# An Extensible and Scalable Knowledge System Architecture for Cancer Research

Daniel Crichton
Dan.Crichton@jpl.nasa.gov
Principal Computer Scientist and Program Manager
Director, Center for Data Science and Technology
Principal Investigator, JPL Informatics Center
NASA Jet Propulsion Laboratory, California Institute of Technology

# NASA Big Data Challenges





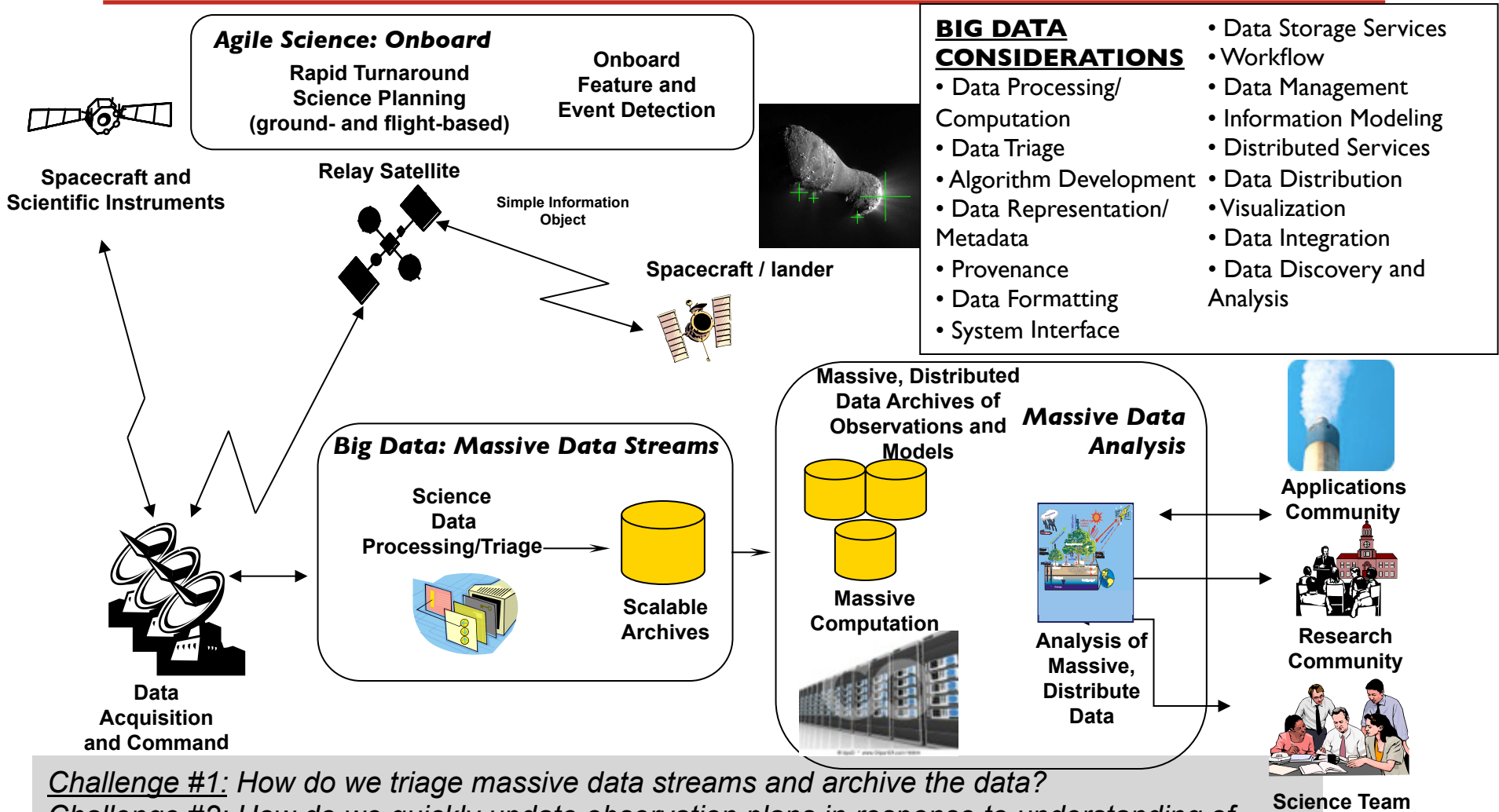- Scientific missions that explore the Earth and solar system return a wealth of data

- Capturing and archiving these data are essential for preserving these data and <u>supporting research and analysis</u>

- Federal research analysis programs at NASA require data be from public archives

# NASA/JPL Big Data Research Areas

**Agile Science: Onboard**
Rapid Turnaround
Science Planning
(ground- and flight-based)

Onboard
Feature and
Event Detection

**Spacecraft and Scientific Instruments**

**Relay Satellite**

Simple Information Object

**Spacecraft / lander**

**BIG DATA CONSIDERATIONS**
- Data Processing/ Computation
- Data Triage
- Algorithm Development
- Data Representation/ Metadata
- Provenance
- Data Formatting
- System Interface
- Data Storage Services
- Workflow
- Data Management
- Information Modeling
- Distributed Services
- Data Distribution
- Visualization
- Data Integration
- Data Discovery and Analysis

**Big Data: Massive Data Streams**
Science Data Processing/Triage

Scalable Archives

**Massive, Distributed Data Archives of Observations and Models**

**Massive Data Analysis**

Massive Computation

Analysis of Massive, Distribute Data

**Applications Community**

**Research Community**

**Science Team**

**Data Acquisition and Command**

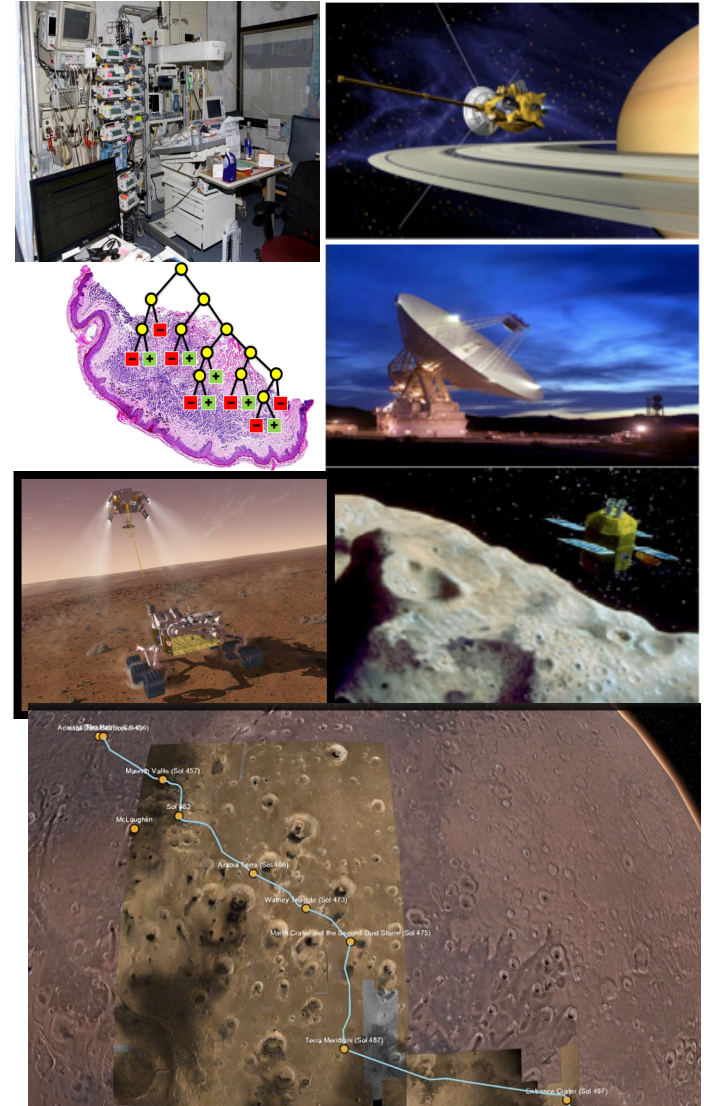*Challenge #1:* How do we triage massive data streams and archive the data?
*Challenge #2:* How do we quickly update observation plans in response to understanding of newly acquired science data, especially for time-limited missions?
*Challenge #3:* How can we use advanced data science methods to systematically derive scientific inferences from massive, distributed science measurements and models?
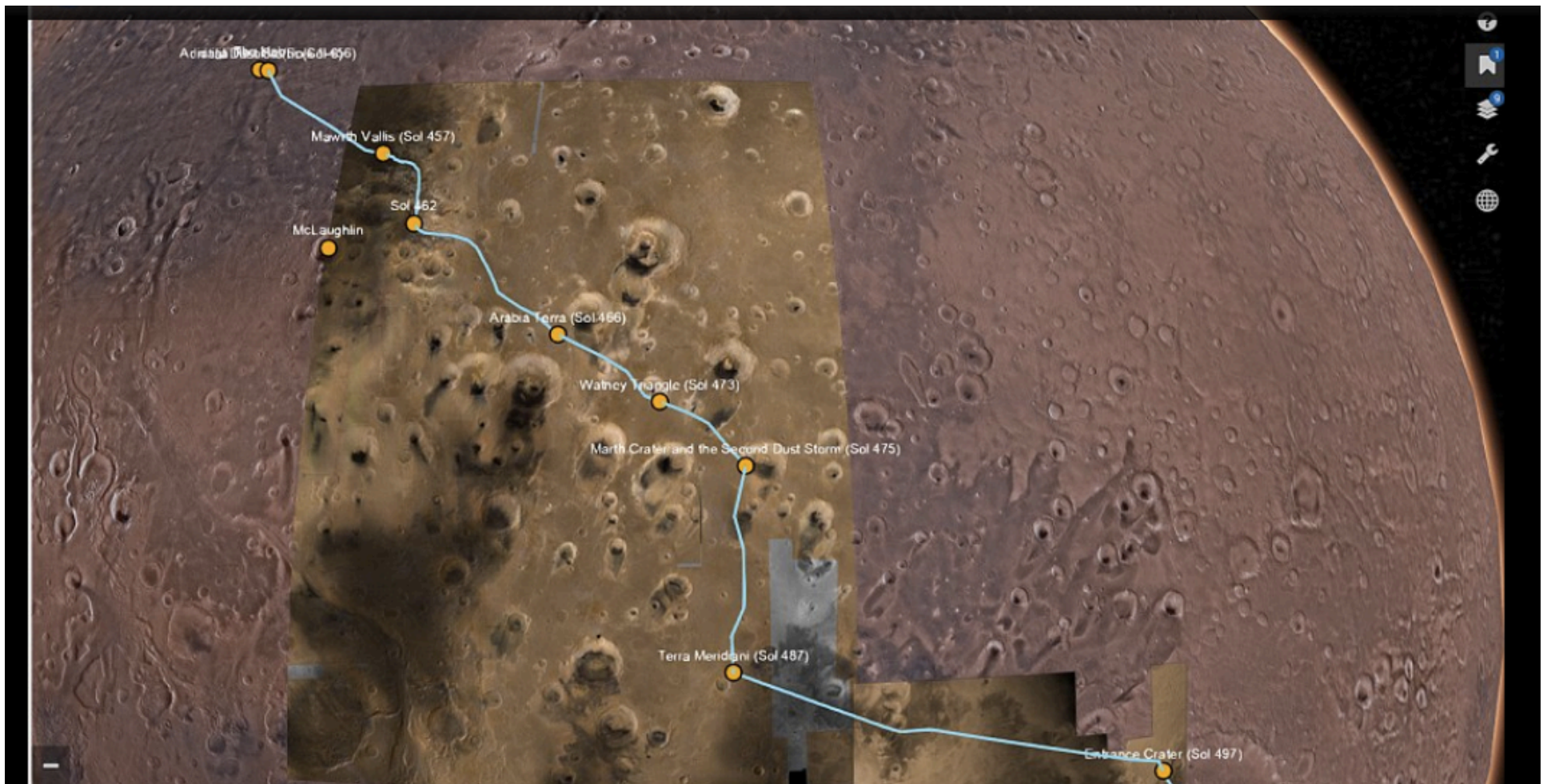
# NASA/JPL Informatics Center:
## Crossing Disciplines to Support Scientific Research



- Development of an advanced Knowledge System to *capture*, *share* and support *reproducible analysis* from the biomarker data results
  - Supporting EDRN program since 2001
  - Supporting MCL program since 2015

- NASA-NCI partnership, leveraging informatics and data science technologies from planetary and Earth science
  - Reproducible, Big Data Systems for exploring the universe
  - 2011 NASA Group Award for "*innovative use of NASA software technologies to support cancer research*"

# Creating Knowledge from Big Data: Exploring Mars Imaging Data with Mars Trek



Derived from about 800 TBs of Imaging Data in the NASA Planetary Data System

# Informatics and Data Science Needs and Capabilities



**National Data Sharing Architectures**



**Common Data Elements & Models**



**Big Data Infrastructure**



**Analytical Data Pipelines**



**Intelligent Data Algorithms**



**Visualization Techniques**

# EDRN – Early Detection Research Network

- Early Detection Research Network (EDRN) is a network of 40+ institutions all performing research geared towards the discovery of cancer biomarkers, which are early indicators of onset of disease
- NCI/NIH funded program
  - Started in ~2000
  - NCI's flagship program
- Informatics efforts cited as a model for biomarker research
- Cross-disciplinary collaboration (FHCRC, JPL, Dartmouth and NCI)

Discovery

Assay Development

**Biomarker Reference Laboratories**

**Biomarker Developmental Laboratories**

**Clinical Validation Centers**

Validation

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.

**EDRN Organizational Structure**

# MCL - Consortium for the Characterization of Screen Detected Lesions



To conduct a comprehensive molecular and cellular characterization of tumor tissue, cell, and microenvironment components to distinguish screen-detected early lesions from interval and symptom-detected cancers.  (Funded 2015)

# Key Informatics Accomplishments in Life Sciences

- Implemented a national, biomarker knowledge system re-using advanced informatics technology developed for planetary data science

- Pioneered the concept of providing access to information about biospecimens across EDRN at a national level (2001)

- Developed a repository for capturing scientific data sets; captured 90 data sets; integrated with the Canary Foundation infrastructure.

- Developed a biomarker database for capturing and annotating EDRN biomarkers; high-quality curation on more than 900 biomarkers

- Developed a public portal that provides dissemination of EDRN information as well as scientific data and results; over 2400 unique visitors a month

- Developing new tools for the Laboratories to support the processing, capture, curation and sharing of data before publications

- Received NASA Award in 2011 for the "innovative use of NASA software technologies to support cancer research" due to significant reuse of capability

- Began leveraging the architecture across multiple NCI programs

# Cancer Biomarker Bioinformatics Workshop

- The NCI and NASA Jet Propulsion Laboratory held a workshop in May 2013 at Caltech to address informatics and data-driven research in cancer biomarkers
  - http://edrn.nci.nih.gov/cancer-bioinformatics-workshop/cancer-biomarker-bioinformatics-workshop-report-may-2013
  - A major outcome focused on data usability, reproducibility of results, methods and algorithms to systematize data analysis, and scalable computing infrastructures.

- Key Recommendations
  - Systematic approaches to the generation, capture, management of data to enable reproducibility
  - Increased emphasis on data curation to promote data reuse
  - Automation of data process/analytics software pipelines
  - Data integration and fusion of data from multiple platforms, studies
  - Scalable data infrastructures and repositories
  - Use of big data tools and bioinformatics techniques to scale data analysis
  - Increased training of scientists in the use of computational tools/methods

# Data and Computational Flow



**Overall Architecture**

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

Instrument

Laboratory
Biorepository
(LabCAS)

Publish Data Sets

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

Specimens

**Data
Distribution
Portal**

Public
Knowledge-base
(eCAS)

**External
Science
Community**

Instrument
Operations

Science Data
Processing

**Analysis Team**
◆ Local algorithmic
processing

Scientific Results

**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**

# Common Data Elements and Information Models

- CDEs provide a common set of data semantics to capture and share data
  - EDRN CDEs
  - MCL CDEs

- Work with CBIIT to reuse standard elements in caDSR for consistency across NCI

- All archived data are captured as information objects (same as planetary)
  - Metadata described using CDEs
  - Data captured and stored in a data repository





Biomarker Ontology Information Model

# LabCAS: Laboratory Catalog and Archive Service

*LabCAS is a new capability under development*

- Provides investigators with a secure, reliable means to capture their **pre-publication** research datasets

- Provides integrated data processing

- Enables investigators and collaborative groups/projects to share data in a secure manner as <u>early as possible</u>

- Scales to support data intensive projects

- Facilitates repeatable data processing pipelines

# Data Capture , Processing and Ingestion

# eCAS: Capture and Sharing of Public Data Sets

- **EDRN has a warehouse of public biomarker data**
  - Uses the CDEs to populate a catalog describing the data sets
  - Supports public release/access to the data
  - Supports peer review of the data by collaborative groups prior to public release
  - Integrated with the rest of the knowledge system
  - Supports reproducibility studies

- **Provides a long term and central capture of biomarker study results for the broad community**

- **Being extended to MCL**



Credit: Sam Hanash (Validation of Protein Markers for Lung Cancer Using CARET Sera and Proteomics Techniques)

# Virtual, Distributed Specimen System

H. Lee Moffitt Cancer Center

University of Texas, San Antonio

Creighton University

University of Colorado

University of Pittsburgh

University of Michigan/Dartmouth College
   (Great Lakes New England Consortium)

Brigham and Womens Hospital

MD Anderson

New York University

UCSD

Center for Disease Control

Johns Hopkins

Duke University

Fred Hutchinson Cancer Research Center

Fox Chase Cancer Research Center



- 51002 - Blood
- 336 - Bone Marrow
- 17618 - Tissue
- 555 - Bronchial Washings
- 12956 - Sputum
- 6523 - Urine



National Data Sharing Infrastructure
Supporting Collaboration In Biomedical Research For EDRN

**JNCI** CANCER SPECTRUM

Content Sources... ▼ go

HOME    SEARCH    BROWSE BY TOPIC    CUSTOMIZE    HELP    FEEDBACK

**Institution: NIH Library** | Sign In as Personal Subscriber | Contact Subscription Administrator at your Institution | FAQ

Go To: Home > JNCI > Archive > Vol. 95, No. 3 > Tenenbaum, pp. 186-187.

# JNCI
*Journal of the National Cancer Institute*

## NEWS

# Serving Up Specimens: NASA-NCI Project Links Databases Across the Country

David Tenenbaum

# Enrichment with Biological Database References and Information



Facilitates capture and sharing of 'omics annotations

Provides connection to the following:
- Protocol
- Scientific Data
- Publications
- Additional Biomarker Resources

# Portal: Dissemination and Access to Knowledge System Data



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

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

# Navigating the Knowledge System: Data Semantically Linked



Biomarker Annotations



Protocols



Biomarker Data Results



Specimens



Linked through Public Portal



Access to download data

# A Virtual, National Integration Biomarkers Knowledge System

# Moving towards data-driven discovery for cancer biomarkers



**"LabCAS"**

**Instrument**

**Laboratory Biorepository**

Publish Data Sets

**Public Biorepository**

**Data Distribution**

**External Science Community**

"eCAS"

Results

**Analysis Team**
- Local algorithms processing

**Bioinformatics Tools**

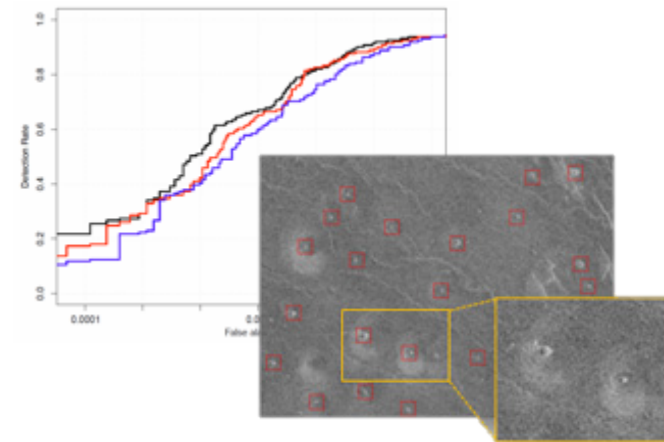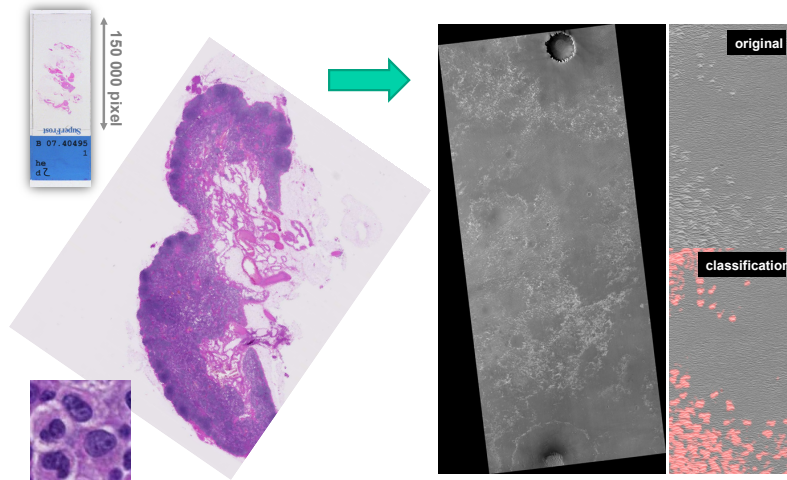**Instrument Operations**

**Science Data Processing**

- Automated pipelines
- Complex Workflows
- Scalable Computational Algorithms
- (genomic, proteomic)
- Automated feature detection
- Automated curation

- Scalable Computational Biology Infrastructures (cloud, HPC, etc)

- On-demand algorithms
- Algorithms
- Data fusion methods
- Machine learning techniques

**Bioinformatics Community**

22

-Cross-cutting Data/Information Architectures -

# Application of Machine Learning Techniques



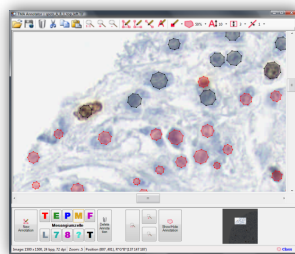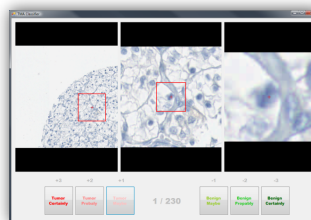**Volcanoes on Venus**



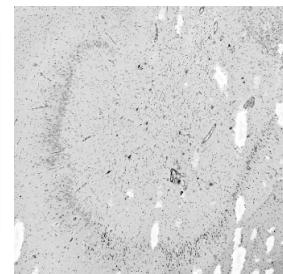## TMA Estimator

Estimate the Staining on a whole spot

## TMA Annotator

Detect nuclei on a whole spot

## TMA Classifier

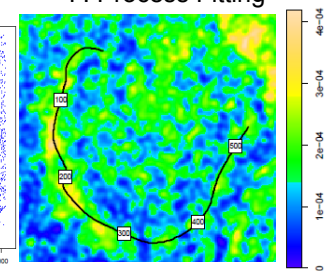**Classify** single nuclei into tumor, non-tumor and stained, not-stained

**Automated Classification**

Original Image
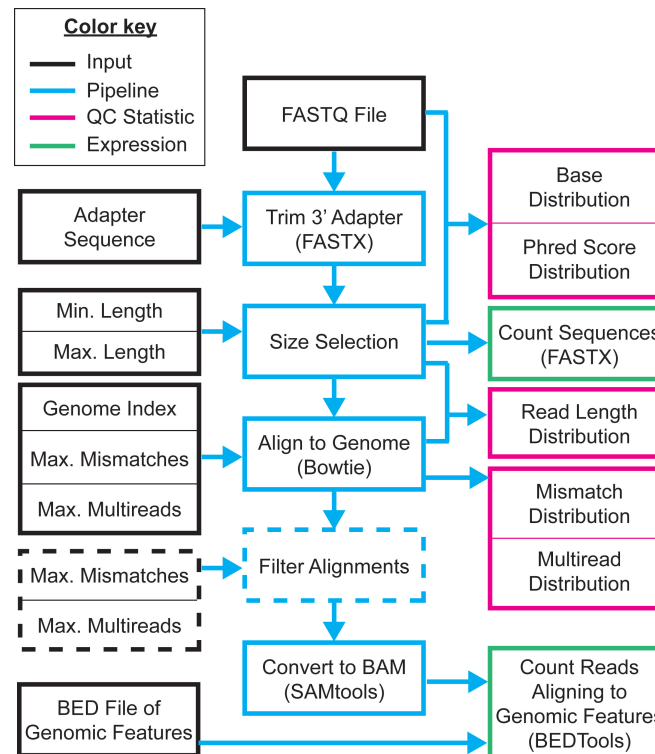
**Discriminative** Object Detection

**Generative** P. Process Fitting

**Feature/Object Detection**
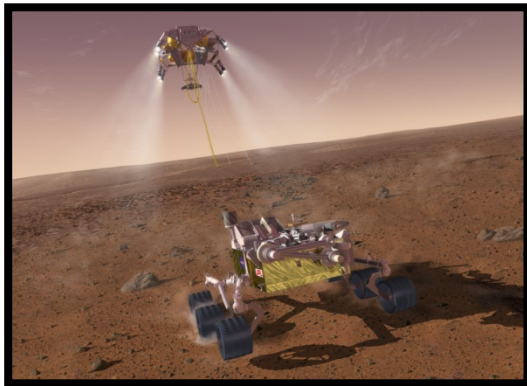
# Potential Collaboration Areas

- Reproducibility experiments

- Analytic Data Pipelines/Computational Methods

- Bioinformatics Tool Integration

- Data Integration

- Data Capture, Curation and Sharing

- Image Archiving and Visualization

- Linking of Distributed Capabilities

- Ontologies, Common Data Elements, etc.



*Courtesy of Josh Campbell and Teresa Wang*

Reproduced from *Wired* magazine

The JPL Informatics Center would like to collaborate with the ITCR Program to explore how data, tools, and methods can be shared to expand the knowledge system and support the EDRN and MCL Programs and other NCI research.

http://twitter.com/edrn_ic

http://www.facebook.com/group.php?gid=56938589930