emerging technology developed at JPL
- A Cloud-based/Cluster-based data platform that performs scalable handling of observational parameters analysis designed to scale horizontally by
  - Leveraging high-performance indexed, temporal, and geospatial search solution
  - Breaks data products into small chunks and stores them in a Cloud-based data store

**Data Volumes Exploding**
- NISAR & SWOT missions coming
- File I/O is slow

**Scalable Store & Compute is Available**
- NoSQL cluster databases
- Parallel compute, in-memory map-reduce
- Bring Compute to Highly-Accessible Data

**Pre-Chunk and Summarize Key Variables**
- Easy statistics instantly (milliseconds)
- Harder statistics on-demand (in seconds)
- Visualize original data (layers) on a map quickly

**A growing collection of data analytic microservices**



NEXUS
Deep Data Platform

EDGE

| Search and Access | Metadata | Analytic |

Data Aggregation Service

Geospatial Metadata Repository

Data Management

Data Access and Distribution

Workflow

Data Analysis

Earth Science Technology Office

# Analytics & Summarization of Stack


Display Variables on Map


Latitude-Time Hovmoller


Plot Aggregate Statistics

**Solr DB Cluster**

| Chunk | Chunk | Chunk |
|---|---|---|
| Chunk | Chunk | Chunk |
| Chunk | Chunk | Chunk |

...

**Fast & Scalable**

| Meta Data | Meta Data | Meta Data | Meta Data |
|---|---|---|---|

...

Metadata (JSON): Dataset and granule metadata, Spatial Bounding Box & Summary Statistics

**Cassandra DB Cluster & Spark In-Memory Parallel Compute!**

Custom Analytics

Subset Variables & Chunk Spatially

Each file contains many high-resolution geolocated arrays

| SMAP | | | | |
|---|---|---|---|---|

| MODIS | | | |
|---|---|---|---|

| GRHSST | | | |
|---|---|---|---|

| JASON | | | |
|---|---|---|---|

**Slow File I/O**    30-Year Time Series of archival HDF & netCDF files (daily or per orbit)

Earth Science Technology Office

# Enable Ocean Science



"The Blob is a result of a high pressure system that has parked itself in the Gulf of Alaska for the past few years that has driven the polar jet stream north into northern Canada and then it plunged rapidly out of northern Canada into the American Midwest and northeast. And so the result was hot dry winters on the west coast, and fierce winters with heavy snow pack in the Midwest." – Bill Patzert, NASA/JPL

# Xtremes Explorer



High Resolution Data Visualization for the Web



Data Analysis Workbench

Aug 02, 2012



Aug 02, 2013



Aug 02, 2014



Aug 02, 2015

# The Notebook

Interact with OceanXtremes using Jupyter Notebook



- **/capabilities**: list of capabilities
- **/chunks**: list data chunks by location, time, and datasets
- **/correlationMap**: Correlation Map
- **/datainbounds**: Matchup operation to fetch values from dataset within geographic bounds
- **/datapoint**: Matchup operation to fetch value at lat/lon point
- **/dailydifferenceaverage**: Daily difference average
- **/latitudeTimeHofMoeller**: Latitude Time Hovmoeller
- **/list**: list available datasets
- **/longitudeLatitudeMap**: Longitude Latitude Map
- **/longitudeTimeHofMoeller**: Longitude Time Hovmoeller
- **/stats**: Statistics (standard deviation, count, min/max, time, mean)

# Data Tiling Scheme

- Pre-processing occur during ETL phase
- Breaking geospatial arrays into small geo-addressable data chunks (or partitions)
- Tile → small → in memory processing
- All spatial indexes are managed by Apache Solr

**Tiling Algorithm**

$c = Number\ of\ tiles\ desired$
$d = Number\ of\ dimensions$
$L_d = Length\ of\ dimension\ d$
$S_d = Step\ size\ for\ dimension\ d$

$$S_d = \left\lfloor \frac{L_d}{\sqrt[d]{c}} + \frac{1}{2} \right\rfloor$$

**MUR Data in 0.01 degrees, Tiles 2.5° x 5°**

$c = 5184$
$d = 2$
$L_{latitude} = 17999$
$L_{longitude} = 36000$

$$S_{latitude} = \left\lfloor \frac{17999}{\sqrt[2]{5184}} + \frac{1}{2} \right\rfloor = \left\lfloor 249.986111111 + \frac{1}{2} \right\rfloor = 250$$

$$S_{longitude} = \left\lfloor \frac{36000}{\sqrt[2]{5184}} + \frac{1}{2} \right\rfloor = \left\lfloor 500 + \frac{1}{2} \right\rfloor = 500$$

# Real-time Ingestion Solution

**A real-time data ingestion system**

1. Data discovery
2. Metadata extraction
3. Data partition (tiles)
4. Pre-compute metrics
5. Register to NEXUS

**Core components**

- Admin
- Containers
- High-performance message broker
- Distributed synchronization service

**Deployed under OpenStack Cloud**

**18 virtual instances**



Real-time Data Ingestion Architecture



High-level Ingestion Workflow

# Investigated Parallel Performance

**Four technologies:**

| | |
|---|---|
| Multicore on single node | 8 core = 8-way parallelism |
| PySpark on YARN scheduler | 8 nodes x 4 cores on each |
| PySpark on Mesos scheduler | 32-way parallelism |
| DPark on Mesos (fastest) | |

**Multiple runs over different numbers of tiles**

- Query for tiles that intersect a user-chosen lat/lon rectangle and time range
- Multiple rectangles:  1, 5, 10, 30, and 90 degree lat/lon boxes

**Vary number of partitions to keep cores busy (> 2-3X)**

- 32, 64, 128, 256  (128 best, 256 saturates)

# Performance Benchmark

## Time-Series Generation Performance



| Environment | 1x1 | 5x5 | 10x10 | 30x30 | 90x90 |
|---|---|---|---|---|---|
| Spark on YARN | 7.562 | 36.663 | 106.101 | 118.678 | 121.306 |
| DPark on Mesos | 5.638 | 29.353 | 96.799 | 103.839 | 107.826 |

**DPark on Mesos is <u>fastest</u> and scales with # of tiles**

- Mesos vs. YARN:  shorter startup time, faster task scheduling
- DPark:  no data movement between python runtime and JVM
- As # of tiles grows, 7 nodes x 4 cores all kept busy.

Earth Science Technology Office

# Climatology Algorithm: Gaussian Interpolation

- Armstrong, E. and J. Vaquez-Cuervo, A New Global Satellite-Based Sea Surface Temperature Climatology, *Geophysical Research Letters* Volume 28, No. 22, Pages 4199-4202, November 15, 2001

- A time/space Gaussian interpolation to generate global sea surface temperature climatology

- The Fortran-based implemented was ported to execute on the Deep Data Computing Cluster

- Python wrapper is being implemented to simplify integration into Xtremes Climatology

- Allow users to rapidly create regional and custom period climatologies for SST, wind etc.



Sea Surface Temperature

Sea Surface Temperature (degree_C)

-2.0E+00    5.4E+00    1.3E+01    2.0E+01    2.8E+01    3.5E+01

Data Min = -2.0E+00, Max = 3.5E+01

SST - 4km
Climatology 2002 - 2016

# Algorithm: Empirical Orthogonal Function (EOF)

- EOF analysis is widely used in meteorology and oceanography to extract dominant modes of behavior in scalar and vector datasets
- EOF is computationally demanding, especially for large high resolution data sets, e.g. MUR SST.
- Typical workstations don't have the capacity to handle "high resolution" and long time series datasets available at PO.DAAC
- Key algorithmic and implementation issue - parallelized/distributed computation of the Singular Value Decomposition (SVD).
- For efficient SVD computation, we have experimented with a Lanczos algorithm that allows computation of only the most dominant eigen/singular values/vectors.
- Codes from an externally maintained package, the ARnoldi PACKage (ARPACK), have been used in prototyping works. These codes/algorithms still need to be realized in an efficient way for our "tiled" data structure.

# Data Sources

| Phenomenon | Dataset | Key Variables | Time Range | Data Mining Operators Needed |
|---|---|---|---|---|
| El Nino genesis, anomaly detection and characterization in different regions (3.4 vs 4). Coastal upwelling | CCMP L4 | Wind | 1987-2015 | Anomaly calculation from fixed or on-the-climatology, Threshold detection. Variance characterization |
| | Integrated Altimeter L4 | SSH | 1992-2013 | |
| | MODIS Aqua/Terra L3 | SST | 2000-present | |
| | AVHRR_OI L4 | SST | 1982-present | |
| | MUR L4 | SST | 2002-present | |
| El Nino and other teleconnections. Regional correlations | CCMP L4 | Wind | 1987-2015 | Cross correlations. Covariabilty and EOFs. |
| | Integrated Altimeter L4 | SSH | 1992-2013 | |
| | MODIS Aqua/Terra L3 | SST | 2000-present | |
| | AVHRR_OI L4 | SST | 1982-present | |
| | MUR L4 | SST | 2002-present | |
| | Aquarius L3 | Salinity | 2011-present | |
| | MODIS Aqua L3 | Chl A | 2002-present | |

# Data Sources

| Phenomenon | Dataset | Key Variables | Time Range | Data Mining Operators Needed |
|---|---|---|---|---|
| Upwelling. Hurricane genesis | CCMP L4 | Wind | 1987-2015 | Divergence and curl. |
| | MODIS Aqua/Terra L3 | SST | 2000-present | |
| | AVHRR_OI L4 | SST | 1982-present | |
| | MUR L4 | SST | 2002-present | |
| Gradients, edges, and eddy detection | MODIS Aqua/Terra L3 | SST | 2000-present | Matched filter (e.g., Sobel operator). First derivatives. |
| | MUR L4 | SST | 2002-present | |
| | MODIS Aqua L3 | Chl A | 2002-present | |
| Trends. Basin scale variability | CCMP L4 | Wind | 1987-2015 | Regression, Polynomial fits. Variance. |
| | Integrated Altimeter L4 | SSH | 1993-2013 | |
| | MUR L4 | SST | 2002-present | |

# Engagements

April 2015: PO.DAAC User Working Group

June 2015/2016: Earth Science Technology Forum

July 2015: ESIP Federation Summer Meeting

October 2015: IEEE Big Data Conference

December 2015: American Geophysical Union Fall Meeting

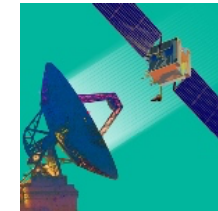January 2016: ESIP Federation Winter Meeting
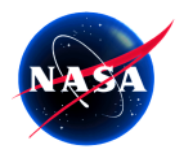
February 2016: Ocean Sciences Meeting

March 2016: Ground System Architectures Workshop

March 2016: PO.DAAC User Working Group

July 2016: ESIP Federation Summer Meeting

October 2016: International Conference on Marine Data and Information Systems – Gdanski, Poland

# Near-term Plan

**Data Services**

- Xtremes Ingester
  - Improve tiling performance and additional tile-level stats
- Xtremes Processors
  - MapReduce framework
  - Automatic detection workflow
- Xtremes Analyzer: Search and metadata capabilities
- Xtremes Speaker: Datacasting feed management
- Docker deployment process

**Science and Algorithms**

- Catalog know anomalies (e.g. El Nino, hurricane, etc)
- Empirical Orthogonal Function (EOF)

**Web Portal**

- More visualizations
- Anomaly search
- User-defined anomaly detection

**Datasets**

- MODIS Terra L3 – SST and Chl A
- CCMP L4 - Wind
- Integrated Altimeter L4 - SSH
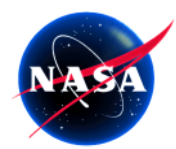- Aquarius L3 - Salinity

# Spark and Resource Management

**Issue:** From our benchmark comparison, we have concluded the common Spark + YARN combination, while it is faster than Hadoop, the bridge to PySpark with YARN don't yield the desired performance.

*PySpark is a python wrapper on Spark, which is implemented in Scala (Java).  Data is being copied between Java memory space to Python memory space.  Python, because of numpy and scipy, is still the leading programming language for scientific programing*

*YARN got popular with Hadoop in the Cloudera distribution.  It works well with Hadoop, but we discovered the scheduling overhead with YARN is less than desirable.*
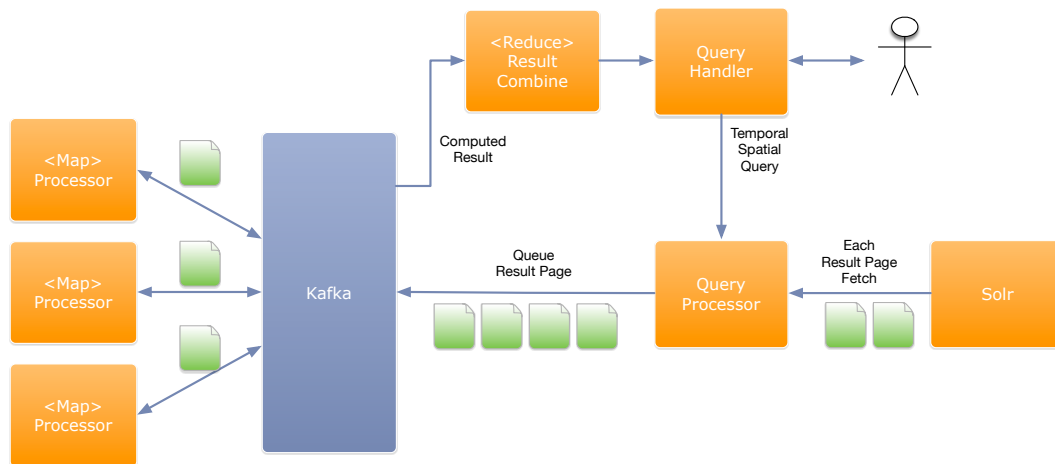
*DPark is pure python implementation of Spark.  Our benchmarking shows DPark + Mesos out performs PySpark + YARN by ~30%.*
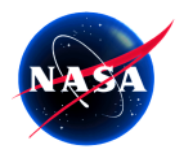
# High-throughput Distributed Processing

**Issue:** Result retrieval from Solr could create huge performance bottleneck. What happen when a temporal-spatial query returns 1M matches. Current implementation fetches all 1M matches before start processing.

*A new high-throughput distributed processing framework is developed to farm jobs for each Solr page fetch. It frees the system from high memory utilization and also increase parallelism, which yields faster response.*

# Questions?