

Data Science: Unlocking Scientific Research from Space to Biology

2018 Jean Golding Institute Showcase University of Bristol

Dan Crichton, Leader, Center for Data Science and Technology July 3, 2018



Overview

>JPL's Unique Mission >The Data Challenge > Autonomy > Data Systems > Data Analytics > Methodology Transfer

JPL's Mission for NASA

Robotic Space Exploration

Earth Science • Mars • Solar System • Astrophysics • Exoplanets • Interplanetary network



Our mission has introduced unique challenges for protecting space system assets and information

JPL: From Caltech students testing rockets to the planets and universe in our lifetimes



Caltech students (1936)



Mars Exploration Rovers (2004 – present)



Missiles (1940s)



Spitzer Space Telescope (2004 – present)



Explorer 1 (1958)



Earth Science (1978 – now)

7/16/2018 JPL is a NASA FFRDC and a Division of Caltech

A History of Firsts...

Surveyor 1, First soft landing on the moon



Voyager 1, First interstellar traveler



Viking, first landing on another planet



Continuous presence on Mars since 1997



Deep space exploration enabled by NASA's Deep Space Network (DSN)







The Data Challenge

Tackling Big Data

- 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.





Data Science Projects at JPL



Planetary Science



Biology



Defense and Intelligence



Earth Science



Medicine



Radio Astronomy





on Moore's law time scales

Understanding of complex phenomena requires complex data!

From data poverty to data glut requires complex da
From data sets to data streams
From static to dynamic, evolving data
From anytime to real-time analysis and discovery
From centralized to distributed resources
From ownership of data to ownership of expertise

Transformation and Synergy: Big Data to Data Science

- All science in the 21st century is becoming cyber-science (aka e-Science) and with this change comes the need for a new scientific methodology
- The challenges we are tackling:
 - Management of large, complex, distributed data sets

 - These challenges are universal
- A great synergy of the computationally enabled science, and the science-driven informatics



Hypothesis-driven science

Data-driven science

Understanding



The two approaches are complementary

A Modern Scientific Discovery Process

their pipelines...)

Databases

Data grids

Data commons

Data Gathering (instruments, sensor networks,

Data Farming:

Storage/Archiving Indexing, Searchability Data Fusion, Interoperability

Data Mining

Pattern or correlation search Clustering analysis, classification Outlier / anomaly searches Hyperdimensional visualization





+feedback

U.S. National Research Council Report:

Frontiers in the Analysis of Massive Data

- Chartered 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 Big Data Lifecycle Model

Emerging Solutions

- Onboard Data
 Analytics
- Onboard Data
 Prioritization
- Flight Computing



Observational Platforms and Flight Computing

Emerging Solutions

- Intelligent Ground
 Stations
- Agile MOS-GDS



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

Ground-based Mission Systems



SMAP (Today): 485 GB/day NI-SAR (2020): 86 TB/day

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

Massive Data Archives and Big Data Analytics



Emerging Solutions

- Data Discovery from Archives
- Distributed Data Analytics
- Advanced Data Science Methods
- Scalable Computation and Storage

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

Autonomy

0



Increasing Computing Capability Onboard

Heading Toward Multicore in Space





Voyager computer

 8,000 instructions/sec and kilobytes of memory



iPhone

- 14 billion instructions/sec and gigabytes of memory

Curiosity (Mars Science Laboratory) Processor: 200 MOPS BAE RAD750



HPSC (NASA STMD / USAF) Processor: 15 GOPS, extensible

Surface Mobility Mars Rover Navigation

Flight Deployed

- **1996 Mars Pathfinder:** obstacle avoidance with structured light
- 2003 Mars Exploration Rover: obstacle avoidance with stereo vision; pose estimation and slip detection with visual odometry; goal tracking
- 2011 Mars Science Laboratory: enhanced obstacle avoidance, visual odometry and goal tracking

Research and Development

Enhanced hazard detection, traversability analysis and motion planning for Mars 2020 and beyond





Terrain classifier



Terrain Classification



Fido



Onboard Analysis *Dust Devils on Mars*

Dust devils are scientific phenomena of a transient nature that occur on Mars

- They occur year-round, with seasonally variable frequency
- They are challenging to reliably capture in images due to their dynamic nature
- Scientists accepted for decades that such phenomena could not be studied in real-time



Spirit Sol 543 (July 13, 2005)

New onboard Mars rover capability (as of 2006)

 Collect images more frequently, analyze onboard to detect events, and only downlink images containing events of interest

Benefit

- < 100% accuracy can dramatically increase science event data returned to Earth
- First notification includes a complete data product



Credit: T. Estlin. B. Bornestein, A. Castano, J. Biesiadecki, L. Neakrase, P. Whelley, R. Greeley, M. Lemmon, R. Castano, S. Chien and MER project team

Data Systems

Data-Driven Capabilities Across the Ground Environment

Intelligent Ground Stations



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

Agile MOS-GDS



Data Analytics and Decision Support



Emerging Solutions

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

Data-Driven Discovery from Archives

Data-Driven Approaches for Deep Space Communication: Detecting Anomalies

Current Inputs: DSN operationally relevant data

Real Time Monitor Data Real Time Operator Logs Track Predicts

Desired Output: Better Fault Detection and Diagnosis

Real-Time, Historically Informed Alerts

Real-Time Insight into Data Points' Criticality and Relationships DSN Software Quality Assessment (SQA) Data Archive

- Relational database
- 10 years of data
- 1.3+ billion records

Credit: Rishi Verma, JPL

Scaling Processing of NISAR in AWS Cloud



NASA Archives: Access to Data*





Planetary Science

Highly distributed/federated Collaborative Information-centric Discipline-specific Growing/evolving Heterogeneous International Standards & Interoperability



Multiple Data Centers



Astronomy



Earth Observation

Multiple Data Centers Heliophysics

* Petascale environment that is moving to an exascale requirement

Growth of Planetary Data Archived from U.S. Solar System Research

U.S. Planetary Data Archives (TBs)



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
 - 4000 different types of data
 - Over 1.4 PB of data
 - International interoperability
 - Distributed federation of US nodes and international archives
- Realized through an international data science platform!





28

Enabling a Metadata Model-Driven Data System

Information System Architecture



Crichton, D. Hughes, J.S.; Hardman, S.; Law, E.; Beebe, R.; Morgan, T.; Grayzeck, E. 7/16/2018 A Scalable Planetary Science Information Architecture for Big Science Data.

IEEE 10th International Conference on e-Science, October 2014.

International Planetary Data Alliance: Collaboration and Access to Data Archives



LADEEMAVENOsiris-RExExoMarsBepiColomboMars 2020Psyche(NASA)(NASA)(NASA)(ESA/Russia)(ESA/JAXA)(NASA)(NASA)



InSight (NASA) JUICE (ESA)

Europa (NASA) Hyabusa-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

3000

3008			75%
50%	60%	10%	65%

4k

3k

2k

1k

Data Analytics

Increasing Need for Data-Driven Analysis



Setup Automated Data Analytic Pipelines (from archives to analytics)

- Support construction of online data analytic capabilities from archives
 - Ad hoc workflows and data pipelines
 - Rapid integration of different methodologies (e.g., feature detection, classification, etc)
 - Establishment of highly scalable databases for analytics
 - Derive additional metadata to support analysis



Visualization and Analytics



Examples: Hydrology and sea level rise



Integration of *multiple* Earth observing remote sensing instruments; comparison against models



Real-time feature extraction and classification in astronomy

Exploring the Moon through Data Analytics



E. Law, S. Malhotra, G. Chang

Mars Trek: The Google Earth of Mars

Curiosity Landing Site





QG

Curiosity landed in Gale Crater on Mars on August 6th, 2012. With a diameter of 154 km and a central peak 5.5 km tall, Gale Crater was chosen as the landing site for the Mars

Science Laboratory Curiosity rover. The choice was based on evidence from orbiting spacecraft that indicate that the crater may have once contrained large amounts of liquid water. The central peak, Mount Sharp, exhibits layered rock deposits rich in sedimentary minerels including clays, sulfates, and saits that require water to form.

egion Information Download for 3D Printer





central peak 5.5 km tall, Gale Crater was chosen as the landing site for the Mars Science Laboratory Curiosity rover. The choice was based on evidence from orbiting spacecraft that indicate that the crater may have once contained large amounts of liquid water. The central peak, Mount Sharp, exhibits layered rock deposits rich in

With a diameter of 154 km and a

sedimentary minerals including clays, sulfates, and saits that require water to form.





Credit: Emily Law, Shan Malhotra

Content-Based Image Classification

- About ~1.3M images from MRO Mission HiRISE instrument
- Previously no way to easily find images with certain landmarks (e.g., craters)
- New Approach:
 - 1) Determine high salience (i.e., distinctive) regions by computing statistical differences between pixel and surrounding context
 - 2) Classify landmarks using *machine learning model* and user *training data*



MOC, June 2000

Salience Map

Crater

Barchan dune



Impact ejecta



Examples of classified landmarks

Ground-based Data Analysis Understanding Underground Water in the California Central Valley



Credit: Tom Farr

Build Analytic Data Infrastructures



Open source and scalable to cloud; 180 billion data points accessible < 1 second

WaterTrek: Interactive Data Analytics for Hydrology



Virtual and Augmented Reality: Virtual Mars

41

50% 1k BOB 25% Methodology Transfer 3000

4k

3k

2k

100%

95

2

65%

80%

60%

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

- Development of an advanced Knowledge System to capture, 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





Data Science Architectures in Cancer Research



Common Data Elements and Information Models

- Provide standard data semantics to capture and share biomarker data
 - EDRNCDEs

- MCL CDEs

 Used as models to drive the knowledge system



Biomarker Ontology Information Model

Data Pipelines and Cancer Research



Scalable computing, common data elements, computational methods





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



Methodology Transfer

From Astrophysics...



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, C**4**8ch.

Feature classification in images



Crowd Sourcing Image Analysis

Lung cancer

ABOUT CLASSIFY TALK COLLECT



00 🗘 🔳



Integrated Knowledge Data Environment



Biomarker Database



Conclusion

- Data Science is changing the paradigm for how we can do science
 - Science is increasingly data-rich and computationally enabled
 - Many new informatics fields are emerging
 - Good successes but more work to do!
- Much of the focus has been on building infrastructures for capturing data
 - Cloud infrastructures and other capabilities have paved the way
- However, great opportunity to consider how analytics is integrated end to end
 - From point of collection to analysis
- New methods are needed to bring together and drive interactive analytics

These are Extraordinary Times





Changing the paradigm for how scientific data is captured, organized, and analyzed

