# Detecting radio-astronomical "Fast Radio Transient Events" via an OODT-based metadata processing pipeline
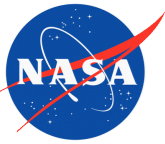
**Chris Mattmann, Andrew Hart , Luca Cinquini**
**David Thompson, Kiri Wagstaff, Shakeh E. Khudikyan**
NASA Jet Propulsion Laboratory,
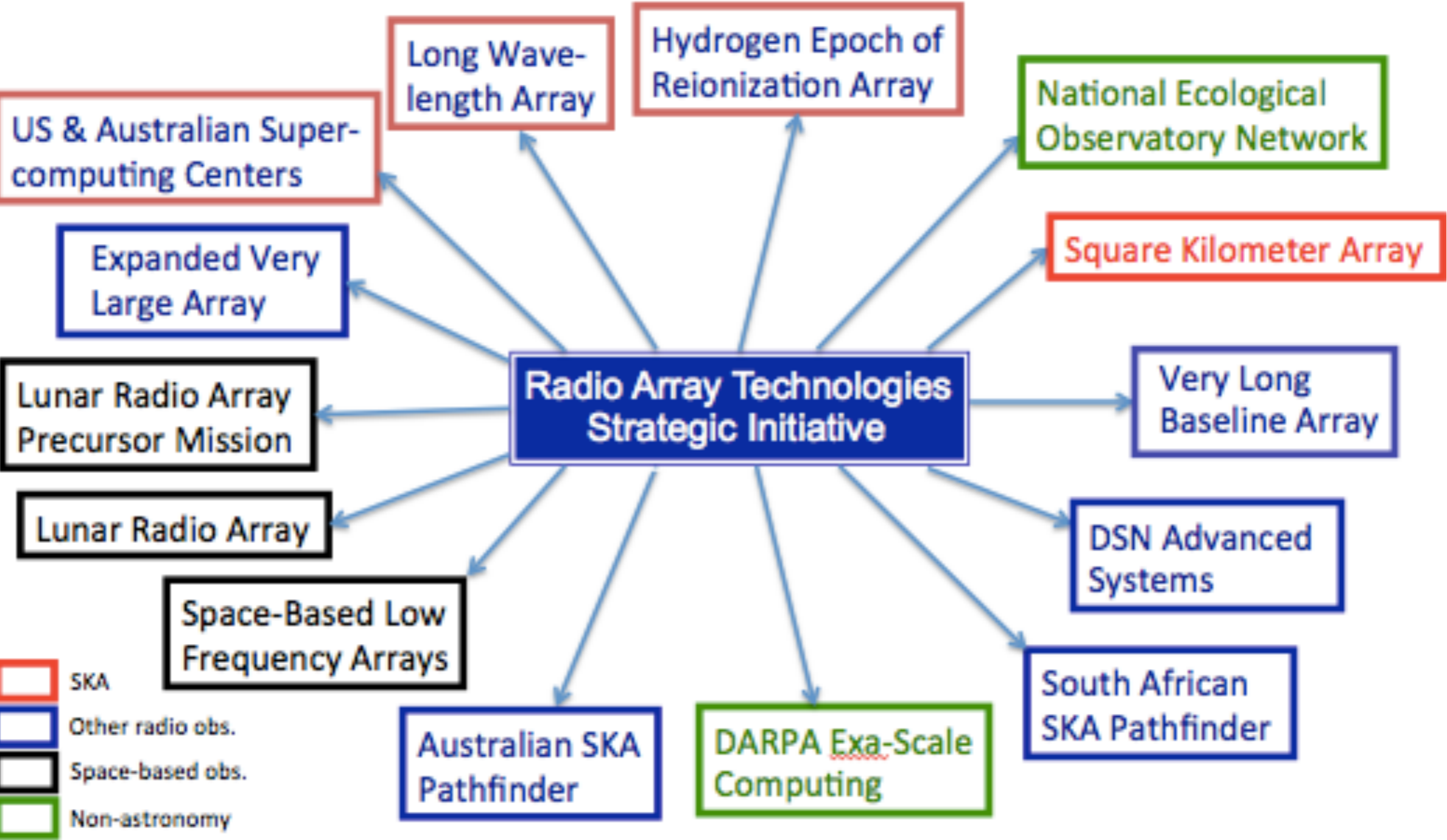California Institute of Technology

- **<u>Initiative Lead: Dayton Jones; Champion: Robert Preston</u>**
- **We will define the necessary data services and underlying substrate to position JPL to compete for and lead "big data" management efforts in astronomy, specifically, SKA, HERA, SKA precursors, and NRAO.**

- **Perform prototyping and deployment to demonstrate JPL's leadership in the "big data" and astronomy space.**

- **Collaborate on Data Products and Algorithms from Adaptive Data Processing task**

- **Establish partnerships with major SKA potential sites and pre-cursor efforts (South Africa, Australia)**

- The Big Picture
  - Astronomy, Earth science, planetary science, life/ physical science all *drowning in data*
  - Fundamental technologies and emerging techniques in archiving and data science
    - Largely center around *open source* communities and related systems
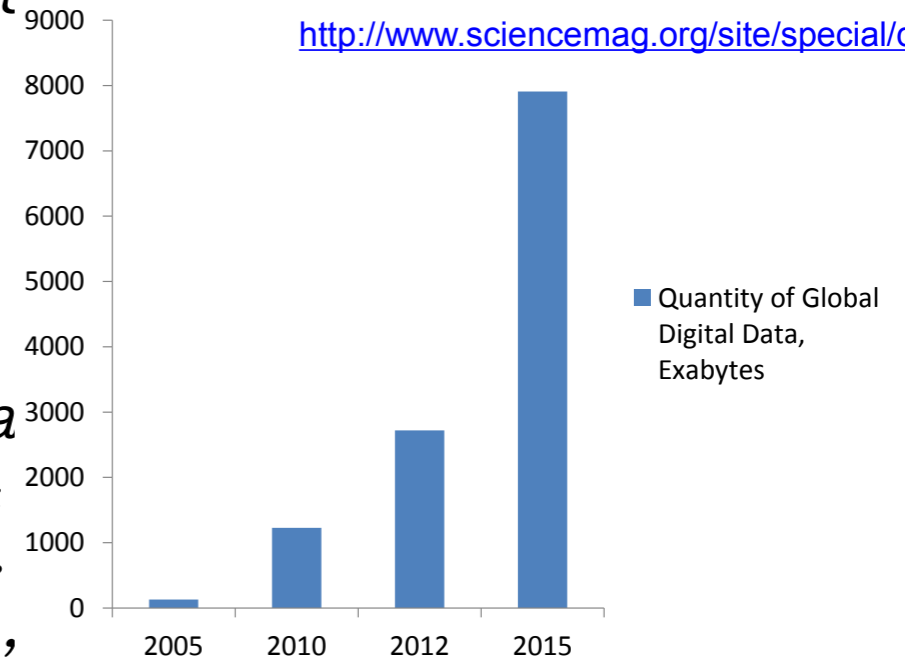- Research challenges (adapted from NSF)
  - *More data is being collected than we can store*
  - *Many data sets are too large to download*
  - *Many data sets are too poorly organized to be useful*
  - *Many data sets are heterogeneous in type, structu*
  - *Data utility is limited by our ability to use it*

- Our Focus: <u>Big Data Archiving</u>
  - *Research methods for integrating intelligent algorithms for data triage, subsetting, summariza*
  - *Construct technologies for smart data movement*
  - *Evaluate cloud computing for storage/processing*
  - *Construct data/metadata translators "Babel Fish'*



http://www.sciencemag.org/site/special/data/

■ Quantity of Global Digital Data, Exabytes

Source: EMC/IDC Digital Universe Study, 2011
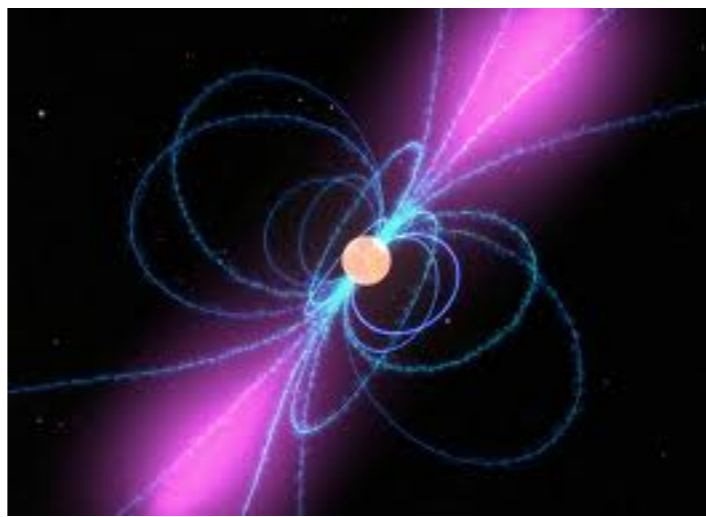
# Some "Big Data" Grand Challenges

- *How do we handle 700 TB/sec of data coming off the wire when we actually have to keep it around?*
  - Required by the Square Kilometre Array - will talk about **data triage** here

- *Joe scientist says I've got an IDL or Matlab algorithm that I <u>will not change</u> and I need to run it on 10 years of data from the Colorado River Basin and store and disseminate the output products*
  - Required by the Western Snow Hydrology project

- *How do we compare petabytes of climate model output data in a variety of formats (HDF, NetCDF, Grib, etc.) with petabytes of remote sensing data to improve climate models for the next IPCC assessment?*
  - Required by the 5th IPCC assessment and the Earth System Grid and NASA

- *How do we catalog all of NASA's current planetary science data?*
  - Required by the NASA Planetary Data System

Image Credit: http://www.jpl.nasa.gov/news/news.cfm?release=2011-295

VFASTR ("VLBA Fast Radio Transients") is a project that aims at detecting short radio pulses (approx. a few milliseconds) from extra-terrestrial sources within the large amounts of data collected by the VLBA ("Very Large Baseline Array")

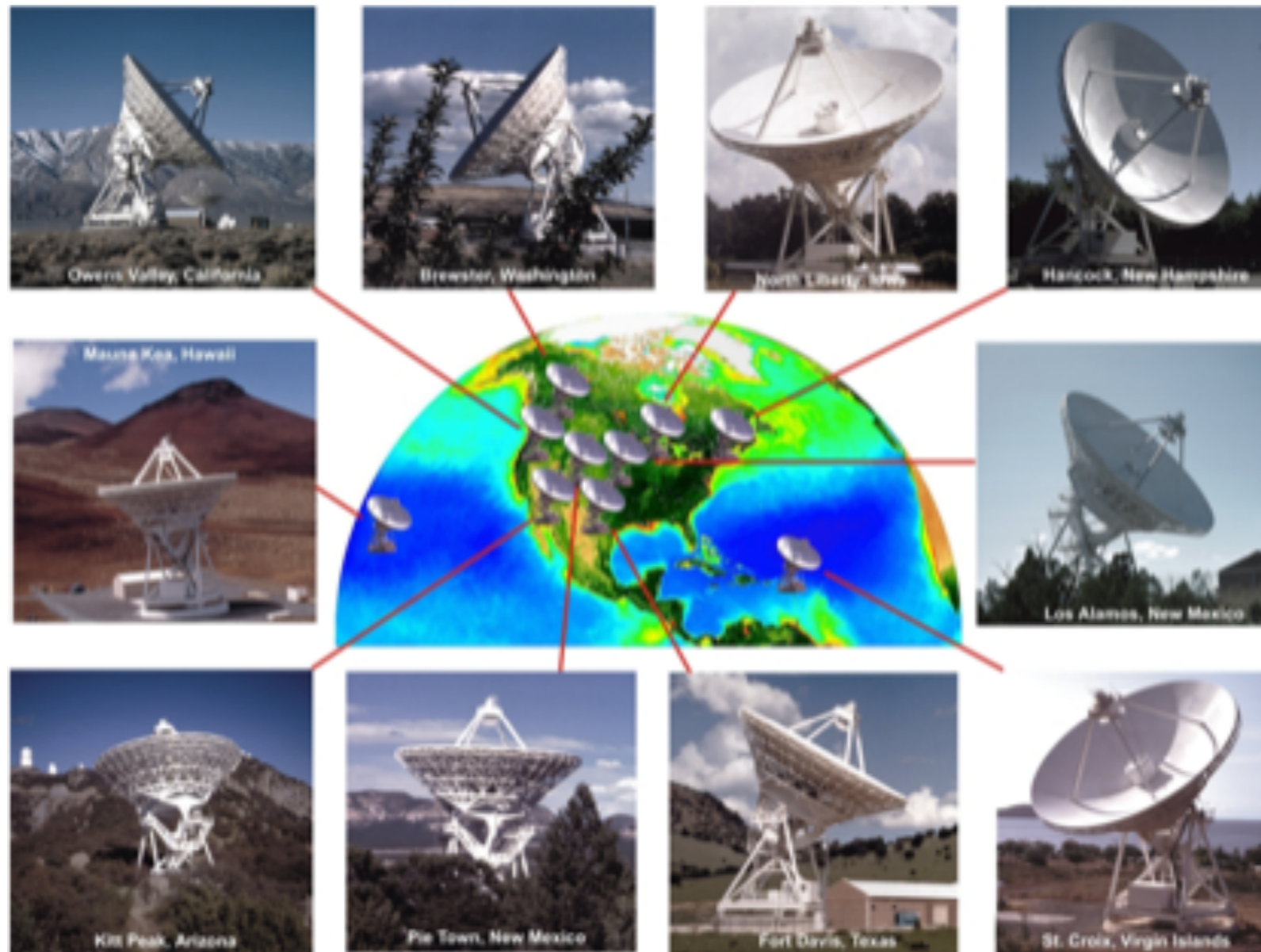Fast Radio Transients may be generated by <u>known</u> and yet <u>unknown</u> sources:

- Pulsars
- Intermittent pulsars
- X-Ray binaries
- Supernovae

- Merging neutron stars
- Annihilating black holes
- ET signals ?
- New deep space objects

VFASTR is one of a new generation of Radio Astronomy experiments that aim at analyzing the "dynamic radio sky" as opposed to mapping and inspecting known static sources

VLBA ("Very Large Baseline Array") is a group of 10 large radio-telescopes (25m diameter) distributed across the U.S.A. from Hawaii to the Virgin Islands.
- No two antennas are within each other's local horizon
- The overall array has a baseline of 800 km => resolution of milliarcsecond

VFASTR employs a <u>commensal</u> (a.k.a. "passive") approach by analyzing data that is collected during normal VLBA operations for other scientific purposes:

- Raw voltages from VLBA antennas are transferred to NRAO, time-correlated, corrected for dispersion through the interstellar medium ("de-dispersion") and separated from instrument noise
- Candidate events are staged on disk and remain available for limited time
- VFASTR team must review tens of candidates daily, archive the promising ones and disregard the others



VFASTR Science Team: Astron (The Netherlands), ICRAR (Australia), JPL (U.S.A.), NRAO (U.S.A.)

# VFASTR Data System Overview

- The software engineering team at JPL has developed an end-to-end data system in support of VFASTR activities with two major <u>goals</u>:
  - ▸ provide a web-based platform for easy and timely review of candidate events by the science team
  - ▸ enable the automatic identification of interesting events by a self-trained machine agent

- The system is composed of three major <u>components</u>:
  - ▸ <u>Data processing pipeline</u>: responsible for data transfer from NRAO and archiving at JPL, and for metadata extraction and cataloging
  - ▸ <u>Web portal</u>: easy accessible application for display of product data and metadata, and selection and tagging of interesting events
  - ▸ <u>Data mining algorithm</u>: analyzes the events pool and tags candidates with characteristics similar to sample interesting events

The VFASTR data processing pipeline was built by using Apache OODT in combination with other Apache and Open Source technologies



OODT: Object Oriented Data Technology: framework for management, discovery and access of distributed data resources. Main features:
• Modularity: eco-system of standalone components that can be deployed in various configurations to fulfill a project specific requirements
• Configurability: each component can be easily configured to invoke alternate out-of-the-box functionality or deployment options
• Extensibility: components can be extended by providing alternate implementations to its core APIs (expressed as Java interfaces) or configuring custom plugins

OODT is used operationally to manage scientific data by several projects in disparate scientific domains:

- Earth Sciences:
  - ▶ NASA satellite missions (SMAP,...) are using OODT components as the base for their data processing pipeline for generation and archiving of products from raw observations
  - ▶ ESGF (Earth System Grid Federation) used OODT to build and publish observational data products in support of climate change research
- Health Sciences: EDRN (Early Detection Research Network) uses OODT to collect, tag and distribute data products to support research in early cancer detection
- Planetary Science: PDS (Planetary Data System) is developing data transformation and delivery services based on OODT as part of its world-wide product access infrastructure
- Radio Astronomy: several projects (ALMA, Haystack,...) are adopting OODT based on successful VASTR example

# VFASTR Data System Architecture

Data products are continuously generated by the VLBA ground and processing system and stored on temporary cache at NRAO.



Data products are transferred to JPL where metadata is extracted, products are made available for review by scientists, sub-selected for further analysis, and tagged.

- Some of the architectural decisions that factored in the data system design were motivated by specific project constraints:
  - ▸ <u>Minimize impact on NRAO resources</u>: because VFASTR is a "guest" project at NRAO, attention had to be paid to limit use of disk storage, network bandwidth and CPU resources
  - ▸ <u>Security</u>: all NRAO resources were exposed as "read-only": no action initiated at JPL could result in any modification of the original products (or compromise the NRAO system)

- Architecture evolved over time as a result of new requirements such as increased data volumes and higher frequency updates
  - ▸ Use of different OODT and Apache components (Lucene vs <u>MySQL data store</u> back-ends, <u>Solr</u> for fast metadata retrieval)
  - ▸ Development of new OODT functionality (<u>RESTful API for metadata updates</u>)

VFASTR data is logically organized into three levels: jobs, scans and events

- <u>Job</u>: a batch of data that is processed at one time, and stored together on physical disk. They are associated with a specific investigator scientist.  Each contains 1-100+ scans.
- <u>Scan</u>: a physical telescope pointing, e.g. a period where the antennas are all directed to a common point on the sky.  They have durations of 1-100+ seconds.
- <u>Event</u>: a time segment that the system thinks is interesting.  Duration is usually about 1-2 seconds.  Most scans have no such events, but some have a dozen or more. The "interesting" part of an event is much shorter: 5-50 milliseconds.



<u>VFASTR Data Product</u>: directory tree containing all data recorded for a single job ("tns_bmp360p2_44")
- Job calibration files
- Scan output files
- Event raw voltages
- Event reconstructed images
- …and other files….
- Approximately 1-100 GB

Telescope signal processing:
- Time correlation
- "De-dispersion" (i.e. corrected for dispersion in interstellar medium)
- "Adaptive excision" (some telescopes are disregarded based on self-learning algorithm)
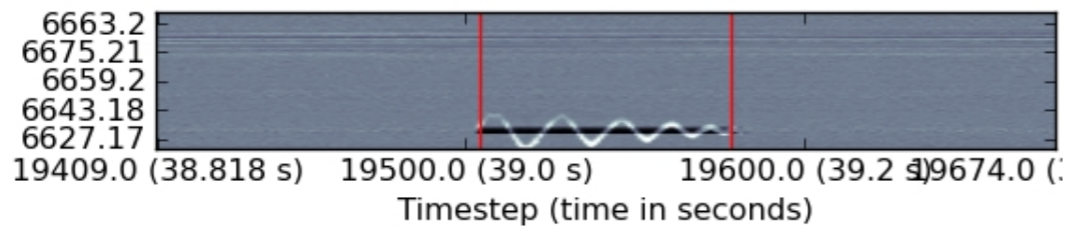
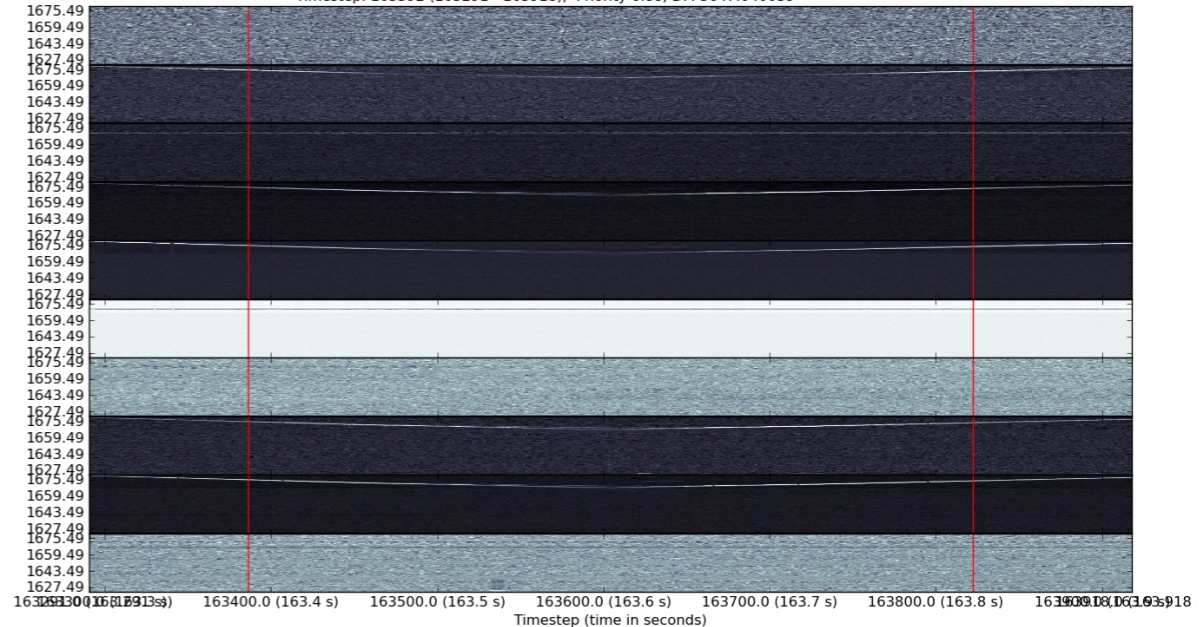- RFI events only detected by one or two antennas, have terrestrial origin

"V-chirp"



Job: G011_02    Scan 8  Antennas: All, Polarization: Sum
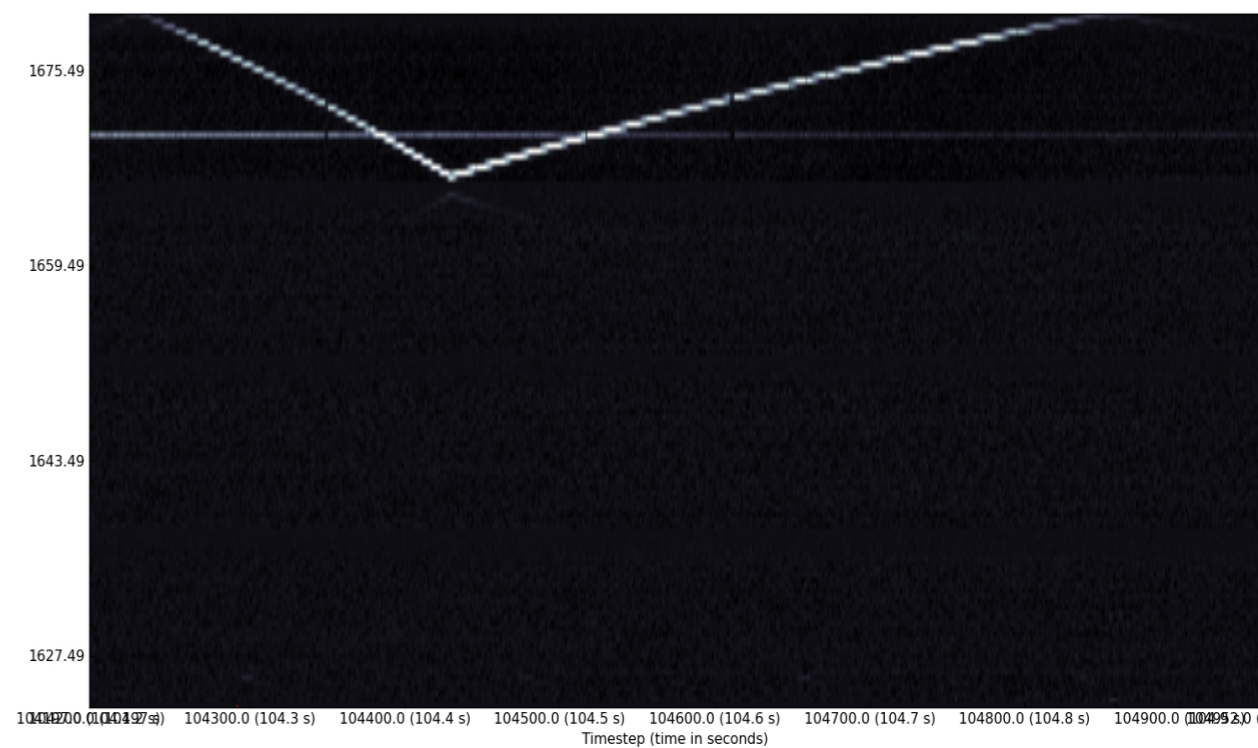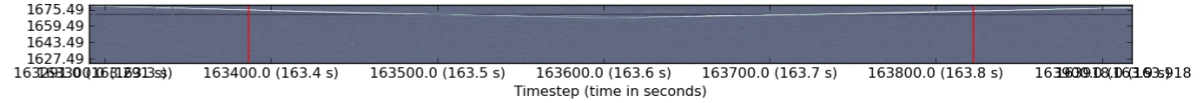Timestep: 19539 (19409 - 19674);  Priority 8.80; DM 4579.450195



Job: inbeam2nd_1    Scan 0  Antennas: All, Polarization: Sum
Timestep: 163392 (163291 - 163918);  Priority 6.88; DM 3647.040039

# Rsync

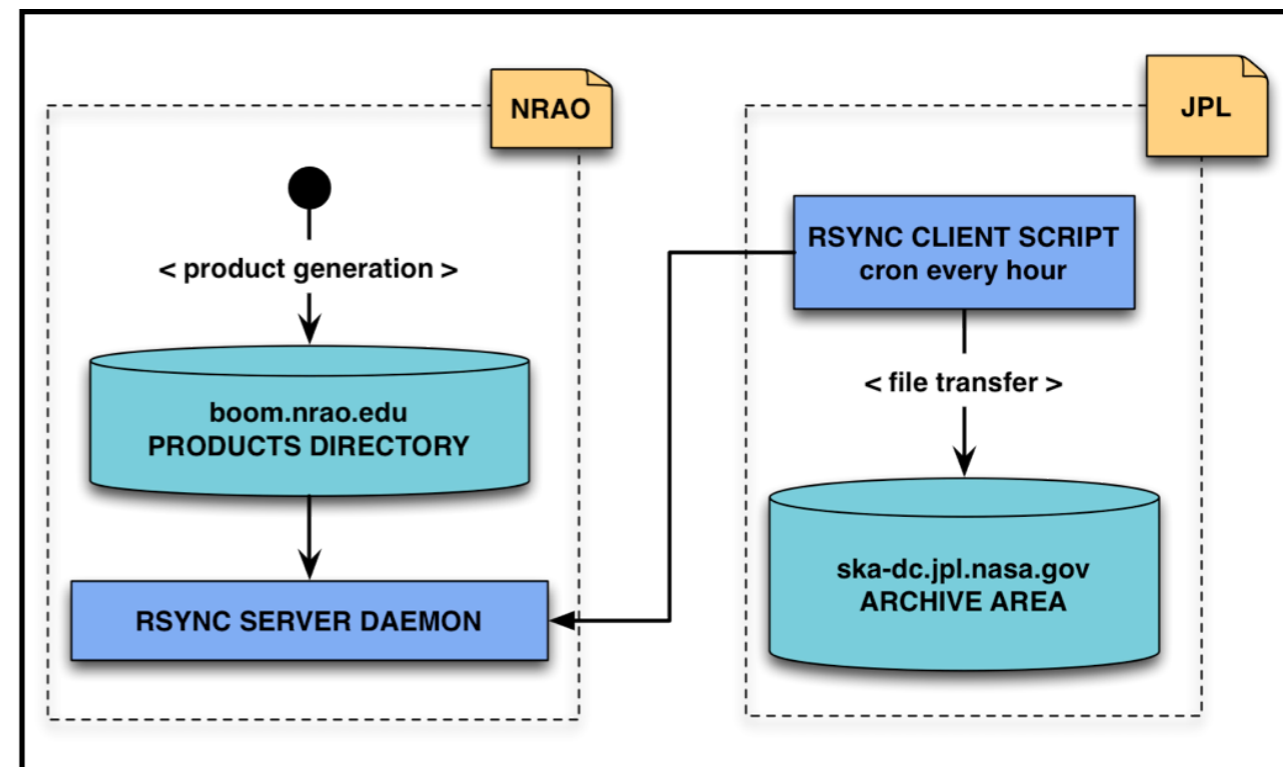Rsync: freely available utility for Unix systems that can be used to synchronize the content of directory trees between two hosts with minimal human intervention.

Features:
- Easy deployment
- Extensive range of configuration options
- High performance: only file changes are transferred ("delta encoding") between sub-sequent invocations, + optional compression
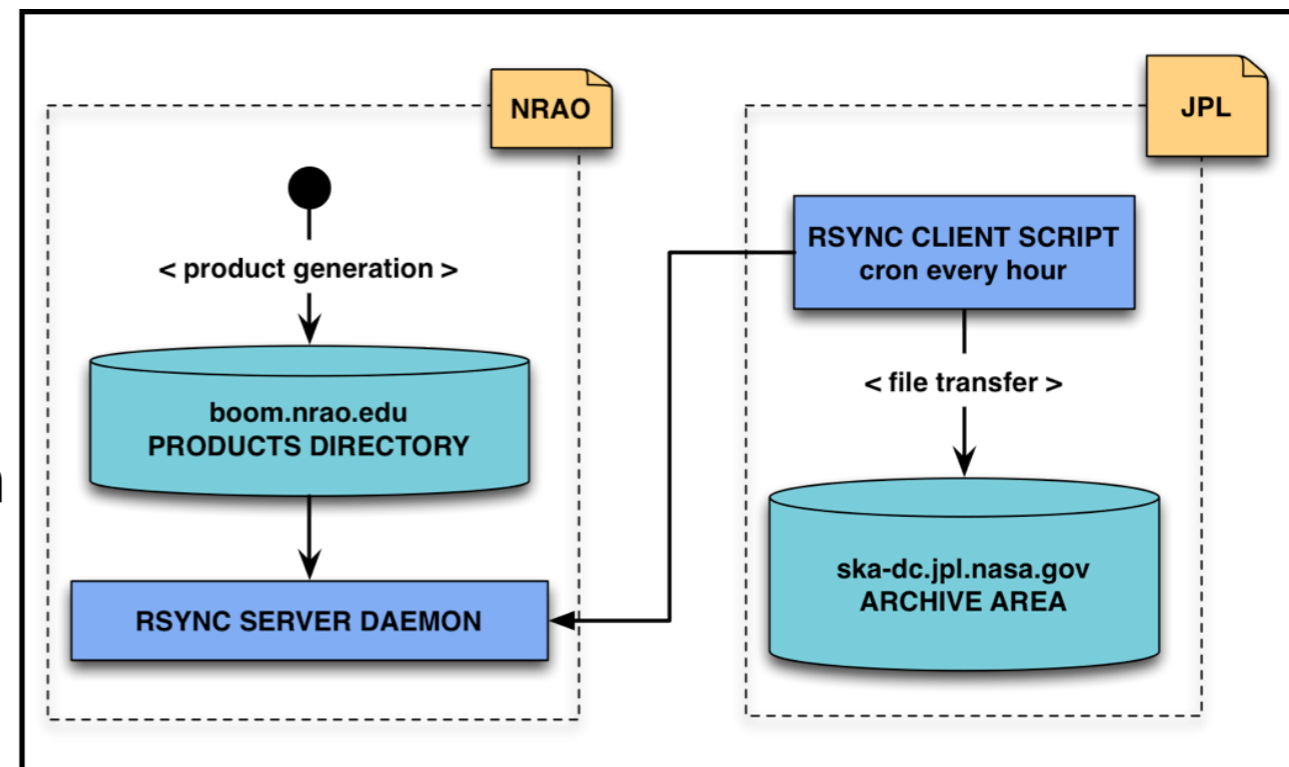- Optional recursion into sub-directories
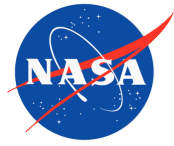- Reliability: "turn-key" toolkit

# Rsync

VFASTR deployment:

- rsync server daemon was deployed at NRAO to make VFASTR products available for download
  - ▶ Configured for read-only mode
  - ▶ Limited to requests coming from JPL IPs
- rsync client running at JPL as system cron job to pull data every hour
  - ▶ Configured to only transfer a subset of the product files (images, output, calibration files)

Measured Data Transfer Rates:

- ~ 2MB/sec between NRAO and JPL
- Approximately 10-20 products per day
- Average volume for transferred product: ~50MB (reduced from 50GB)
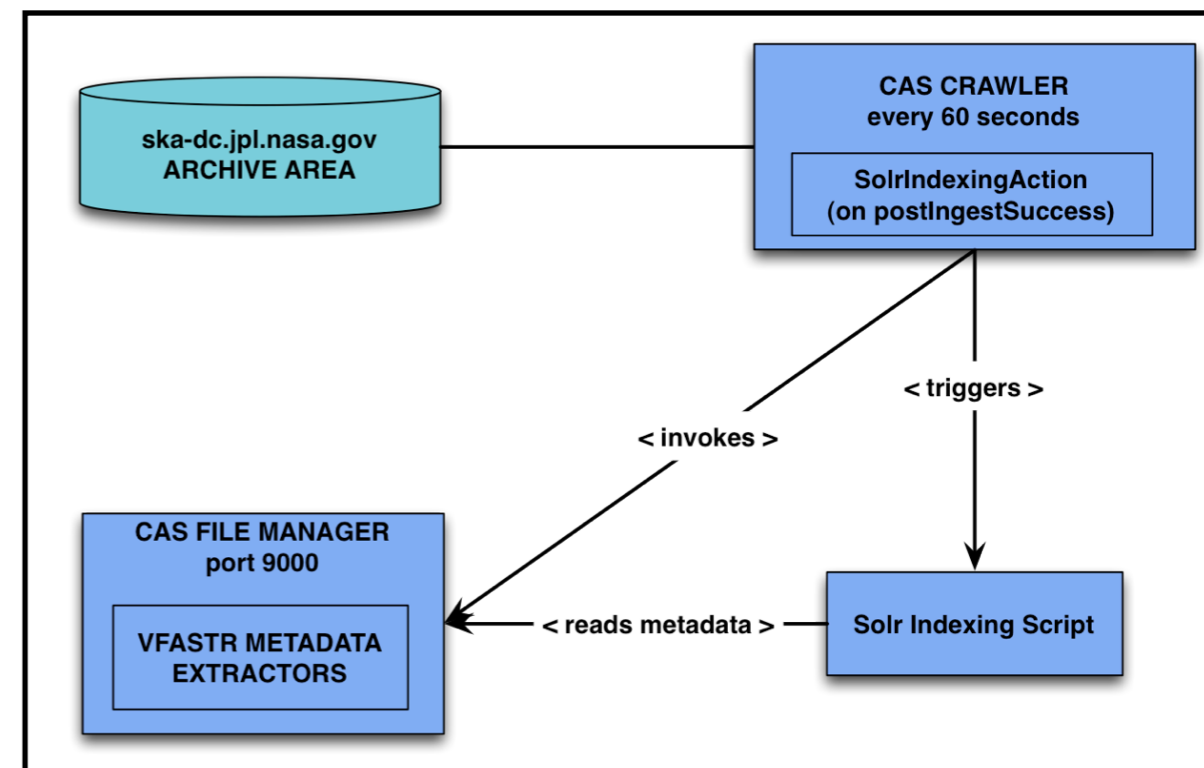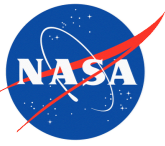- Can transfer all (reduced) daily products in a few minutes!

The CAS Crawler is an OODT component that can be used to list the contents of a staging area and submit products for ingestion to the CAS File manager. Typically used for automatic detection of new products transferred from a remote source.

VFASTR deployment:
- Run as daemon every 300 seconds
- In-place archiving of products (no movement)
- Preconditions:
  ‣ Product must be complete
  ‣ Product must be no older than 10 days
  ‣ Product must not exist in catalog already
- Post-ingest action on success:
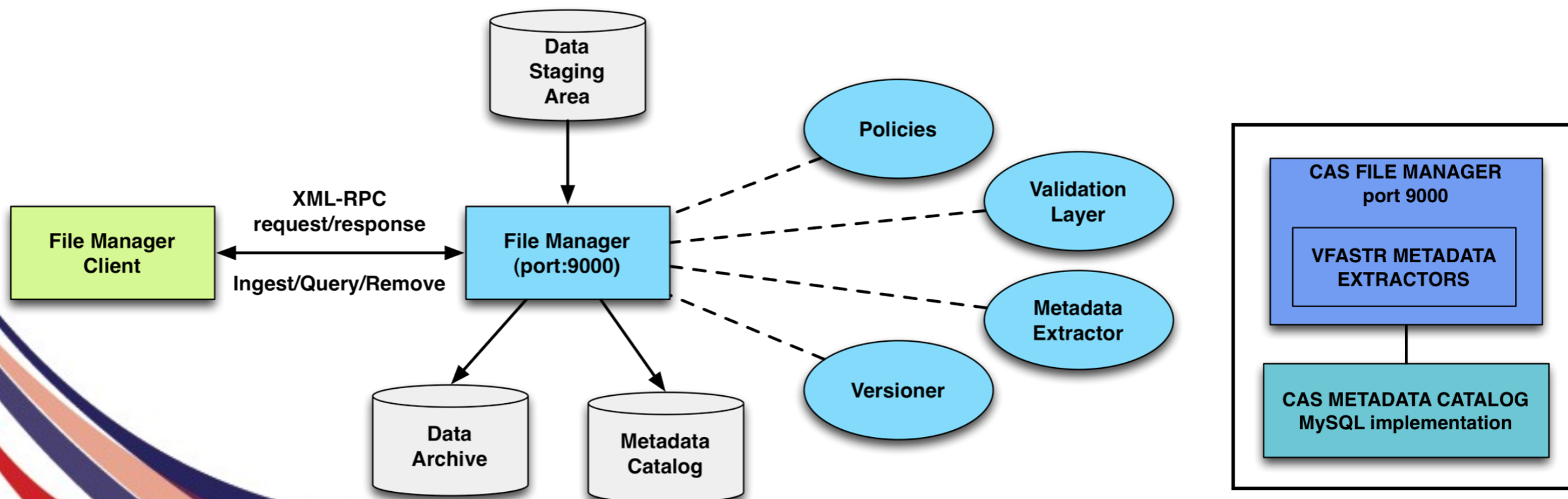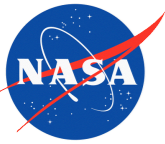  ‣ Trigger metadata harvesting by Solr script

The CAS File Manager is an OODT service for cataloging, archiving and delivery of data products (files and directories) and associated metadata. It is used as core data management component in most OODT-based data systems.
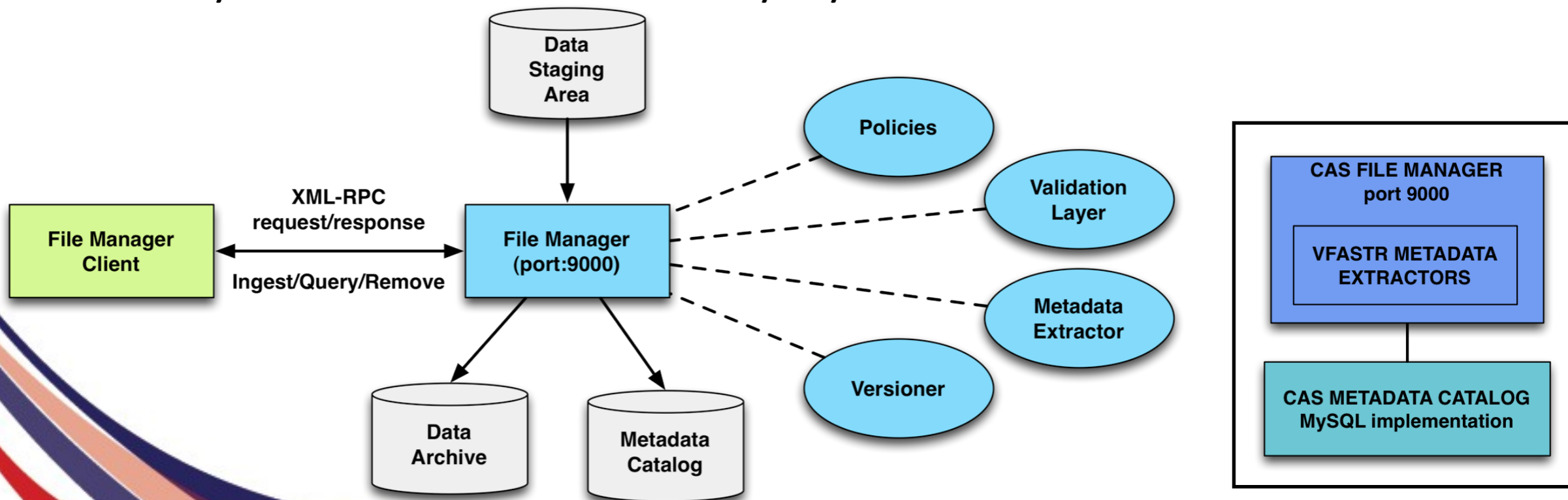
VFASTR Deployment:

• Policy Files: define a single VFASTR metadata type to capture ALL information associated with a single product (job, scans and events)

▸ The full metadata for a product can be retrieved by a client with one request

▸ Metadata keys must be named dynamically to capture job-scan-event references

▸ Example: key=EventStartDateTime_s6 values=2013-01-12T15:48:21.800-0800, 2013-01-12T15:48:22.830-0800 (scan 6 contains 2 events)
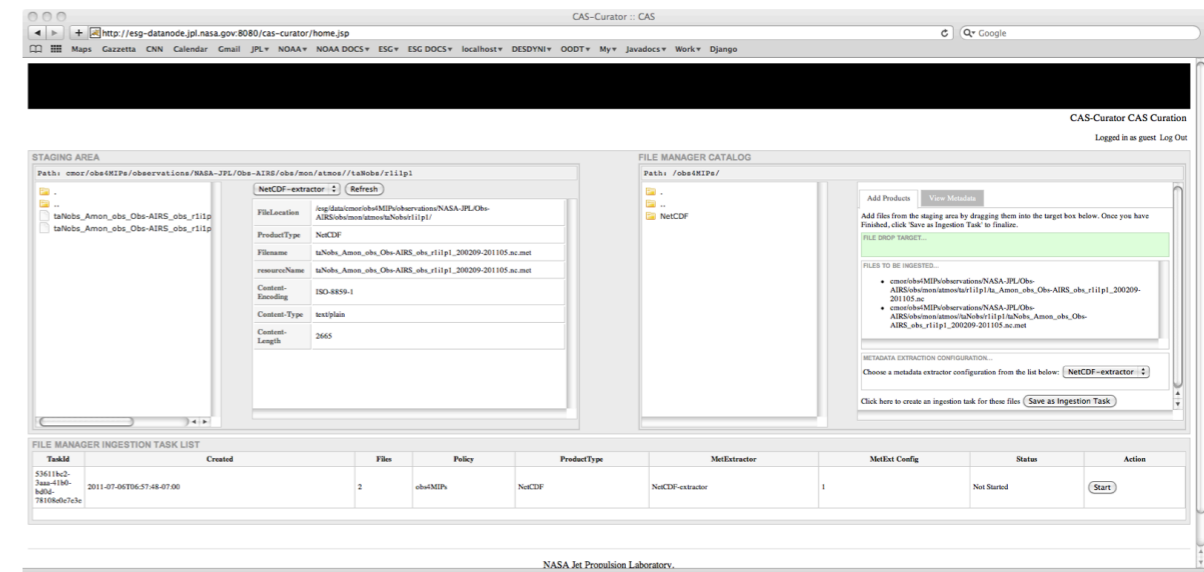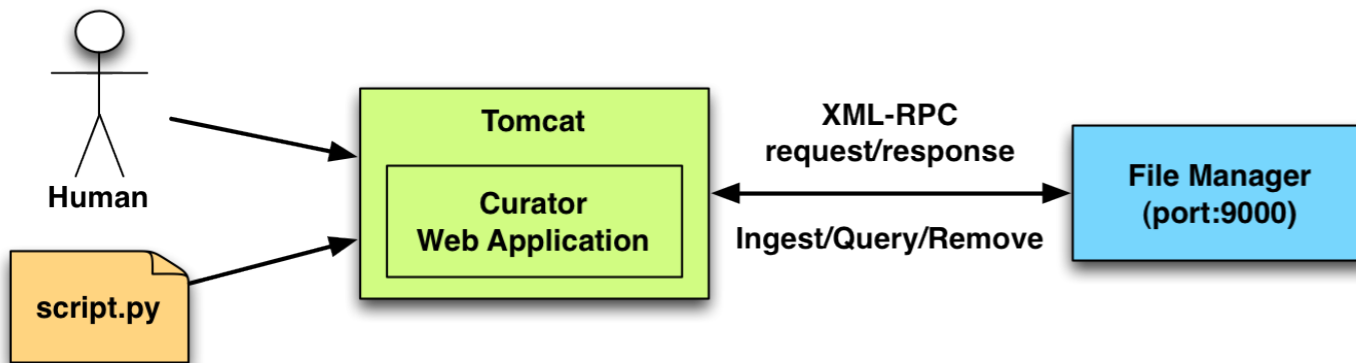
- <u>Validation Layer</u>: no validation applied as metadata fields are not known a-priori
  - ▸Back-end catalog implementations had to be extended to allow for optional "lenient" behavior for ingested metadata
- <u>Metadata Extractors</u>: custom metadata extractors written to parse information for job, scans, events from directory structure, calibration and output files, and to assign detection images to the events that generated them
- <u>Metadata Catalog</u>: used both Lucene and MySQL back-ends
  - ▸Switched to MySQL to support high-frequency updates
  - ▸Lucene File Manager implementation now fixed to support high frequencies
- <u>Data Transfer Protocol</u>: archive products in place
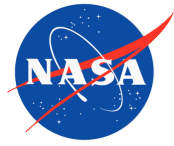  - ▸Otherwise they would be re-transferred by rsync

The CAS Curator is a web application for interacting with the File Manager (i.e. web-based client for File Manager service):
- Submit data product ingestion jobs
- Inspect, add and update product metadata ("curation")



Features: provides two interfaces for interacting with the File Manager:
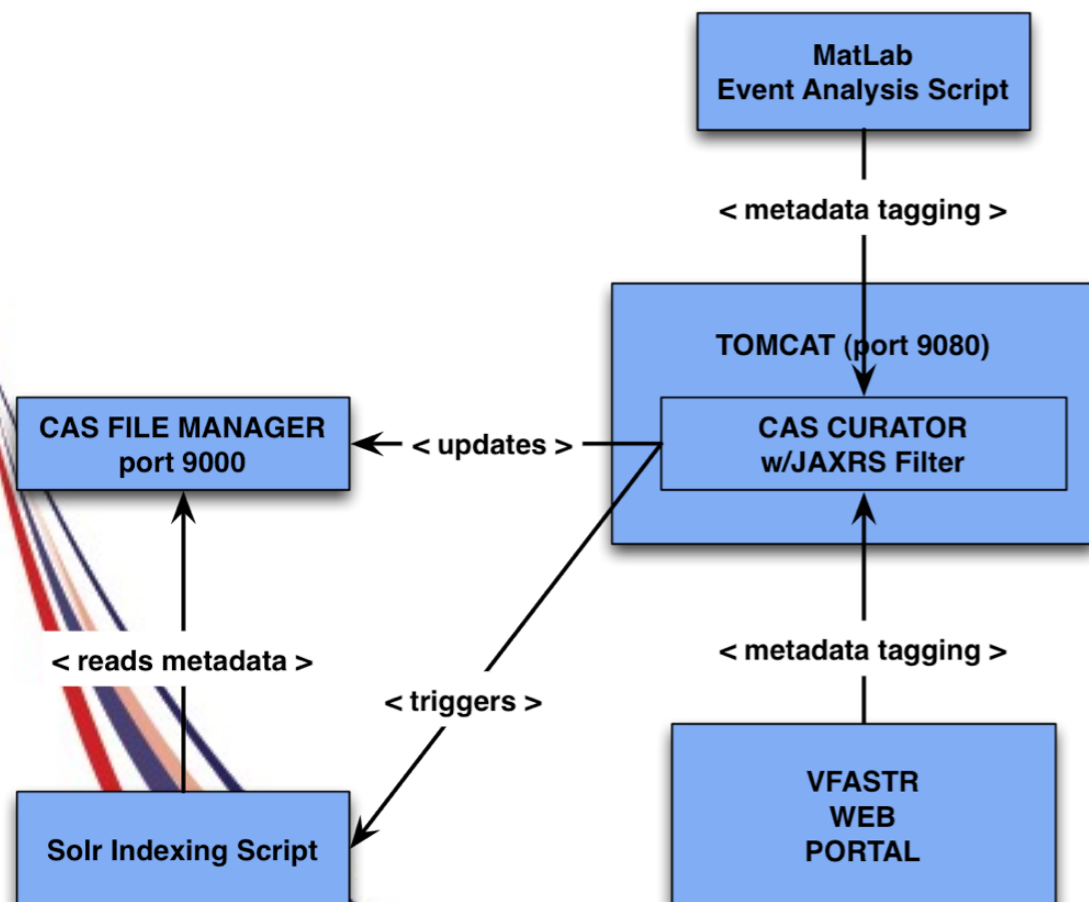- Web User Interface: used by humans to manually interact with the system
  ‣ Drag-and-drop selection of files from the staging area
  ‣ Selection of metadata extractor, versioner from available pool
  ‣ Submission of job for bulk ingestion to File Manager
  ‣ Widget for display and update product metadata

- Web RESTful API: used by programs and scripts for machine-machine interaction
  - ▶ Based on Apache JAX-RS project (project for RESTful web services)
  - ▶ Allows to annotate existing products with enhanced metadata
  - ▶ Example HTTP/POST invocation:
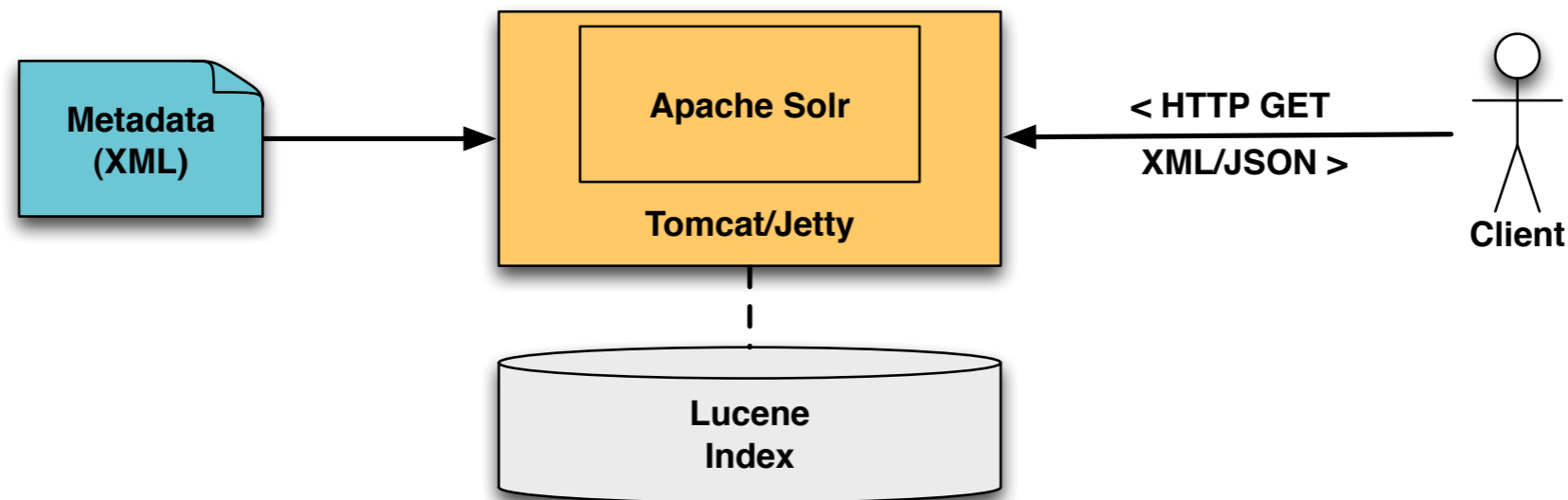  - ▶ curl --data "id=<product_id>&metadata.<name>=<value>" http://<hostname>/ curator/services/metadata/update



VFASTR deployment:
- REST API used by Web Portal and MatLab script to tag interesting events
- Updated product metadata submitted to FM via XML/RPC request
- ▶ VFASTR Curator was wired with JAXRS ResponseHandler (servlet filter invoked before response is sent back to client) to invoke the script for updating the Solr metadata
- metadata.event_s0_e1="pulsar|machine|date"

# Apache Solr

Solr is a high-performance web-enabled search engine built on top of Lucene
- Used in many e-commerce web sites
- Free text searches (w/ stemming, stop words, ...)
- Faceted searches (w/ facet counts)
- Other features: highlighting, word completion,...
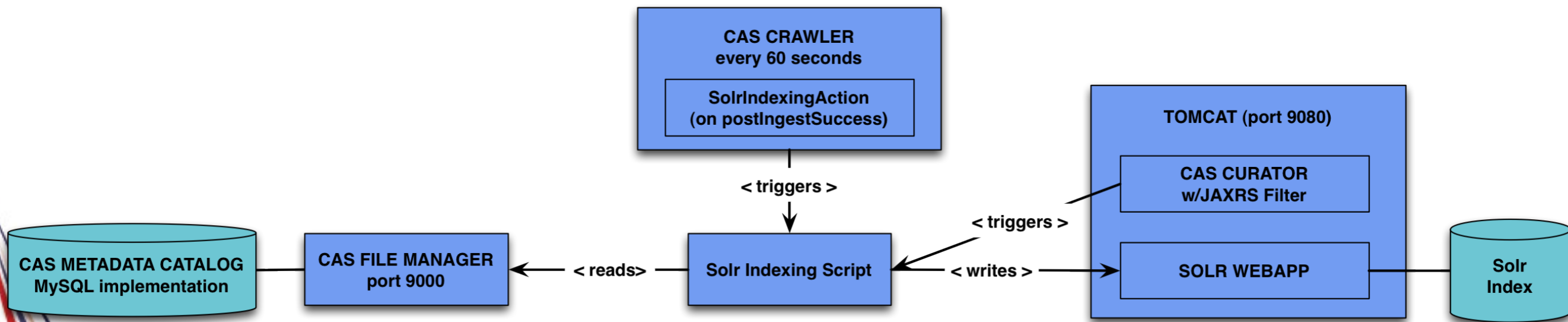- Flat metadata model: (key, value+) pairs

Metadata
(XML)

Apache Solr

Tomcat/Jetty

< HTTP GET
XML/JSON >

Client

Lucene
Index

Scalability: Solr includes features to scale to 10-100 M of records:
- Multiple Cores to partition records into distinct Indexes
- Multiple Shards to distribute the query across complementary Indexes
- Replicated Indexes for high availability and low latency

VFASTR deployment: Solr is used to enable high performance metadata querying by clients: Web Portal and MatLab script

- Solr web application deployed within same Tomcat container as CAS Curator
- Python Indexing script harvests metadata from CAS Catalog to Solr Index
  - ▶Triggered by CAS Crawler when a product is first ingested
  - ▶Triggered by CAS Curator when the product metadata is updated
- VFASTR Solr schema specifies name and data type for all metadata fields
  - ▶"Type=job/scan/event" field used to discriminate among records



Examples of VFASTR queries to Solr Index:
- ▶List of latest products by date
- ▶Full metadata for a given job, scan or event
- ▶All events that were assigned a given tag
- ▶All tags assigned to all events

# What is it?

- Web-based view of the metadata associated with nightly observations

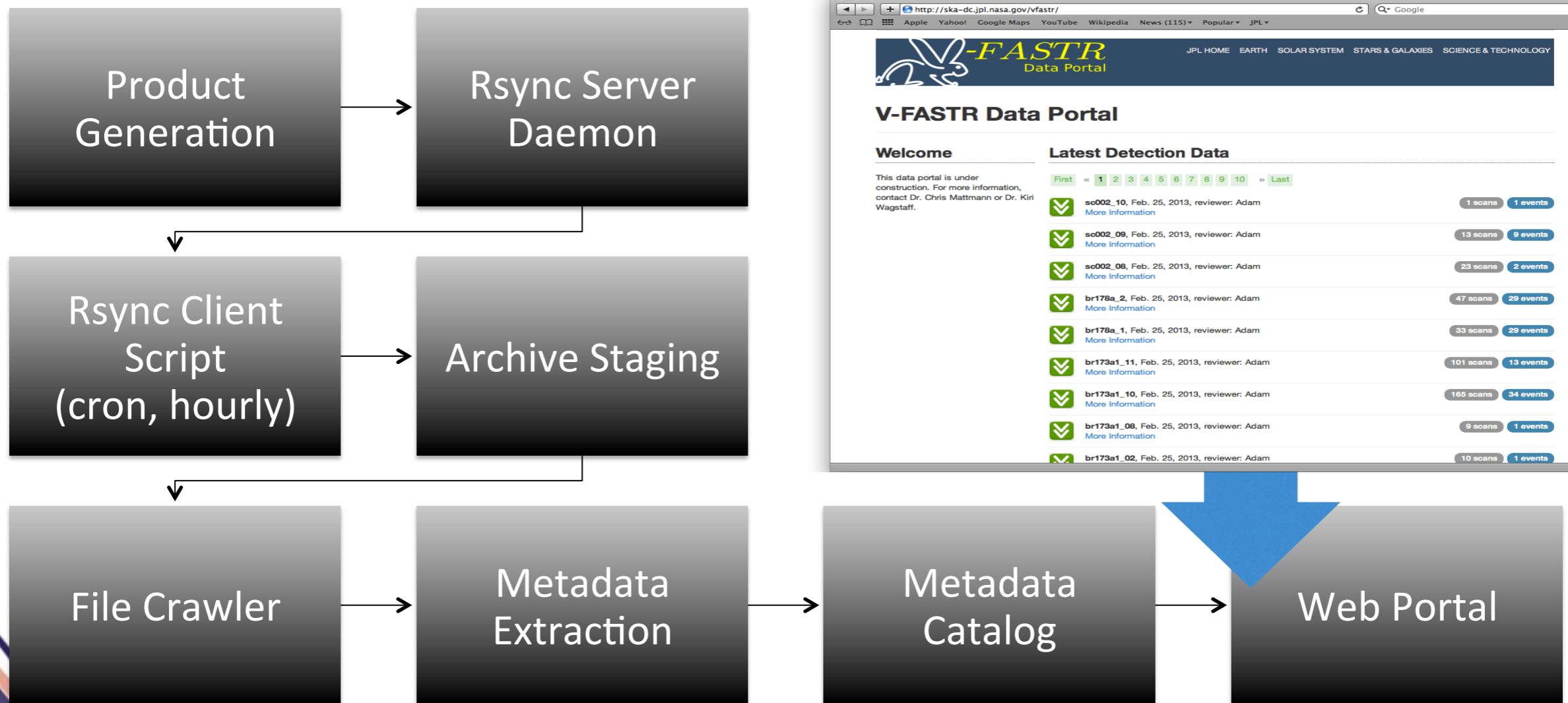- Collaborative environment for review by V-FASTR science team

# Why does it exist?

- Provide distributed science team with convenient access to metadata

# How does it fit?

- Focal point for end-user access to the data pipeline

# What is it built with?

- Some of the technologies behind the data portal:

http://ska-dc.jpl.nasa.gov/vfastr/

Q▾ Google

Apple  Yahoo!  Google Maps  YouTube  Wikipedia  News (115)▾  Popular▾  JPL▾

# V-FASTR
## Data Portal

JPL HOME   EARTH   SOLAR SYSTEM   STARS & GALAXIES   SCIENCE & TECHNOLOGY

# V-FASTR Data Portal

## Welcome

This data portal is under construction. For more information, contact Dr. Chris Mattmann or Dr. Kiri Wagstaff.

## Latest Detection Data

First  «  **1**  2  3  4  5  6  7  8  9  10  »  Last

**sc002_10**, Feb. 25, 2013, reviewer: Adam
More Information
1 scans | 1 events

**sc002_09**, Feb. 25, 2013, reviewer: Adam
More Information
13 scans | 9 events

**sc002_08**, Feb. 25, 2013, reviewer: Adam
More Information
23 scans | 2 events

**br178a_2**, Feb. 25, 2013, reviewer: Adam
More Information
47 scans | 29 events

**br178a_1**, Feb. 25, 2013, reviewer: Adam
More Information
33 scans | 29 events

**br173a1_11**, Feb. 25, 2013, reviewer: Adam
More Information
101 scans | 13 events

**br173a1_10**, Feb. 25, 2013, reviewer: Adam
More Information
165 scans | 34 events

**br173a1_08**, Feb. 25, 2013, reviewer: Adam
More Information
9 scans | 1 events

VFASTR Data Portal

http://ska-dc.jpl.nasa.gov/vfastr/scan/13d2b64f-0918-48ff-bbde-2dbcf7e4ccbf/2

Google

Apple    Yahoo!    Google Maps    YouTube    Wikipedia    News (115) ▾    Popular ▾    JPL ▾

# V-FASTR
## Data Portal

JPL HOME    EARTH    SOLAR SYSTEM    STARS & GALAXIES    SCIENCE & TECHNOLOGY

Home → Job: tns_hiresah_7 → Scan: 2

# Scan Viewer

## Metadata

| | |
|---|---|
| Job: | tns_hiresah_7 |
| Scan: | 2 |
| Pointing Source: | J0038+4137 |
| Scan Date: | 2013.02.02 |
| Scan Start: | 21:58:20 |
| Event Count: | 4 |

## Events in this Scan (4)

**Event 1** - *Start time*: 21:58:32, *Priority*: 6.531, *Median DM*: 0.000
Details

**Event 2** - *Start time*: 21:58:40, *Priority*: 7.202, *Median DM*: 0.000
Details

**Event 3** - *Start time*: 21:59:13, *Priority*: 7.055, *Median DM*: 0.000
Details

**Event 4** - *Start time*: 21:59:16, *Priority*: 6.693, *Median DM*: 0.000
Details

SKA Data Center - VFASTR

NASA Jet Propulsion Laboratory

Created with Apache OODT Balance 0.3-SNAPSHOT

Privacy

Responsible Official:
Chris A. Mattmann
Site Contact:
Andrew F. Hart

USA.gov
Government Made Easy

http://ska-dc.jpl.nasa.gov/vfastr/event/fc6c9247-9455-4d7f-9ce4-b2a6037ca6d6/2/1

Apple  Yahoo!  Google Maps  YouTube  Wikipedia  News (126) ▾  Popular ▾  JPL ▾

# V-FASTR Data Portal

JPL HOME    EARTH    SOLAR SYSTEM    STARS & GALAXIES    SCIENCE & TECHNOLOGY

Home → Job: tns_sc002_09 → Scan: 2 → Event: 1

# Event Viewer

## Metadata

| | |
|---|---|
| Job Id: | tns_sc002_09 |
| Scan Id: | 2 |
| Event Id: | 1 |
| Median DM: | 0.000 |

## Timing

| | |
|---|---|
| Start Date/Time: | 2013.02.08 11:33:42 |
| Start Time Step: | 0000000001010 |
| End Date/Time: | 2013.02.08. 11:33:43 |
| End Time Step: | 0000000001011 |

## Imagery

**Detection**
Download Full-size

**Dedispersion**
Download Full-size

**Detection Summed**
Download Full-size

## Detection Imagery
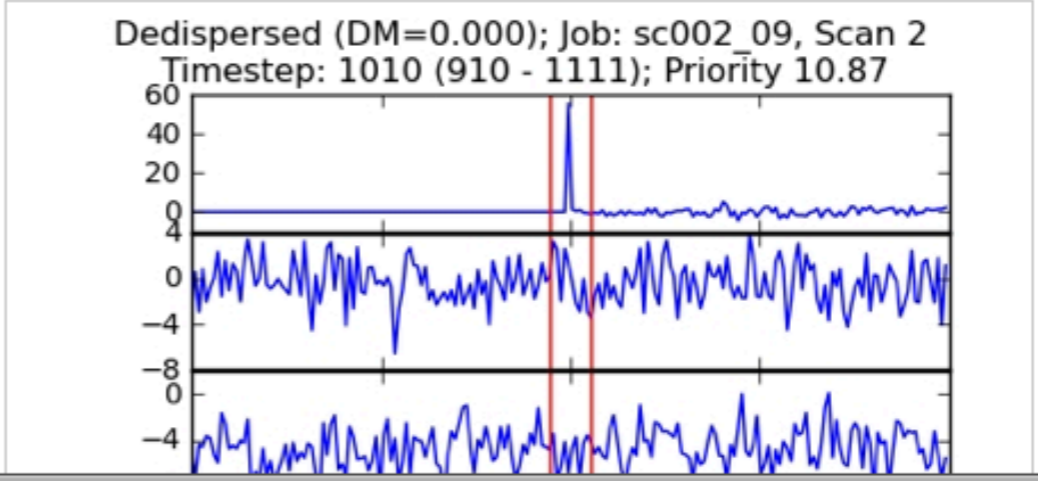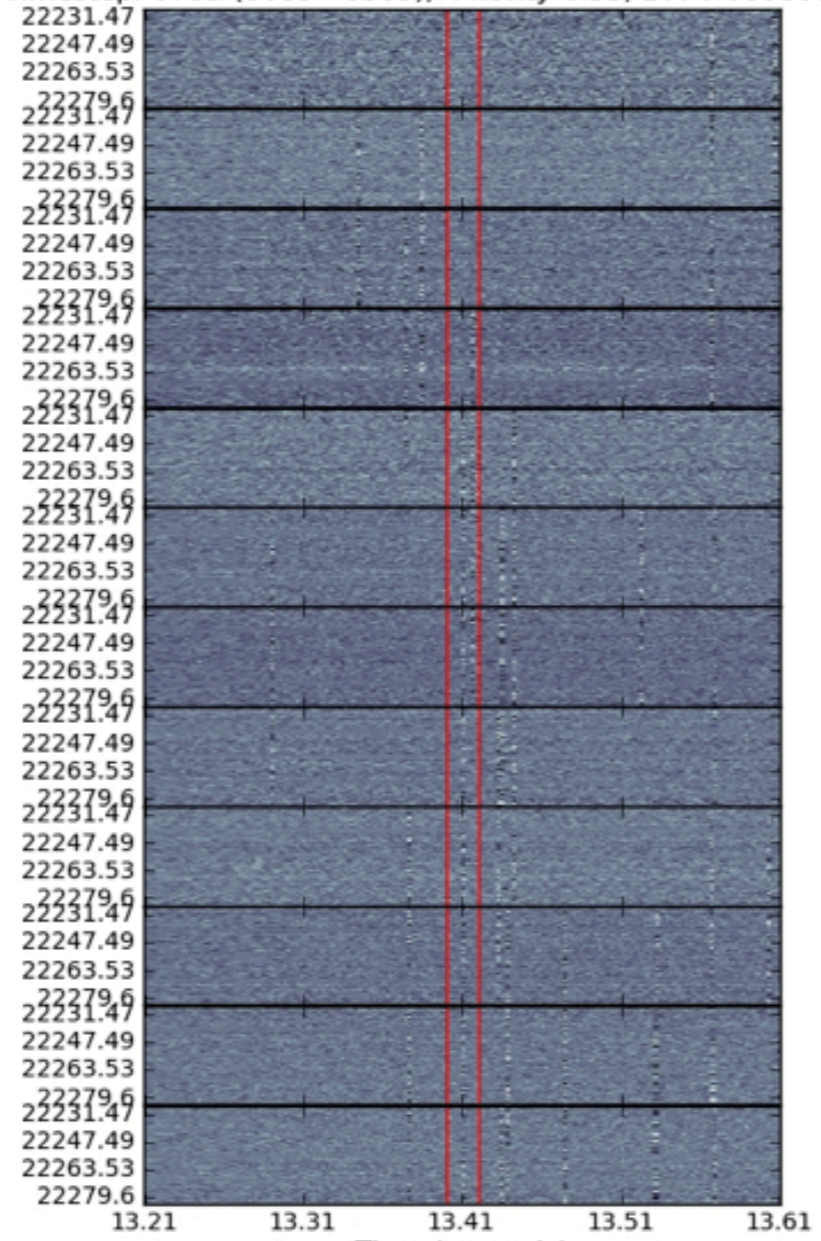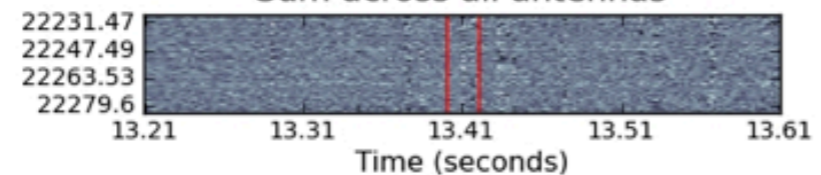
Q view zoomed sum across antennas    ⌃ view fullsize

Job: sc002_09   Scan 2  Antennas: All, Polarization: Sum
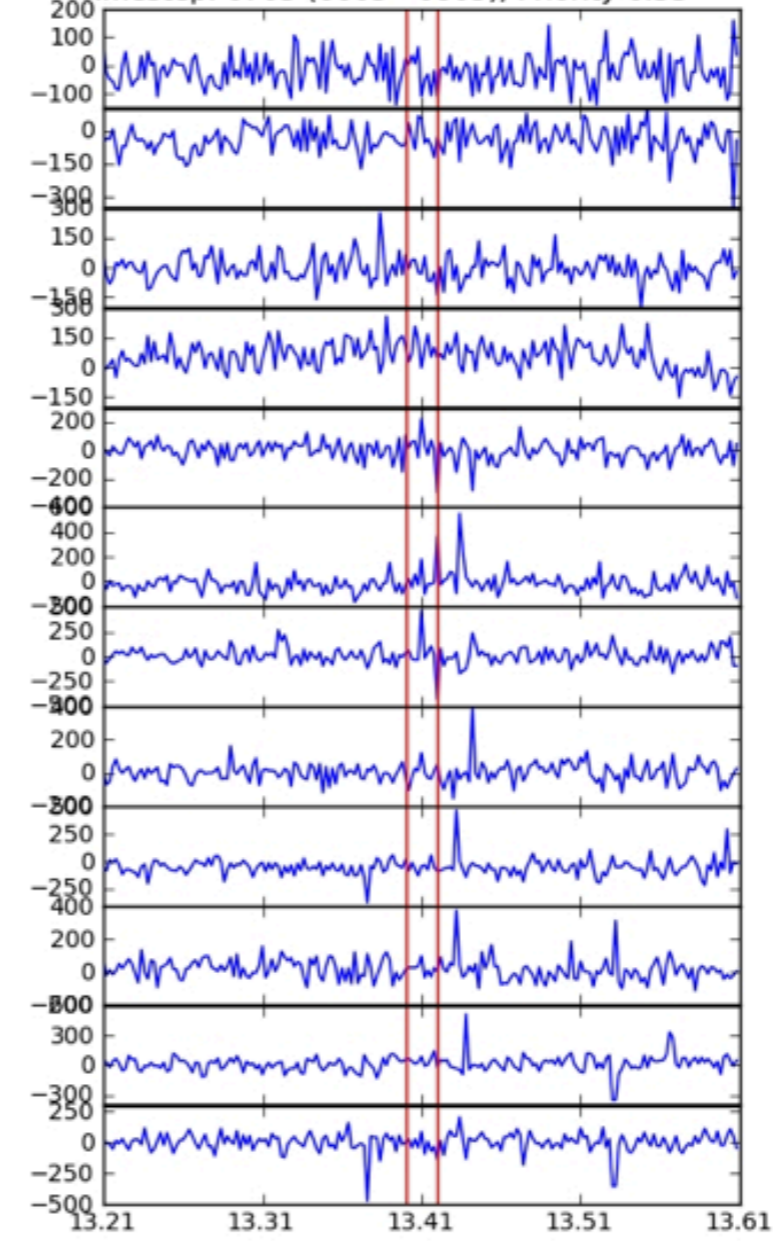Timestep: 1010 (910 - 1111); Priority 10.87; DM 0.000000

## Dedispersion Imagery

⌃ view fullsize

Dedispersed (DM=0.000); Job: sc002_09, Scan 2
Timestep: 1010 (910 - 1111); Priority 10.87

# What are tags?

- Descriptive metadata associated with an event

- Enable classification and filtering

- Serve as training for AI (now)

- Serve as guide for what to archive (soon)

http://ska-dc.jpl.nasa.gov/vfastr/review-job/ec373a23-d3b5-4099-a8c7-fcaa835d0e68/tns_br178a_

Google

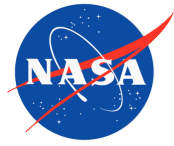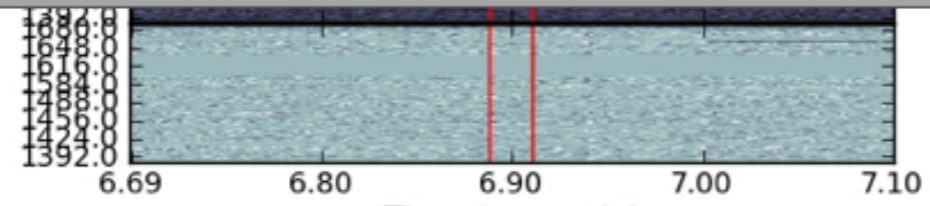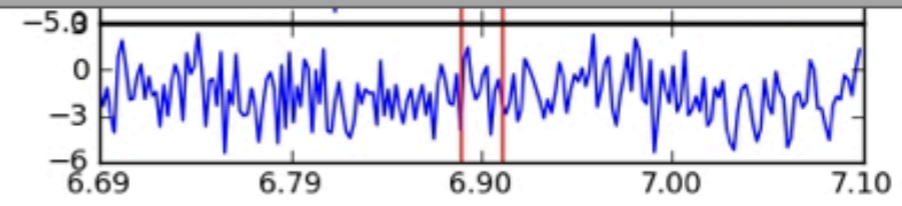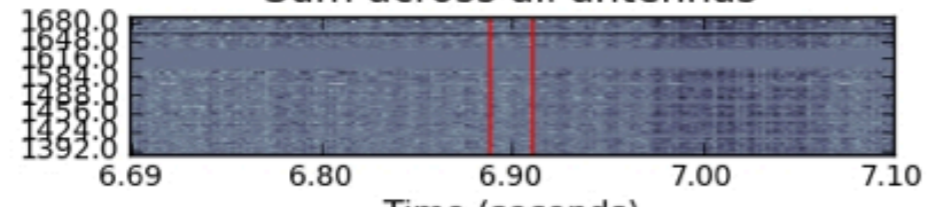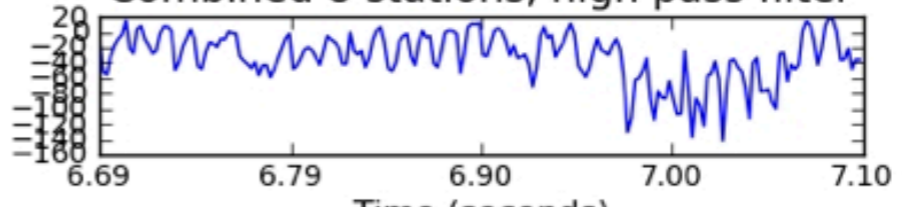Apple Yahoo! Google Maps YouTube Wikipedia News (115)▾ Popular▾ JPL▾

W

Sum across all antennas

Time (seconds)

Combined 8 stations, high pass filter

Time (seconds)

## Event Tags

Human generated tags look like ( this ), whereas machine generated tags look like ( this ). Hover over a tag with your mouse to see details.

**Tags for this Event:**

pulsar

[ r ]   [ Add Tag ]

RFI
data_dropout
salt_and_pepper
uninteresting
interesting
pulsar

Ev                    0:25,  *Priority*: 6.586,  *Median DM*: 2.760

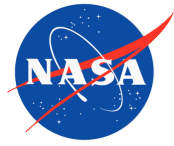**D**         **ery**                    **Dedispersion Imagery**

🔍 view zoomed sum across antennas    ⌃ view fullsize                    ⌃ view fullsize

Job: br178a_2   Scan 0  Antennas: All, Polarization: Sum          Dedispersed (DM=2.760); Job: br178a_2, Scan 0
Timestep: 36578 (36478 - 36680);  Priority 6.59; DM 2.760000          Timestep: 36578 (36478 - 36678); Priority 6.59
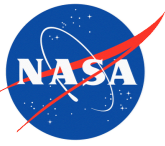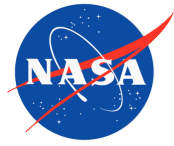
# What can users do?

- Event tagging (now)

  - Nightly classification, viewable by all users

  - Used as training input to the automatic candidate detection

- Job archiving (soon)

  - Initiate archival of job on NRAO hardware based upon the contents of the tags and/or other metrics
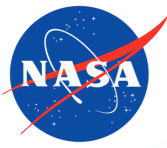
# What have users thought?

- Interface with all event imagery on-screen is an improvement over command-line methods (more efficient)

- Organization of the interface should support rapid evaluation of an entire job (minimize clicks)

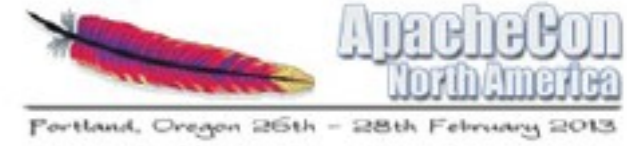- Improved accessibility of the information a big plus (mobile access)

# What is next?

- Provide more comprehensive job/event search capability

  - Facet by tag values, index metadata

- Continued efficiency improvements

  - Bulk tagging of all events in a job

- Implement front-end capability to initiate back-end archive process

# Questions ?