# Data Science Methodology Transfer: Space Science to Biomedicine

Dan Crichton
PI, JPL Informatics Center
Leader, Center for Data Science & Technology
Manager, JPL Data Science Office

Rich Doyle
Manager, Information and Data Science Program Office
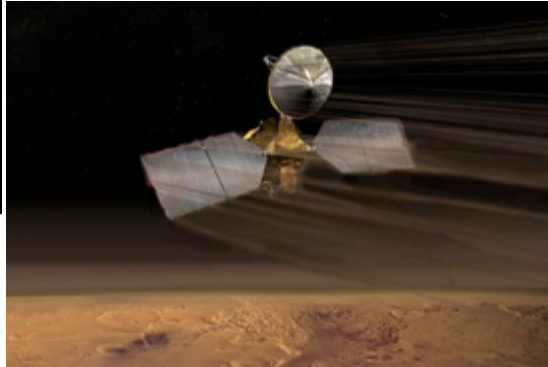Manager, HPSC

June 27, 2019

# Data Lifecycle Model for NASA Missions
## From Onboard Computing to Scalable Data Analytics

**Emerging Solutions**
- *Next-Generation Flight Computing*
- *Onboard Data Analytics*



*Observational Platforms and Flight Computing*

**Scaling Pressures Expose the Need for an Integrated End-to-End Data and Computational Architecture**

**Emerging Solutions**
- *Intelligent Ground Stations*
- *Agile Mission Operations*



*Ground-based Mission Systems*

**Emerging Solutions**
- *Data-Driven Discovery from Archives*
- *Scalable Computation and Storage*



*Interactive Analytics and Visualization and Decision Support*
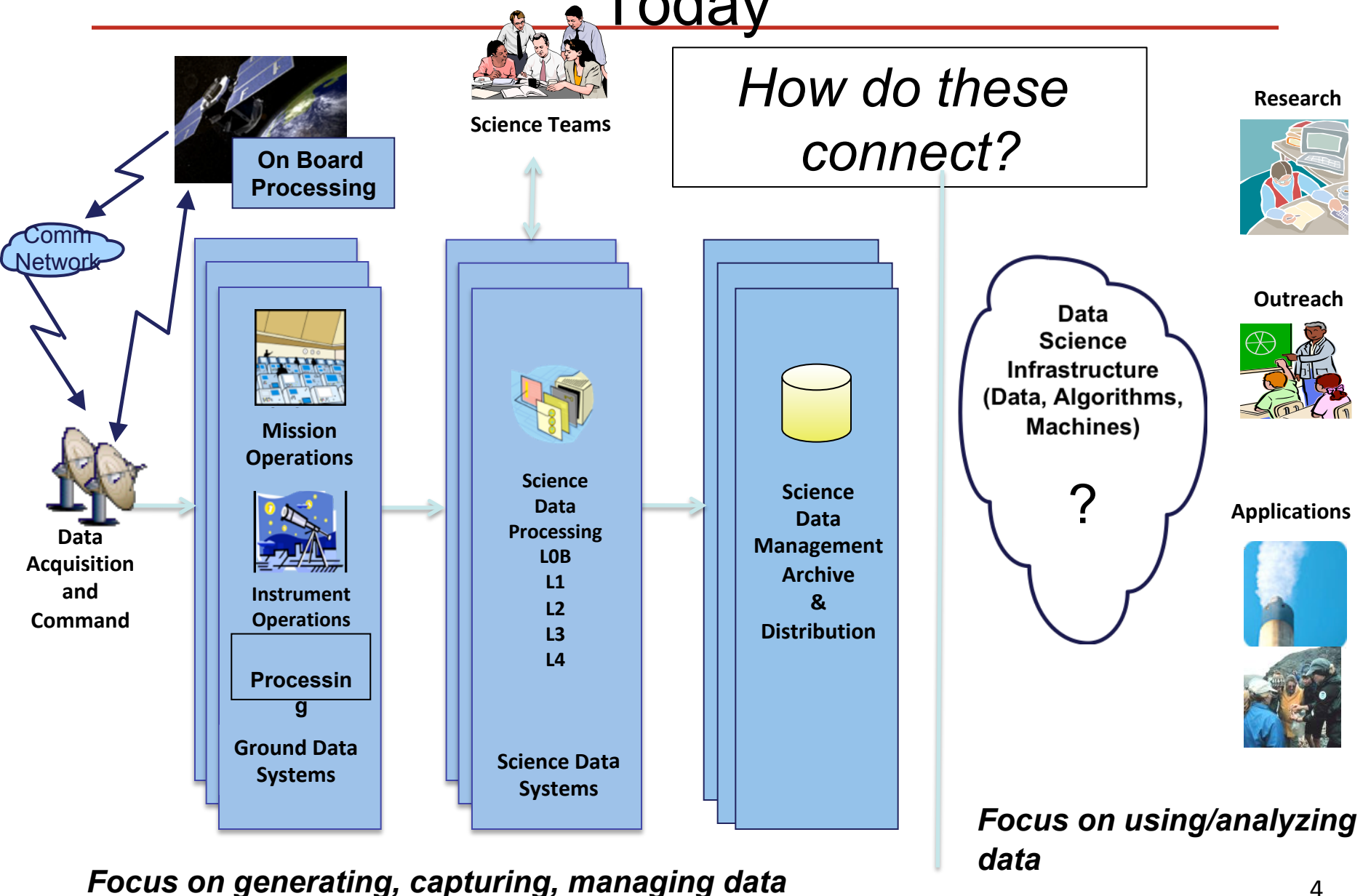
# NASA Data Archives

- **NASA captures robust scientific archives of data to support long-term analysis**
  - A requirement on every mission
  - Adheres to specific standards (both in terms of structure + description/metadata)
  - Mandates public access after a specified time period

- **The capture of well-curated, organized data collections is the basis for enabling data analysis**
  - This is a critical precursor step to analysis
  - The scientific community uses these collections as the basis for studies (often grant funded)
  - The development of robust architectures and systems to support data management and data analysis services are emerging

- *Evolving from research preservation and stewardship towards driving analysis from these archives…*
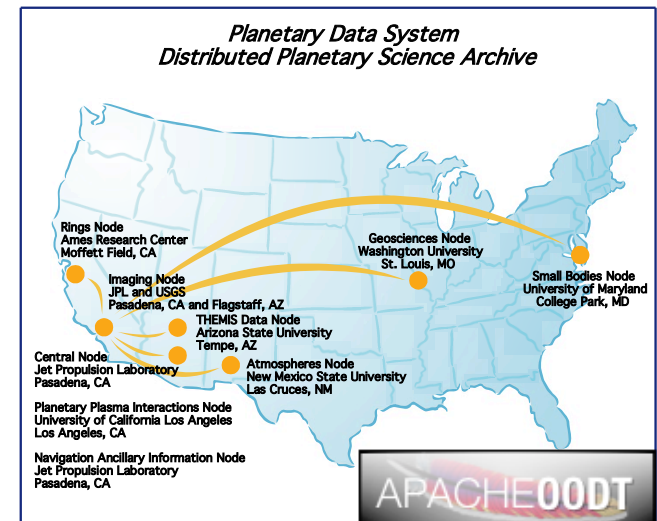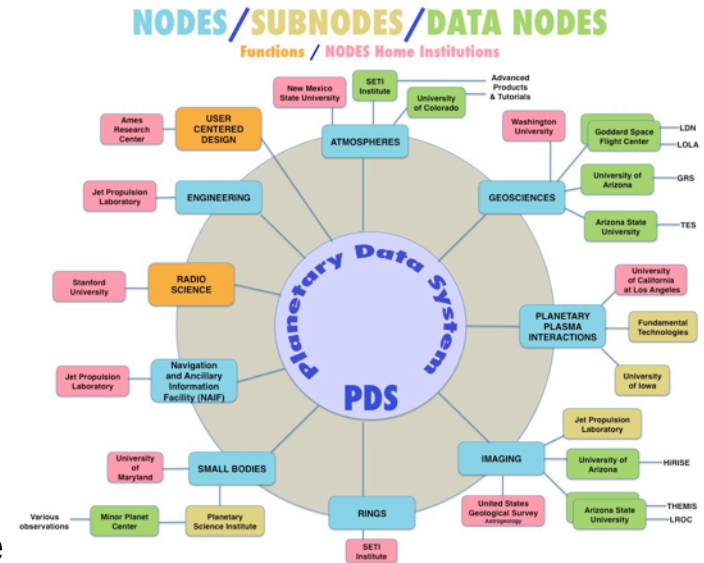
# Planetary Science Data Analytics Support Today



**On Board Processing**

**Science Teams**

*How do these connect?*

**Research**

**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**

Data Science Infrastructure (Data, Algorithms, Machines)

**?**

**Outreach**

**Applications**

*Focus on using/analyzing data*

*Focus on generating, capturing, managing data*

# Expanding to Data-Driven Analytics to Enable Science



Formulation

On Demand Algorithms

Machine Learning/ Deep Learning

Scalable Data Infrastructures

Visualization

On-Demand, Interactive Data Analytics

Research/ Knowledge

Applications

*Today*

**Data Archiving and Stewardship of Massive Data**

Data Integration

Other Data Archives & Models

*Future*

**Data Analytics**

Decision Support

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.

# Planetary Data System

- **Purpose:** 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.

- **Infrastructure:** 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 ontology
  - Approximately 1.7 PB of data
  - About 4000 different types of data and 40M data products
  - International adoption
  - NASA's de facto archive for all planetary data



NODES / SUBNODES / DATA NODES
Functions / NODES Home Institutions



Planetary Data System
Distributed Planetary Science Archive

Rings Node
Ames Research Center
Moffett Field, CA

Geosciences Node
Washington University
St. Louis, MO

Small Bodies Node
University of Maryland
College Park, MD

Imaging Node
JPL and USGS
Pasadena, CA and Flagstaff, AZ

THEMIS Data Node
Arizona State University
Tempe, AZ

Central Node
Jet Propulsion Laboratory
Pasadena, CA

Atmospheres Node
New Mexico State University
Las Cruces, NM

Planetary Plasma Interactions Node
University of California Los Angeles
Los Angeles, CA

Navigation Ancillary Information Node
Jet Propulsion Laboratory
Pasadena, CA

APACHE OODT

# International Collaboration on PDS4 Through IPDA



LADEE
(NASA)

MAVEN
(NASA)

Osiris-REx
(NASA)

ExoMars
(ESA/Russia)

BepiColombo
(ESA/JAXA)

Mars 2020
(NASA)

Psyche
(NASA)



InSight
(NASA)

JUICE
(ESA)

Europa
(NASA)

Hyabussa-2
(JAXA)

Chandrayaan-2
(ISRO)

Lucy
NASA

Endorsed by the **International Planetary Data Alliance** in July 2012 –
https://planetarydata.org/documents/steering-committee/ipda-endorsements-recommendations-and-actions

# Planetary Image Archiving

# Mars Trek: The Google Earth of Mars

# NCI/JPL Informatics Collaboration: Crossing Disciplines to Support Scientific Research



- Development of an advanced Knowledge System to *capture*, process, *share* 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 based on JPL open source technologies and architectures





APACHE OODT

# Biomarkers Knowledge Environment

- Integrate diverse research into an online data environment
  - Integrate data as opposed to re-managing data
  - Be agnostic to data formats and structure

- Provide a well architected data management environment that captures, integrates, and shares data for biomarker research including:
  - Biomarkers
  - Biospecimens
  - Validation Study Information
  - Protocol Information
  - Study Results and Data Sets
  - Publications
  - Artifacts supporting collaboration
  - Curated metadata

- Support diverse community needs

# Two Key Programs



**Organizational Structure of the MCL Consortium**

MCL Steering Committee

MCL Laboratories: JHU, MDACC, Stanford, UCLA, UCSF, UVM, Vanderbilt

NCI Programs: EDRN, TMEN, NCATS, TCGA, PROSPR, BCSC, BETRNet

Coordination and Data Management Group (CDMG) (Dartmouth)

Informatics Center (JPL)

NCI CBIIT

Disease Sites: Breast, Lung, Pancreas, Prostate

EDRN Organizational Structure

Discovery

Assay Development

Validation

Biomarker Reference Laboratories

Biomarker Developmental Laboratories

Clinical Validation Centers

Network Consulting Team — Chair: Larry Norton, M.D.

Steering Committee — Chair: Ian Thompson — Co-Chair: Joshua Labaer

Data Management and Coordinating Center — Director: Ziding Feng, Ph.D.

Develop cross-cutting informatics capabilities to support the capture, curation, management, distribution, and analysis of diverse data and results

# Capture − Process − Search − Access − Analyze - Visualize



Capture data and metadata

Process reproducibly

Search Metadata

Access

Analyze with algorithms, pipelines

Visualize

# MCL Integrated Data Ecosystem Overview

# Portal: Dissemination and Access to Biomarker Data



- Gateway to information
- Information managed both within and outside the knowledge system
- Initial starting point for community to get to research data
- Google-like search to access the wealth of data
- Multi-level Security

http://cancer.gov/edrn
http://mcl.jpl.nasa.gov

15

# Navigating the Knowledge System:
# Data Semantically Linked



Biomarker Annotations



Protocols



Biomarker Data Results



Specimens



Linked through
Public Portal



Access to download data

# Data Models and Elements

Data generated in the course of research from any source







Instruments, etc.

Data linked across the Data Commons



Data is linked and available for integration of different analytical methods and tools to drive data-driven discovery.

Addition of data sets over time – any data type generated

*Drives a consistent data architecture across the consortium.*

# LabCAS: Capturing and Sharing Science Data Analysis Data and Pipelines

- A secure, reliable means to capture, process and manage data

- Plug in analytical methods

    - Repeatable data processing pipelines

- Integrate visualization

# Data Capture, Processing and Ingestion

# Data Pipelines and the Knowledge Environment



## "LabCAS"

**Instrument**

**RNASeq**

**Publish Data Sets with Common Data Elements**

*Laboratory Data Repository*

*Science Data Processing Algorithm*

**Instrument**

**Imaging**

**Publish Data Sets with Common Data Elements**

*Laboratory Data Repository*

*Science Data Processing Algorithm*

**Instrument**

**Pathology**

**Publish Data Sets with Common Data Elements**

*Laboratory Data Repository*
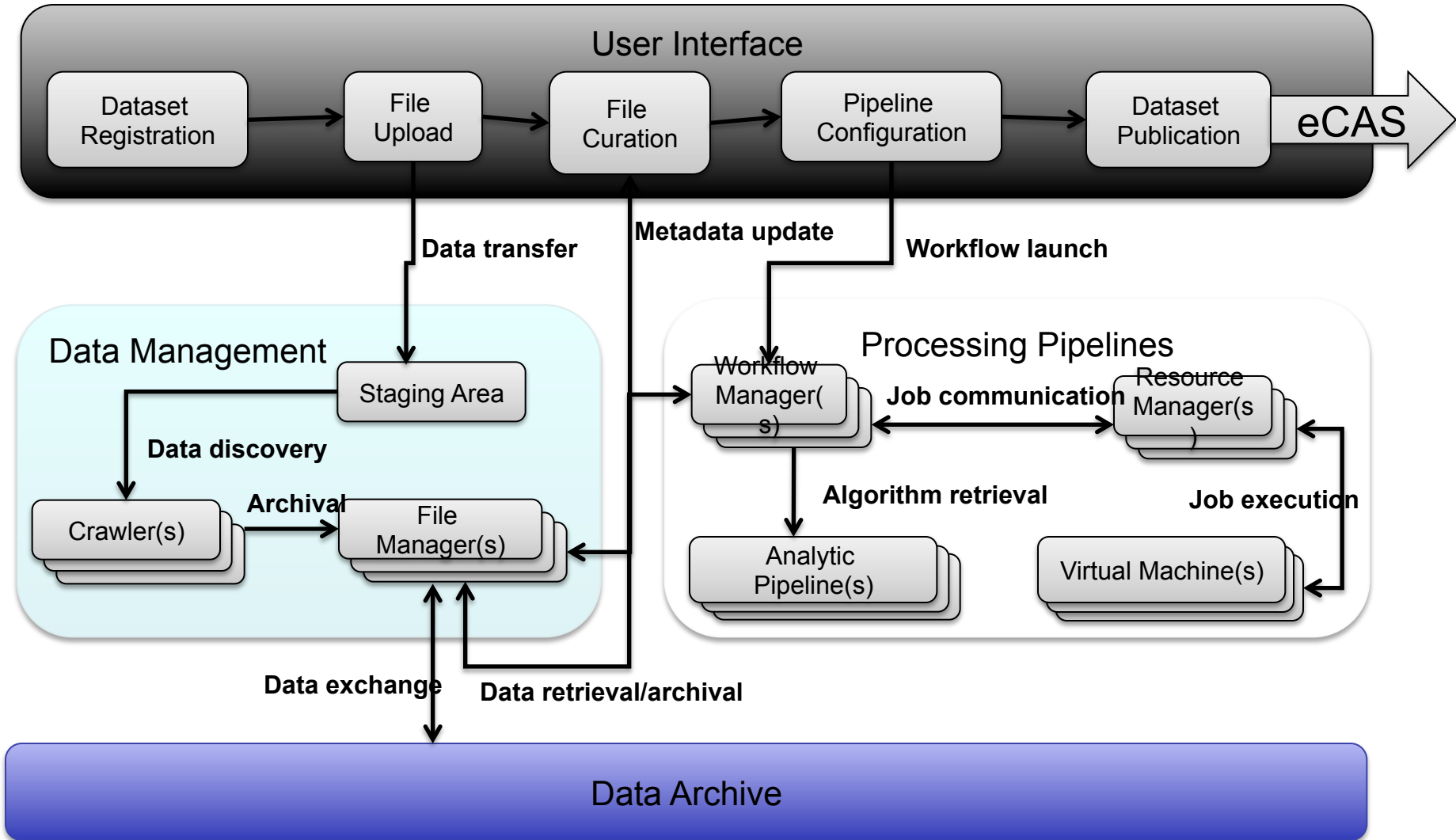
*Science Data Processing Algorithm*

Integrated Knowledge Database

Biomarkers

Image Data

Pathology Data

Genomics Data

Protocols

Sites

## Other Resources

NCI Imaging Data

caBio

Other Resources

Scalable computing, common data elements, computational methods

# Scalable architecture for ML on AWS

- Uses Docker Swarm, Apache OODT workflows (from NASA/JPL), RabbitMQ messaging
- Can scale/auto-scale to any number of EC2 nodes

# Benchmarking on AWS

- Executing genomics workflow on single EC2 server (r5.2xlarge, 8 CPUs, 64 GB memory)

- Measuring time to execute 100 workflows vs # of containers, # of processes/container

- Can scale horizontally until processes start to compete for resources

# Lung Adenocarcinoma Gene Expressions

- In collaboration with Dr. Chris Amos, Yafang Lee (Dartmouth)
- 10 TB of data for full study comparing gene expression profiles between smoking and non-smoking patients.
- Integrated Dartmouth analysis tools with 99% accuracy into LabCAS pipelines

# Smart-3SEQ

- In collaboration with Robert West, Joseph Foley, Sujen Vennam @ Stanford School of Medicine

- Smart-3SEQ: new method for quantifying gene expressions in small RNA samples including single cells (see: https://github.com/jwfoley/3SEQtools)

- Smart-3SEQ pipeline is composed of 3 steps:
  - FastQ generation and alteration (w/ Illumina bcl2fastq)
  - Gene alignment (w/ Samtools, STAR, UMI-dedup)
  - Read counting (w/ Bioconductors

| FastQ generation | → | Gene alignment | → | Read counting |

# DNA-Sequencing

- In collaboration with Olivier Harismendy at UCSD

  – See https://github.com/bcbio, http://bcbio.io

- Using bcbio ("Bue Collar Bioinformatics"): Python framework and community tools for analysis of biological data (variant calling, RNA-SEQ and small RNA pipelines)

- Experimented with running sample pipelines (big and small)

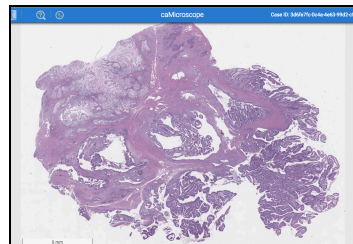| mutect2 | → | freebayes | → | vardict | → | varscan |

# Multi-dimensional Integrated Data

Collaboration with Sandy Borowsky/UC Davis and Anirban Maitra/ MD Anderson

• Developed a pathology archive for MCL

• Dr. Maitra goal: multi-dimensional data – IPMN ppts
Collaboration with Radka Stoyanova and Alan Pollack

• MAST (Mapped Active Surveillance Trial). Longitudinal multivariate data (mpMRI, pathology and gene expression) is obtained from patients on Active Surveillance for prostate cancer.

*Path Viz Tool Credit:*
*caMicroscope - Dr. Joel Saltz Lab*
*ITCR Program*

# Support for Image Archiving

# Integrate different analytical tools
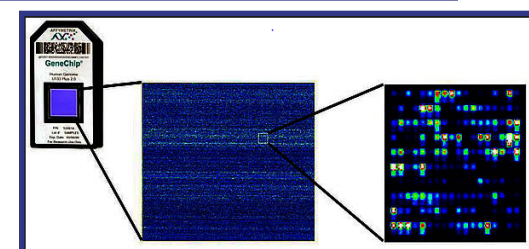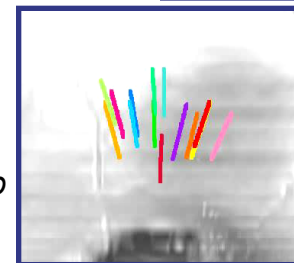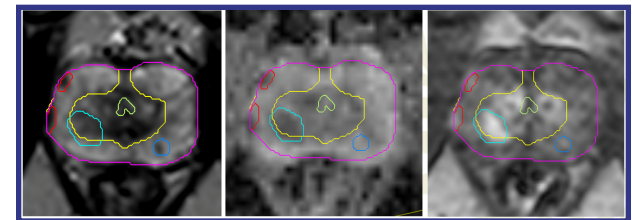
- Framework to plug-in many open-source viewers (client and server side)
  - caMicroscope
  - OHIF DICOM Viewer
  - 3D Slicer
- Capture additional metadata
- User specified organization of data
- Explore VR and other approaches for multii-dimensional imaging









3D rendered isosurfaces: section of lung tissue

# Amazon Workspaces

– Remote desktops running on the AWS Cloud

– Support analysis tools that must be installed on client machines, to access local data:

- QuPath for use by Sandy Borowsky's group (UC Davis)

- 3D Slicer for use by Radka Stoyanava's group (Univ. Of Miami)

# AWS Workspaces: QuPath

UC Davis Pathology data



UCSD DNA-Sequencing data

# AWS Workspaces: 3dSlicer



Moffitt lung data

Univ. of Miami prostate data

# Feature classification in images

# Crowd Sourcing Image Analysis



Adaptation of Zoonverse
Reduce False Positives from traditional ML
Enable radiologists seed unsupervised clusters
Drive towards ML pipelines

# What's in Place Today

- A national, biomarker knowledge system serving multiple programs

- A biomarker data infrastructure consisting of ~1000 biomarkers, ~200 protocols, 1500 publications, 100 TB data

- Tools for laboratories to support the processing, capture, curation and sharing of data before publications

- Pilot projects in imaging, scalable workflows, data integration, etc

- Support for data-driven approaches for data discovery and analysis

- Common portals to access the knowledge environment

# NIH Data Science Strategic Plan*

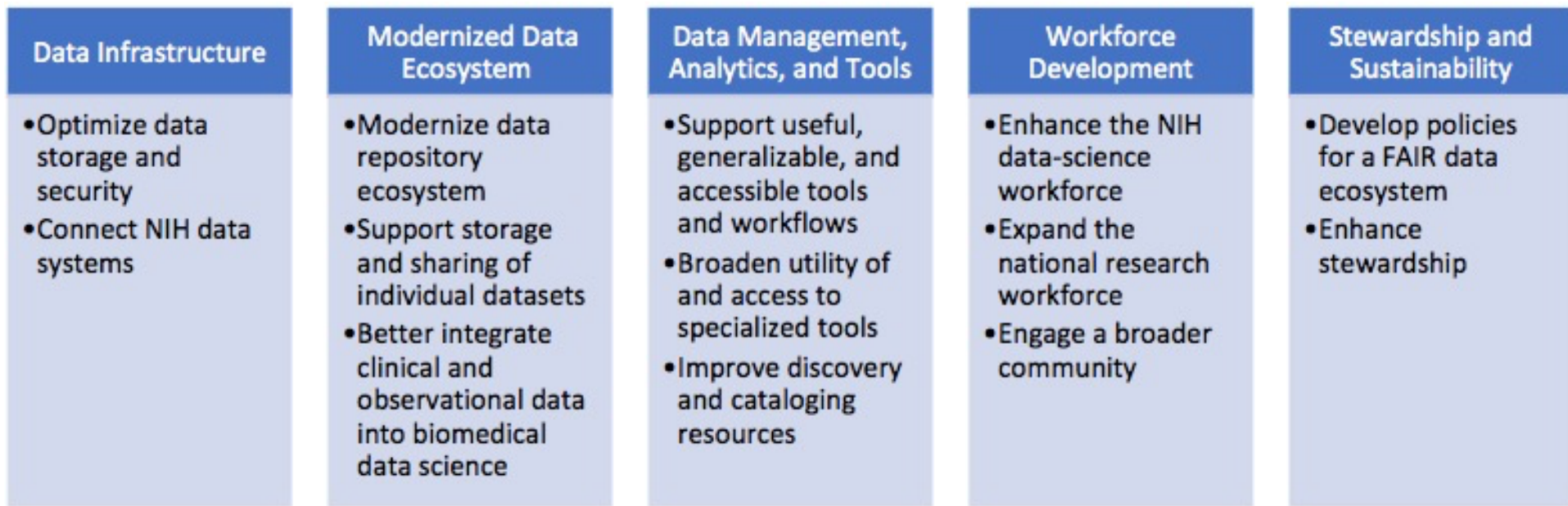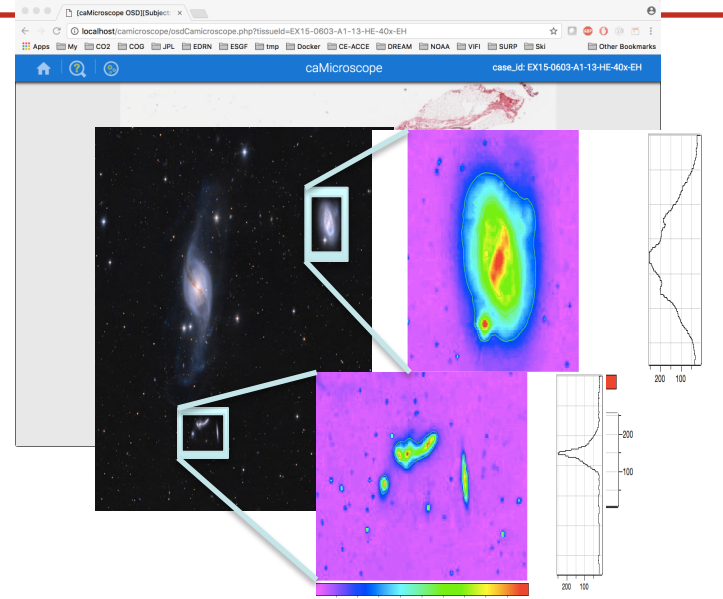| Data Infrastructure | Modernized Data Ecosystem | Data Management, Analytics, and Tools | Workforce Development | Stewardship and Sustainability |
|---|---|---|---|---|
| • Optimize data storage and security<br>• Connect NIH data systems | • Modernize data repository ecosystem<br>• Support storage and sharing of individual datasets<br>• Better integrate clinical and observational data into biomedical data science | • Support useful, generalizable, and accessible tools and workflows<br>• Broaden utility of and access to specialized tools<br>• Improve discovery and cataloging resources | • Enhance the NIH data-science workforce<br>• Expand the national research workforce<br>• Engage a broader community | • Develop policies for a FAIR data ecosystem<br>• Enhance stewardship |

**Figure 2.** NIH Strategic Plan for Data Science: Overview of Goals and Objectives
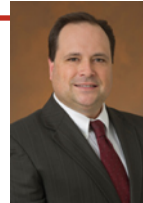
\* By 2025, the total amount of genomics data alone is expected to equal or exceed totals from the three other major producers of large amounts of data: astronomy, YouTube, and Twitter.

# Future Directions

- Systematize the capture of data end-to-end into a knowledge environment

- Integrate data-driven techniques and tools such as machine learning as part of the end-to-end knowledge environment

- Enable collaborative analysis that scales

- Unify consortium enterprises and data

- Enable science through an explicit and well architected data science strategy and platform

# JPL Informatics Center Data Science Team


Dan Crichton
NASA/JPL
Principal Investigator

Data Architecture →
•Luca Cinquini NASA/JPL
•Asitang Mishra NASA/JPL

Machine Learning Visualization →
Ashish Mahabal
Caltech
•Alphan Altinok
NASA JPL
•Santiago Lombeyda
Caltech

Portal System Engineering →
Sean Kelly    David Liu
NASA/JPL

Bioinformatics Biomarker Curation/CDEs →
Kristen Anton Maureen Colbert
University of North Carolina

Project Management Data Coordination →
Heather Kincaid
NASA/JPL

System Admin Cloud Computing →
Paul Zimdars
NASA/JPL
•Susan Neely NASA JPL
•Rojeh Yaghoobi NASA/JPL

# Backup