

JPL-Caltech Virtual Summer School

Big Data Analytics

September 2 – 12, 2014

Richard Doyle, *Information and Data Science Program Office*

Daniel Crichton, *Center for Data Science and Technology*

NASA Jet Propulsion Laboratory

Welcome and Introductory Remarks



Sponsored by the Jet Propulsion Laboratory and by the Keck Institute for Space Studies, California Institute of Technology





Objectives

- Covers topics in the area of **data-intensive science** ranging from programming to computational methods for analyzing massive data
- Oriented towards computational science rather than computer science
- Familiarizes students with concepts and support learning through hands-on labs
- Delivered entirely as a MOOC

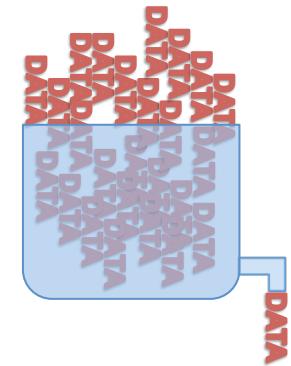


Introduction

What is Big Data? What to Do About It?

The CHALLENGE: Big Data

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

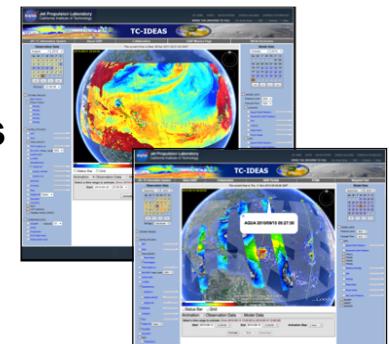


The SOLUTION: Data Science

- Scalable architectural approaches, techniques, software and algorithms which alter the paradigm by which data is collected, managed and analyzed

The RELEVANCE:

- Addressing the challenge of Big Data is on the critical path for many organizations
 - the size and distribution of data sets and predictive models continues to burgeon
 - core objectives such as providing the basis of an analytic result are becoming compromised



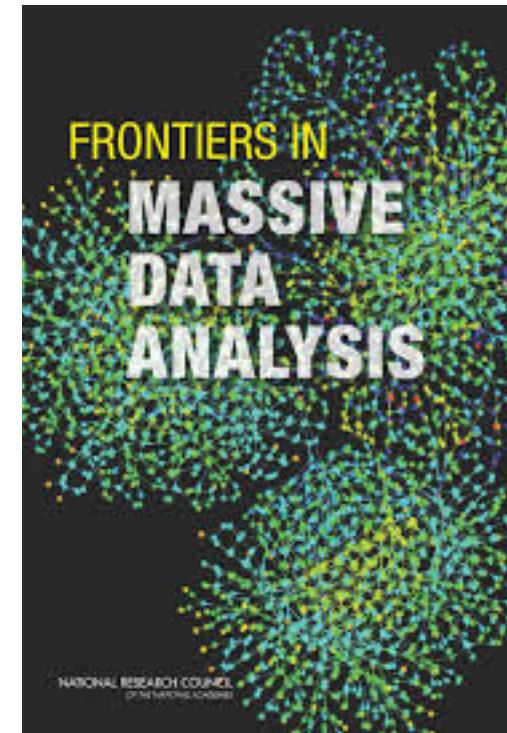


NASA
Jet Propulsion Laboratory
California Institute of Technology

NRC Report

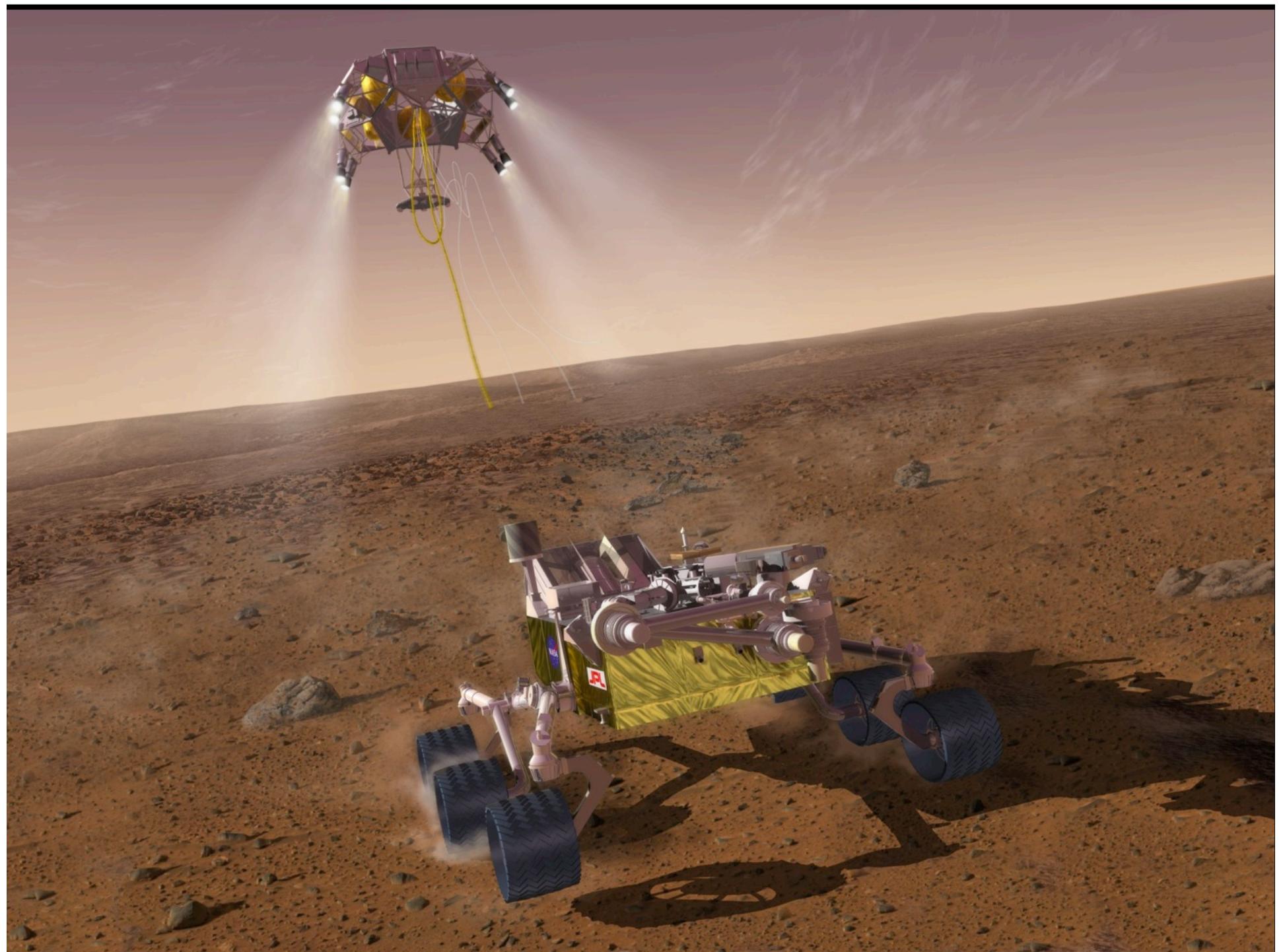
Frontiers in the Analysis of Massive Data

- Chartered in 2010 by the National Research Council
- Chaired by Michael Jordan, Berkeley, AMP Lab (Algorithms, Machines, People)
- Co-author: Dan Crichton, JPL
- Consideration of the architecture for big data management and analysis
- Importance of systematizing the analysis of data
- Need for end-to-end lifecycle: from point of capture to analysis
- Integration of multiple discipline experts
- Application of novel statistical and machine learning approaches for data discovery



**Published in
2013**

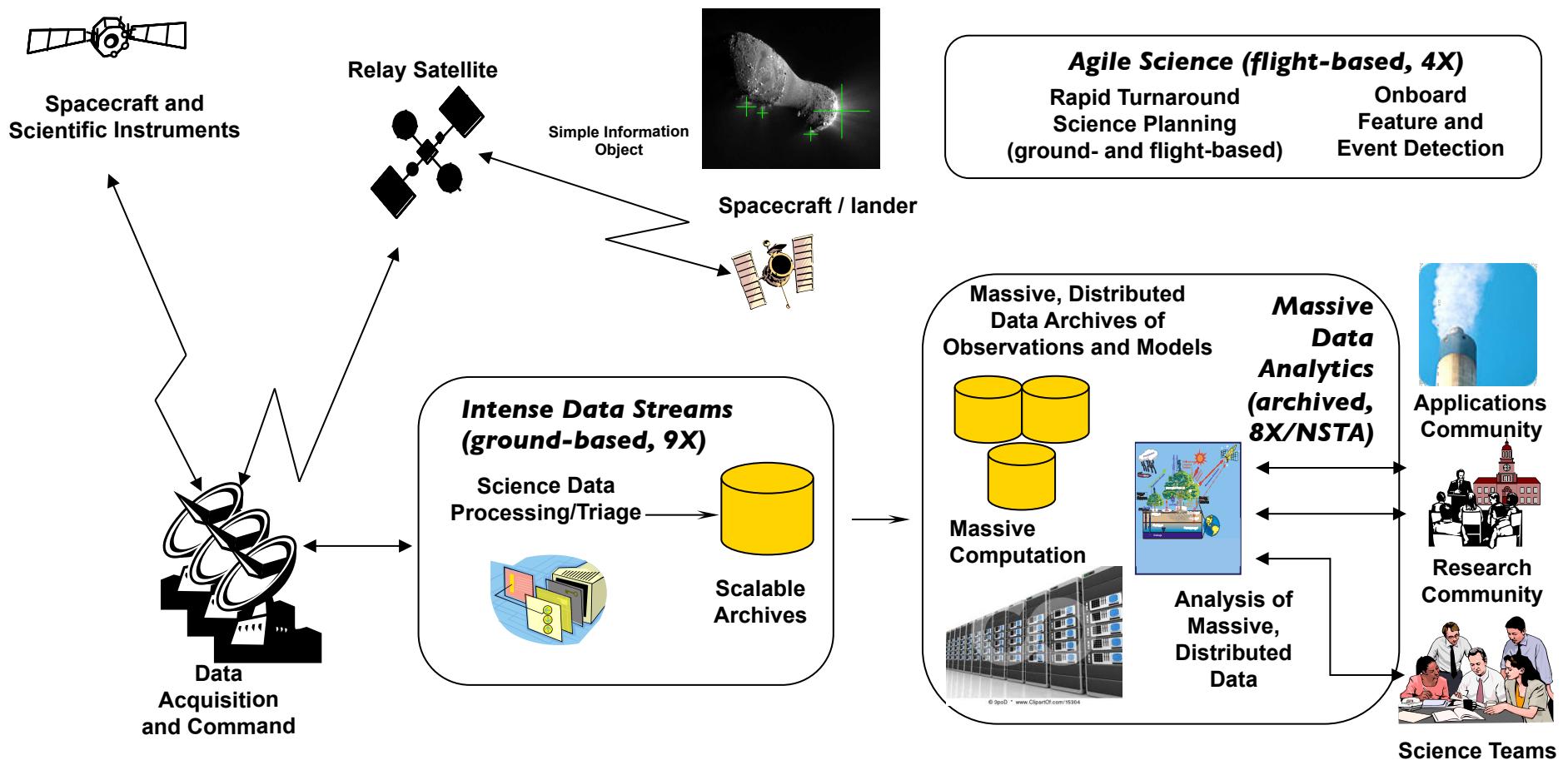
- A Major Shift from Compute-Intensive to Data-Intensive -





NASA
Jet Propulsion Laboratory
California Institute of Technology

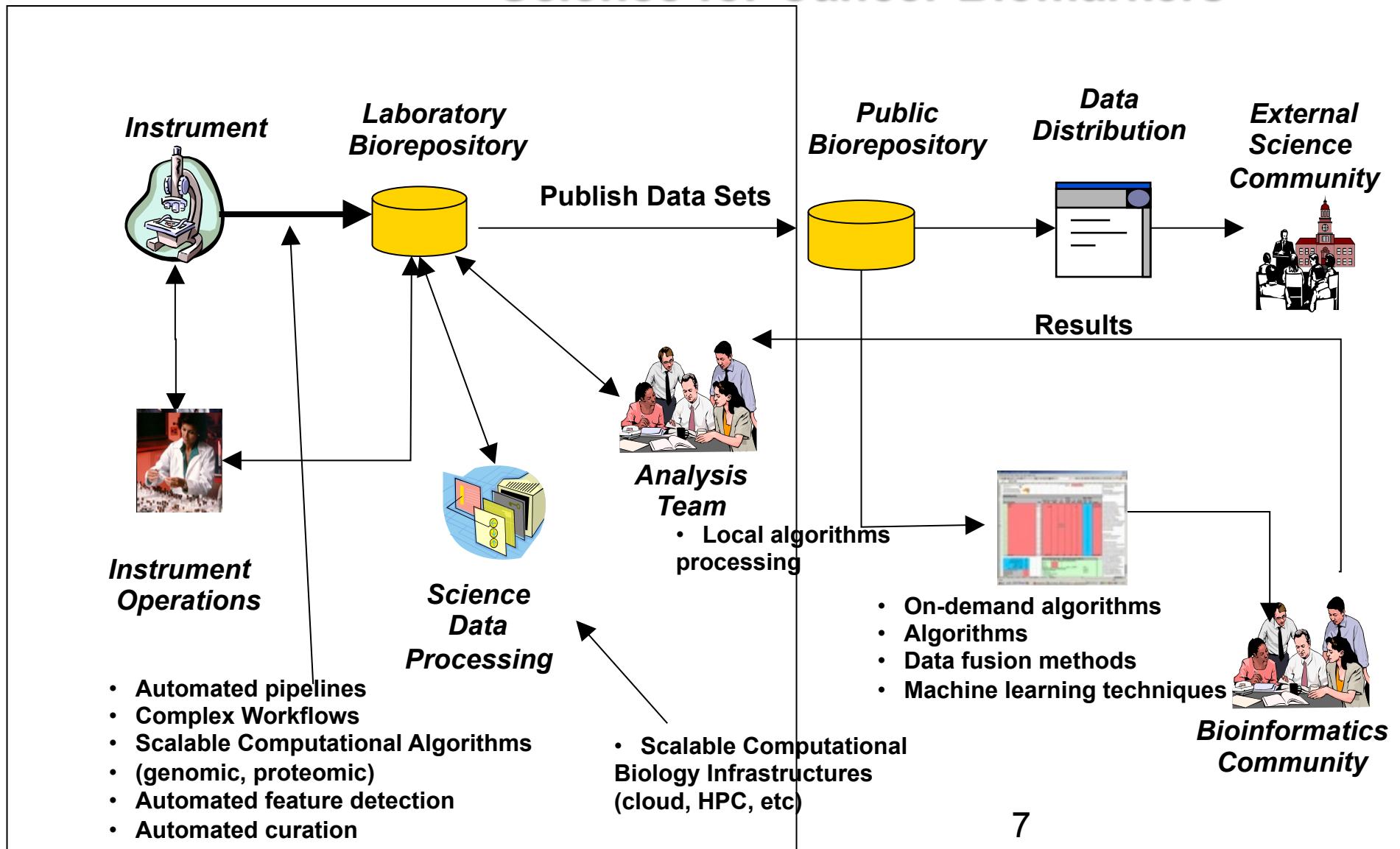
Data Lifecycle Model for Space Missions





NASA
Jet Propulsion Laboratory
California Institute of Technology

NCI Early Detection Research Network: Moving towards Data-Driven Science for Cancer Biomarkers

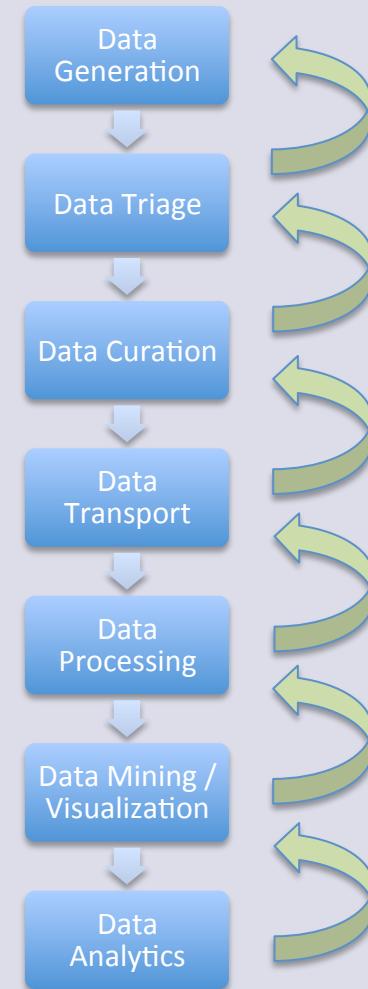




The Need for Data Science

The Data Lifecycle: Too Fast, Too Much, Not Enough

- Perform original processing at the sensor / instrument*
- Make choices at the collection point about which data to keep*
- Anticipate the need to work across multiple data sources*
- Improve resource efficiencies to enable moving the most data*
- Increase computing availability at the data to generate products*
- Develop and apply analysis techniques to enable data understanding*
- Create analytics services effective across massive, distributed data*



Data Science is the focused research to develop principled techniques and scalable architectures to address challenges across the entire Data Lifecycle



Common Challenges in Massive, Distributed, Heterogeneous Data

- 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
- Developing mechanisms for identifying and extracting interesting features / patterns
- Developing methodologies for reconciling predictive models vs. measurements
- Methods for visualizing massive data
- Providing a trusted basis for actionable results of data analytics



Jet Propulsion Laboratory
California Institute of Technology

OODT: Object-Oriented Data Technologies

An Open Source Data Science Framework

1 Ob 2 Or 3 D 4 T
Object Oriented Data Technology

Catalogs, archives, metadata, & more
Data grid framework for transparent search and discovery of disparate science resources

Apache SOFTWARE FOUNDATION

- An open source data science framework
 - Developed at NASA/JPL
 - Top Level Project at the Apache Software Foundation (2011)
 - Used across multiple NASA centers (JPL, GSFC, LaRC)
 - Used across multiple agencies (NASA, NIH, NSF, DARPA, NOAA)
 - Integrates with an information architecture – metadata and ontology (e.g., earth science, biomedicine, etc.)
 - Significantly reduces cost and increases performance of data processing and management
 - Integrates with data analytics
- Applied to multiple Earth Science missions
 - Seawinds, OCO-2, SMAP, NPP Sounder Peate, JPSS
 - CARVE, Airborne Snow Observatory
- Applied to Earth science, planetary science, astronomy, biomedicine, defense

<http://oodt.apache.org>



POC: Dan Crichton, Chris Mattmann

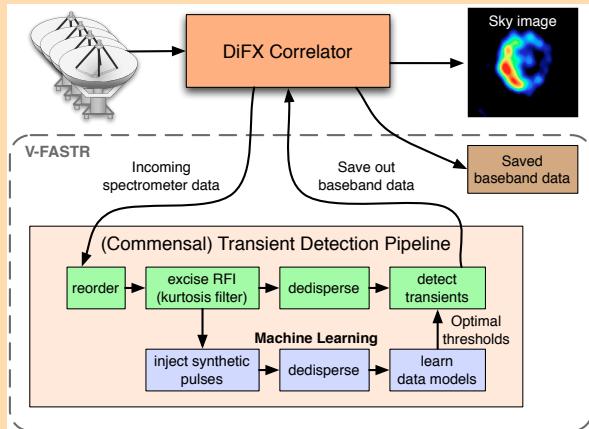


NASA
Jet Propulsion Laboratory
California Institute of Technology

Triage, Analysis, and Understanding of Massive Data

- Detection: fast identification of signals of interest (triage)

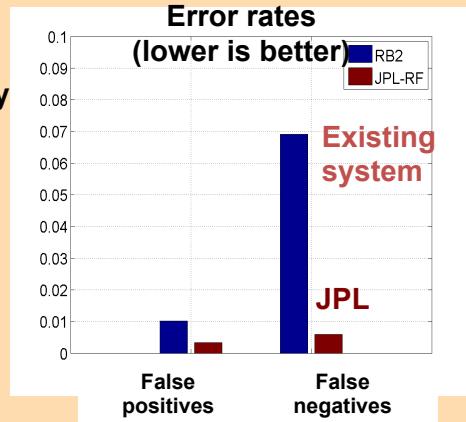
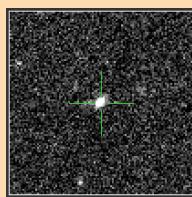
Radio astronomy:
V-FASTR
realtime system
at the VLBA



- Classification: online, real-time source type classification

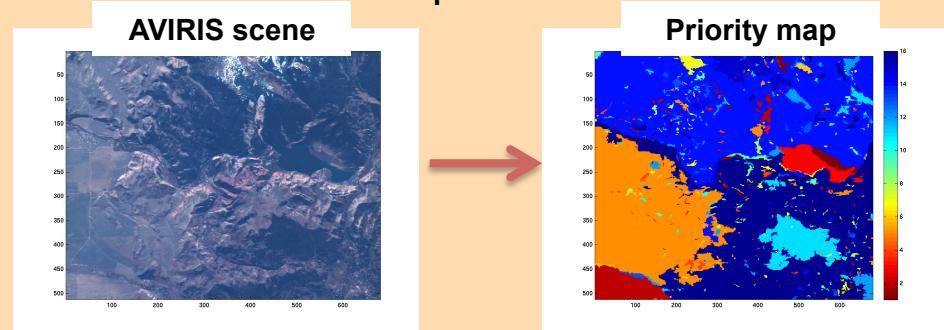
Optical astronomy:
Reducing false positives for
the Palomar Transient Factory

Real or spurious?



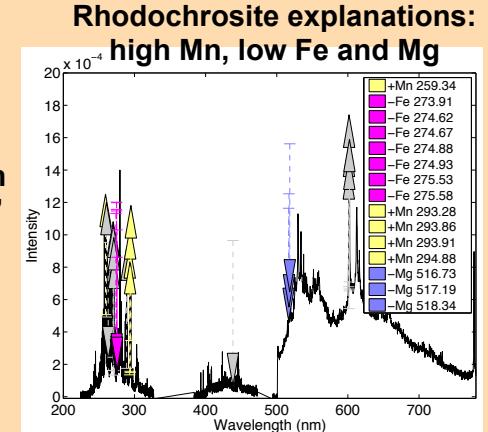
- Prioritization: use triage decisions to inform adaptive data compression

Earth science:
Onboard content-sensitive data compression



- Understanding: generate human-understandable explanations for decisions

Planetary science:
Anomaly detection in ChemCam emission spectra from Mars, with content-sensitive “explanations” indicated with arrows (higher than expected vs. lower than expected)



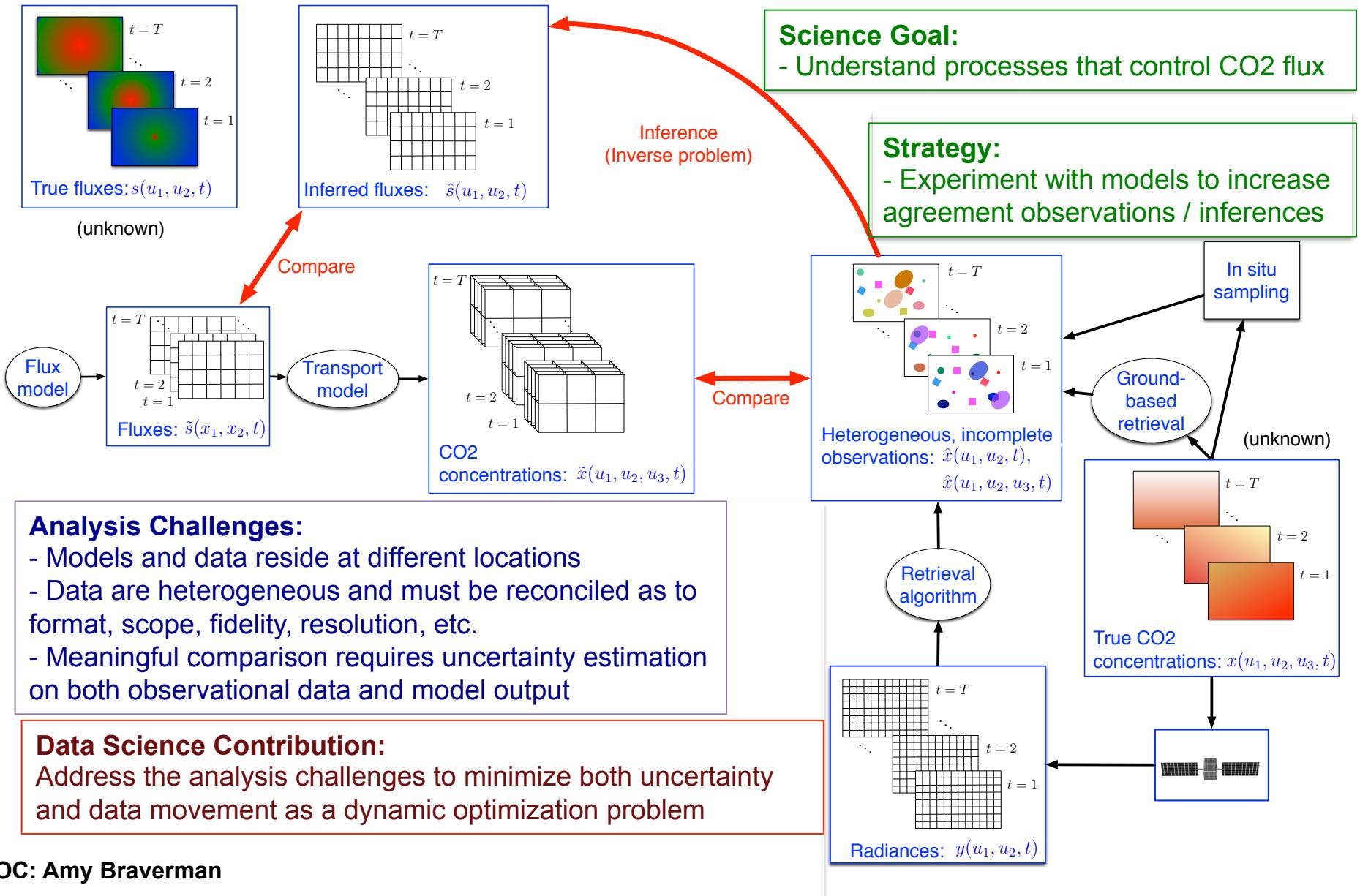
Credits: Kiri Wagstaff, Umaa Rebbapragada, David Thompson, Benyang Tang, Hua Xie



NASA
Jet Propulsion Laboratory
California Institute of Technology

Distributed Data Analytics

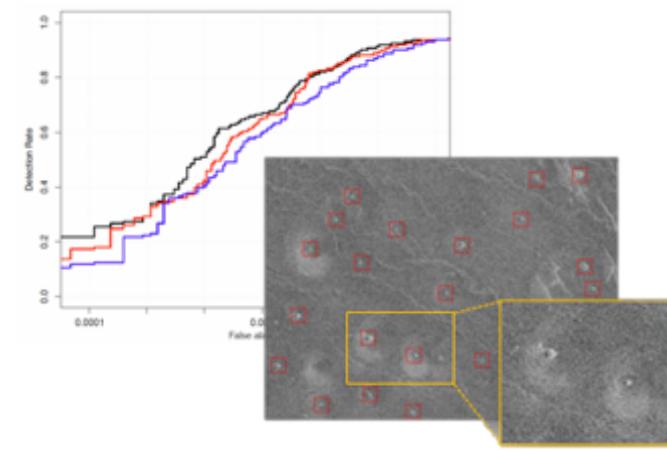
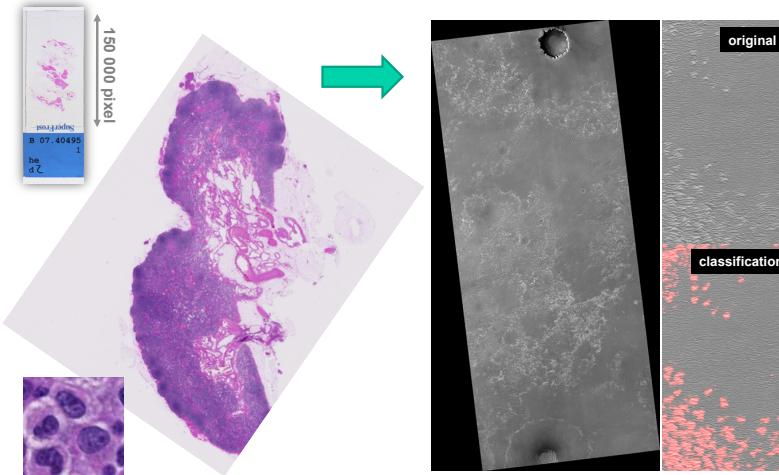
Model-to-Data Reconciliation for Carbon Cycles



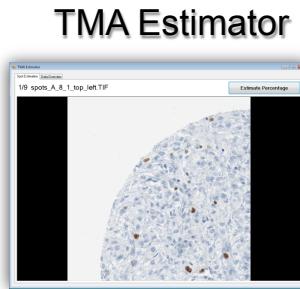


NASA
Jet Propulsion Laboratory
California Institute of Technology

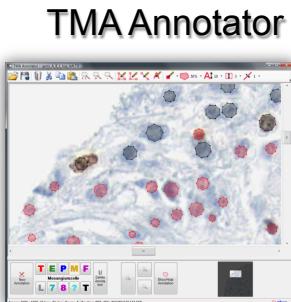
Application of Machine Learning Methodologies to Cancer Biomarker Research



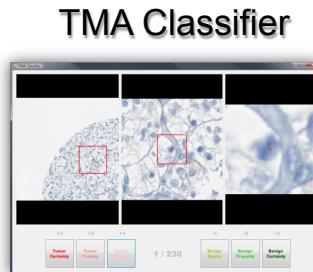
Volcanoes on Venus



Estimate the Staining
on a whole spot



Detect nuclei on a
whole spot

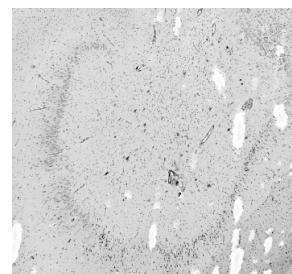


Classify single nuclei into
tumor, non-tumor and
stained, not-stained

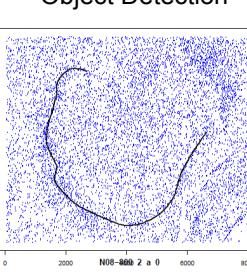
Automated Classification

POC: Thomas Fuchs

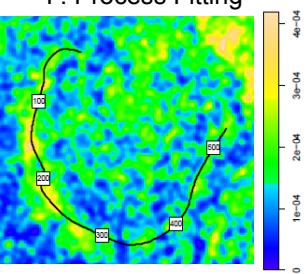
Original Image



Discriminative
Object Detection



Generative
P. Process Fitting



Feature/Object Detection



Summary

- **Data Science** is a growing area that requires new thinking across the data lifecycle, in software/system architectures, in the application of intelligent algorithms, in addressing uncertainty
- JPL and Caltech have both established Centers to respond to this important growing area
- Caltech and JPL are also partnering to maximize the value of combining our collective expertise
- JPL is working with NASA, other government agencies, academia and industry to bring together solutions