

Introduction to Data Science and Its Applications

Presentation for the SDSC Summer Institute
August 3, 2021 (Delivered remotely due to COVID-19.)

İlkay ALTINTAŞ, Ph.D.

Chief Data Science Officer & Division Director of Cyberinfrastructure Research, Education and Development, **San Diego Supercomputer Center**

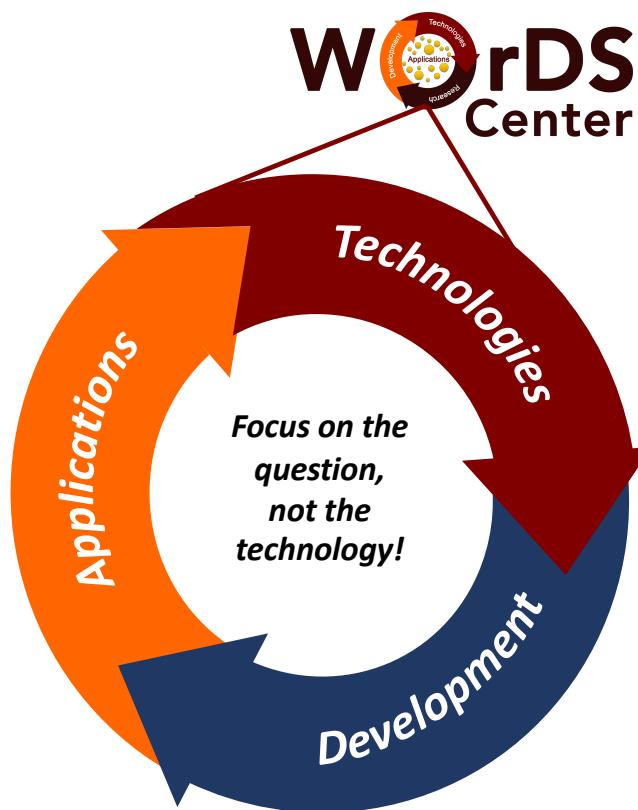
Founding Fellow, **Halıcıoğlu Data Science Institute**

Founding Director, **Workflows for Data Science Center of Excellence**

Founding Director, **WIFIRe Lab**

Part 1:

Data-Driven Problem Solving



Workflows for Data Science Center of Excellence at SDSC

<http://WorDS.sdsc.edu>

Mission:

Methodology and tool development
to enable collaborative workflow-driven science
and create solution architectures
on top of big data and advanced computing platforms.

Common Theme...

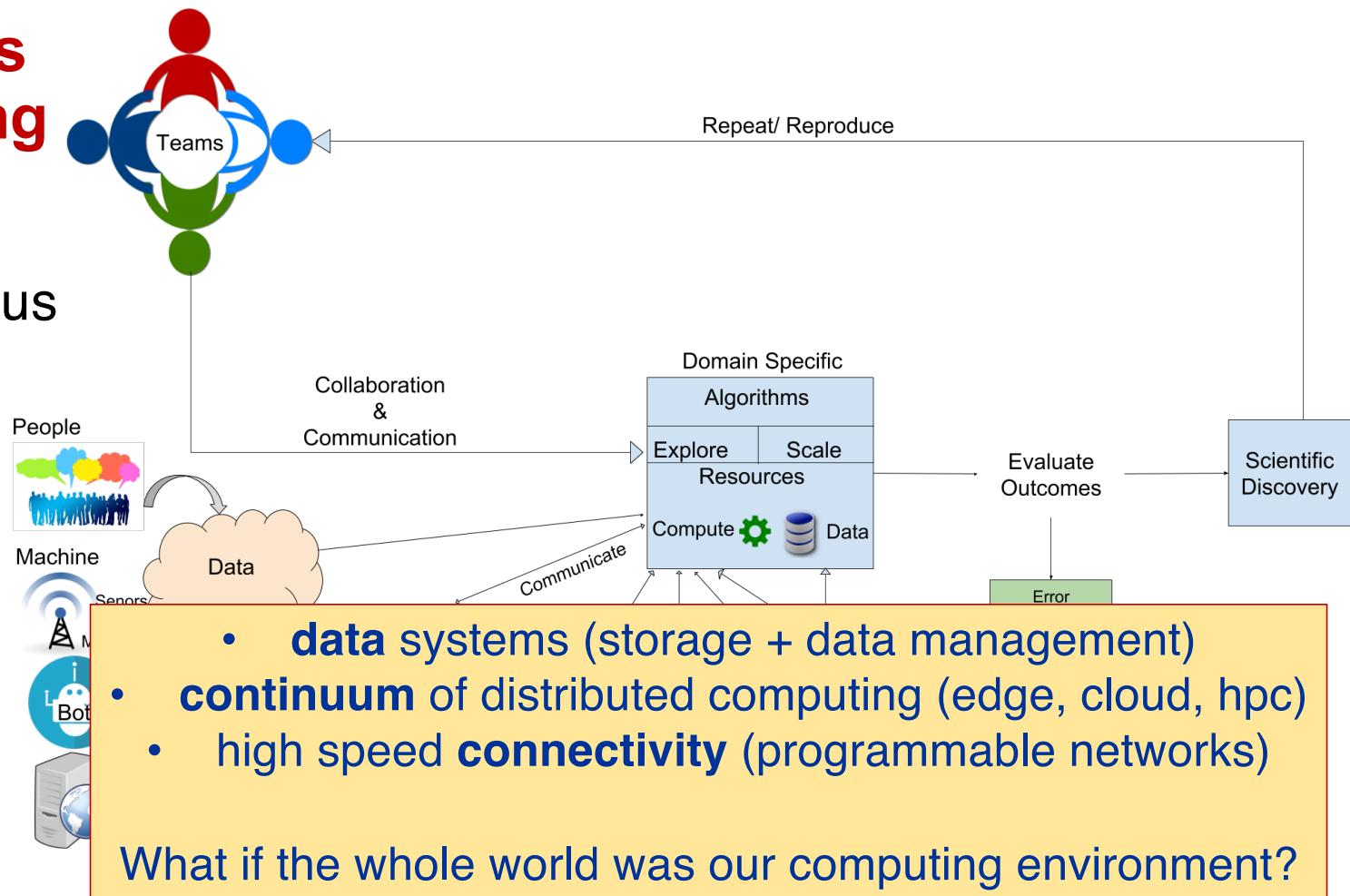
**“Big” Data, Computational
Science, Data Science, Cyberinfrastructure,**

and Their Applications

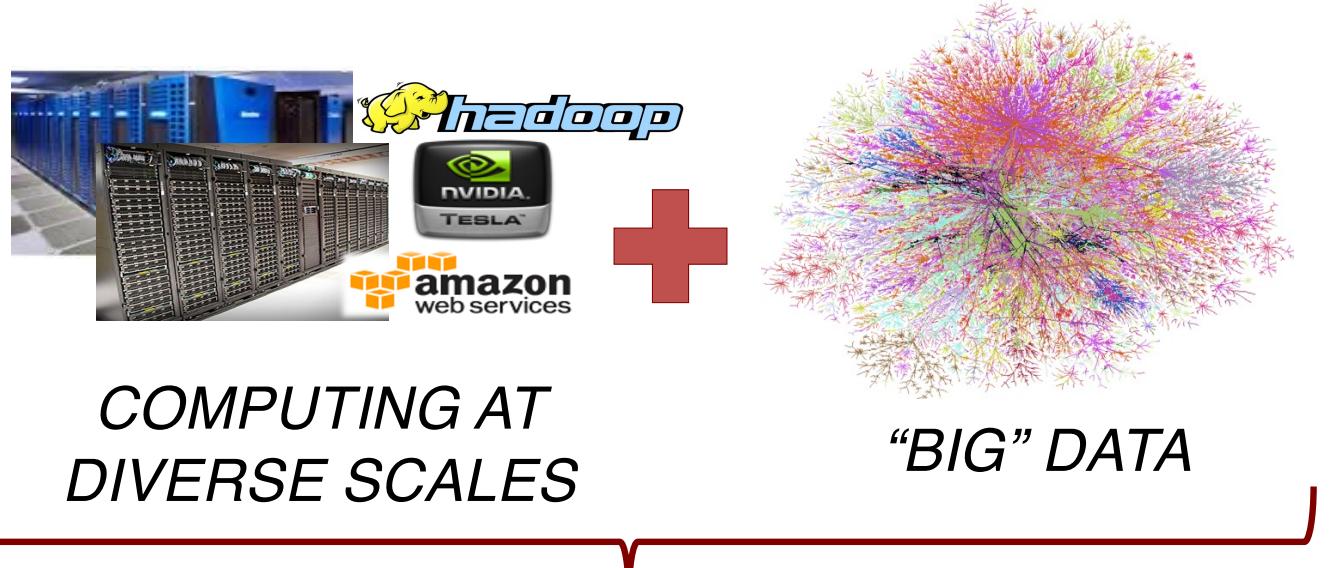
i.e., the problems we are solving

The problems we are solving are ...

- data-driven
- heterogeneous
- collaborative



**Big Data
combined with
Scalable
Computing
can be
very valuable.**



Smart Manufacturing

Computer-Aided Drug Discovery



Personalized Precision Medicine

Smart Cities



Smart Grid and Energy Management

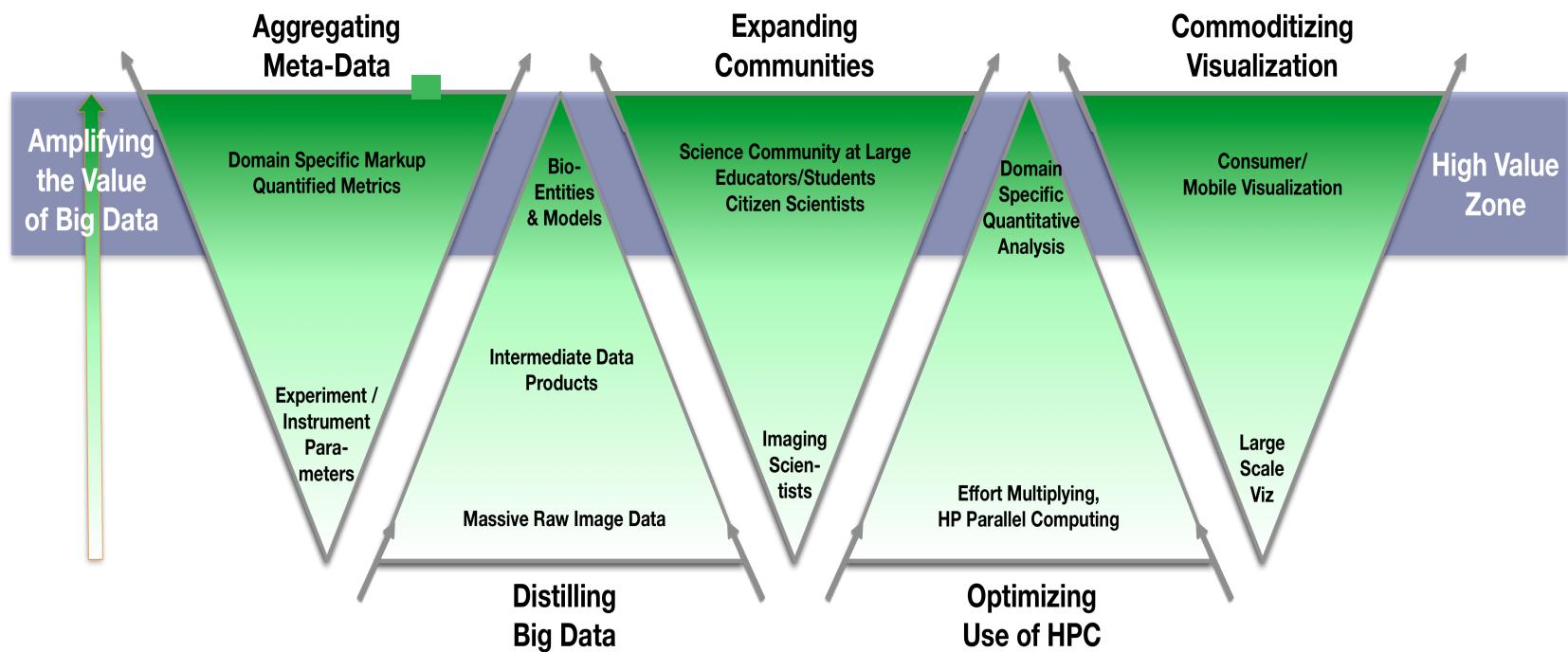
Disaster Resilience and Response



Precision Education

How do we amplify the value of Big Data?

Data → Knowledge → Action



How do we find the connections and impactful answers to our questions, i.e., turn it into benefit?



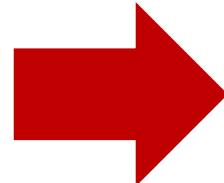
“We are drowning in
information and starving
for knowledge”

– John Naisbitt

Source: Megatrends, 1982

Going from Data to Discovery to Impact

Amplifying the
Value of Data
Related to X



Benefit Y for
Business,
Science,
Society, ...

We need to focus on the problems to solve.

Hidden Technical Debt in Machine Learning Systems

D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips
{dsculley, gholt, dgg, edavydov, toddphillips}@google.com
Google, Inc.

Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-François Crespo, Dan Dennison
{ebner, vchaudhary, mwyong, jcrespo, dennison}@google.com
Google, Inc.

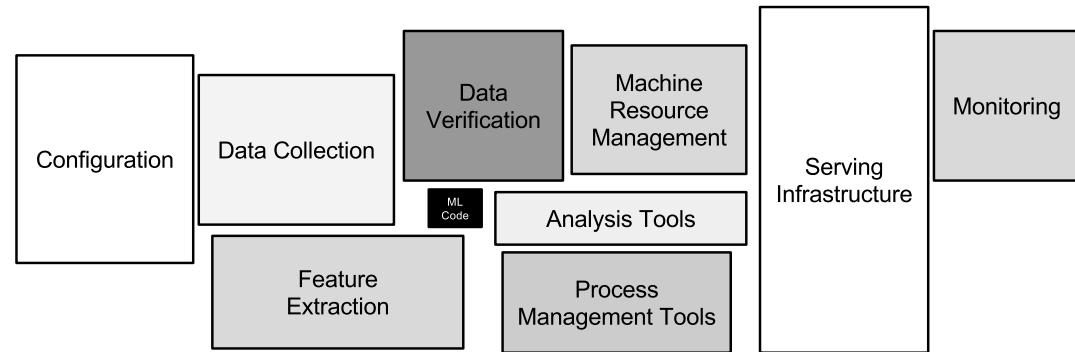


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

The problem is:
solutions are
complicated!

- Heterogenous systems and infrastructure
- Data management
- Machine learning, statistics and analytical methods
- Scalable process management
- Dynamic coordination and resource optimization
- Skilled interdisciplinary team
- Collaborative culture and communication tools

Research Goal:

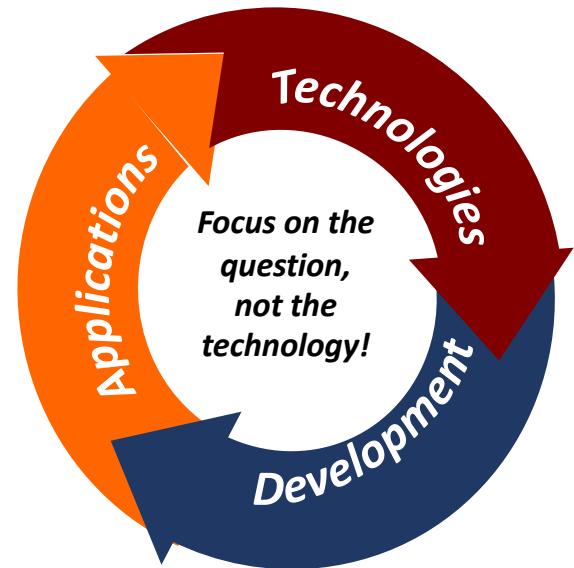
Create a Data Science Solution Ecosystem that Enables Needs and Best Practices

- data-driven
- scalable
- dynamic
- process-driven
- heterogeneous
- accountable
- reproducible
- interactive
- multidisciplinary
- collaborative

Team Data Science

**How can I get smart people
to collaborate and
communicate?**

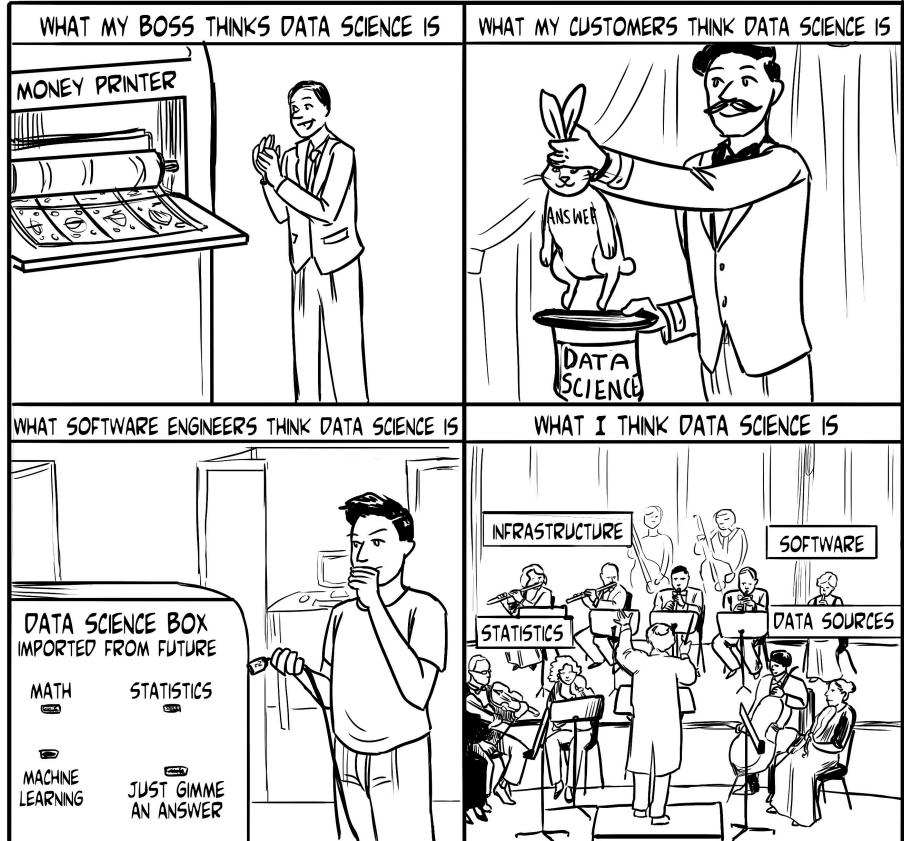
...to utilize data and computing to
generate insights and solve a question.



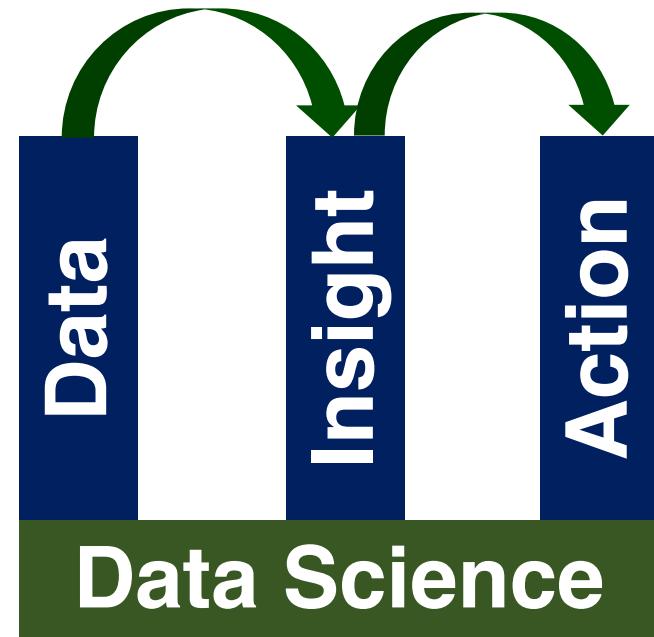
Part 2:

What is Data Science?

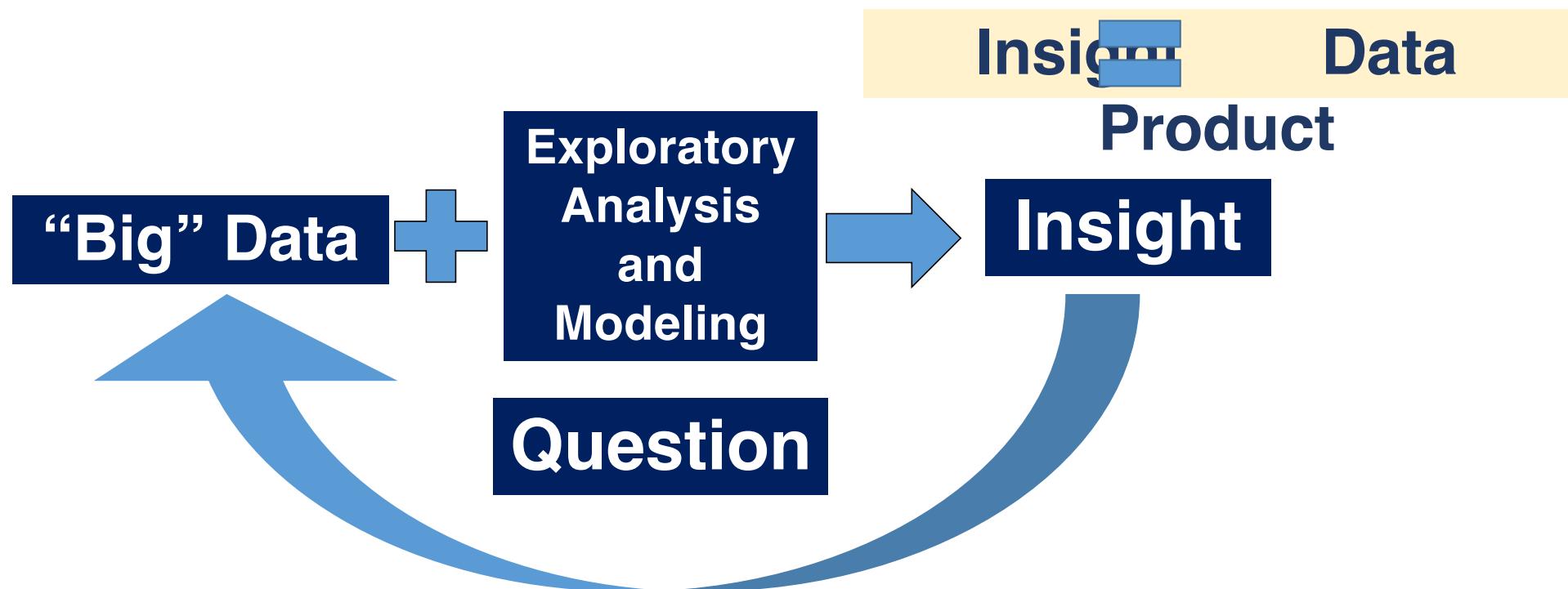
So what is data science?



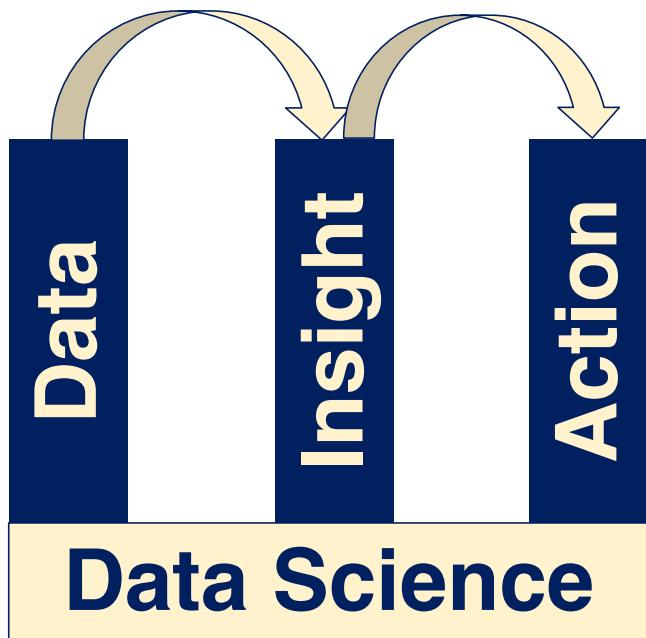
Ultimate Goal of Data Science



How does successful data science happen?



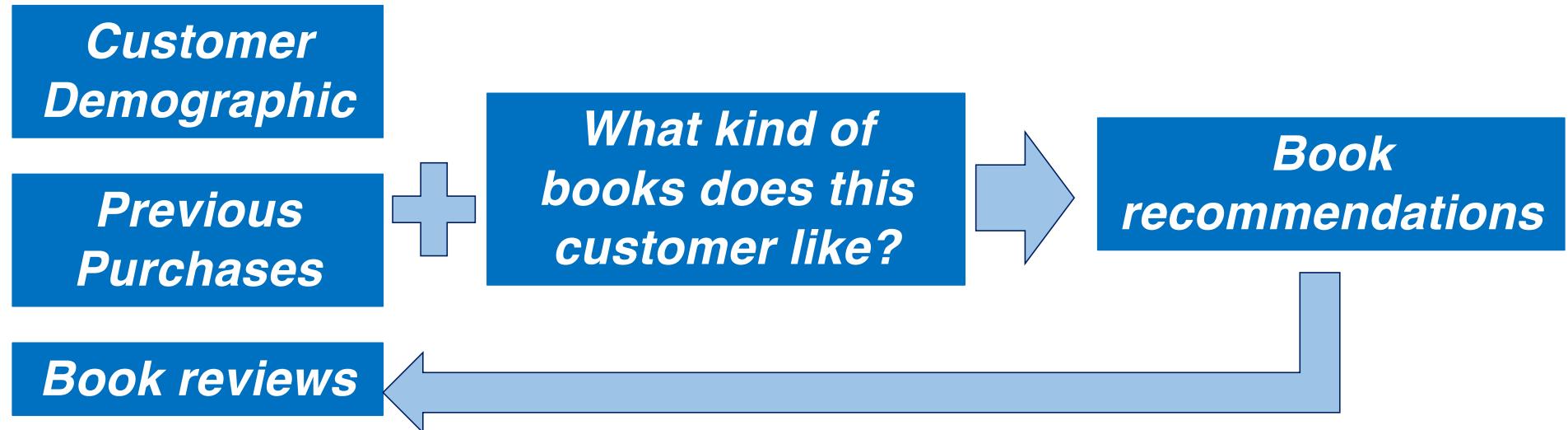
Data Products



“... a product that facilitates an end goal through the use of data”

-- DJ Patil, Former U.S. Chief Data Scientist
in Data Jujitsu

Book Recommendations

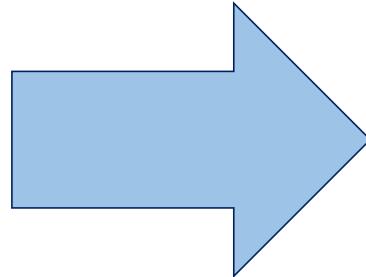


Find Potential Audience for a

*Model of
customer's book*



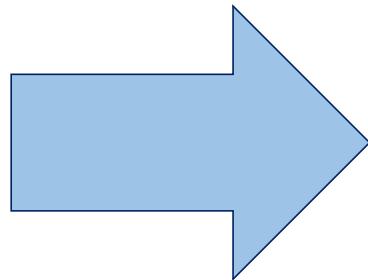
*New book
information*



*Who is likely to
like this book?*

Market a New Book

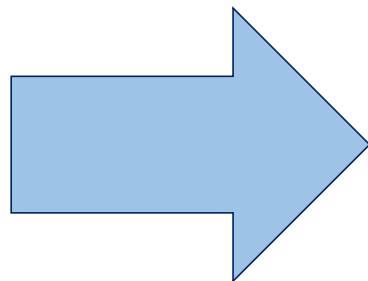
*Who is likely to
like this book?*



*Action to market
the book to the
right audience*

Market a New Book

Who is likely to like this book?



Action to market the book to the right audience

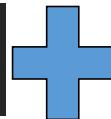
Insight



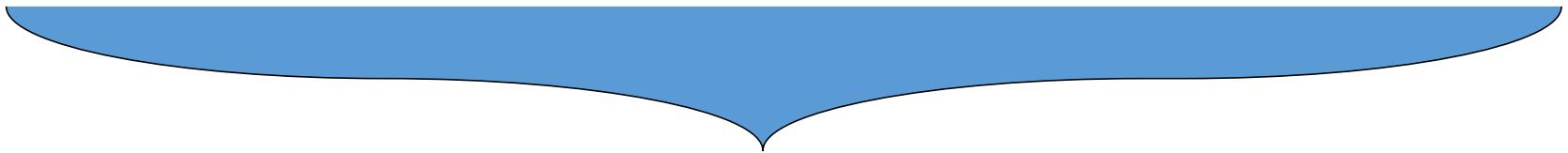
Action

Actionable Information

Historical data



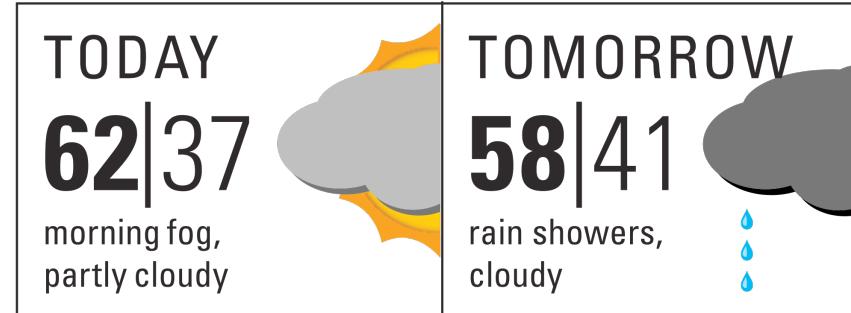
Near real-time data



Prediction

A large blue downward-pointing arrow originates from the bottom of the input boxes and points directly to the prediction box below.

Prediction



Action



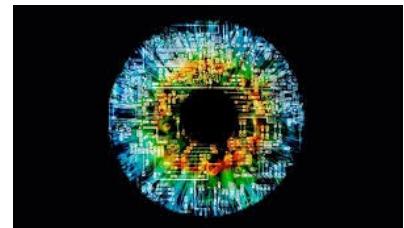
Systems and models
that help us to
understand data
in order to
gain insights and
make predictions
leading to action for impact.

Data Science is “IMPACT” Science!

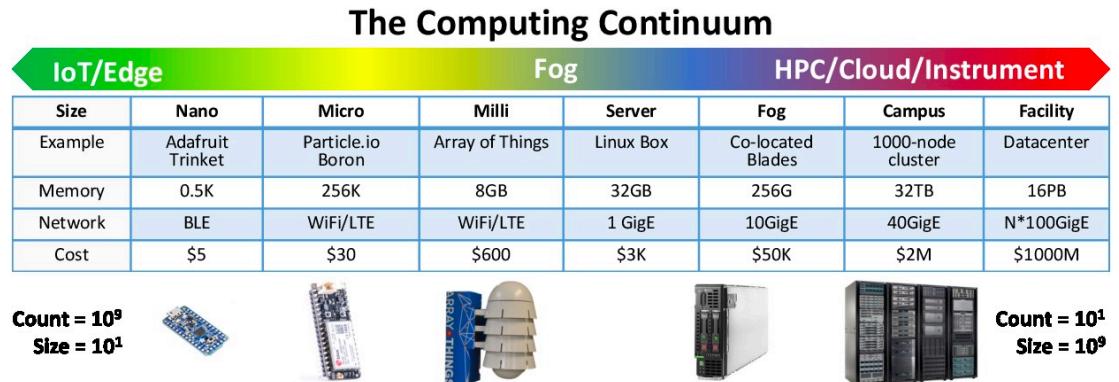
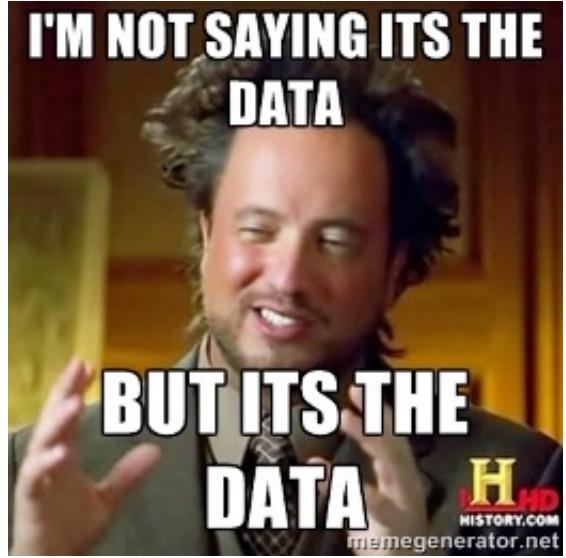
Part 3:

Big Data Science Today

The Disruptors:



Why are we talking about Data Science now?

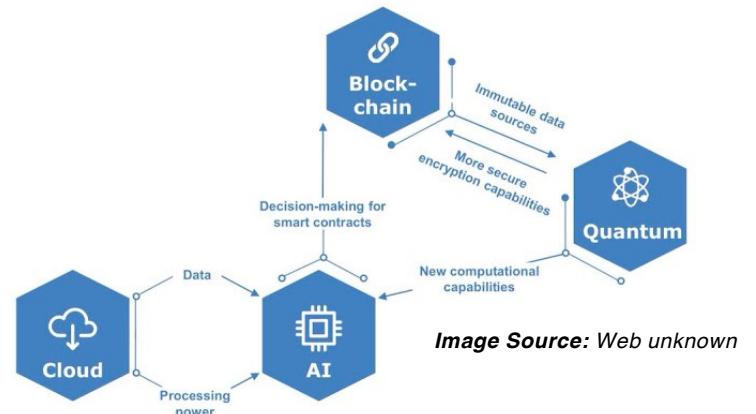
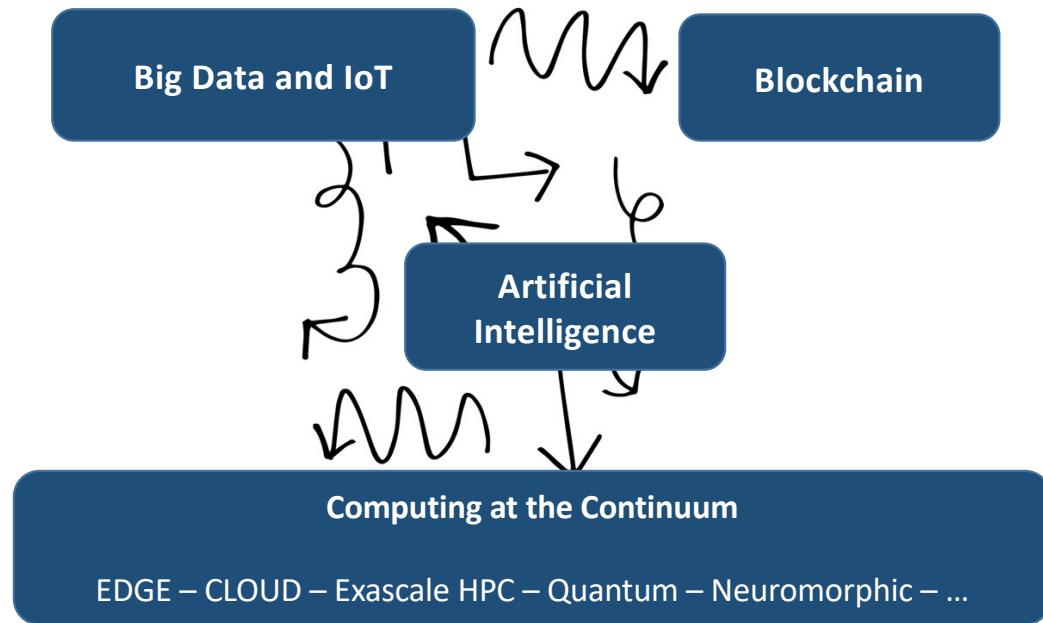


Big Data

Computing at the Continuum

Changed the way we solve problems!

New kids on the block...



How Much Data Is Big Data?



What happens in an 2016 INTERNET MINUTE?

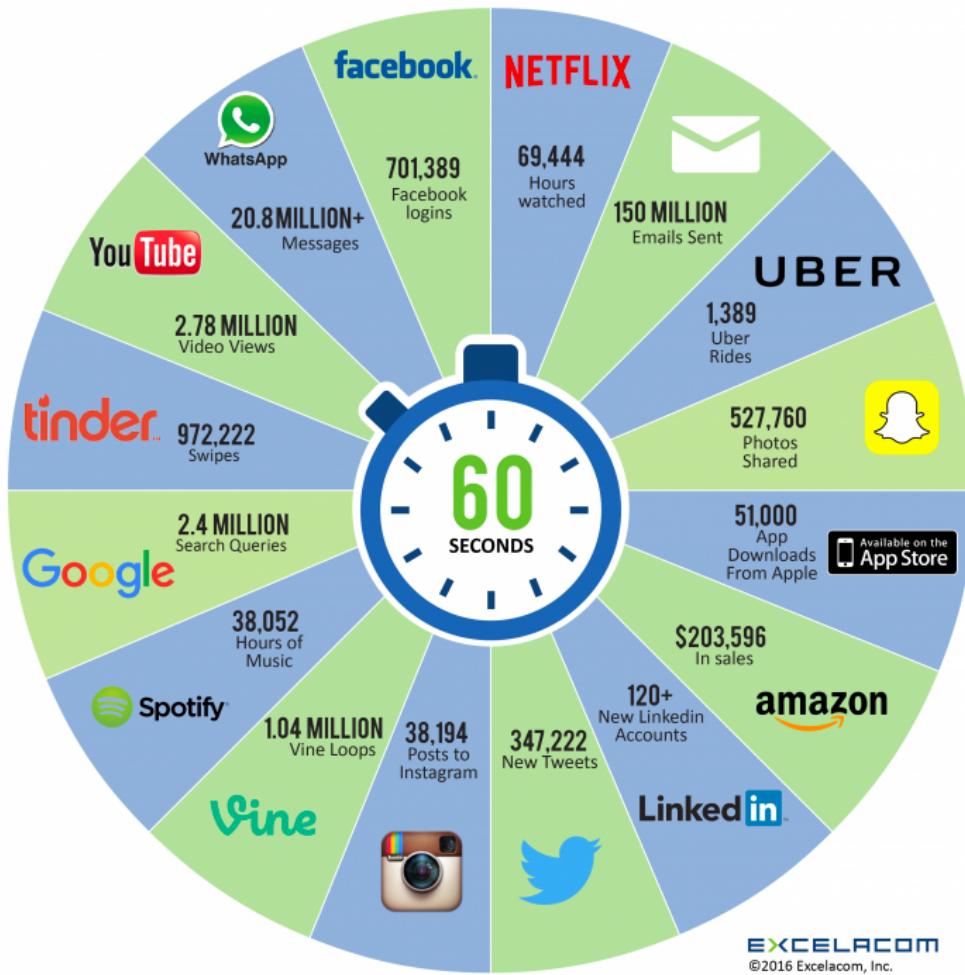
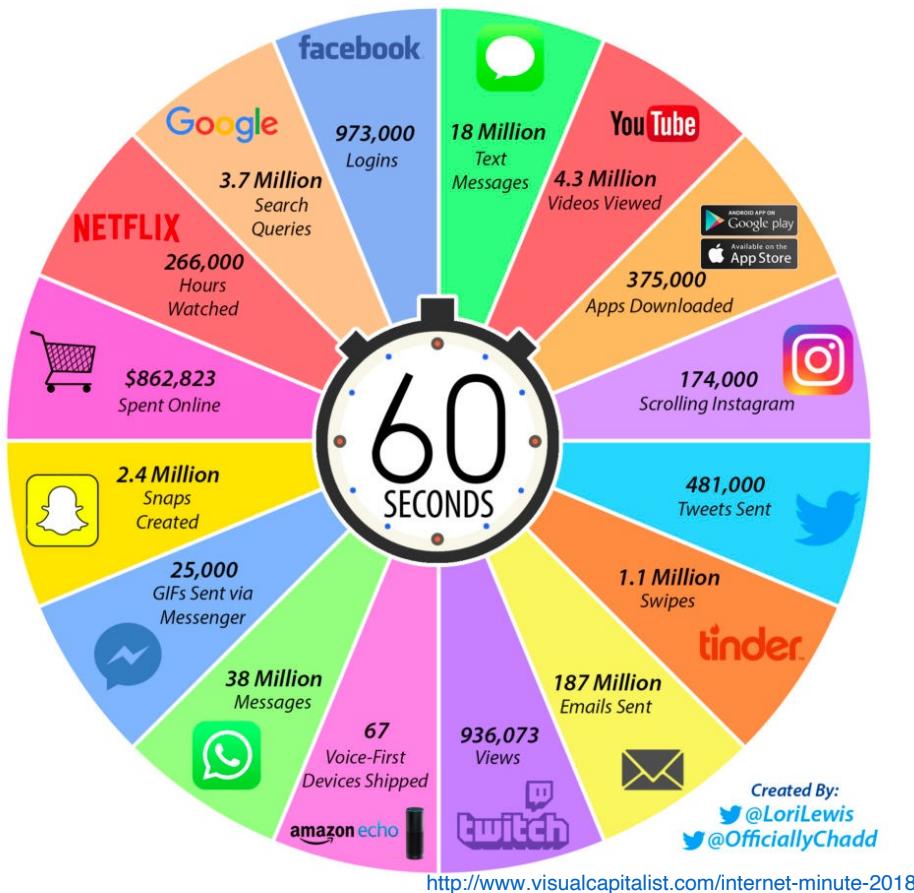
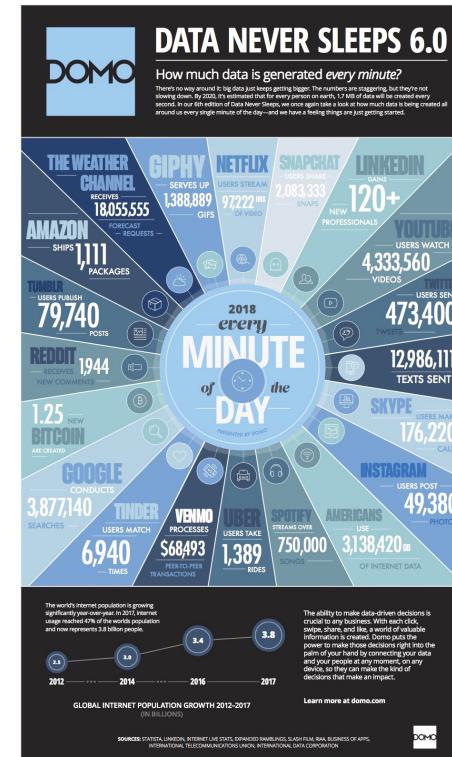


Image Source: <http://www.marketwatch.com/story/one-chart-shows-everything-that-happens-on-the-internet-in-just-one-minute-2016-04-26>

2018 This Is What Happens In An Internet Minute



Data Generated Every Minute on Internet

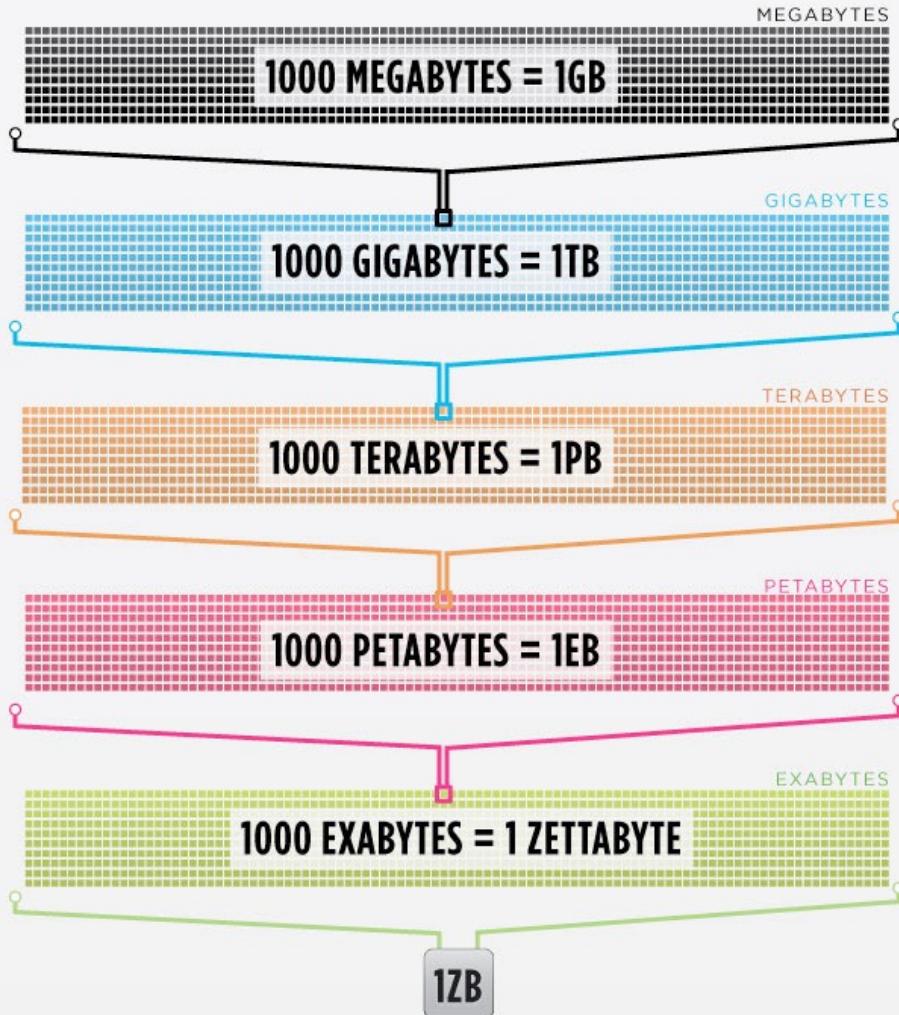


<https://www.adweek.com/digital/domio-data-never-sleeps-5/>

Scientific Data Management and Analysis

- HPWREN: hpwren.ucsd.edu
 - 30 TB of data annually
- MODIS: modis.gsfc.nasa.gov
 - 219 TB of data annually
- Precision Medicine
 - 4 EB (10^{18} bytes) of data in 2016 (www.fastcompany.com)
- LIGO, Deep Space Network, Protein Data Bank, ...

But how much data are we talking about?



100 MBs ≈ couple of volumes of Encyclopedias

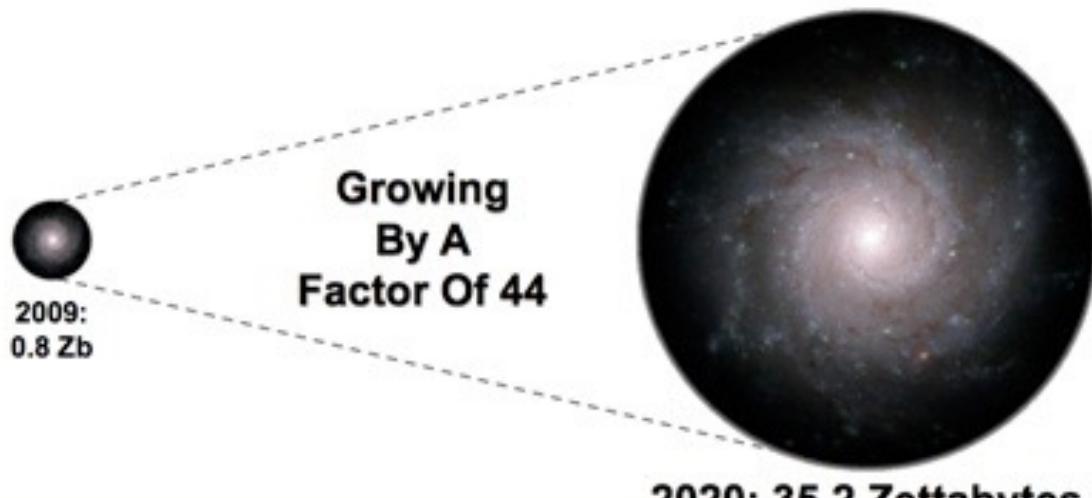
A DVD ≈ 5 GBs

1 TB ≈ 300 hours of good quality video

LHC ≈ 15 PBs a year

D (iltintas@ucsd.edu)

The Digital Universe 2009-2020



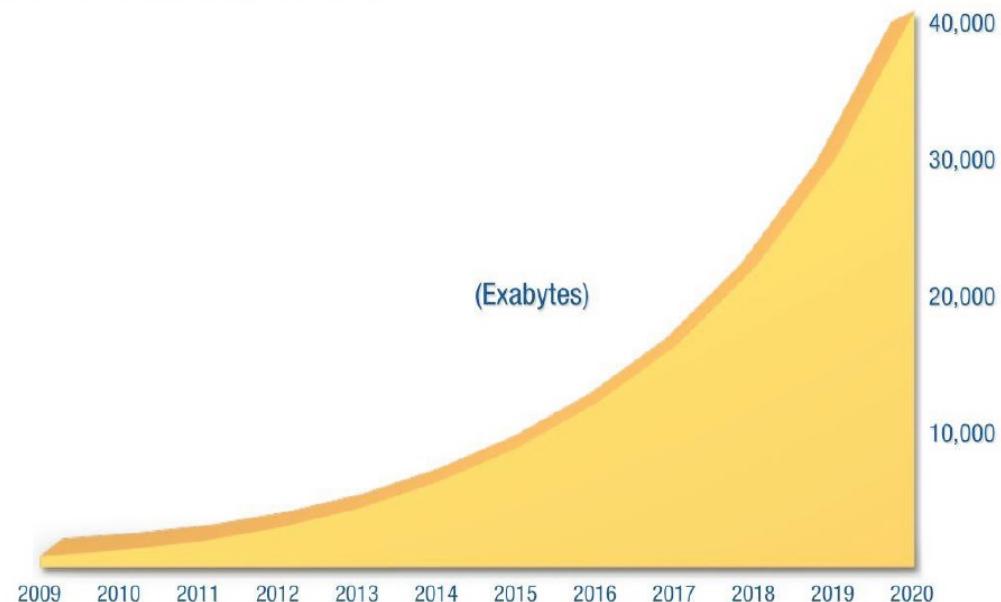
Source: IDC Digital Universe Study, sponsored by EMC, May 2010

©2010 EMC Corporation. All rights reserved.



Exponential data growth!

The Digital Universe: 50-fold Growth from the Beginning of 2010 to the End of 2020



This IDC graph predicts exponential growth of data from around 3 zettabytes in 2013 to approximately 40 zettabytes by 2020. An exabyte equals 1,000,000,000,000,000 bytes and 1,000 exabytes equals one zettabyte. Source: IDC's Digital Universe Study, December 2012, <http://www.emc.com/collateral/analyst-reports/idc-the-digital-universe-in-2020.pdf>.

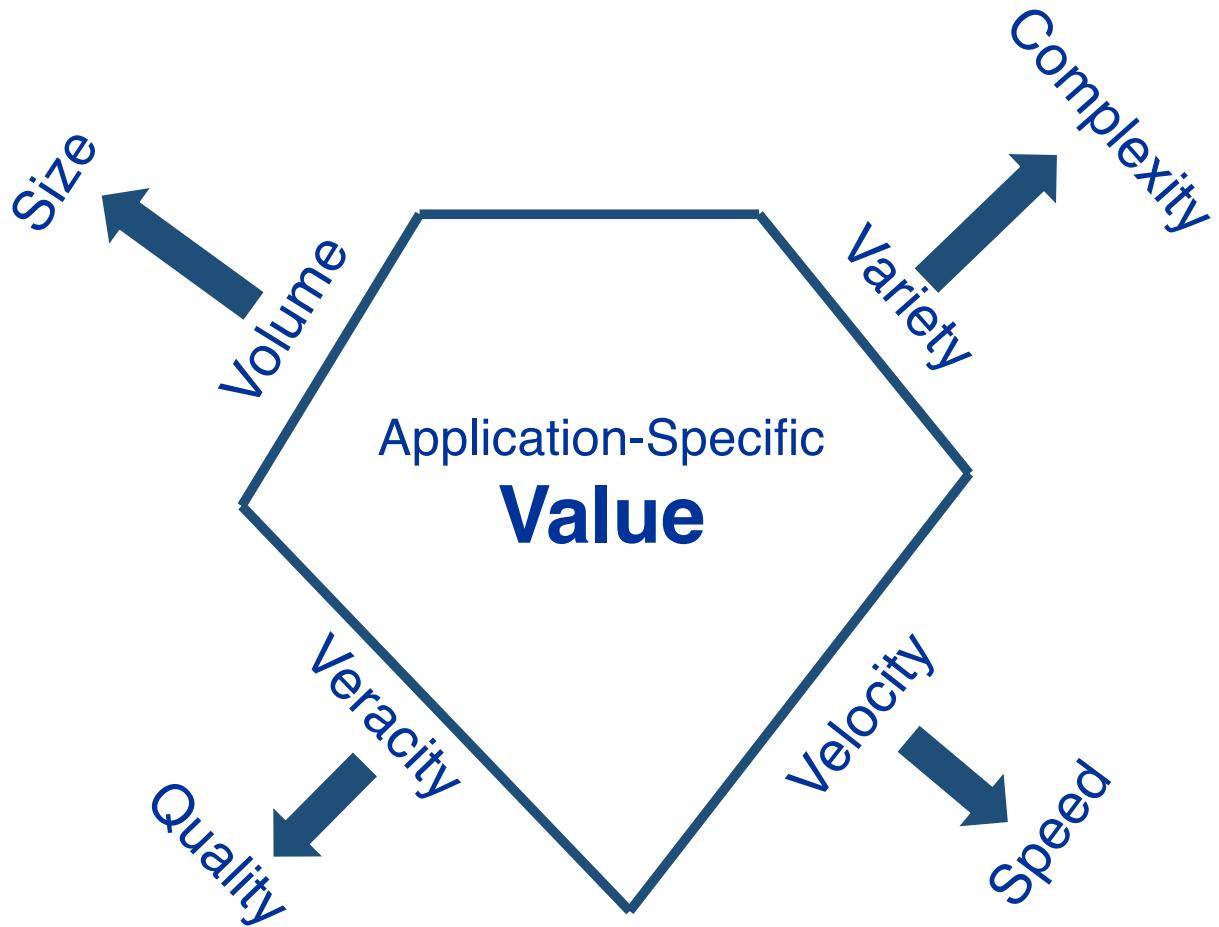
Data Deluge

“We are drowning in information and starving for knowledge”

– John Naisbitt

Source: Megatrends, 1982





What problem are these disruptors helping to solve?

New Opportunities for AI-Driven Approaches and Cyberinfrastructure

- Closing the loop between observation, experimentation and simulation
- Dynamic data-
- Real-time data
- Reactive systems
- Coupling the environment and AI
- New software and tools for AI-driven computing
- Intelligent security, integrity and privacy practices

All require getting value from data as fast as possible!

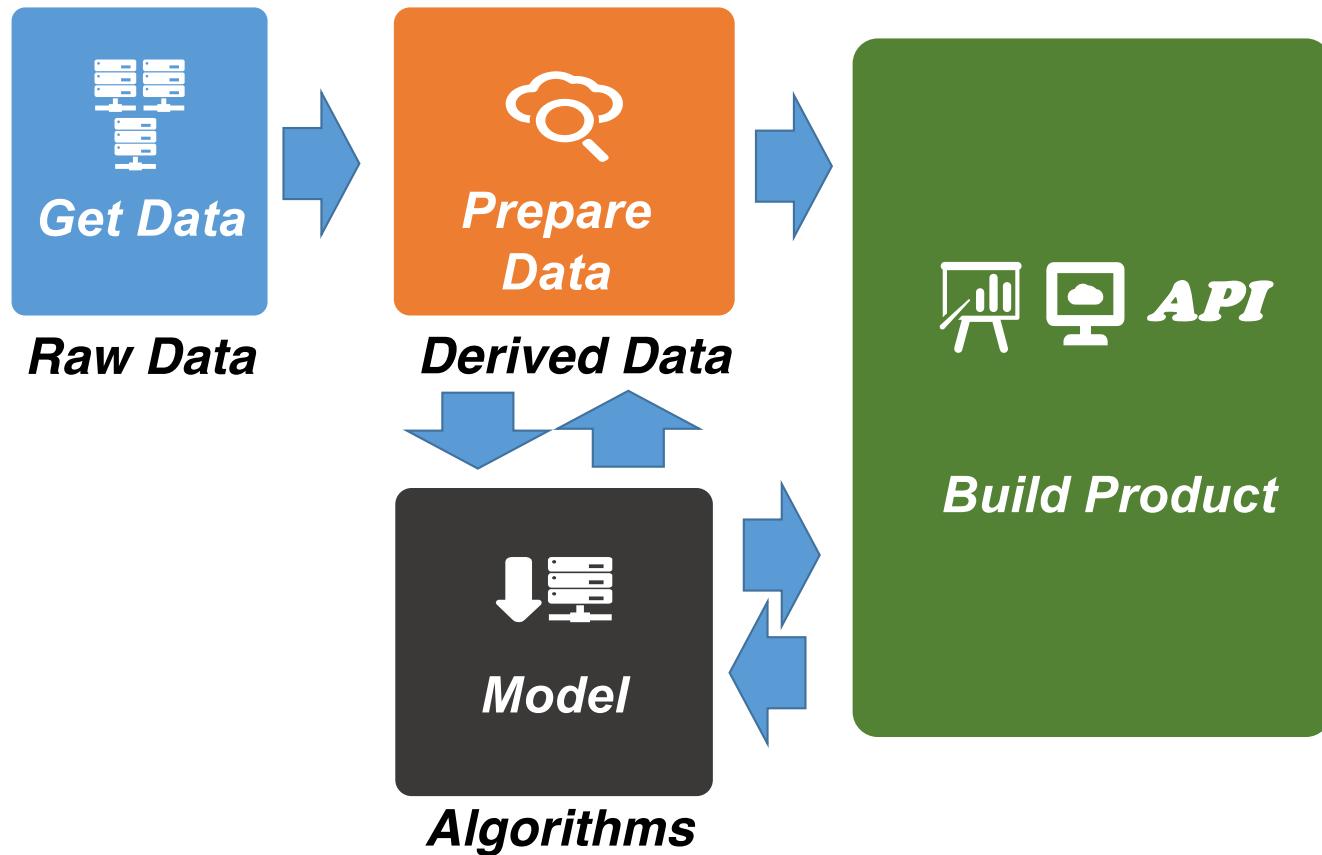
How do we find the connections?



Part 4:

Scalable Data Science Process

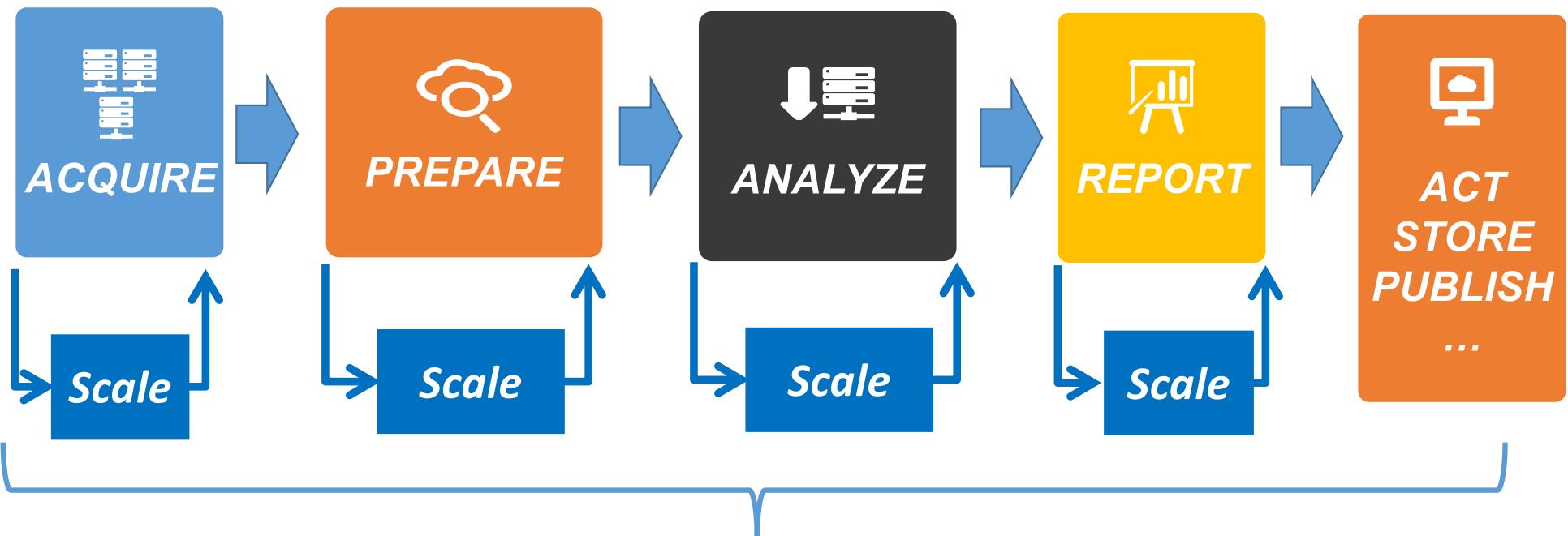
Going from raw data to a model using data science...



- ***Visual Dashboards***
- ***Web Interfaces***
- ***Programming Interfaces***
- ***Robotics platforms***

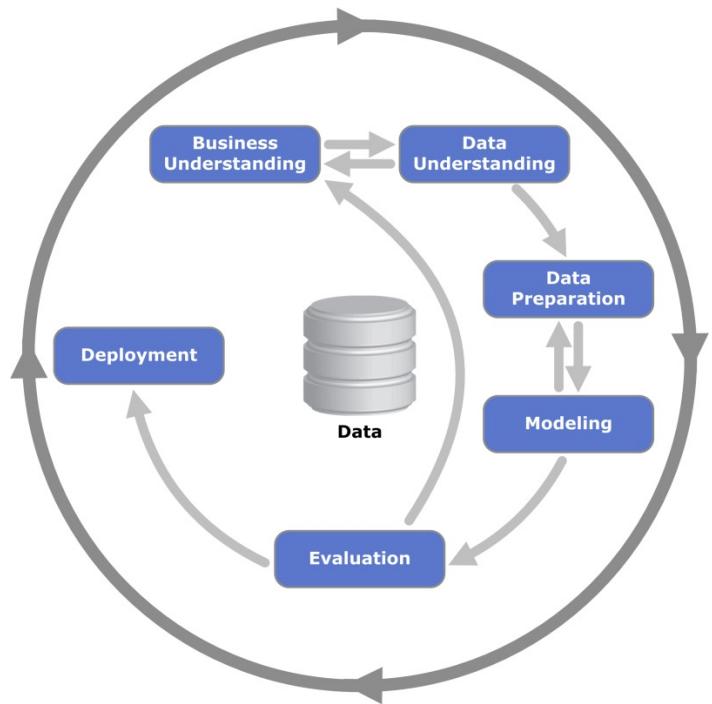
Data Engineering

Computational Data Science

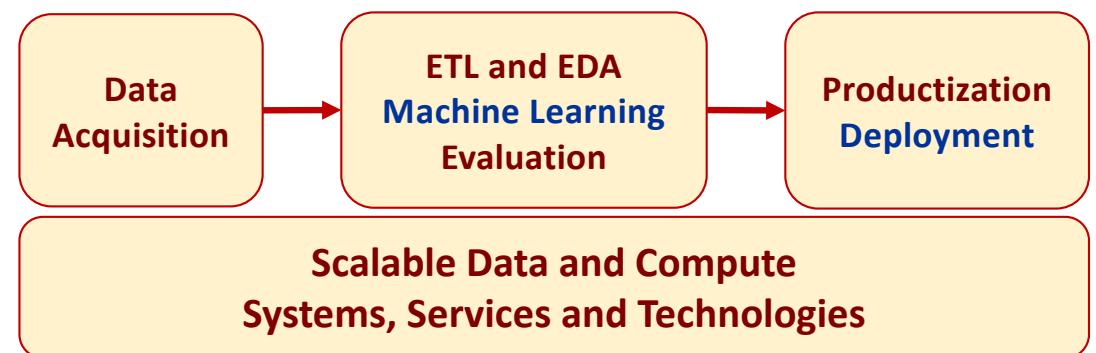


Continuous Iteration, Integration, Programmability, Measurement and Scalability

Scalable Data Science Process



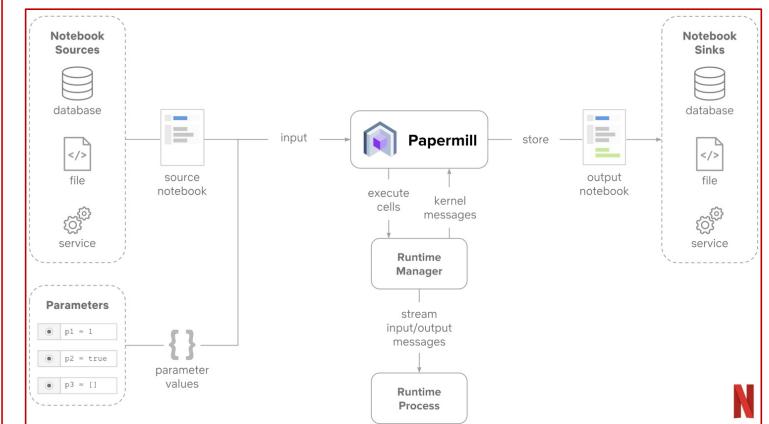
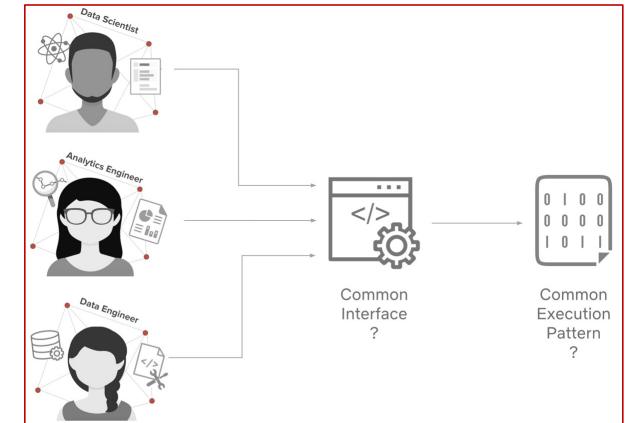
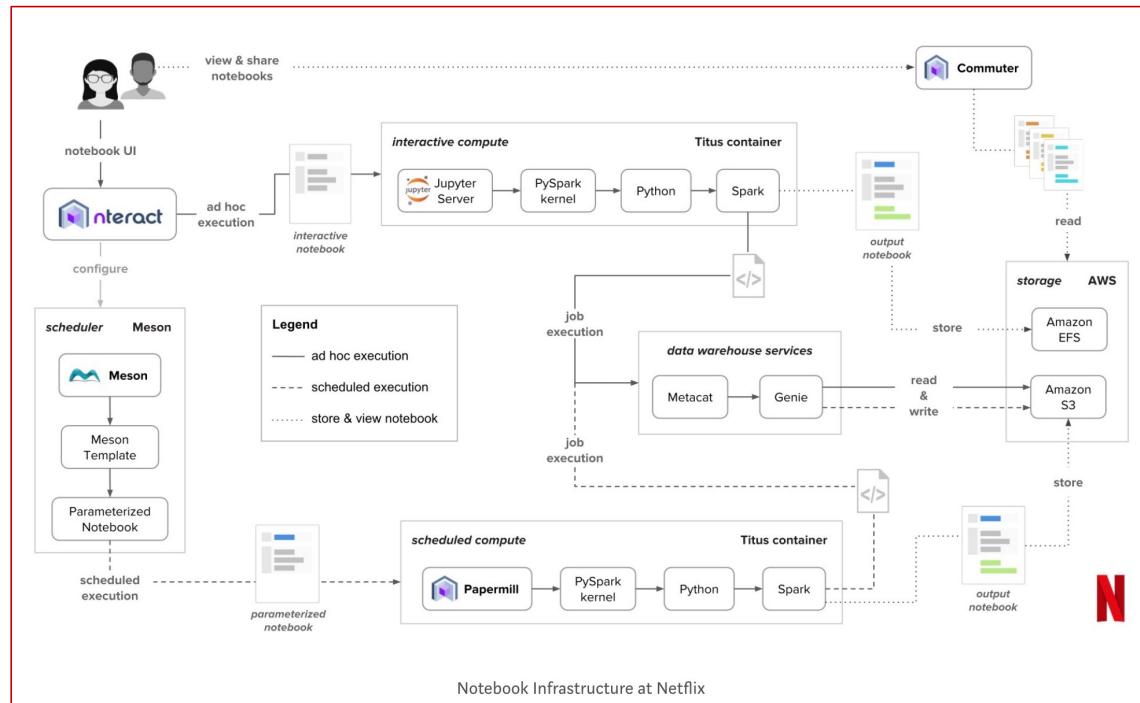
together with



CRISP-DM: Cross-industry standard process for data mining

Example: Notebook Innovation at Netflix

- <https://medium.com/netflix-techblog/notebook-innovation-591ee3221233>
- <https://medium.com/netflix-techblog/scheduling-notebooks-348e6c14cf6>



Pieces of the Solution

- Stakeholders
- Datasets and data lifecycle
- Compliance requirements
- Defined actions
- Analytical methods
- Technical infrastructure

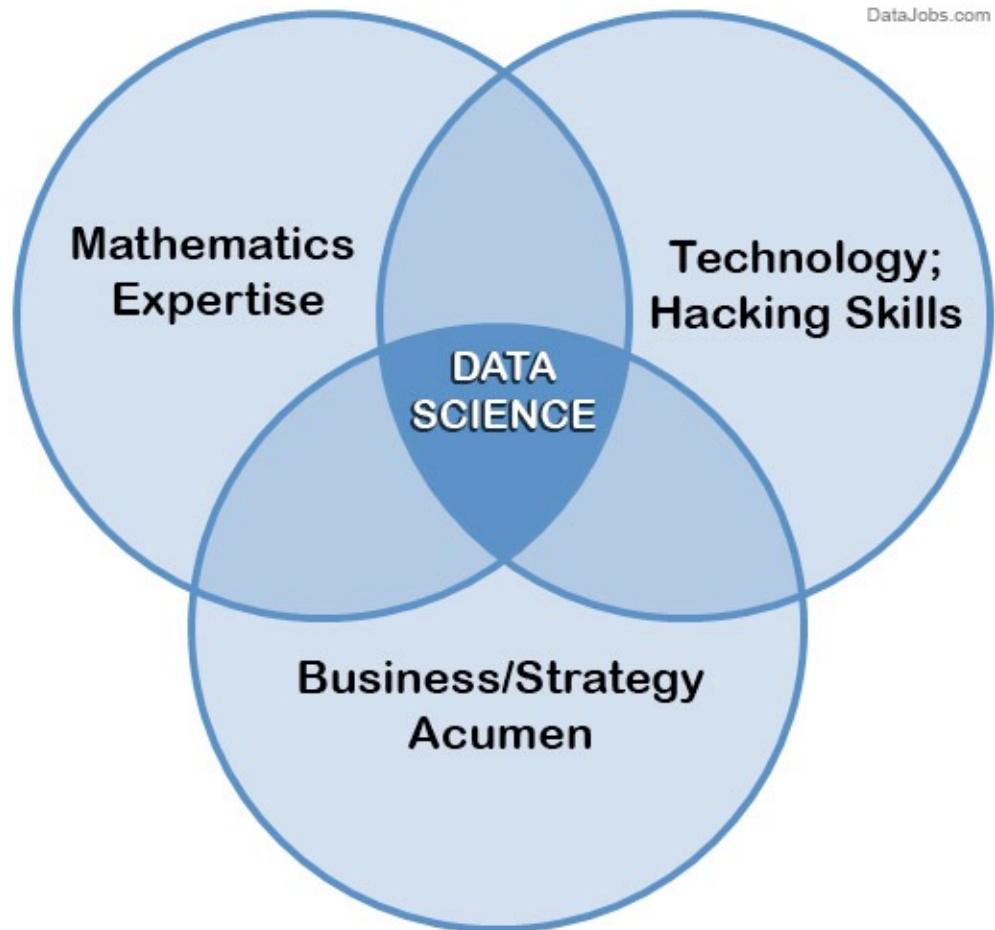


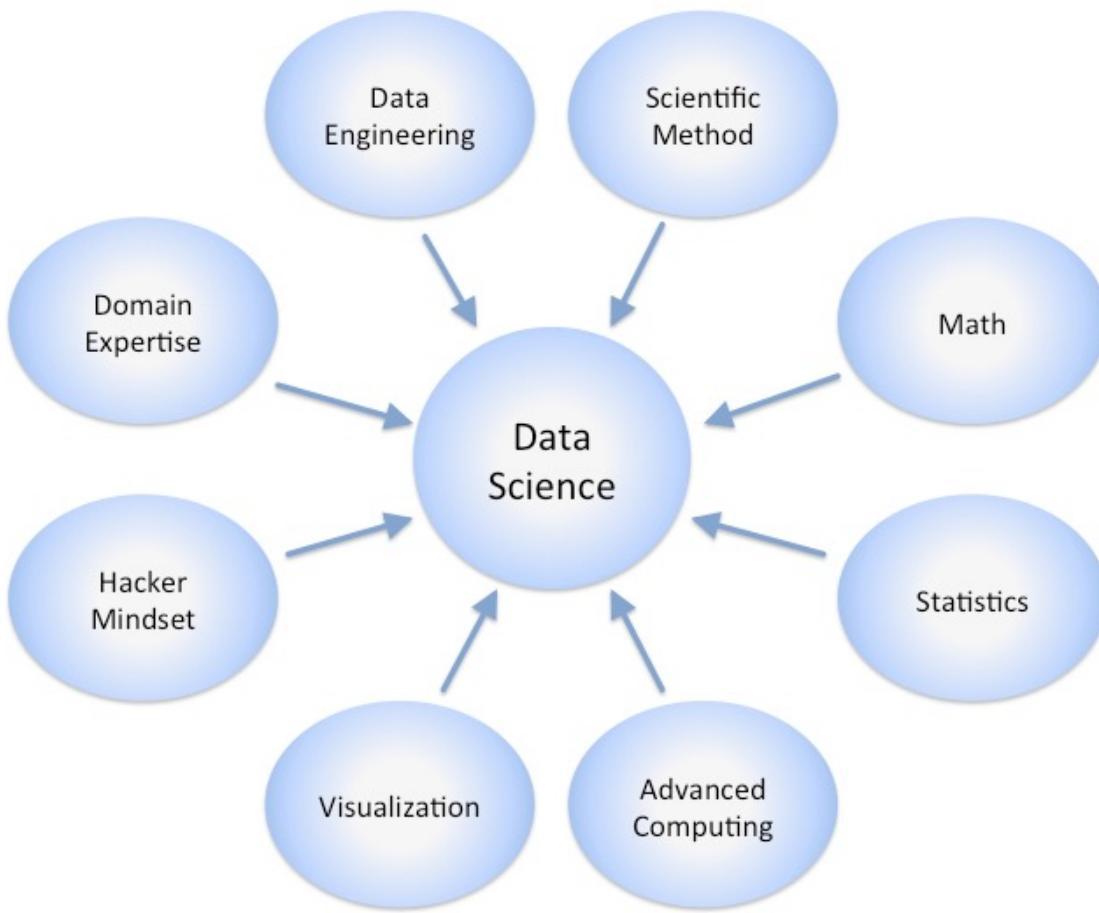
Part 5:

Team Data Science

Modern Data Science Skills

- Programming in Python
- Statistics
- Machine Learning
- Scalable Big Data Analysis





MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ★ Machine learning
- ★ Statistical modeling
- ★ Experiment design
- ★ Bayesian inference
- ★ Supervised learning: decision trees, random forests, logistic regression
- ★ Unsupervised learning: clustering, dimensionality reduction
- ★ Optimization: gradient descent and variants



PROGRAMMING & DATABASE

- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing package e.g. R
- ★ Databases SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with xaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

- ★ Able to engage with senior management
- ★ Story telling skills
- ★ Translate data-driven insights into decisions and actions
- ★ Visual art design
- ★ R packages like ggplot or lattice
- ★ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21th century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ★ Machine learning
- ★ Statistical modeling
- ★ Experiment design
- ★ Bayesian inference
- ★ Supervised learning: decision trees, random forests, logistic regression
- ★ Unsupervised learning: clustering, dimensionality reduction
- ★ Optimization: gradient descent and variants



PROGRAMMING & DATABASE

- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing packages, e.g., R
- ★ Databases: SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with xaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

- ★ Able to engage with senior management
- ★ Story telling skills
- ★ Translate data-driven insights into decisions and actions
- ★ Visual art design
- ★ R packages like ggplot or lattice
- ★ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

MarketingDistillery.com is a group of practitioners in the area of e-commerce marketing. Our fields of expertise include: marketing strategy and optimization: customer tracking and on-site analytics; predictive analytics and econometrics; data warehousing and big data systems; marketing channel insights in Paid Search, SEO, Social, CRM and brand.

Marketing DISTILLERY
© Krzysztof Zawadzki

SDSC SAN DIEGO
SUPERCOMPUTER CENTER

İlkay Altıntaş, PhD (iagtintas@ucsd.edu)

Marketing DISTILLERY
© Krzysztof Zawadzki

UC San Diego
HALİCIOĞLU DATA SCIENCE INSTITUTE



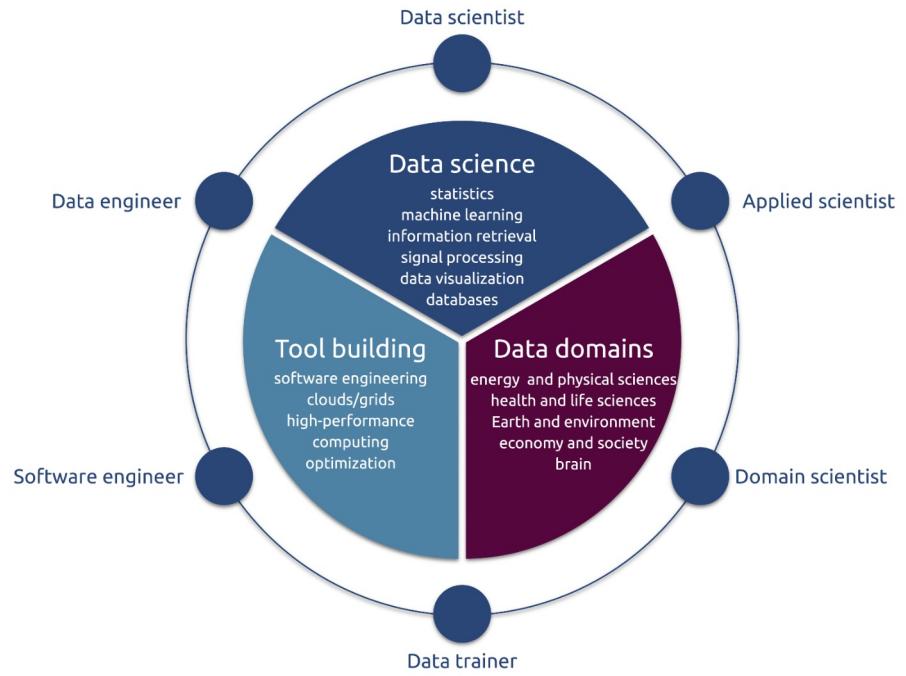
Are data scientists unicorns?

*Data science is
team sport!*

Data Science is “WE” Science!

Data Science Team

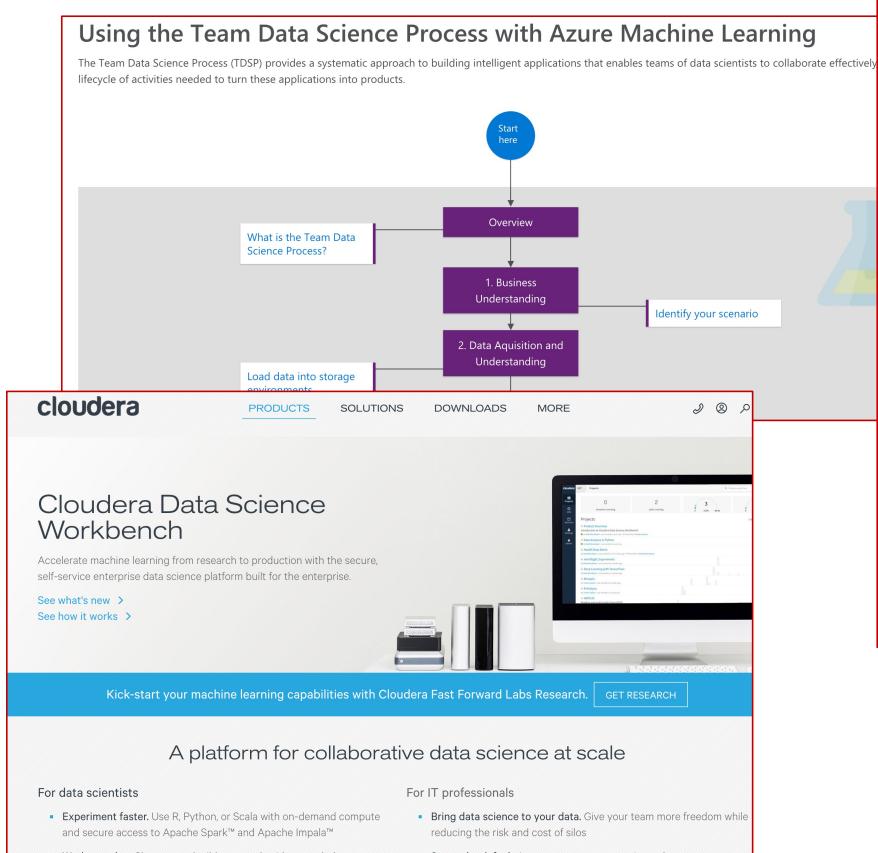
- Data engineer
- Data analyst
- Methods expert
- Scalability and operations expert
- Business manager
- Business analyst
- Scientist
- Visualization and dashboard developer
- Solution architect
- Story teller/coordinator
- Project manager



The data science ecosystem: activities and actors
<https://medium.com/@balazskegl/the-data-science-ecosystem-678459ba6013>

**Expertise and skills often overlap,
but nobody has it all!**

Some examples of team data science practice...



Using the Team Data Science Process with Azure Machine Learning

The Team Data Science Process (TDSP) provides a systematic approach to building intelligent applications that enables teams of data scientists to collaborate effectively throughout the lifecycle of activities needed to turn these applications into products.

```
graph TD; Start((Start here)) --> Overview[Overview]; Overview --> Business[1. Business Understanding]; Business --> Data[2. Data Acquisition and Understanding]; Data --> Identify[Identify your scenario]; Load[Load data into storage] --> Business;
```

cloudera

Cloudera Data Science Workbench

Accelerate machine learning from research to production with the secure, self-service enterprise data science platform built for the enterprise.

[See what's new >](#) [See how it works >](#)

Kick-start your machine learning capabilities with Cloudera Fast Forward Labs Research. [GET RESEARCH](#)

A platform for collaborative data science at scale

For data scientists

- Experiment faster. Use R, Python, or Scala with on-demand compute and secure access to Apache Spark™ and Apache Impala™.
- Work together. Share code with researchers with your whole team.

For IT professionals

- Bring data science to your data. Give your team more freedom while reducing the risk and cost of silos.
- Create by default. Leverage common compute and infrastructure.



IBM Cloud Blog Why IBM Products Solutions Garage Pricing Blogs Docs Support

Garage

The IBM Garage expands to Data Science Insights

Schedule a visit to a Garage

July 12, 2017 | Written by:

Categorized: Garage | Service

Share this post:

[f](#) [in](#) [t](#)

In the IBM Bluemix Garage locations, the innovative id

Jupyter

Install About Us Community Documentation

jupyterhub

A multi-user version of the notebook designed for companies, classrooms and research labs

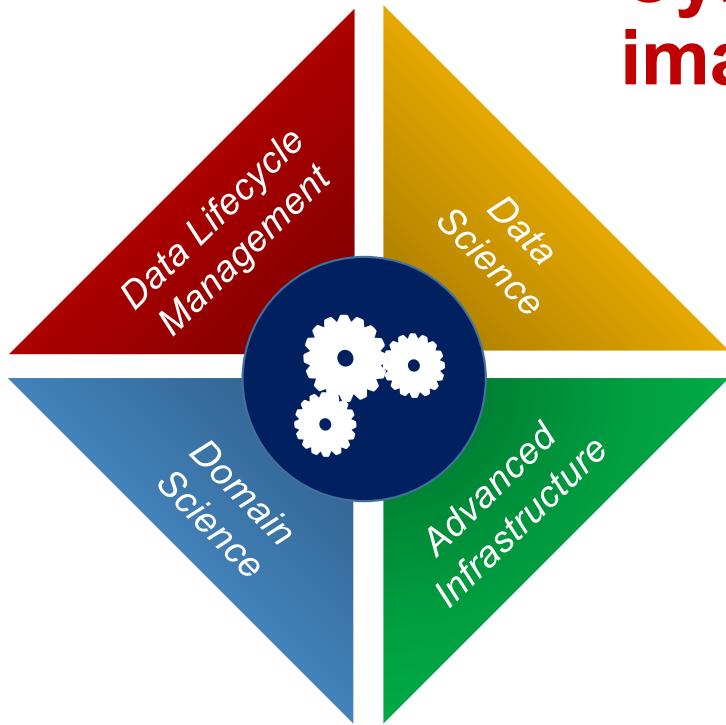
What is JupyterHub?

JupyterHub brings the power of notebooks to groups of users. It gives users access to computational environments and resources without burdening tasks. Users - including students, researchers, and data scientists - can get their work done in their own workspaces on shared resources which can be administrators.

JupyterHub runs in the cloud or on your own hardware, and makes it possible to serve a pre-configured data science environment to any user in the suitable for small and large teams, academic courses, and large-scale infrastructure.

Part 6:

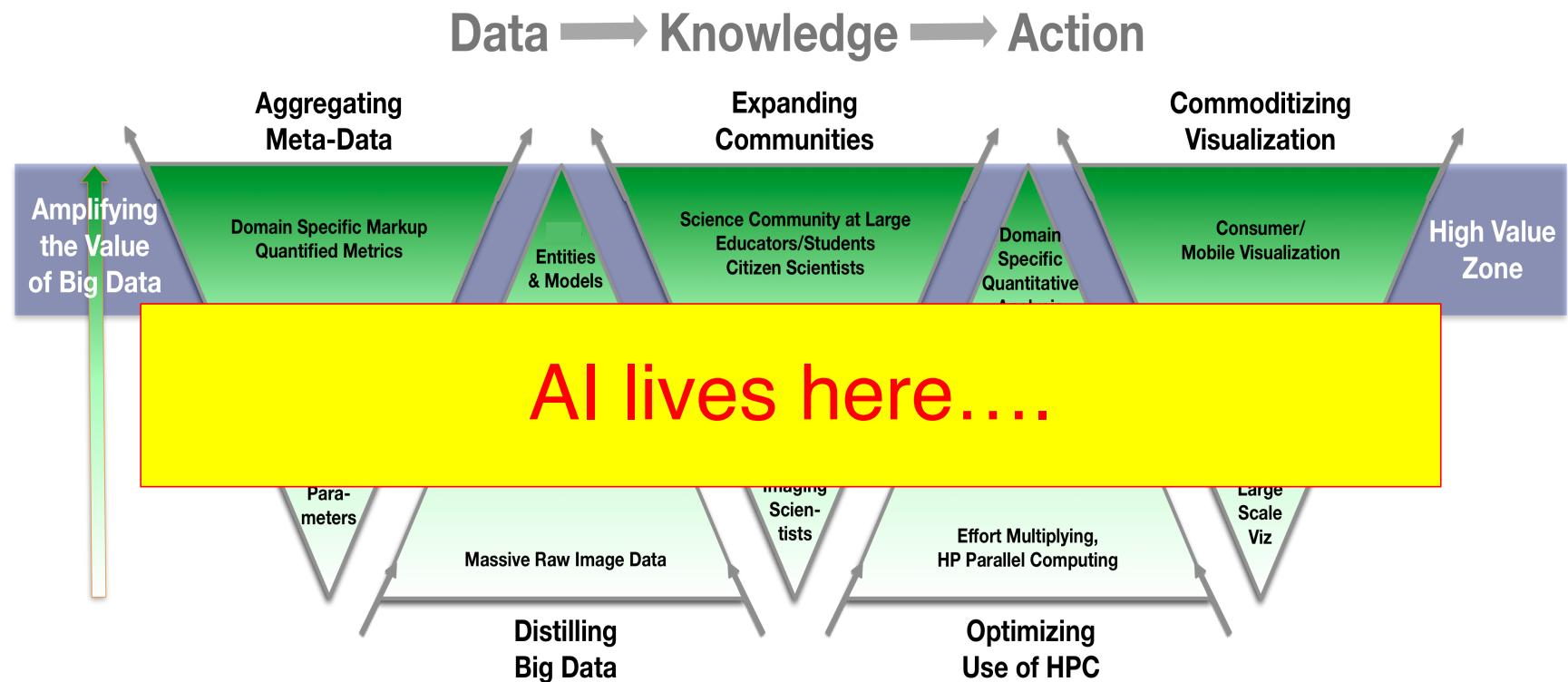
Continuous AI Integration at the Digital Continuum



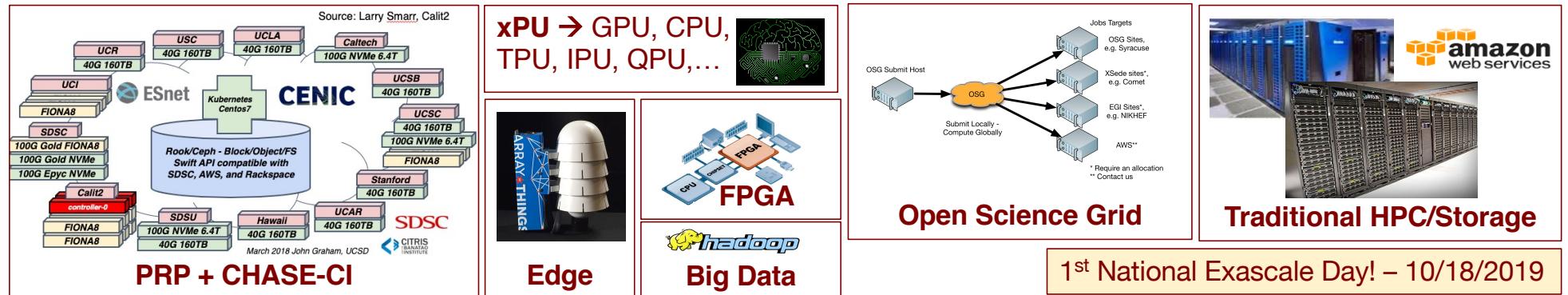
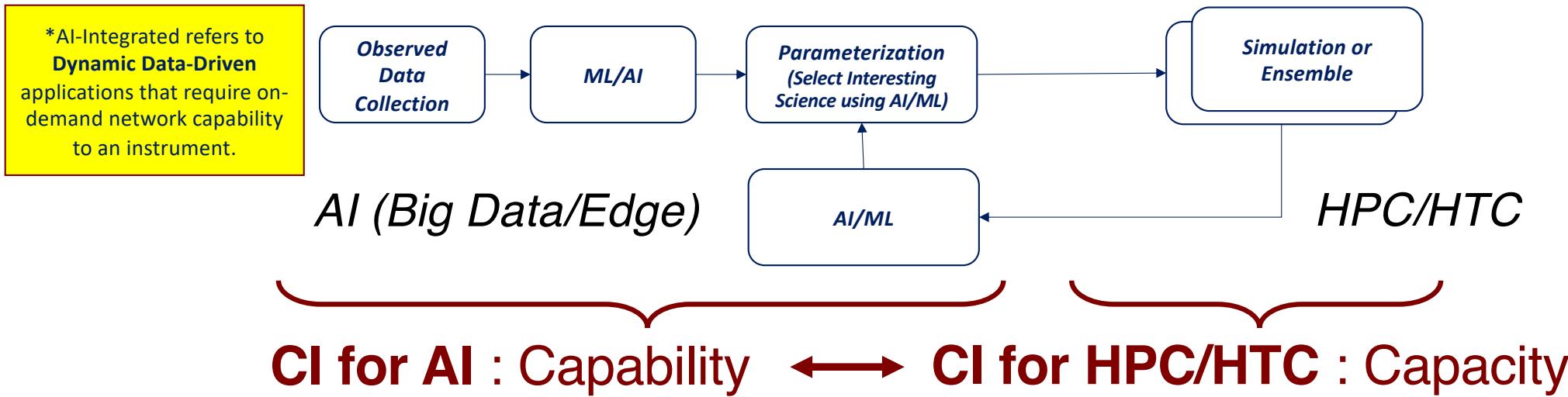
Cyberinfrastructure to support imaging research:

- Heterogenous systems
- Data management
- Data-driven methods
- Scalable tools for dynamic coordination and resource optimization
- Skilled interdisciplinary workforce
- Collaboration tools that enable groups to converge

Methodologies, standards, tools and resources are required to amplify value of image data.



A Typical Heterogeneous AI-Integrated* Workflow



CI Requirements for Convergence



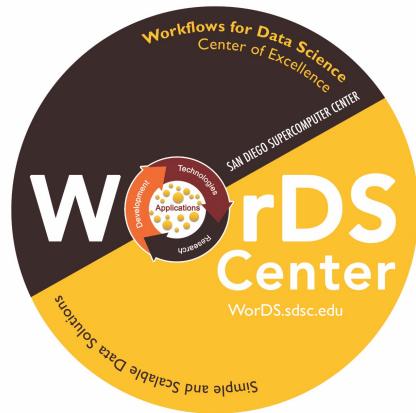
Dynamic composability matters.

Systems are only useful if groups can integrate them into applications.



TEAMWORK

Tools that enhance teamwork need to be coupled with AI systems.



Questions?

The presented work is collaborative work with many wonderful individuals, and parts of it are funded by NSF, DOE, NIH, UC San Diego and various industry partners.

Contact: *İlkay Altintas, Ph.D.*
Email: ialtintas@ucsd.edu

SDSC SAN DIEGO
SUPERCOMPUTER CENTER

<https://words.sdsc.edu/publications>

We are hiring! -- <https://words.sdsc.edu/careers>

NBCR

it²



NIH

National Institutes of Health

NSF



U.S. DEPARTMENT OF
ENERGY

Office of
Science