# Big Data and Data Science at NASA/JPL: Methodology Transfer From Space Science to Biomedicine

## Daniel Crichton

*Leader, Center for Data Science and Technology*
*Proj. Manager, Planetary Data System Engineering*
*Prog. Manager, Data Science Office*

## Dr. Richard Doyle

*Prog. Manager, Information and Data Science*
*Proj. Manager, High Performance Spaceflight Computing*

*leaving the
safe harbor
to explore
uncharted waters*

Jet Propulsion Laboratory
California Institute of Technology

*March 6, 2017*

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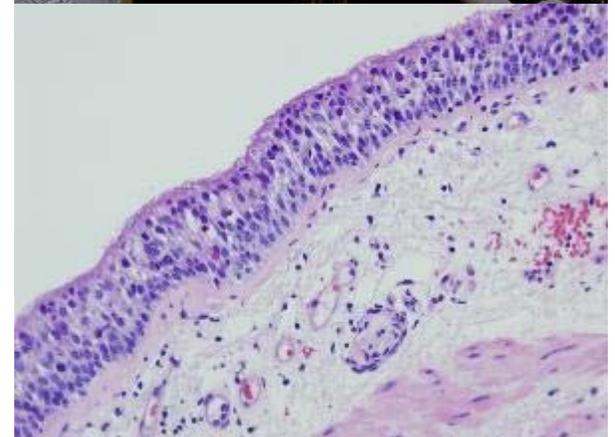**California Institute of Technology**

Jet Propulsion Laboratory
California Institute of Technology

- JPL is involved in the research and development of technologies, methodologies in science, mission operations, engineering, and other non-NASA applications.
  - Includes onboard computing to scalable archives to analytics

- JPL and Caltech formed a joint initiative in Data Science and Technology to support fundamental research all the way to operational systems.
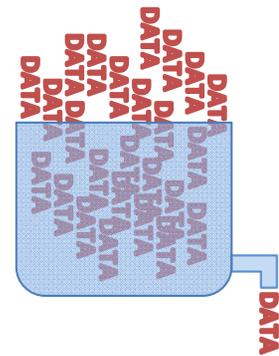  - Methodology transfer across applications is a major goal.

# Terms: Big Data and Data Science

## Big Data

- When needs for data collection, processing, management and analysis go beyond the capacity and capability of available methods and software systems

## Data Science

- Scalable architectural approaches, techniques, software and algorithms which alter the paradigm by which data is collected, managed and analyzed
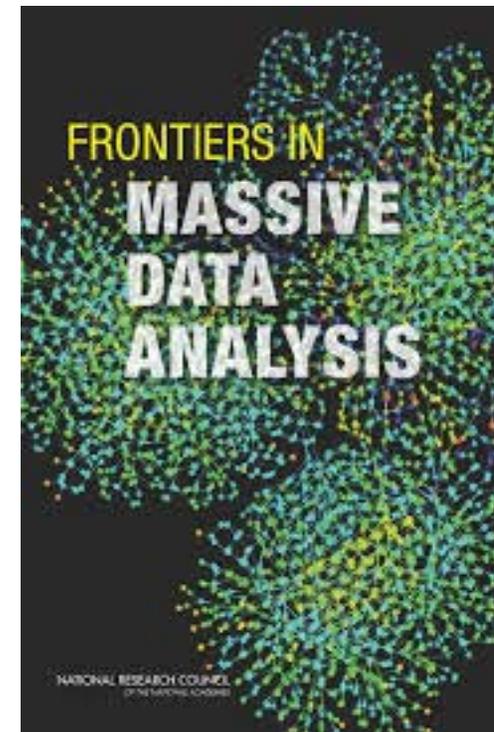
# U.S. National Research Council Report:
## *Frontiers in the Analysis of Massive Data*

- Chartered in 2010 by the U.S. National Research Council, National Academies
- Chaired by Michael Jordan, Berkeley, AMP Lab (Algorithms, Machines, People)
- NASA/JPL served on the committee covering systems architecture for big data management and analysis
- **Importance of more systematic approaches for analysis of data**
- **Need for end-to-end data lifecycle: from point of capture to analysis**
- **Integration of multiple discipline experts**
- Application of novel statistical and machine learning approaches for data discovery
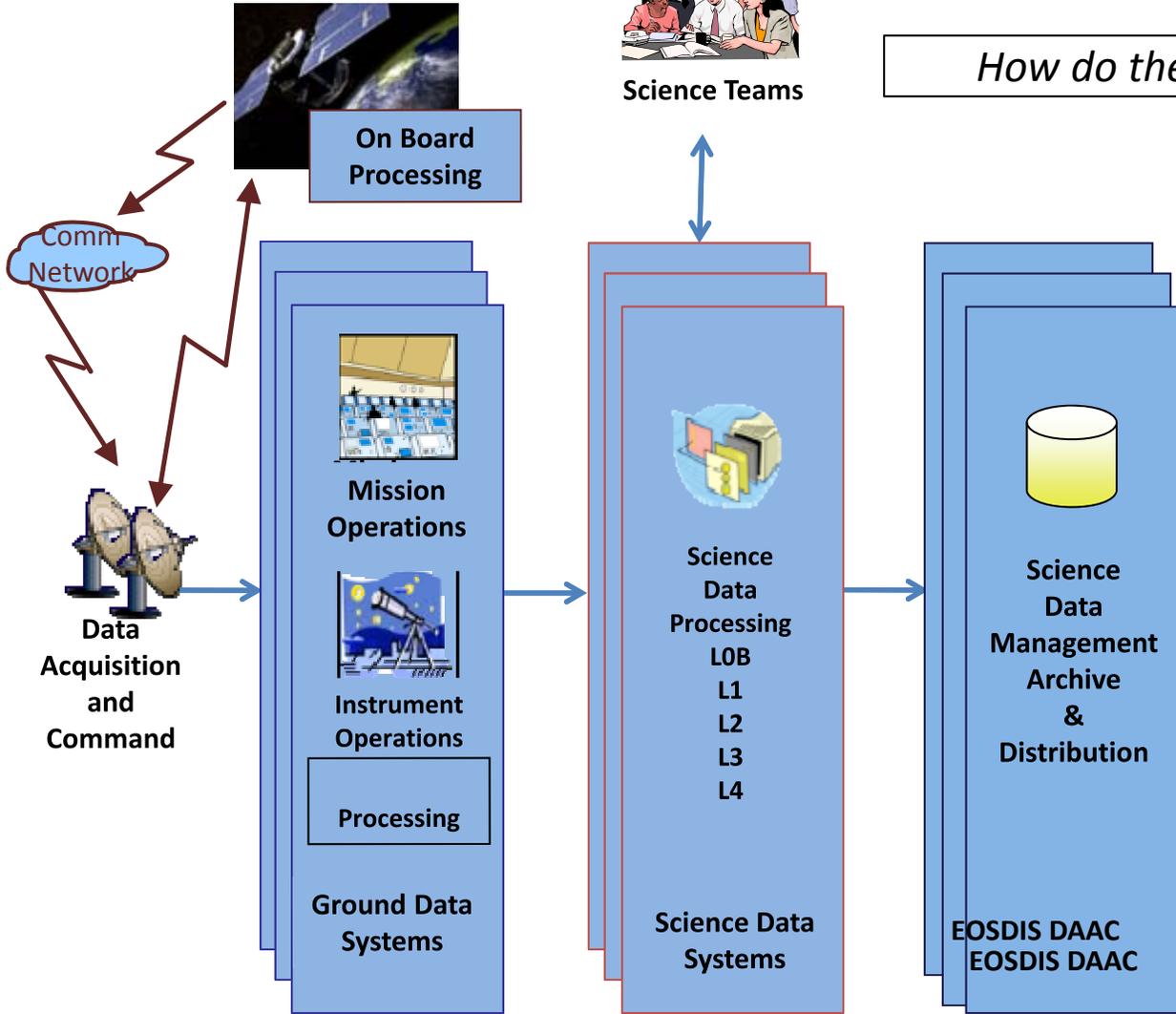
2013

# NASA Science and Big Data Today

Jet Propulsion Laboratory
California Institute of Technology

How do these connect?

**Science Teams**

**On Board Processing**

**Comm Network**

**Data Acquisition and Command**

**Mission Operations**

**Instrument Operations**

**Processing**

**Ground Data Systems**

**Science Data Processing**
L0B
L1
L2
L3
L4

**Science Data Systems**

**Science Data Management Archive & Distribution**

**EOSDIS DAAC**
**EOSDIS DAAC**

**Big Data Infrastructure (Data, Algorithms, Machines)**

?

**Research**

**Outreach**

**Applications**

*Focus on generating, capturing, managing big data*

*Focus on using/analyzing big data*

5

# JPL Data Science Working Group

- Established in 2014 to explore big data use cases and challenges in science and to make a recommendation to JPL senior management.
  - Launched internal investments: planetary science (onboard agile science), earth science (distributed data analytics), and astronomy (machine learning and data collection methods).
  - Engaged cross disciplinary expertise (science, computer science – systems and machine learning, statistics, program management)
  - Partnered with Caltech to bring in research perspectives.

- In November 2016 a chartered Data Science WG reporting to JPL's Leadership Management Council (LMC), chaired by Deputy Director Larry James, was established in data science covering all aspects of the Lab operations.

![Jet Propulsion Laboratory, California Institute of Technology]

# JPL Data Science Strategy
## *Guiding Principles*



**Data Lifecycle**

**Agile Science – Onboard Analysis**

Challenge:
Too much data, too fast; cannot transport data efficiently enough

Future Solutions:
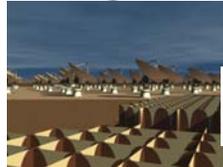Onboard computation and data science

**Extreme Data Volumes – Data Triage**

Challenge: Data collection capacity at the instrument outstrips data transport and data storage capacity

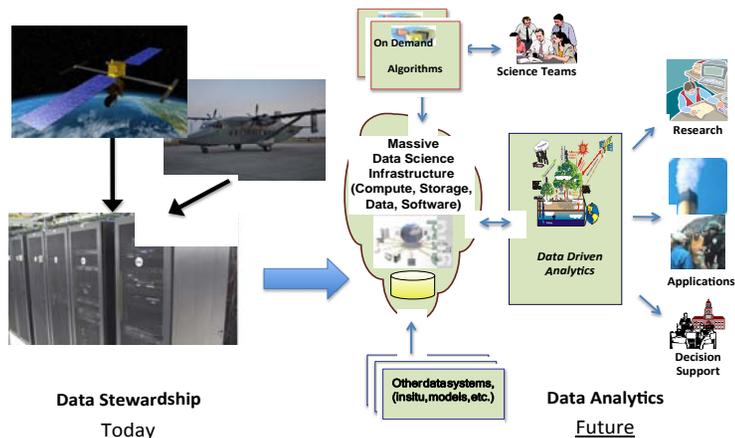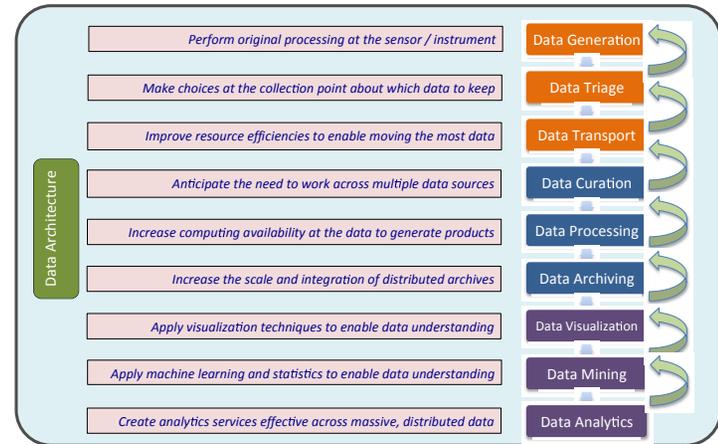Future Solutions: Dynamic architectures to scale data processing and triage exascale data streams

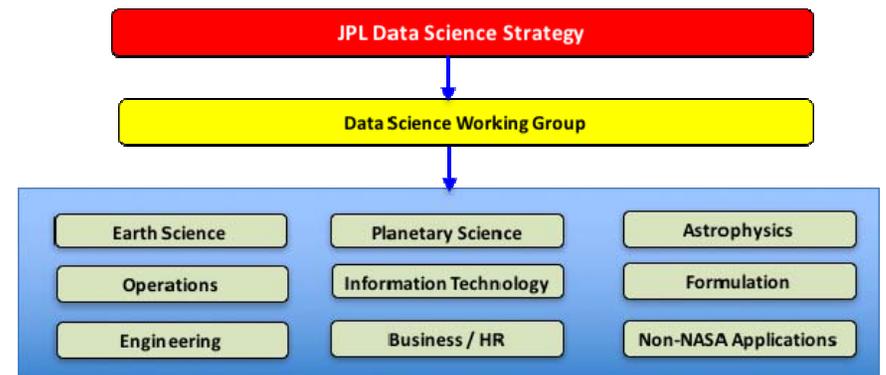**Distributed Data Analytics**

Challenge: Data distributed in massive archives; many different types of measurements

Future Solutions: Distributed data analytics; uncertainty quantification

**Data Architecture**

| Principle | Stage |
|-----------|-------|
| Perform original processing at the sensor / instrument | Data Generation |
| Make choices at the collection point about which data to keep | Data Triage |
| Improve resource efficiencies to enable moving the most data | Data Transport |
| Anticipate the need to work across multiple data sources | Data Curation |
| Increase computing availability at the data to generate products | Data Processing |
| Increase the scale and integration of distributed archives | Data Archiving |
| Apply visualization techniques to enable data understanding | Data Visualization |
| Apply machine learning and statistics to enable data understanding | Data Mining |
| Create analytics services effective across massive, distributed data | Data Analytics |

**Data Ecosystem**

On Demand Algorithms

Science Teams

Massive Data Science Infrastructure (Compute, Storage, Data, Software)

Other data systems, (insitu, models, etc.)

Data Driven Analytics

Research

Applications

Decision Support

Data Stewardship
Today

Data Analytics
Future

**Cross-Cutting**

JPL Data Science Strategy

Data Science Working Group

| Earth Science | Planetary Science | Astrophysics |
| Operations | Information Technology | Formulation |
| Engineering | Business / HR | Non-NASA Applications |

# Data Lifecycle Model
## *for NASA Space Missions*

Jet Propulsion Laboratory
California Institute of Technology

**Emerging Solutions**
- *Onboard Data Products*
- *Onboard Data Prioritization*
- *Flight Computing*

*(1) Too much data, too fast; cannot transport data efficiently enough to store*

Observational Platforms /Flight Computing

Massive Data Archives and Big Data Analytics

**Emerging Solutions**
- *Low-Power Digital Signal Processing*
- *Data Triage*
- *Exa-scale Computing*

**Emerging Solutions**
- *Distributed Data Analytics*
- *Advanced Data Science Methods*
- *Scalable Computation and Storage*

*(2) Data collection capacity at the instrument continually outstrips data transport (downlink) capacity*

*(3) Data distributed in massive archives; many different types of measurements and observations*

Ground-based Mission Systems

# Cross-Cutting Capabilities

**International Data Archive and Sharing Architectures**
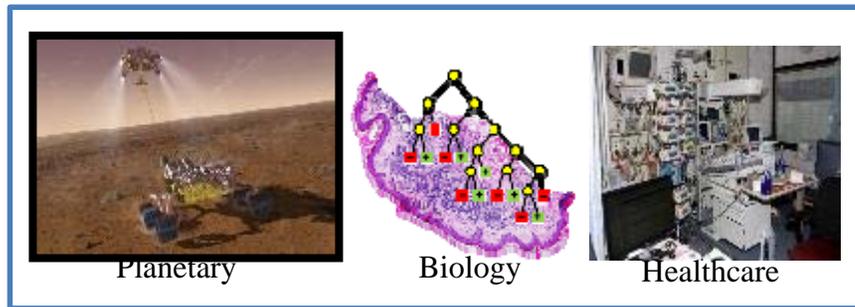


**Big Data Infrastructures**

(from open source to cloud computing and scalable compute infrastructures)



**Intelligent Data Algorithms**

(Machine Learning, Deep Learning)



**Common Data Elements & Information Models**

(discipline and common)



**Analytical Data Pipelines**



Planetary    Biology    Healthcare

*Great Opportunities for Methodology Transfer and Collaboration*    9



**Visualization Techniques**

# Future of Data Science at NASA
## *Enabling a Big Data Research Environment*

Jet Propulsion Laboratory
California Institute of Technology

On Board Processing

Comm Network

Data Acquisition and Command

Mission Operations

Instrument Operations

Processing

Ground Data Systems

Instrument Data Systems

Airborne Data

NASA Data Archives

On Demand Algorithms

Science Teams

Big Data Infrastructure (Data, Algorithms, Machines)

Other data systems (in situ, other agency, etc.)

*Data Analytics Centers*

(Water, Ocean, CO2, Extreme Events, Mars, etc.)

Research

Applications

Decision Support

**Data Capture**

**Data Analysis**

Reducing Data Wrangling: "There is a major need for the development of software components… that link high-level data analysis-specifications with low-level distributed systems architectures."
*Frontiers in the Analysis of Massive Data*, National Research Council, 2013.

# Opportunities and Use Case Across the Ground Environment

Jet Propulsion Laboratory
California Institute of Technology

## Intelligent Ground Stations

**Emerging Solutions**
- *Anomaly Detection*
- *Combining DSN & Mission Data*
- *Attention Focusing*
- *Controlling False Positives*

## Intelligent Archives and Knowledge-bases

**Emerging Solutions**
- *Automated Machine Learning - Feature Extraction*
- *Intelligent Search*
- *Learning over time*
- *Integration of disparate data*

**Technologies: Machine Learning, Deep Learning, Intelligent Search, Data Integration, Interactive Visualization and Analytics**

## Intelligent MOS-GDS

**Emerging Solutions**
- *Anomaly Interpretation*
- *Dashboard for Time Series Data*
- *Time-Scalable Decision Support*
- *Operator Training*

## Data Analytics and Decision Support

**Emerging Solutions**
- *Interactive Data Analytics*
- *Cost Analysis of Computation*
- *Uncertainty Quantification*
- *Error Detection in Data Collection*

# 2015-2016 AIST Big Data Study

- Study led by JPL for the NASA Advanced Information Systems Technology Program (under Mike Little)

- Mapped technology and data needs against the mission-science data lifecycle

- Focuses on expansion from data stewardship to data use across the vast data ecosystem (satellite, airborne, in situ)

- Basis for 2015 IEEE Big Data workshop on Data and Computational Challenges in Earth Science Research

- Key input for 2016 ROSES AIST call (per Mike Little, NASA PM for AIST)

Data and Computational Science Technologies for Earth Science Research

# AIST Big Data Study: 10 Year Capability Needs in Big Data

| System | 2015 | 2025 | NASA Applications |
|---|---|---|---|
| Observational Platforms | Limited onboard computation including data triage and data reduction. Investments in new flight computing technologies for extreme environments. | Increase onboard autonomy and enable data triage (machine learning techniques) to support more capable instruments. Support reliable onboard processing in extreme environments to enable new exploration missions. | Onboard computation across all types of platforms; flight computing capabilities deployed for extreme environments; data triage for satellites and spacecraft. |
| Ground-based Mission Systems | Rigid data processing pipelines; limited real-time event/feature detection. Support for 500 TB missions. | Increase computational processing capabilities for mission (100x); Enable ad hoc workflows and reduction of data; Enable realtime triage/ML techniques, event and feature detection. Support 100 PB scale missions. | Future mission computational challenges; high bandwidth data volumes; more agile airborne, cube sat, multi-sensor campaigns; increase automated event detection across mission lifecycle. |
| Massive Data Archives | Support for 10 PB of archival data; limited automated event and feature detection. | Support exascale archives; automated event and feature detection/ML techniques; virtually integrated, distributed archives. | Turn archives into knowledge-bases to improve data discovery. Leverage massively scalable virtual data storage. |
| Distributed Data Analytics | Limited analytics services; generally tightly coupled to specific data centers; limited cross-archive/data center, cross-agency integration; limited capabilities in data fusion; statistical uncertainty; provenance of the results. | Computational techniques (ML, statistical methods) integrated into mission-science lifecycle; Integrated data, HPC, algorithms across archives; Support for cross product data fusion; capture of statistical uncertainty; virtual missions; specialized Analytics Centers. | Automated data analysis methods; integration of data across spacecraft, remote sensors, satellite, airborne, and ground-based sensors; systematic approaches to addressing uncertainty; complex scientific questions. |

*Derived from AIST Big Data Study & NASA Office of the Chief Technologist TA-11 Roadmap (2015)*

# Planetary Data System

- <u>Purpose:</u> To collect, archive and make accessible digital data and documentation produced from NASA's exploration of the solar system from the 1960s to the present.

- <u>Infrastructure:</u> A highly distributed infrastructure with planetary science data repositories implemented at major government labs and academic institutions
    - System driven by a well defined planetary science information model
    - Over 1 PB of data
    - Movement towards international interoperability
    - Distributed federation of US nodes and international archives

- Being realized through PDS4

# (Some) Big Data Challenges in Planetary Science

- Variety of planetary science disciplines, moving targets, and data
- Volume of data returned from missions including provenance
- Federation of disciplines and international interoperability

- These factors can affect choices in:
    - Data Consistency
    - Data Storage
    - Computation
    - Movement of Data
    - Data Discovery
    - Data Distribution

*Ultimately, having a planetary science information architectural strategy that can scale to support the size, distribution, and heterogeneity of the data is critical*

**A well formed model that drives the software is something that many groups have struggled with!**

15

# PDS4: International Adoption of an Open Planetary Approach


LADEE (NASA)


InSight (NASA)


BepiColombo (ESA/JAXA)


MAVEN (NASA)


Osiris-REx (NASA)


ExoMars (ESA/Russia)


JUICE (ESA)


Mars 2020 (NASA)

**…also Hayabussa-2, Chandrayaan-2**

16

Endorsed by the **International Planetary Data Alliance** in July 2012 – https://planetarydata.org/documents/steering-committee/ipda-endorsements-recommendations-and-actions
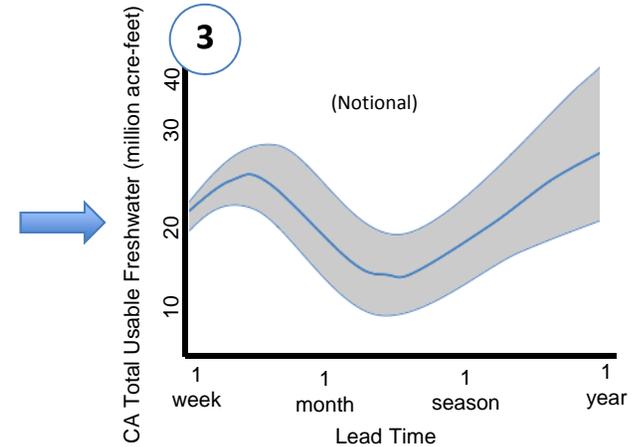
E. Law, S. Malhotra, G. Chang

Built on PDS4

# Western States Water Mission – Understanding Water Availability



**1** Observations

**2** Coupled and Validated Computer Models

Pacific Northwest

Great Basin

California

Upper Colorado

Lower Colorado

**3** Estimates with Uncertainties

(Notional)

e.g., CA Total Usable Freshwater (million acre-feet)

1 week — 1 month — 1 season — 1 year

Lead Time

**4** DEPARTMENT OF WATER RESOURCES

western states water council

(Prospective customers)

U.S. DEPARTMENT OF THE INTERIOR
BUREAU OF RECLAMATION

Colorado River Basin

**Stakeholders and Customers**

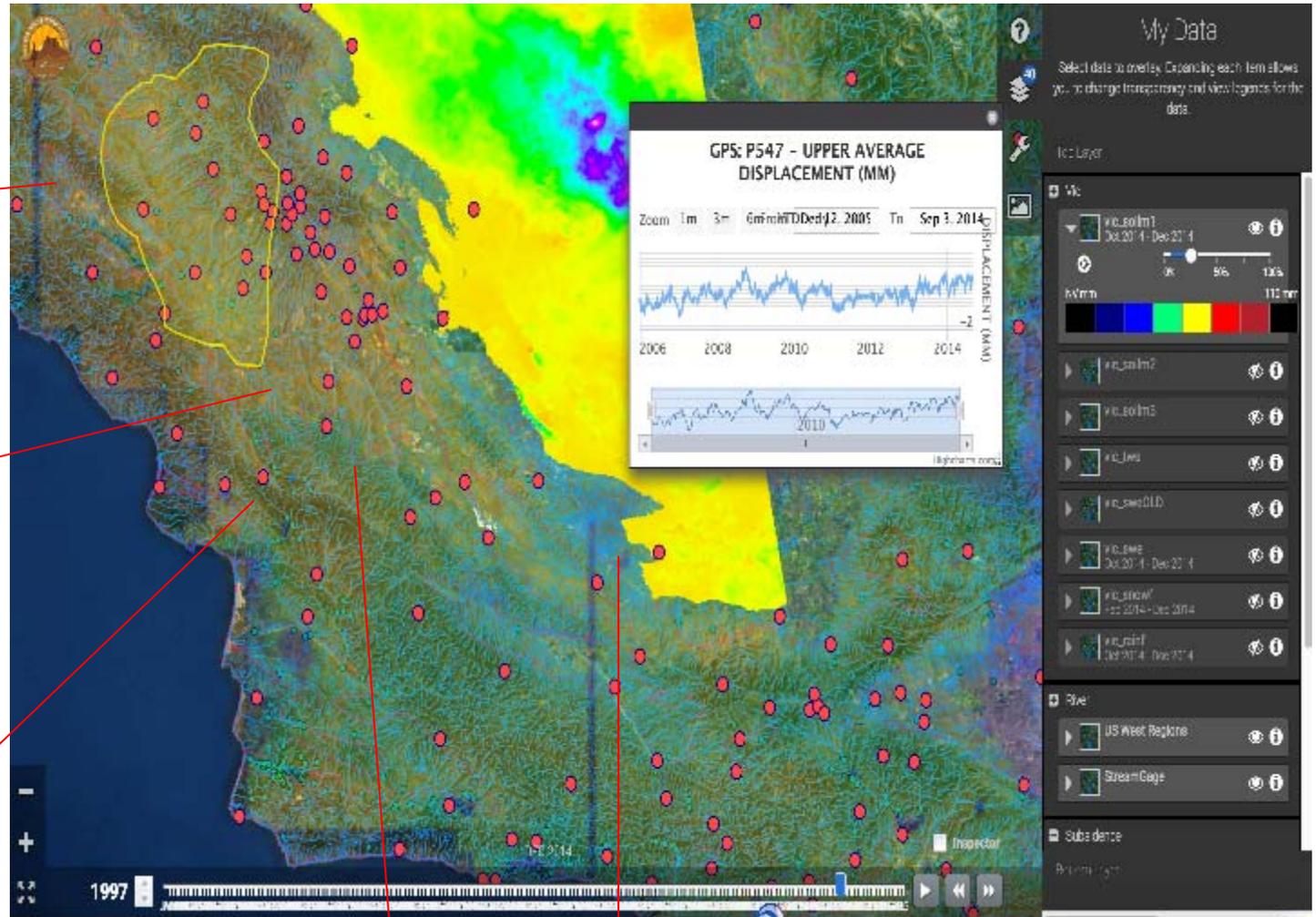**Western States Water Mission (WSWM): A Science/Data Science Collaboration**

Jet Propulsion Laboratory
California Institute of Technology

Decision Support

Research

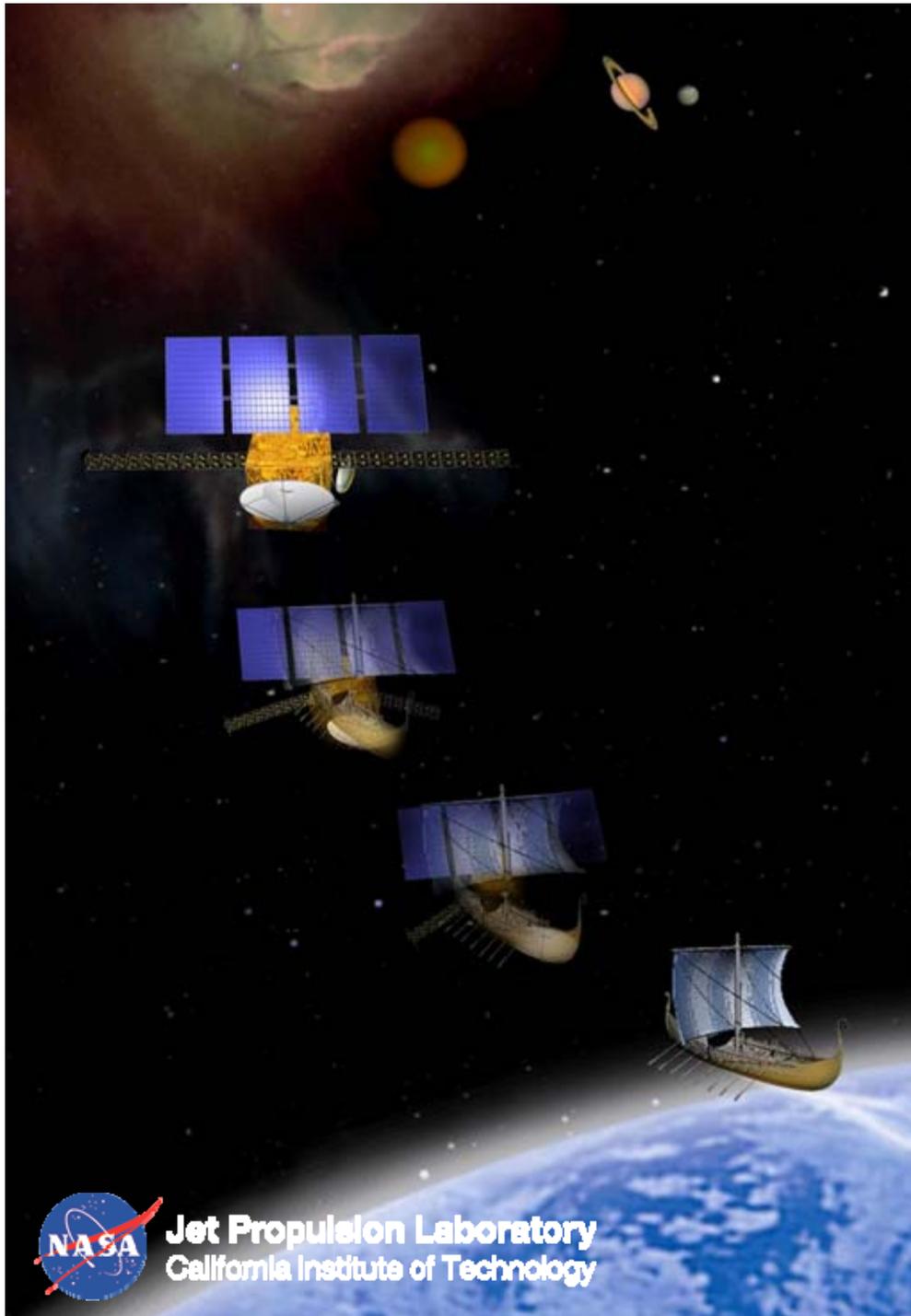Applications

Visualization

(Web-Based Interface)

Standard Reports | Ad Hoc Queries and Custom Reports

Data Analytics

Single-Month Estimates | Short and Long-Term Trends

Models

Snow-Water Equivalent | Surface Water | Ground Water

Observations

Input-Forcing (e.g., GPM) | For Data Assimilation (e.g., MODSCAG)

Data Science Infrastructure (Tools, Services, Methods for Massive Data Analysis)

A Scalable Data Processing System for Hydrological Science

19

# WaterTrek

Fusing In-situ, Air-borne, Space-borne and model generated data using visualization and a big data analytics engine

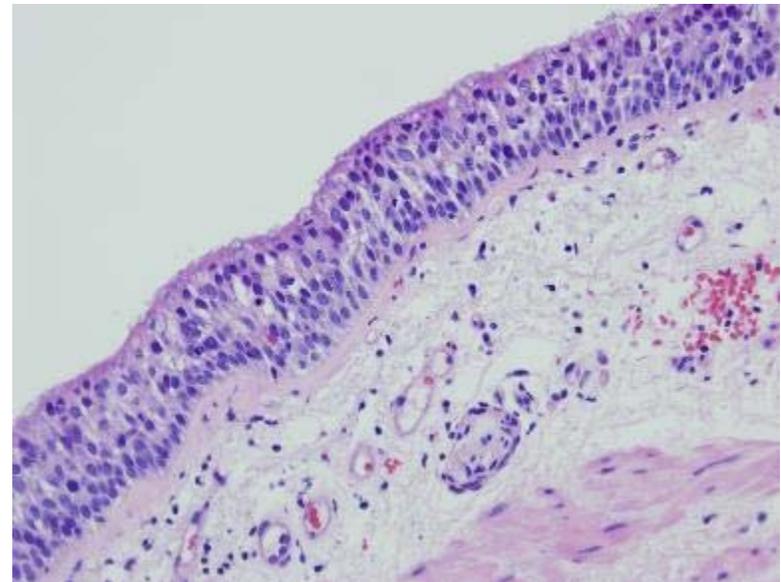# Methodology Transfer in Data Science from Planetary & Earth to Biomedicine

# NASA/JPL Informatics Center: Crossing Disciplines to Support Scientific Research

- Development of an advanced Knowledge System to *capture*, *shar* and support *reproducible analysis* for biomarker research

  - Genomics, Proteomics, Imaging, etc data types of data



- NASA-NCI partnership, leveraging informatics and data science technologies from planetary and Earth science

  - Reproducible, Big Data Systems for exploring the universe

  - Software and data science methodology transfer

  - Presented informatics collaboration at a congressional briefing in October 2015

# Early Detection Research Network: Finding Cancer Biomarkers

Jet Propulsion Laboratory
California Institute of Technology

- **A comprehensive infrastructure to support biomarker data management across EDRN's distributed cancer centers**
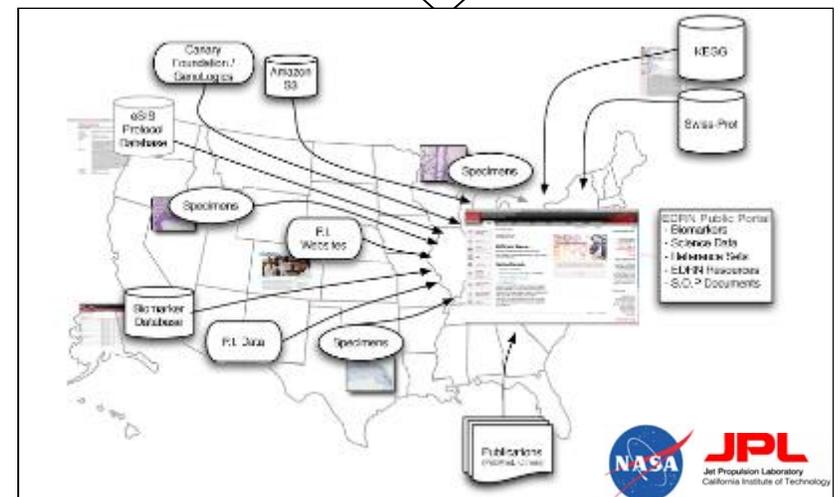  - A national data sharing architecture
  - Data Integration
  - Information model for cancer biomarkers following the PDS4 approach
  - Development of data analytic pipelines
  - Shared open source software capabilities

- **Integration of data across the EDRN (biomarkers, specimens, protocols, biomarker data, publications) including:**
  - Data from over 100 research labs; multiple organs
  - 800+ data elements
  - 900+ biomarkers captured
  - 200+ protocols of study
  - 1500+ publications
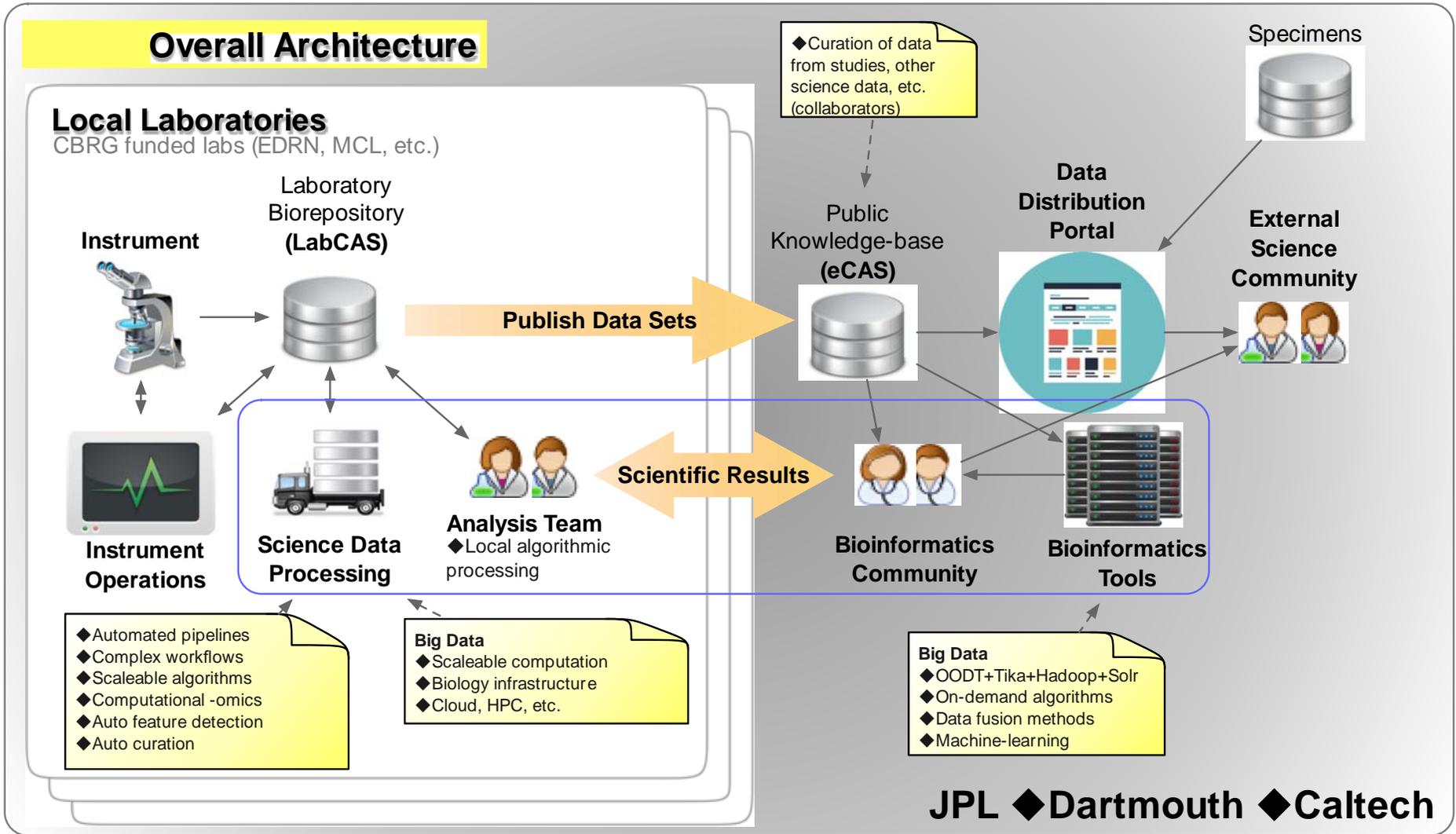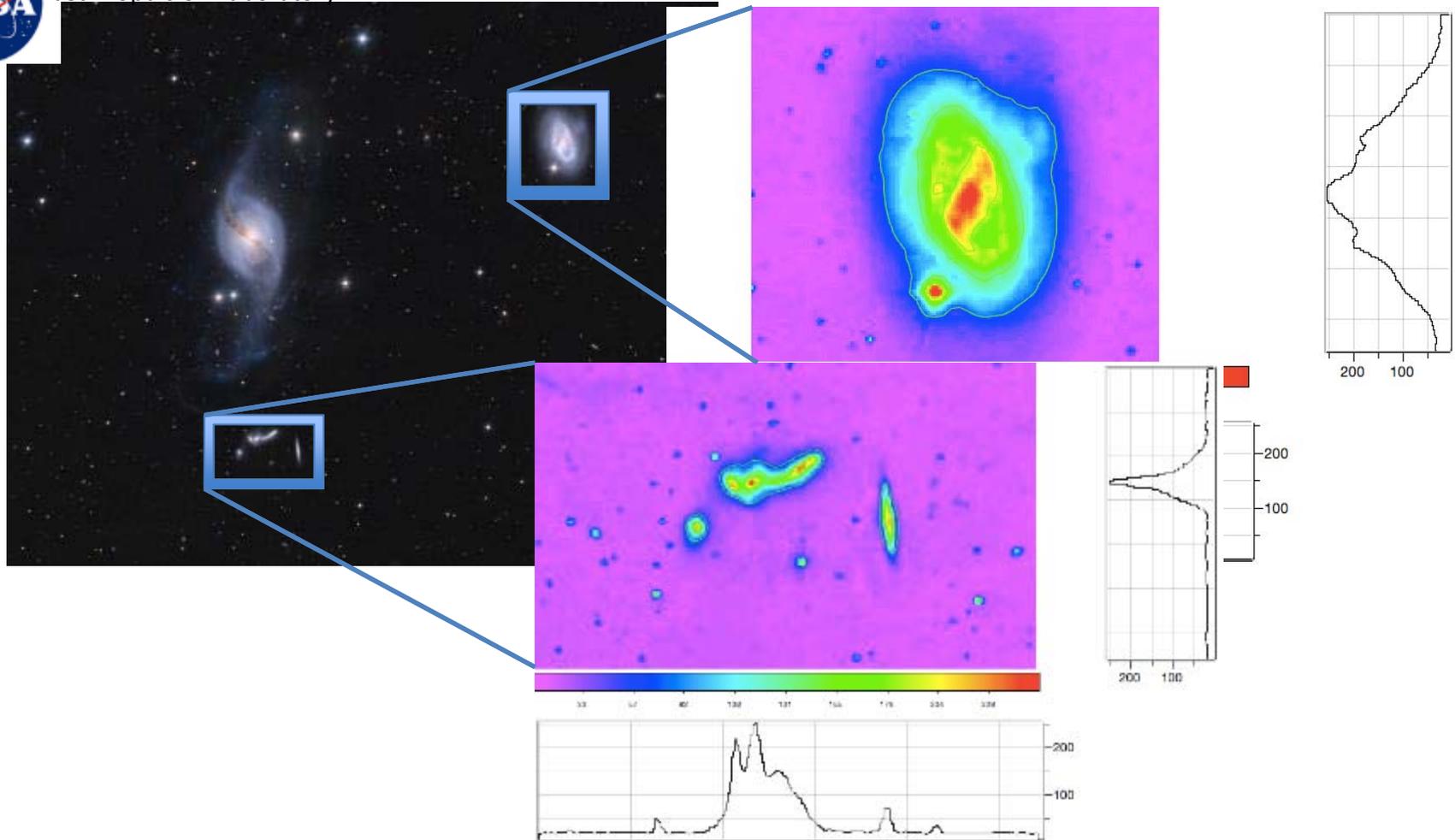  - Multiple terabytes of data from biomarker studies



http://cancer.gov/edrn

# Example of Data Science Capabilities in Cancer Research from NASA

**Jet Propulsion Laboratory**
California Institute of Technology



**Overall Architecture**

**Local Laboratories**
CBRG funded labs (EDRN, MCL, etc.)

- Instrument
- Laboratory Biorepository (LabCAS)
- Instrument Operations
- Science Data Processing
- Analysis Team ◆Local algorithmic processing

**Publish Data Sets**

**Scientific Results**

◆Curation of data from studies, other science data, etc. (collaborators)

- Public Knowledge-base (eCAS)
- Data Distribution Portal
- Specimens
- External Science Community
- Bioinformatics Community
- Bioinformatics Tools

◆Automated pipelines
◆Complex workflows
◆Scaleable algorithms
◆Computational -omics
◆Auto feature detection
◆Auto curation

**Big Data**
◆Scaleable computation
◆Biology infrastructure
◆Cloud, HPC, etc.

**Big Data**
◆OODT+Tika+Hadoop+Solr
◆On-demand algorithms
◆Data fusion methods
◆Machine-learning

**JPL ◆Dartmouth ◆Caltech**

Description: Detecting objects from astronomical measurements by evaluating light measurements in pixels using intelligent software algorithms.
Image Credit: Catalina Sky Survey (CSS), of the Lunar and Planetary Laboratory, University of Arizona, and Catalina Realtime Transient Survey (CRTS), Center for Data-Driven Discovery, Caltech.
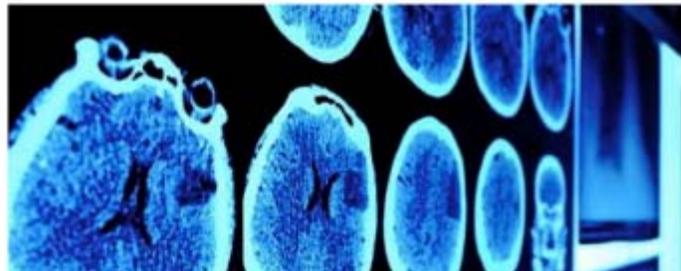
25

Description: Detecting objects from oncology images using intelligent software algorithms transferred to and from space science.
Image Credit: EDRN Lung Specimen Pathology image example, University of Colorado

**INNOVATION**

# 10 ways tech is improving cancer research

New advances in cancer diagnosis and treatment leverage
and even NASA tools to help detect and beat the disease.

By Alison DeNisco 🐦 | September 22, 2016, 6:31 AM PST

Yes.

💬 0     f 114     in 21     🐦     ≡

## 2. NASA: Using space technology to find cancer markers

A NASA machine learning algorithm that
identifies similarities between galaxies will
now analyze tissue samples for signs of
cancer. Earlier this month, NASA's Jet
Propulsion Laboratory and the National
Cancer Institute renewed a research
partnership through 2021 to collect research
on these biomarkers into one searchable
network. This way, physicians can compare,
for example, a CT scan with an archive of
similar images to search for early signs of
cancer, based on a patient's demographics.
Ultimately, this could translate into new techniques for early diagnosis of
cancer or cancer risk.

Dozens of institutions, including Dartmouth College's Geisel School of
Medicine, Harvard Medical School's Massachusetts General Hospital, and
Stanford's NIST Genome-Scale Measurements Group have joined the
network. It is similar to NASA's Planetary Data System, in which all can
share information.

**More about Innovation**

→ When your driverless car
crashes, who will be
responsible? The answer
remains unclear

→ GE makes $1.4B bet on 3D
printing, acquires two firms
to boost additive
manufacturing

→ IoT helping Tassie oyster
farmers avoid unnecessary
closures

→ Subscribe to TechRepublic's
Next Big Thing newsletter.

Sep 22, 2016

Jet Propulsion Laboratory
California Institute of Technology

# Other Partnerships



## Searching deep and dark: Building a Google for the less visible parts of the web

January 8, 2017 8:33pm EST

A geographical map depicting hotbeds of dark web activity related to illegal products. Larger circles indicate more activity. Christian Mattmann, CC

DARPA/Memex, C. Mattmann, JPL



SPAWAR/Data Science for C4CSI,
L. Deforrest, JPL



### Earth System Grid Federation
An open source effort providing a robust, distributed data and computation platform, enabling world wide access to Peta/Exa-scale scientific data.

Learn more ►

DOE/ESGF, L. Cinquini, JPL



## UNC CHARLOTTE WINS $4 MILLION NSF GRANT FOR BIG DATA RESEARCH

Search News and Features
Archive News and Features
UNC Charlotte wins $4 million NSF grant for Big Data research
Faculty Spotlights
Student Spotlights
Research

Tuesday, September 13, 2016

*Ashit Talukder*

The National Science Foundation has awarded a $4 million grant to UNC Charlotte researchers to develop a multidisciplinary research program called Virtual Information Fabric Infrastructure (VIFI) that will create new ways to manage, use and share Big Data and analytic results.

Ashit Talukder, director of the Charlotte Data Visualization Center and the Bank of America Endowed Chair in Information Technology in the College of Computing and Informatics, is the principal investigator for the grant. The award was made under the NSF-CISE/ACI-Data Infrastructure Building Blocks (DIBBS) solicitation.

"Under this large-scale research program, a novel Virtual Information Fabric Infrastructure (VIFI) will be created, allowing scientists to search, access, manipulate and evaluate fragmented, distributed data in the information 'fabric' (the infrastructure to facilitate data sharing) without directly accessing or moving large amounts of data," said Talukder.

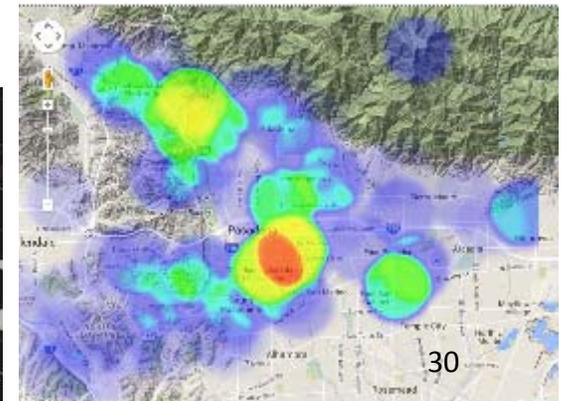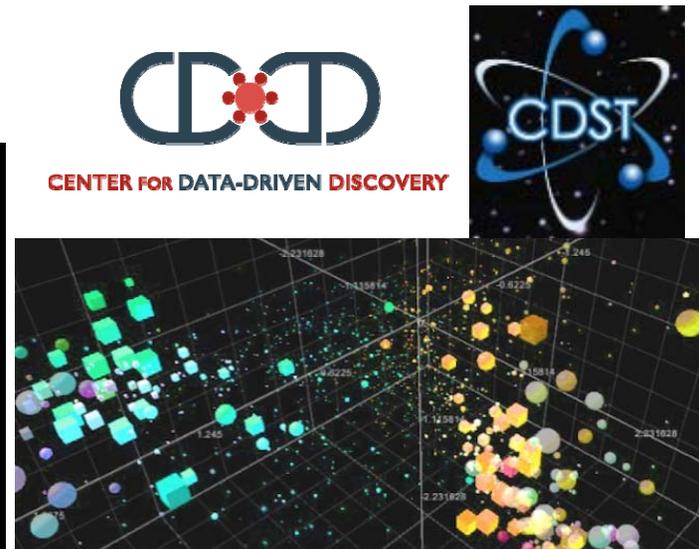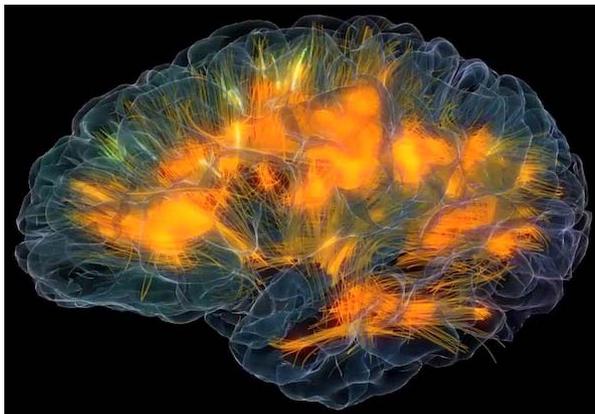NSF/DIBBS, A Talukder, UNC, G. Djorgovski, Caltech,
D. Crichton, JPL

# Driving Forward

# *Caltech-JPL*
# Partnership in Data Science

- Center for Data-Driven Discovery on campus/Center for Data Science and Technology at JPL

- From basic research to deployed systems ~10 collaborations
  - Leveraged funding from JPL to Caltech; from Caltech to JPL

- Virtual Summer School (2014) has seen over 25,000 students
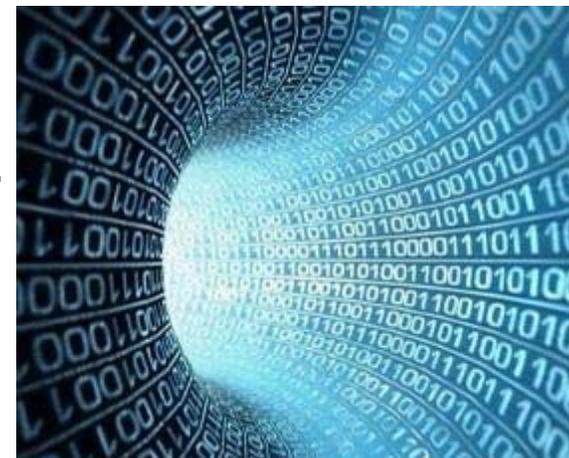
# Example University Partnerships

# Recommendations

- Use the Mission-Science Data Lifecycle to organize Big Data at NASA.
  - From flight computing to data analytics.

- Enable use and data analytics for the community.
  - Promote data ecosystems for sharing data.
  - Support international partnerships.

- Explore opportunities for methodology transfer.
  - Across SMD
  - With other agencies
  - Focused around open source

- Establish multi-disciplinary teams between science/discipline experts, computer science/data science.


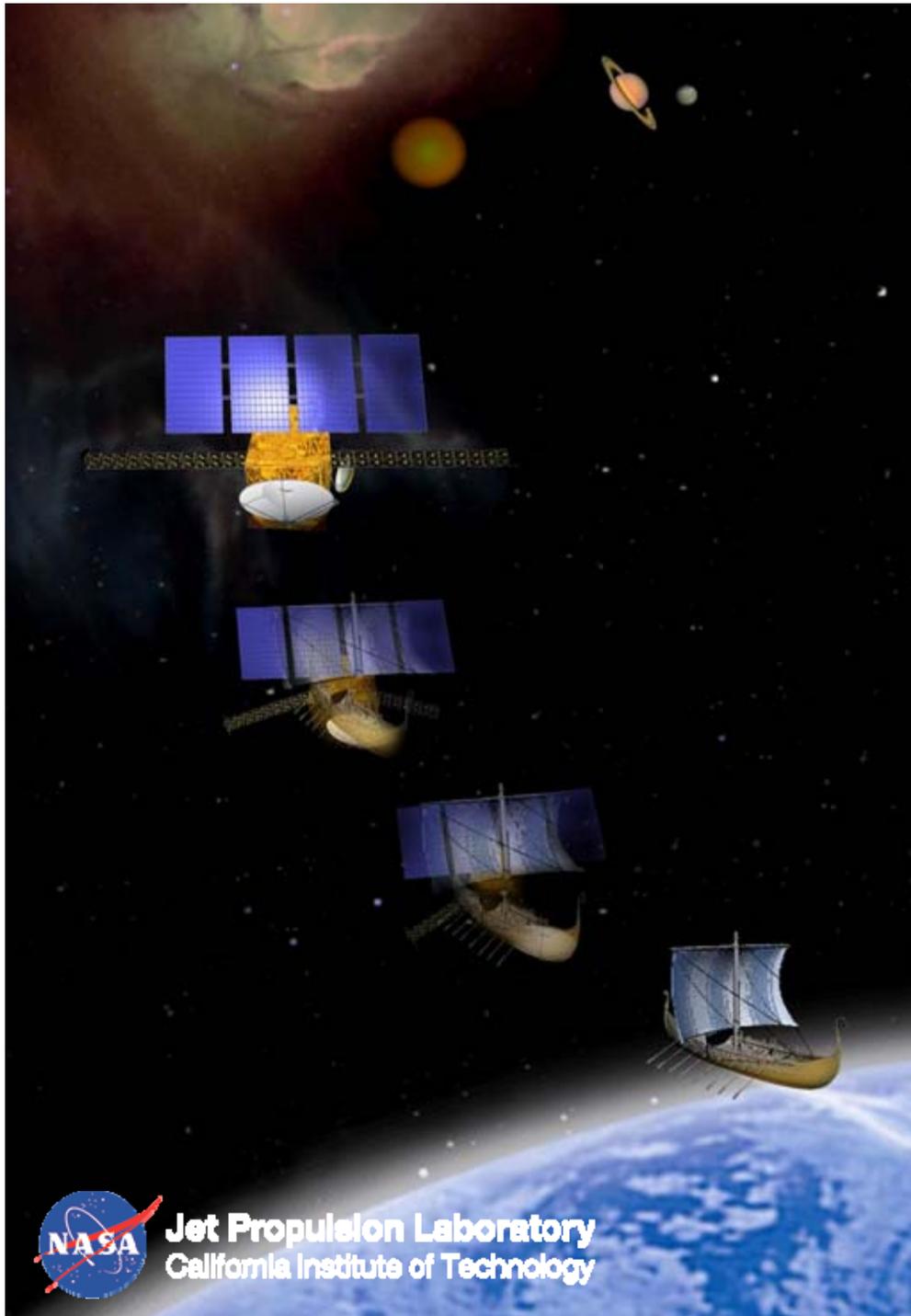
*What do we do with all this data?*



*This is looking like a black hole – but wait, there's light at the end of the tunnel!*

32

- Frontiers on Massive Data Analysis, NRC, 2013

- NASA OCT Technology Roadmap, NASA, 2015

- NASA AIST Big Data Study, NASA/JPL 2016

- IEEE Big Data Conference, Data and Computational Science Big Data Challenges for Earth Science Research, IEEE, 2015

- IEEE Big Data Conference, Data and Computational Science Big Data Challenges for Earth and Planetary Science Research, IEEE, 2016

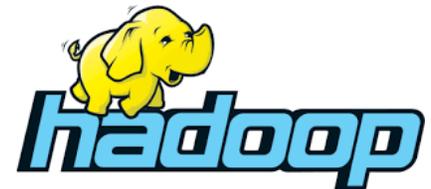- Planetary Science Informatics and Data Analytics Conference, April 2018

**Questions?**

# The Role of Open Source in Big Data Infrastructures

- Open source is an excellent vehicle for collaborations in big data across the science community
  - Great opportunities for sharing software frameworks and tools

- JPL has been involved in the Apache Software Foundation for several years and helped launch Apache in Science.
  - JPLers are *committers* on several Apache projects

# Common Big Data Challenges

- Defining the data lifecycle for different domains in science, engineering, business

- Capturing well-architected and curated data repositories

- Enabling access and integration of highly distributed, heterogeneous data

- Developing novel statistical approaches for data preparation, integration and fusion

- Supporting analysis and computation across highly distributed data environments and silos

- Developing mechanisms for identifying and extracting interesting features and patterns

- Developing methodologies for validating and comparing predictive models vs. measurements

- Methods for visualizing massive data

SPACE TECHNOLOGY RESEARCH GRANTS PROGRAM, Feb 2017