

National Aeronautics and Space Administration



Data Science @ OCIO/NASA

Office of the Chief Information Officer

Brian Thomas
June 29, 2016

www.nasa.gov





Outline

- What is Data Science?
- Data Science activities at OCIO/NASA
- Big Data Issues for Data Science
- Summary



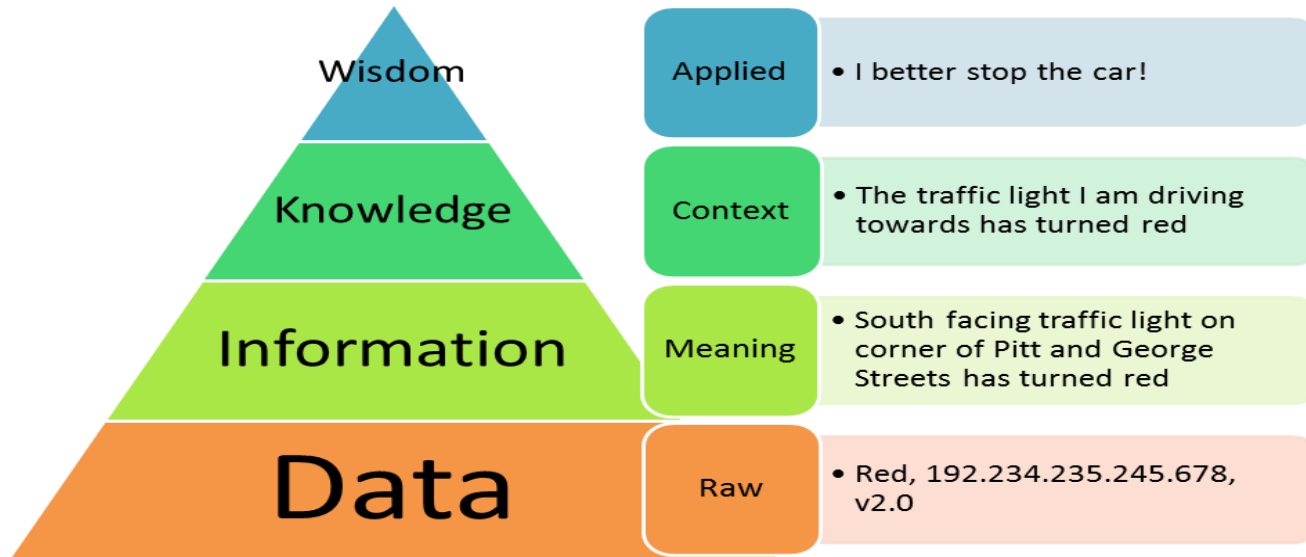
Data Science? What's that?

Wikipedia :

*“Data Science is an interdisciplinary field about processes and systems to **extract knowledge** or insights **from data** in various forms, either structured or unstructured.”*



Data -> Knowledge (and beyond!)





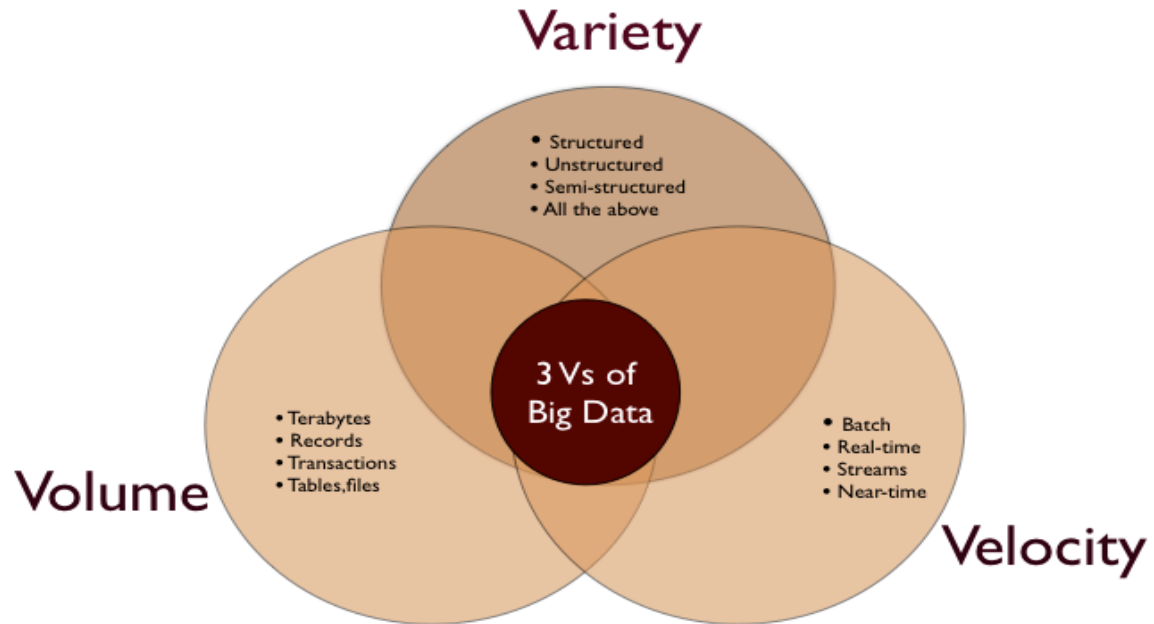
“Problem Datasets”?



© marketoonist.com

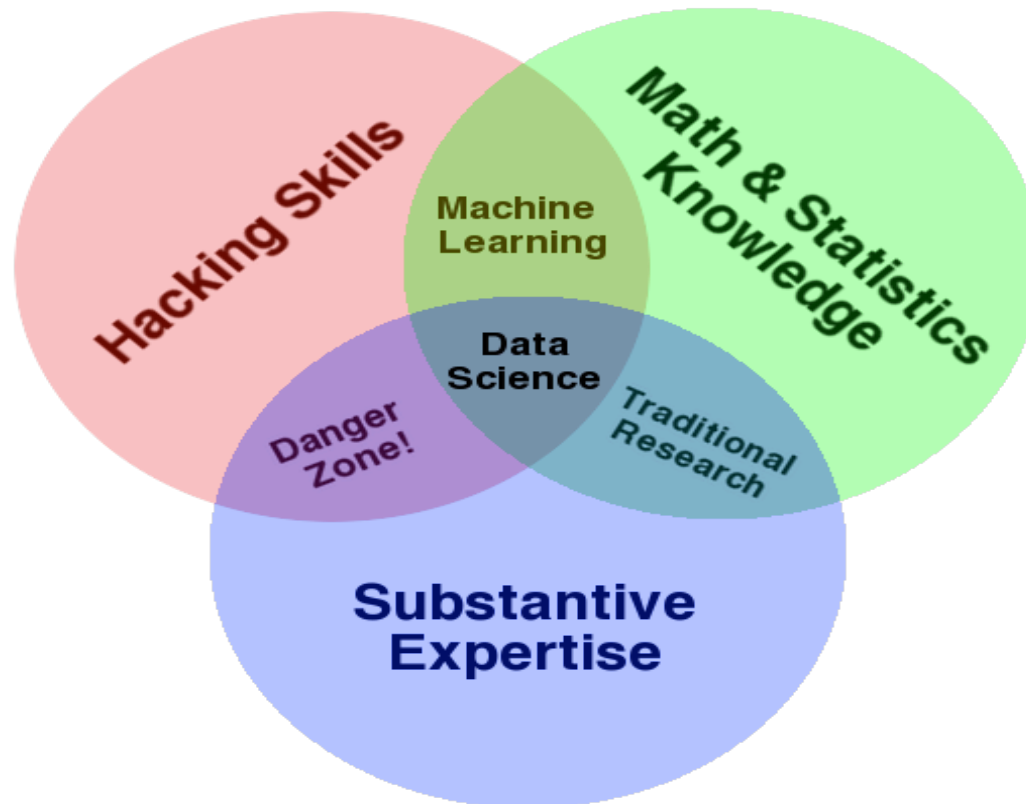


“Problem Datasets”? - Characteristics



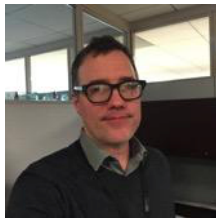


Data Scientists: Distinguishing Technical Skills





Data Science @ NASA/CIO : the “Data Team”



Brian Thomas
Agency Data Scientist



Nick Skytland
NASA Data Evangelist



Sandeep Shetye
Chief Information Architect

*Plus a team of talented Data Scientists, Software Engineers and Interns
(located at various NASA campuses)*



ARCH

Information Architecture

The development and use of models, policies, rules or standards that govern how information is collected, stored, arranged, integrated, and put to use in data systems and in organizations.

Standards and interoperability
IA Reference Architecture
Data Mining
Data Contract Language
Mission Support Tech Consulting

Big Data

NIAM.NASA.Gov

Tech & Innovation Division--What We Do:



DATA

Data Science

The collection, management and analysis of data in order to produce information and drive decision making.

Data Strategy
Data Governance
Data Lifecycle Management
Data Analytics Lab
Data Fellows Program
Data Stewards



OPEN

Open Innovation

The development of new innovation frameworks and techniques, and the development and delivery of machine-readable instructions to access, arrange, and apply data.

Agency Open Data Mgmt
Digital Strategy Reporting
Space Apps Challenge
Innovation Incubator
Data Innovation Pipeline
Women in Data
Open.NASA, Github/NASA
Data.NASA, Code.NASA
API.NASA



DIGITAL

Digital Integration

The delivery of enhanced digital capability to support data interoperability and accessibility to enable data insights and discovery.

Data inventory/Registry
Tagging/discoverability
Data usability/APIs
Computing/Coding
Mission-focused tech applications



TECH

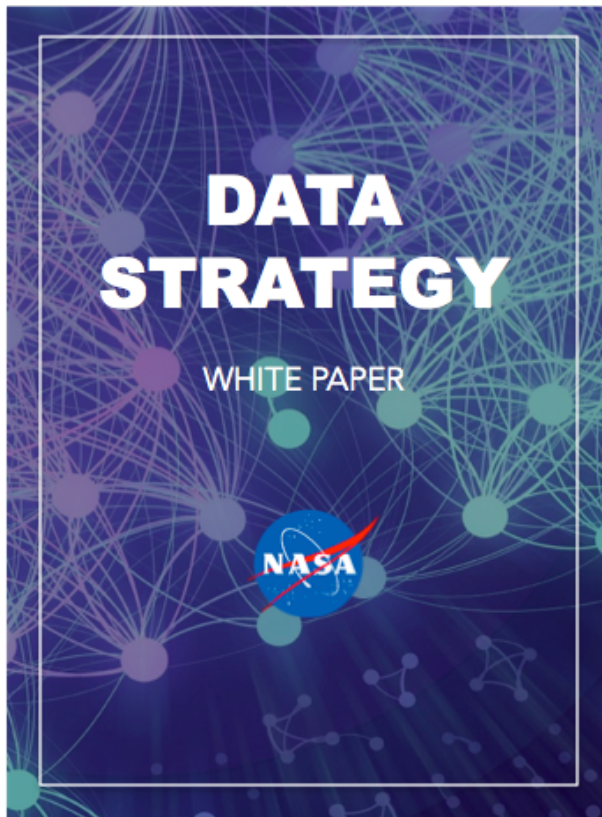
Tech Infusion

Research, prototype and assess the creation, modifications and usage of processes or tools to solve a problem, generate or perform a specific function to meet stakeholder needs

Data Centric Architecture
Internet of Things
Software as a Service
Virtual Desktop Infra
Collaboration (Google Apps/Secure Dropbox, NASATube)



Our Plan



- Data Management
- Unified Data Lifecycle
- Data Governance
- Data Analytics Lab
- Data Fellows Program
- Data Stewards



Our Collaborators

Current and Past (OCIO Data Team) Project Partners include:

NASA Missions/Centers

- Ames Research Center
- Marshall Space Flight Center
- Johnson Space Center

Agency

- Office of Chief Information Officer
- Office of General Council
- Office of Human Capital Management
- Office of the Chief Engineer
- Office of the Chief Scientist

Academia

- University of Maryland UC



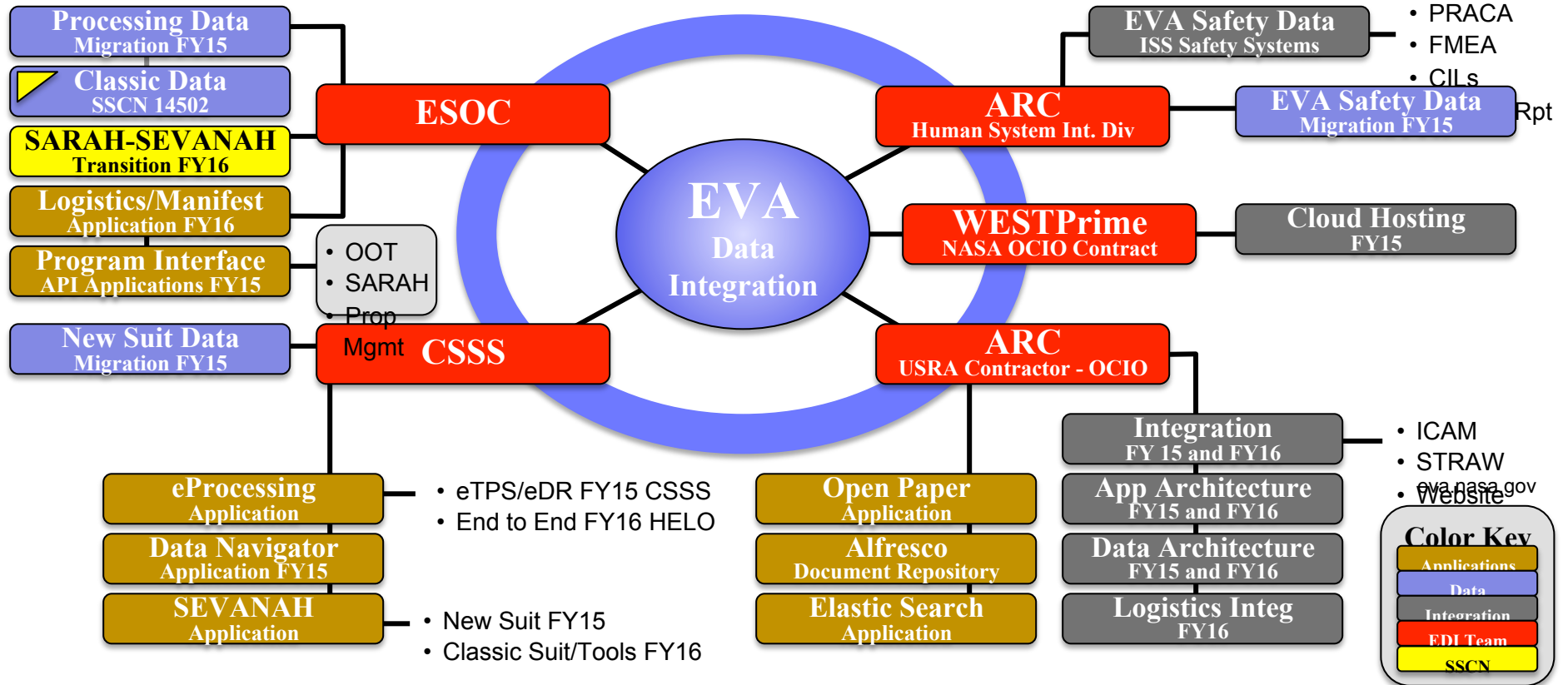
OCIO Data Science Summary Projects Roster

Prototypes/Research into methods/tools/platforms in:

- **Data Analytics/Deep Learning**
 - » ExMC-telemedicine, Mapping Data Universe, Security
- **Search Analytics**
 - » Helping optimize strategies for agency document search
- **Document Classification/Content Tagging**
 - » Security, Open Data, Aerospace, Publications, Records Management
- **Quantum Computing**
 - » API for analytics
- **Financial Analytics**
 - » OCIO, CFO
- **Data Integration/Advanced Data Management**
 - » ISS, MBSE, Orion, EVA, Publications, Records
 - » Microservices (OCIO, Open Data)

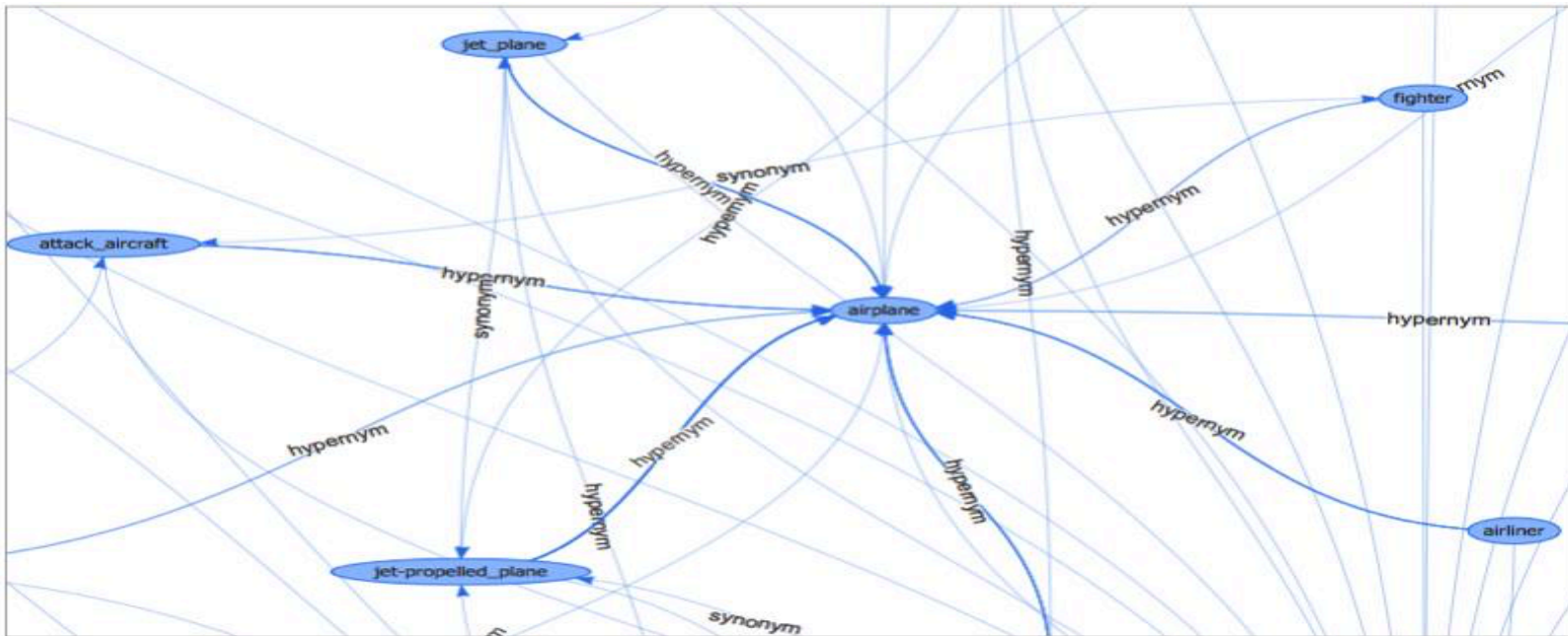


Sampling of Current Projects



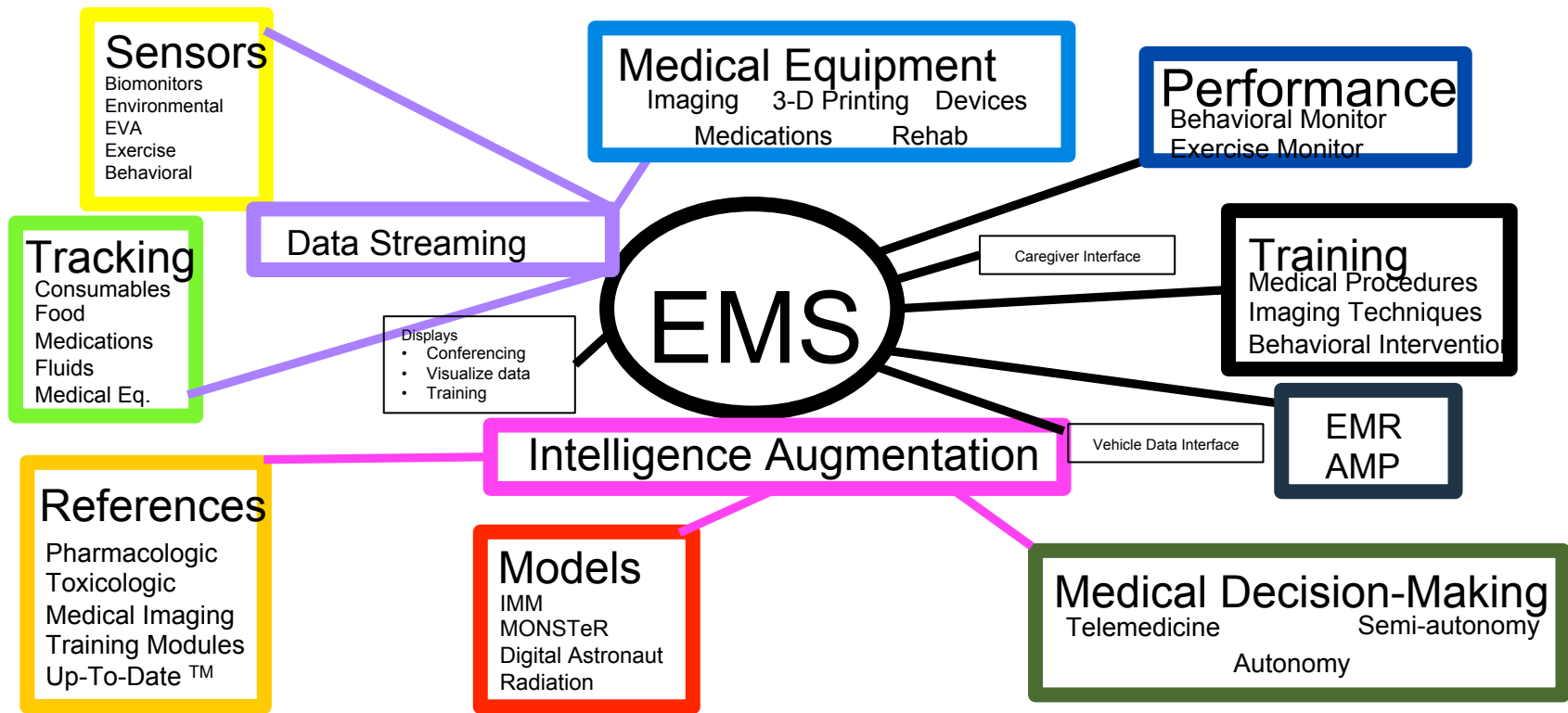


Sampling of Current Projects



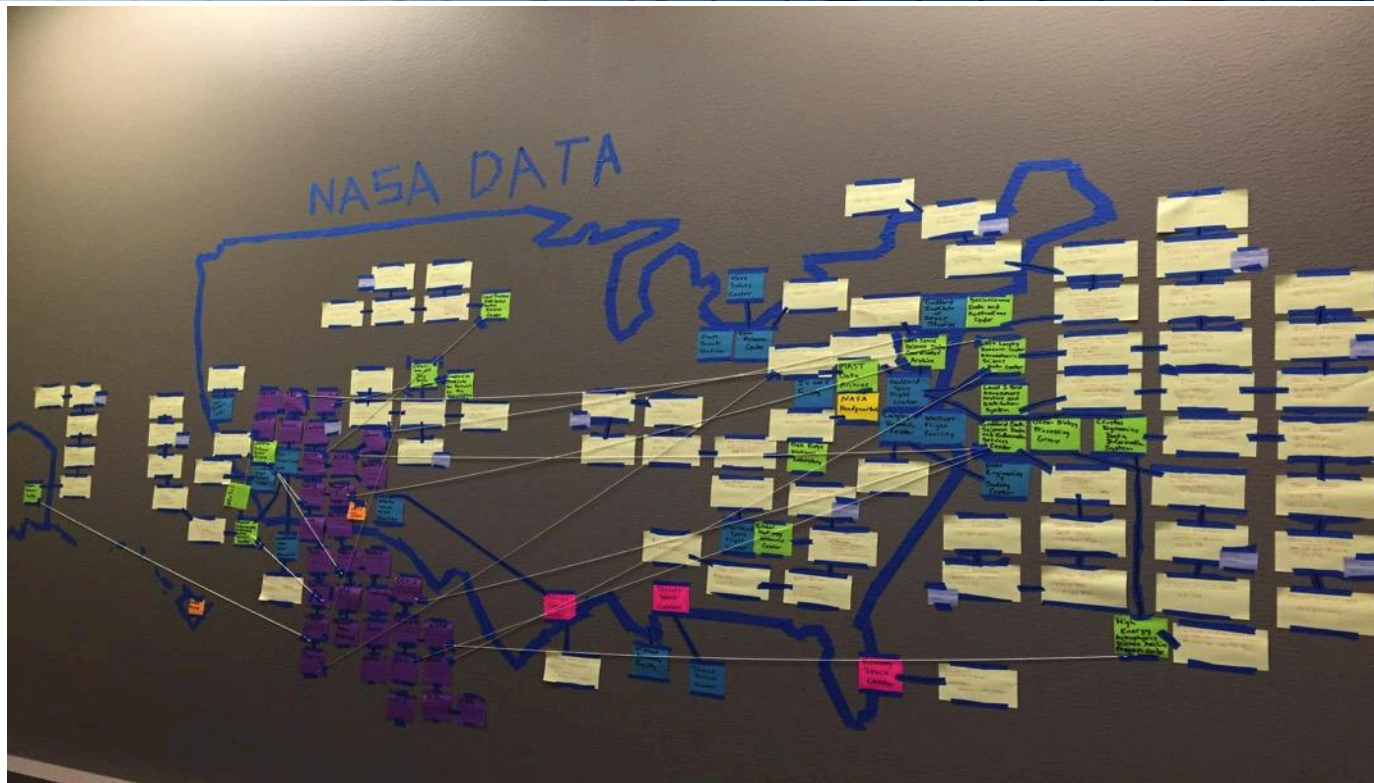


Sampling of Current Projects





Sampling of Current Projects





Sampling of Current Projects

Google Analytics

Add smoother?
Linear or smoothed
loss

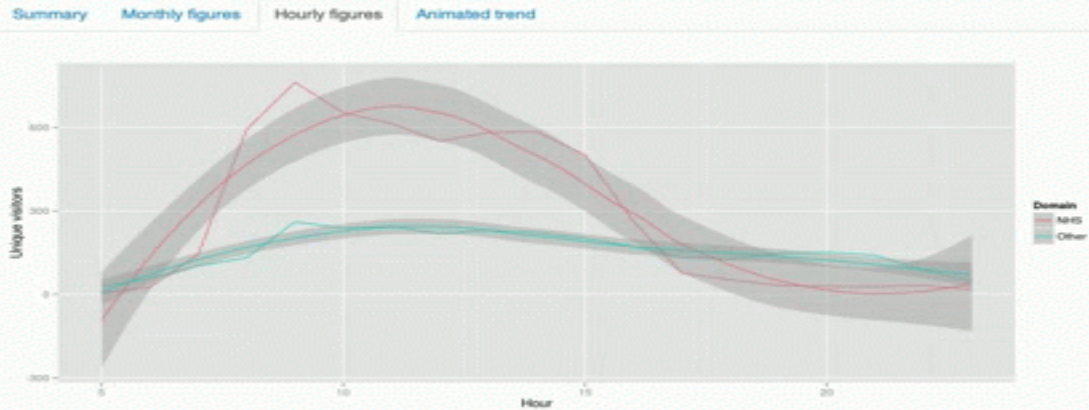
Hours of interest- minimum
0 25

Hours of interest- maximum
0 25

Show NHS and other domain?
 NHS users
 Other

Output required
 Visitors
 Bounce rate
 Time on site

Choose subdomain
nhs.uk





Sampling of Current Projects

Building Microservices

- **Providing easier data access** (files, databases) (NASA aggregate dictionary: <http://nasa-dictionary.herokuapp.com>)
- **Creating higher value data** (mash-ups, APOD : <http://api.nasa.gov/planetary/apod>)
- **Providing platform-independent capability** (<http://nasa-tagging.herokuapp.com>)



Observed Trends

- **NASA has all types of “V’s” problems**
 - » “Variety” is the hardest to solve from a technical standpoint
- **NASA implementation of data analytics is uneven**
 - » Science doing the best with their data
 - » Business units are doing the poorest
- **Need for much greater Data Sharing**
 - » Microservices (and other solutions) promise an easy way to overcome silos but will create new challenges for NASA's information architecture.
 - » Variety problem will become acute
 - » Adequate Governance/Access/Security are critical
 - » Provenance is also important



Big Data Issues for Data Science

- **Network Throughput**
 - » As we enable greater exchange of data between former silos
- **“Code Shipping” ability**
 - » we will never be able to move all of the data
- **Cost effective cloud computing**
 - » Microservices, by their nature, are likely to arise from small teams
 - » PAAS which will support ‘off the shelf/ready to go’ data science
- **Infrastructure support for Governance/Quality**
 - » Who can access, original or derived data?
 - » Provenance of creation? transforms?
- **Discovery of Data at the Agency**
 - » Data in all areas still rapidly growing
 - » Combining datasets is powerful -- more important than ever to facilitate the finding of data across NASA
 - » Enabling better collaboration and return on data



Summary

- **NASA has been doing “Data Science” for a long time**
 - » Many different missions/centers are engaged in Data Science activities
- **Uneven distribution of experience/uptake**
 - » Science areas @NASA are doing the best job with Big Data, Business units least well
- **Hard issues to solve with the “variety” problem**
- **Data Silos and Data Provenance are pressing issues**
- **Big Data will often require shipping analysis to the data**
- **Need to support 21st century environment for NASA**

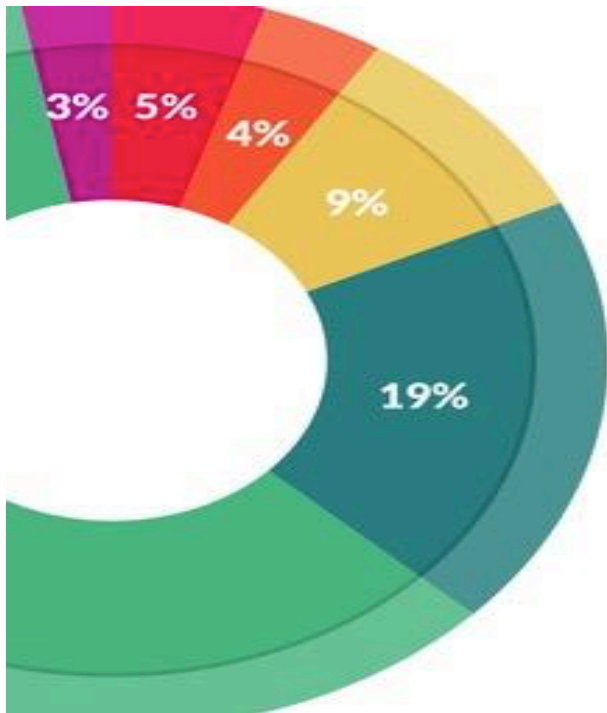


Spare Slides

Spare Slides



Why APIs?



What data scientists spend the most

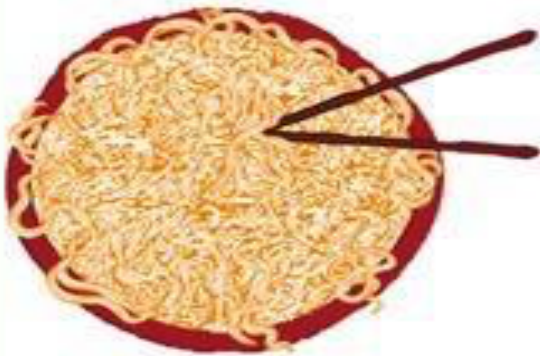
- *Building training sets: 3%*
- *Cleaning and organizing data: 60%*
- *Collecting data sets; 19%*
- *Mining data for patterns: 9%*
- *Refining algorithms: 4%*
- *Other: 5%*



Sampling of Current Projects

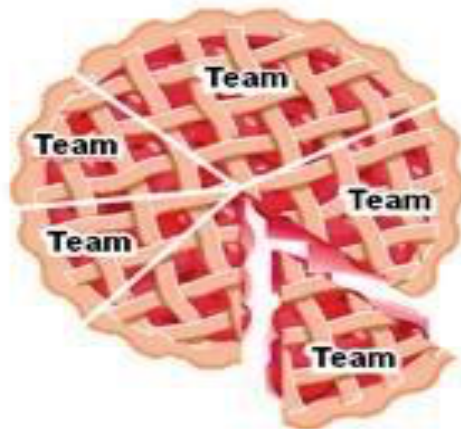
1990s and earlier

Pre-SOA (monolithic)
Tight coupling



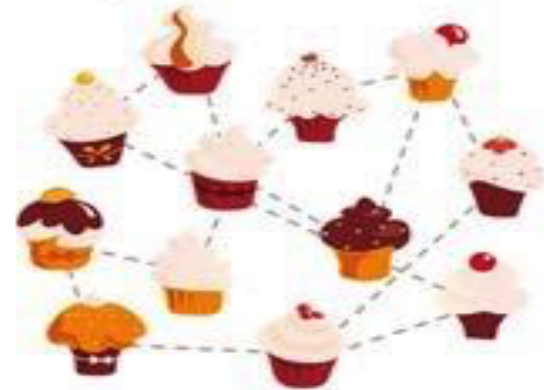
2000s

Traditional SOA
Looser coupling



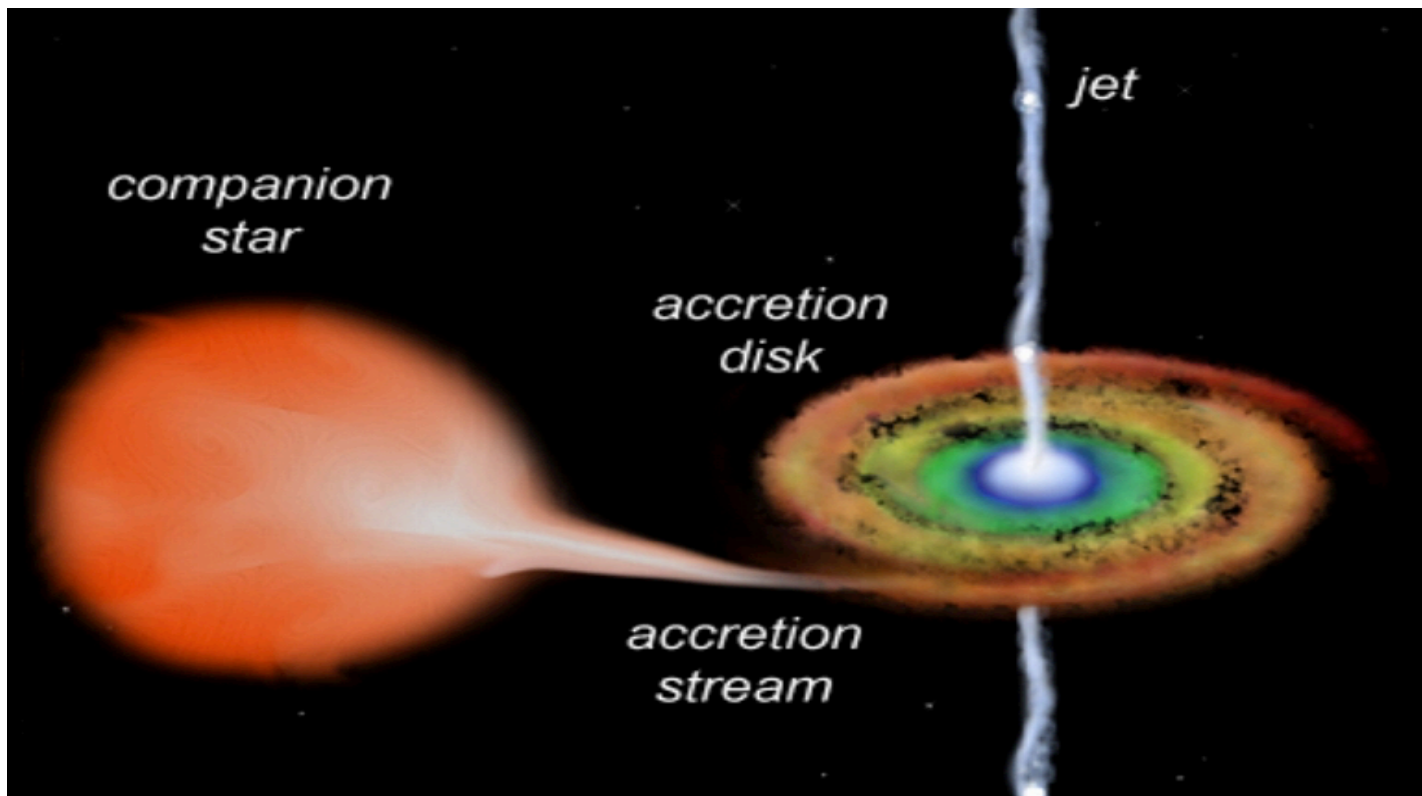
2010s

Microservices
Decoupled





Field of Study: High Energy Astronomy





Finding, Munging & Analyzing the Data

Space-based Data



Ground-based Data



Useable Data

Alt (km)	00a-00a	00a-00a	EXTENSION (1/Min)			18-04	YEAR 1980-1990		
			00a-20a	20a-10a	10a-00a		00a-00a	00a-00a	00a-00a
6.5	3.2	2.1	2.1	2.0	2.4	4.2	14.3	13.0	
7.5	4.0	2.5	1.8	2.4	4.0	6.5	11.5	10.5	
8.5	4.9	2.6	1.8	2.4	2.0	9.3	10.0	8.6	
9.5	4.3	2.0	1.6	1.9	2.0	4.8	7.2	6.7	
10.5	3.3	2.0	1.8	1.5	2.0	4.8	7.8	7.4	
11.5	2.9	2.0	1.9	1.4	2.3	3.9	4.3	4.3	
12.5	3.0	2.0	1.8	1.5	1.8	2.8	4.9	5.3	
13.5	3.2	2.6	1.7	1.5	1.8	2.4	4.0	4.6	
14.5	3.7	2.7	1.6	1.6	1.8	2.3	3.6	4.0	
15.5	3.0	3.0	1.6	1.0	1.8	2.2	3.4	3.6	
16.5	3.0	3.4	1.8	1.4	1.6	2.2	3.2	3.0	
17.5	3.5	3.5	2.2	1.6	1.7	2.4	3.0	2.4	
18.5	3.2	3.4	2.0	2.0	2.1	2.7	2.5	2.0	
19.5	2.8	2.9	2.1	2.0	2.0	2.7	1.9	1.2	
20.5	1.8	2.5	2.0	2.4	2.2	2.2	1.3	0.8	
21.5	1.1	1.7	2.8	2.3	2.1	1.6	0.8	0.5	
22.5	0.4	1.2	2.0	2.1	2.2	0.9	0.5	0.4	
23.5	0.4	0.8	1.6	2.0	2.2	0.4	0.4	0.3	
24.5	0.2	0.5	1.1	2.1	1.7	0.4	0.2	0.2	
25.5	0.1	0.3	0.8	1.8	1.5	0.3	0.2	0.2	

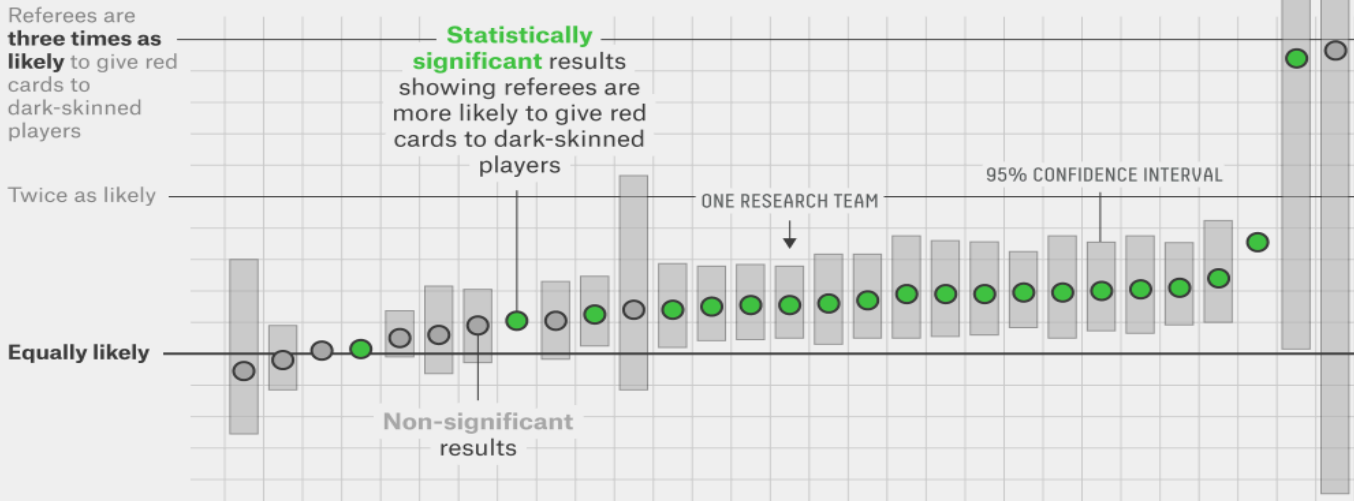




Example of the Danger Zone “Work”

Same Data, Different Conclusions

Twenty-nine research teams were given the same set of soccer data and asked to determine if referees are more likely to give red cards to dark-skinned players. Each team used a different statistical method, and each found a different relationship between skin color and red cards.



FIVETHIRTYEIGHT

SOURCE: BRIAN NOSEK ET AL.

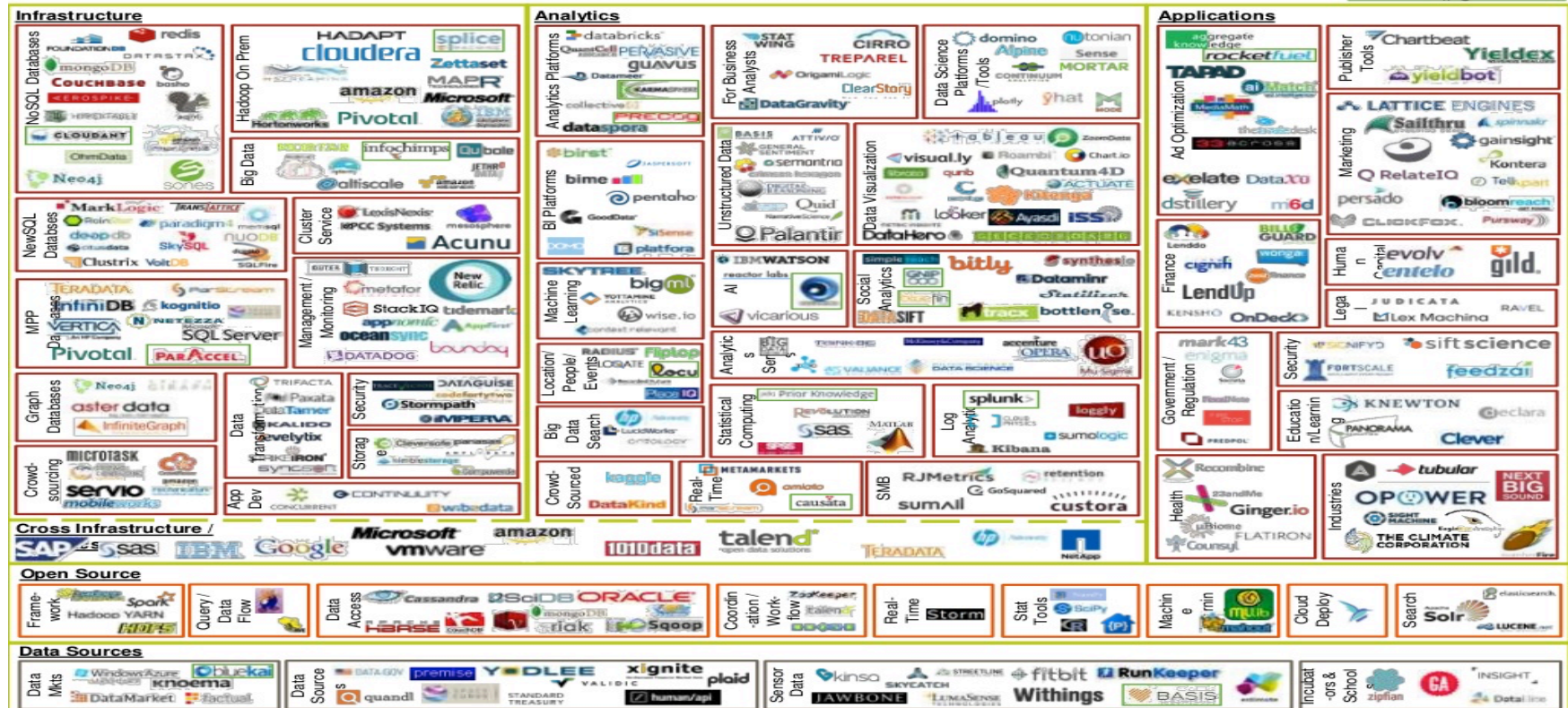
<http://fivethirtyeight.com/features/science-isnt-broken>



Tools & Technologies Big Data/Data Science Landscape

BIG DATA LANDSCAPE, VERSION 3.0

Exited: Acquisition or IPO



© Matt Turck (@mattturck), Sutian Dong (@sutiandong) & FirstMark Capital (@firstmarkcap)



Optimizing Processes

Legacy Subsystem

NBL Traing Tool Tracking (NB)	Shipping System (SH)	Inventory Apps. (IC)	Physical Inventory (PH)
GSE Calib & Per. Maint (CL)	As-Built Bill of Materials (BM)	Receiving (RE)	Property Mgmt (PM)
TPS Routing and WO (WO)	Configuration Mgmt (PM)	Bill of Materials (BM)	Document Maint. (DM)
			Event Bill of Materials (ER)
	Anthropometric (AN)	Fabrication Apps. (FB)	Engineering Apps (EN)
Discrepancy Reporting (DR)	Procedure Deviation (PD)	Test Readiness Rev (TR)	H/W Limited Life (LL)
			Quality Records (QR)
	Work Breakdn Struc (WB)	Purchase Orders (PO)	Purchase Requisition (PR)
		Application Utilities (AU)	Menu System Maint (ME)
		Data Cleanup Utilities (DC)	SQLs to Menu Dr. Rpts (SQ)

28
FUNCTION

New Common Functions

- Physical Asset Service (ship, tool, inventory)
- Logical Asset Service (Items, BOM, Event planning)
- Engineering and Fabrication Services
- Quality Management Service (procedure deviation, ...)
- Purchasing Service
- Administrative Services (groovy obj manipulation, report)
- Authorization Administration Service

EVA Functional Groups

- Logistics
- Property Management
- Configuration Management
- Production Control
- Design Drafting
- Design Engineering
- Quality Assurance
- Purchasing
- CSSS PO&C
- Admin PP&B

10
FUNCTION