# Data Science at NIH: Opportunities and Challenges

Philip E. Bourne, PhD, FACMI

Associate Director for Data Science

National Institutes of Health

## Council of Councils

September 1, 2015

# Agenda

- Drivers of change

- Scientific motivators

- The NIH response

    – The Office of Data Science

    – The Big Data to Knowledge (BD2K) initiative

- The 3-legged stool strategy

    – Infrastructure

    – Policies
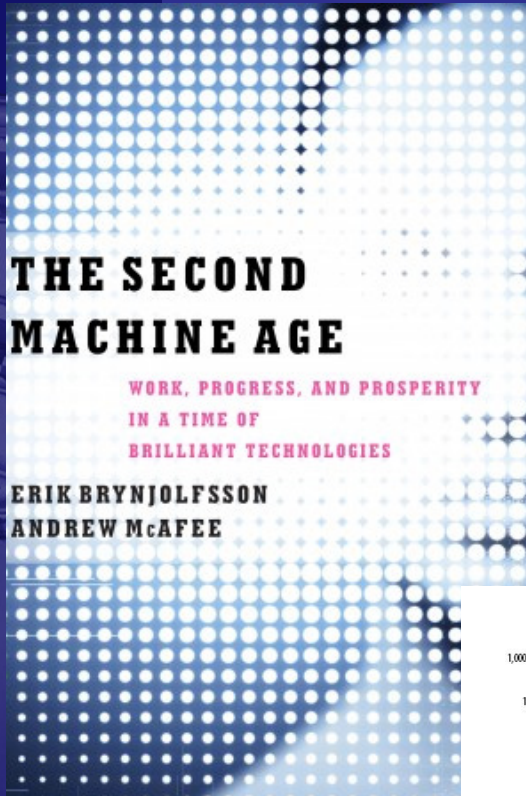
    – Communities

- Diversity

- Your thoughts here?

# Agenda

- (Drivers of change)
- Scientific motivators
- The NIH response
  - The Office of Data Science
  - The Big Data to Knowledge (BD2K) initiative
- The 3-legged stool strategy
  - Infrastructure
  - Policies
  - Communities
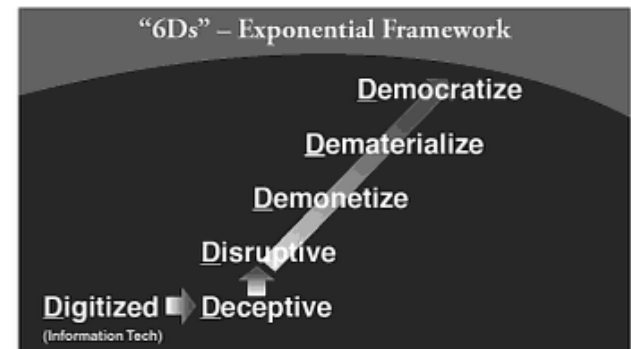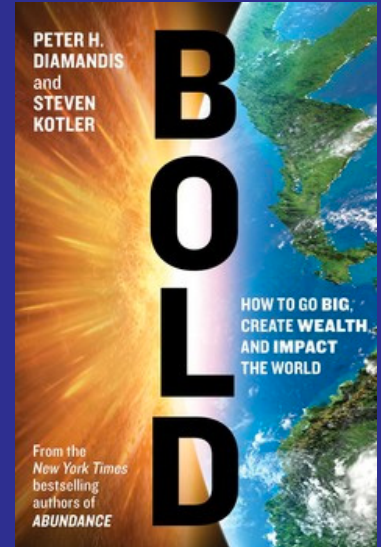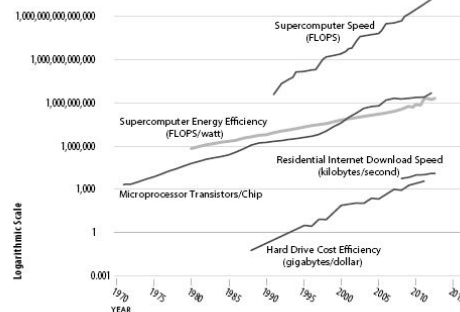- Diversity
- Your thoughts here?

# We are at a Point of Deception …

- Evidence:
  - Google car
  - 3D printers
  - Waze
  - Robotics
  - Sensors



FIGURE 3.3 The Many Dimensions of Moore's Law



"6Ds" – Exponential Framework

Democratize
Dematerialize
Demonetize
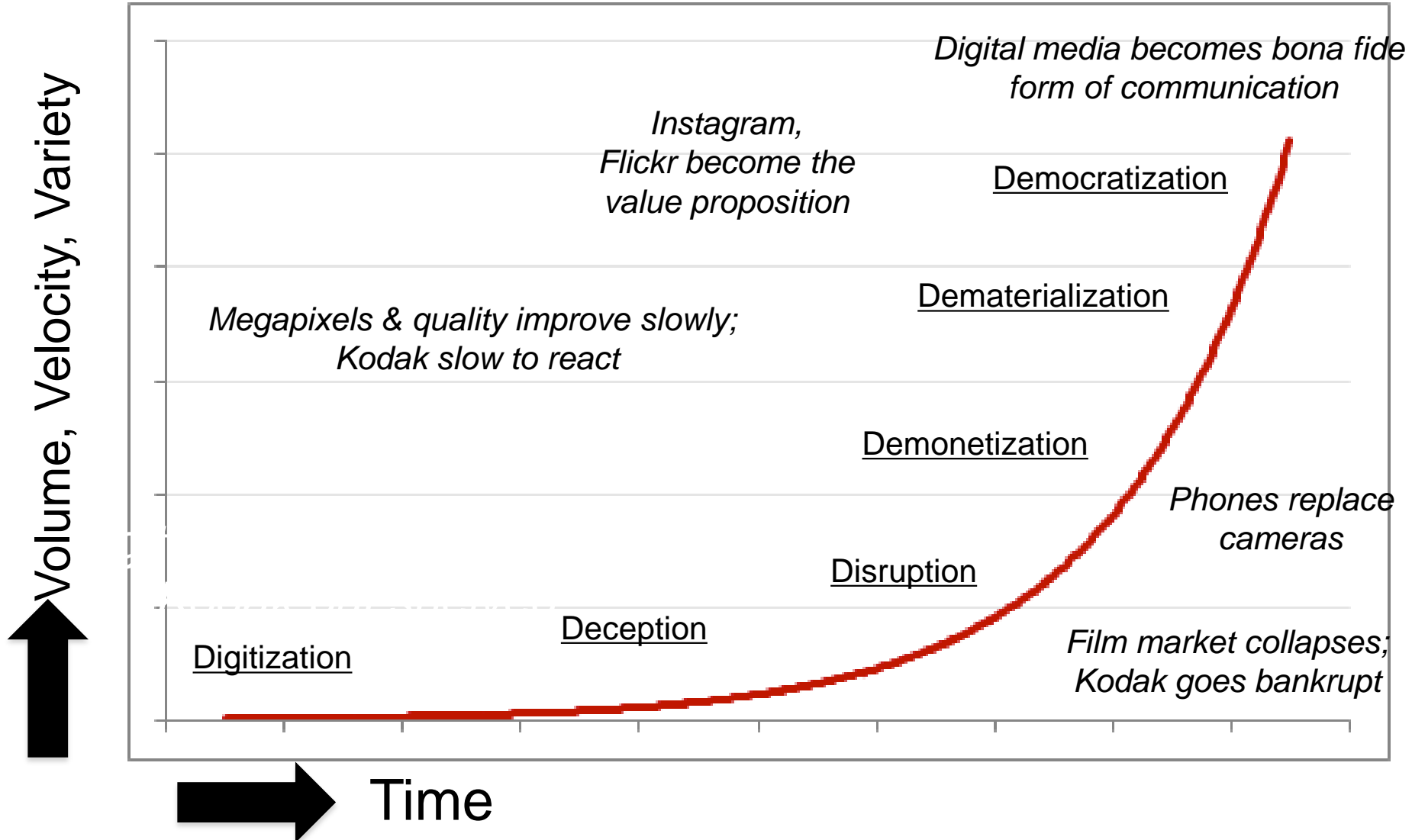Disruptive
Digitized ▶ Deceptive
(Information Tech)

The 6 Ds of Exponentials: Digitalization, Deception, Disruption, Demonetization, Dematerialization, and Democratization
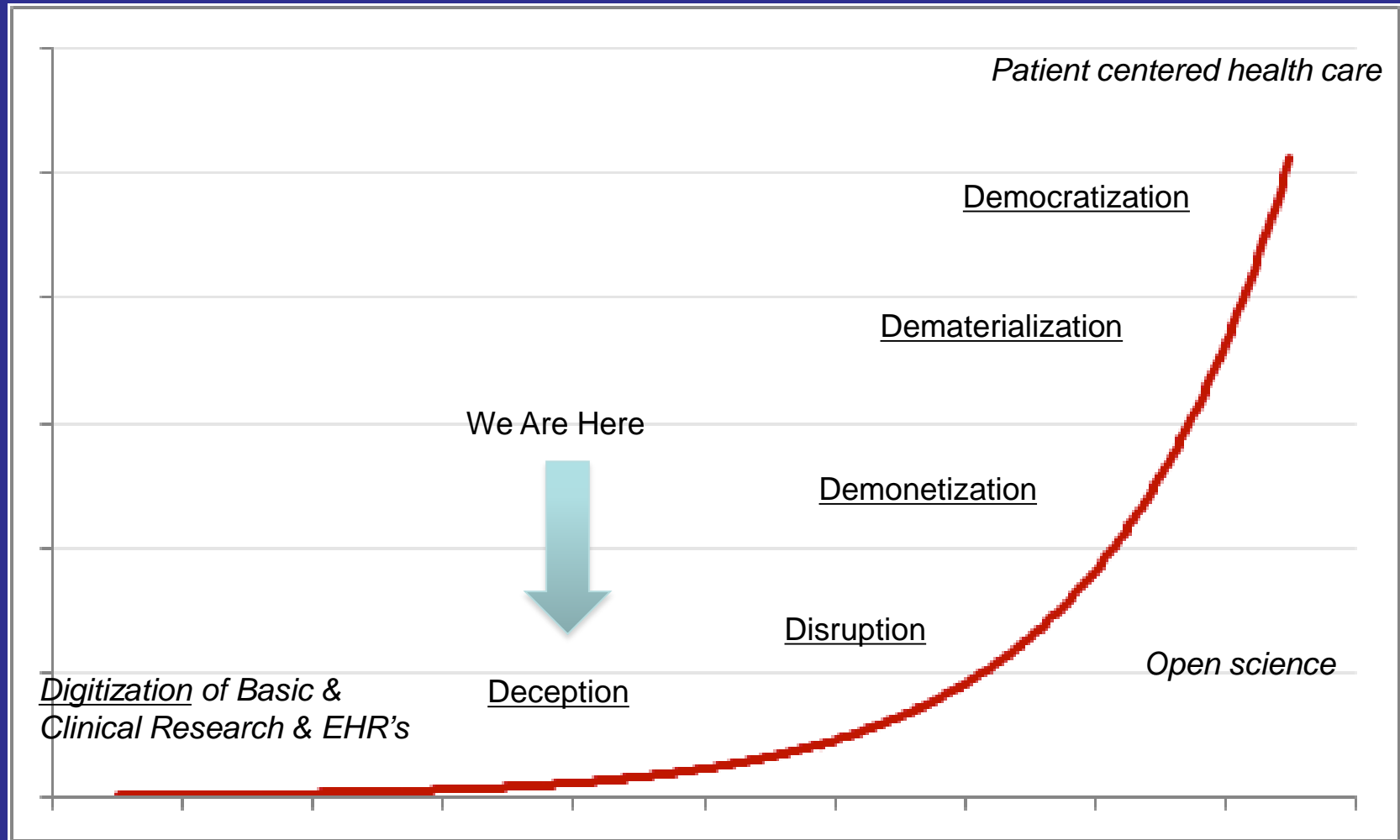
Source: Peter H. Diamandis, www.abundancehub.com

From: The Second Machine Age: Work, Progress, and Prosperity in a Time of Brilliant Technologies by Erik Brynjolfsson & Andrew McAfee
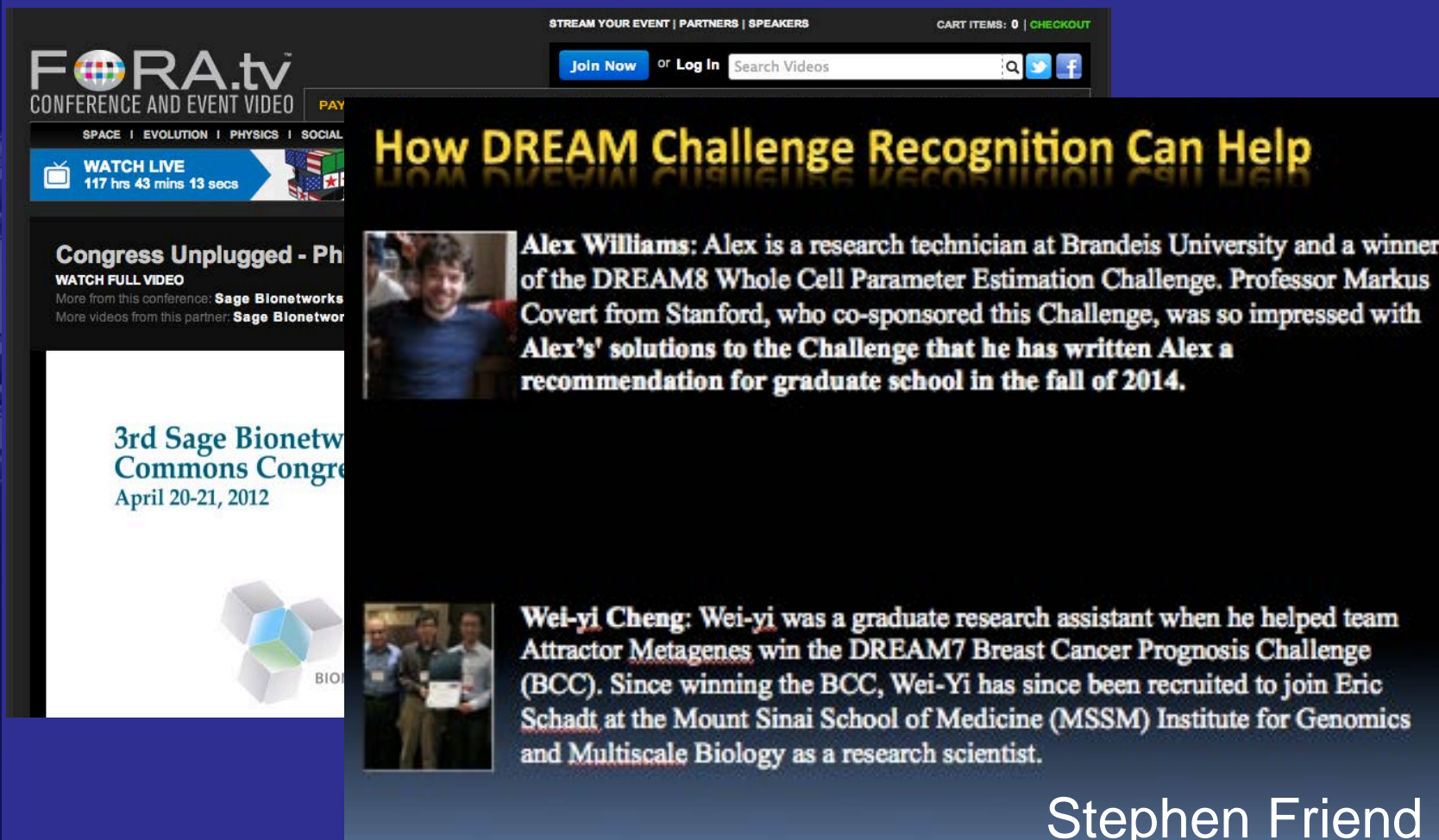
# Example - Photography

Volume, Velocity, Variety → (y-axis)

Time → (x-axis)

Digitization

Deception

Disruption

Demonetization

Dematerialization

Democratization

*Megapixels & quality improve slowly; Kodak slow to react*

*Instagram, Flickr become the value proposition*

*Digital media becomes bona fide form of communication*

*Phones replace cameras*

*Film market collapses; Kodak goes bankrupt*

# We Are At a Point of Deception
# The 6D Exponential Framework



Patient centered health care

Democratization

Dematerialization

We Are Here

Demonetization

Disruption

Open science

Deception

Digitization of Basic &
Clinical Research & EHR's

# Another Driver of Change



**How DREAM Challenge Recognition Can Help**

Alex Williams: Alex is a research technician at Brandeis University and a winner of the DREAM8 Whole Cell Parameter Estimation Challenge. Professor Markus Covert from Stanford, who co-sponsored this Challenge, was so impressed with Alex's' solutions to the Challenge that he has written Alex a recommendation for graduate school in the fall of 2014.

Wei-yi Cheng: Wei-yi was a graduate research assistant when he helped team Attractor Metagenes win the DREAM7 Breast Cancer Prognosis Challenge (BCC). Since winning the BCC, Wei-Yi has since been recruited to join Eric Schadt at the Mount Sinai School of Medicine (MSSM) Institute for Genomics and Multiscale Biology as a research scientist.

Stephen Friend

http://fora.tv/2012/04/20 Congress_Unplugged_Phil_Bourne

Let's Make Gender Diversity in Data Science a Priority Right from the Start
2015 Berman & Bourne *PLOS Biology* 13(7): e1002206

# Agenda

- Drivers of change
- Scientific motivators
- The NIH response
  - The Office of Data Science
  - The Big Data to Knowledge (BD2K) initiative
- The 3-legged stool strategy
  - Infrastructure
  - Policies
  - Communities
- Some future activities
- Your thoughts here?

NIH

"And that's why we're here today. Because something called precision medicine … gives us one of the greatest opportunities for new medical breakthroughs that we have ever seen."

President Barack Obama
January 30, 2015

# Precision Medicine Initiative

- **National Research Cohort**
  - >1 million U.S. volunteers
  - Numerous existing cohorts (many funded by NIH)
  - New volunteers
- Participants will be centrally involved in design and implementation of the cohort
- They will be able to share genomic data, lifestyle information, biological samples – all linked to their electronic health records

# An Example of That Promise: Comorbidity Network for 6.2M Danes Over 14.9 Years

# Agenda

- Drivers of change

- Scientific motivators

- The NIH response
  - The Office of Data Science
  - The Big Data to Knowledge (BD2K) initiative

- The 3-legged stool strategy
  - Infrastructure
  - Policies
  - Communities

- Some future activities

- Your thoughts here?

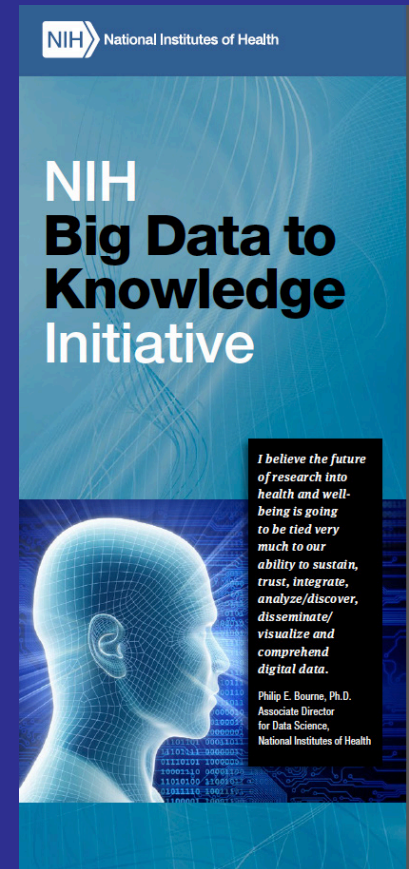# NIH Office of Data Science Mission Statement



To use data science to foster an open *digital ecosystem* that will accelerate **efficient, cost-effective** biomedical research

*to enhance health, lengthen life, and reduce illness and disability*

Goals expanded from recommendations in the June 2012 DIWG and BRWWG reports.

# The BD2K Program is Central to the Mission



Planned – Black; Available – Green

# BD2K Awards

Legend

△ Centers

◆ DDICC

⬠ LINCS

★ Training

● Secondary site (color matches parent initiative)

# Example: BD2K Center
## *Working Across Strategic Areas*

**Strategic Areas**

**Sustainability**

**Workforce Development & Diversity**

**Discovery & Innovation**

**Policy & Process**

**Leadership**

Research Objects in the Commons

Over 100 Public Lectures
Collaboration with a Minority Institution

Voxel Wide Genome Scanning
MRI standardization

Genomic Data Sharing
Policy

185 Institutions Involved

# Agenda

- Drivers of change

- Scientific motivators

- The NIH response
  - The Office of Data Science
  - The Big Data to Knowledge (BD2K) initiative

- The 3-legged stool strategy
  - Infrastructure
  - Policies
  - Communities
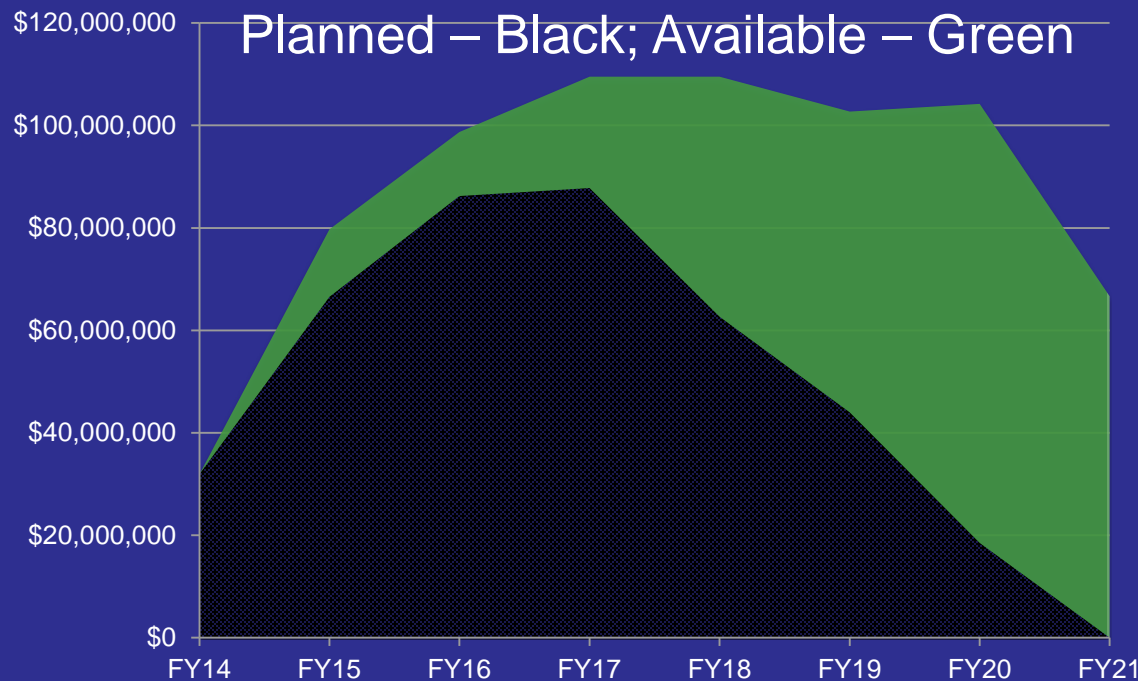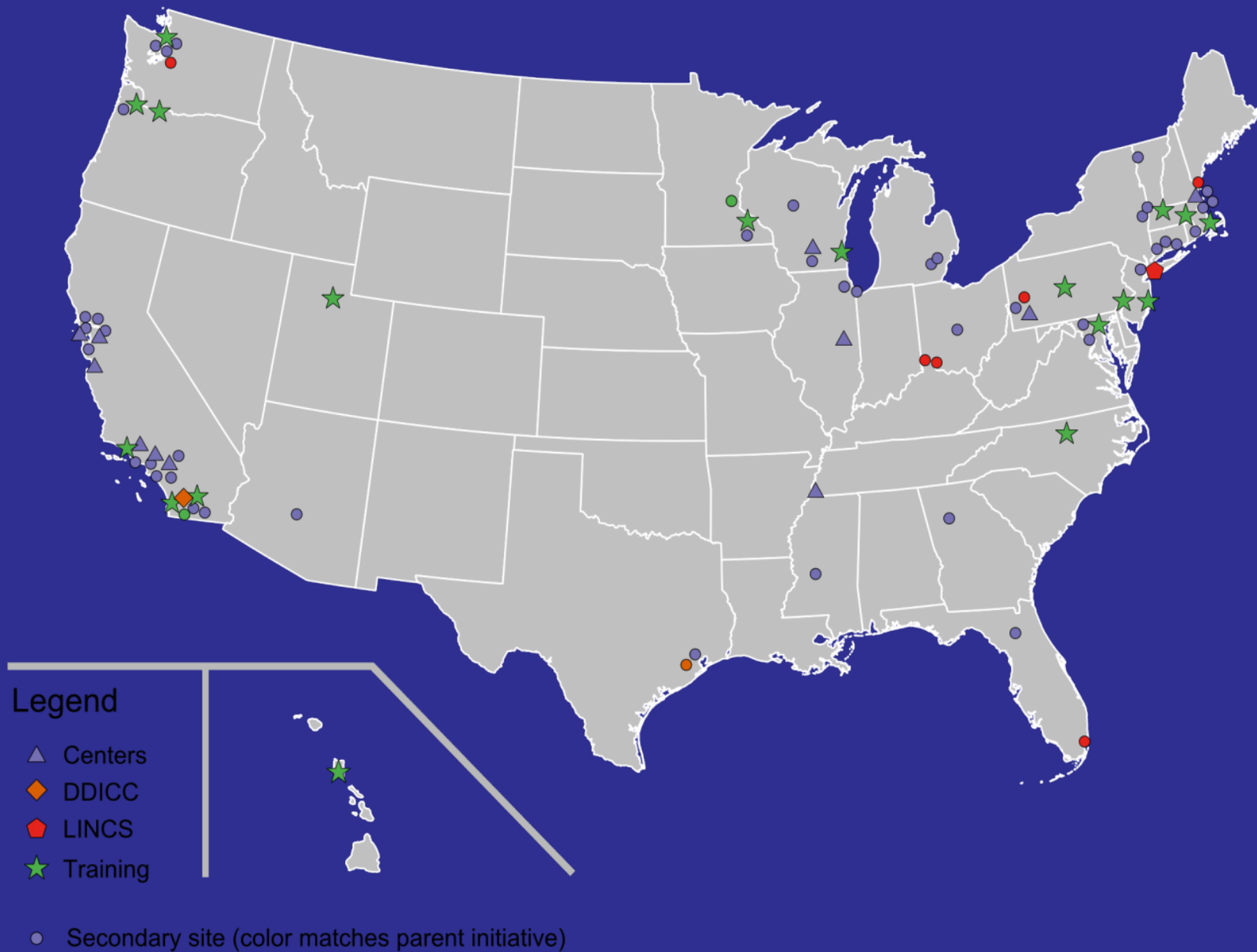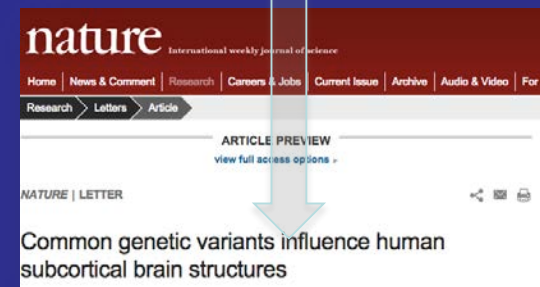
- Some future activities

- Your thoughts here?

# Elements of Our Strategy

Communities

Policies

Infrastructure

- Intersection:
  - Sustainability
  - Efficiency
  - Collaboration
  - Training

# Elements of Our Strategy

Communities

Policies

**Virtuous Research Cycle**

- Intersection:
  - Sustainability
  - Efficiency
  - Collaboration
  - Training

Infrastructure

# Consider Each Component Starting with the Infrastructure…

# Infrastructure - The Commons

# The Commons

The Commons

Digital Objects
(with UIDs)

Search
(indexed metadata)

Computing
Platform

Vivien Bonazzi
George Komatsoulis

# The Commons: Compute Platforms

The Commons
Conceptual Framework

Public Cloud Platforms

Super Computing (HPC) Platforms

Other Platforms ?

- Google, AWS (Amazon)
- Microsoft (Azure), IBM, other?

- Traditionally low access by NIH

- In house compute solutions
- Private clouds, HPC
  - Pharma
  - The Broad
  - Bionimbus

# The Commons:
## *Business Model*



[George Komatsoulis]

Communities     Policies

- Intersection:
  - Sustainability
  - Efficiency
  - Collaboration
  - Training

Infrastructure

# And Now Communities …

# Communities: Example Activities

- Visioning workshop convened 9/3/14

- Launched BD2K ($32M)

  – 12 Centers of data excellence

  – Data Discovery Index Coordination Consortium (DDICC)

  – Training awards

- First successful consortia meeting 11/3-4

- Workshops to inform future funding

  – Software indexing and discoverability

  – Gaming

Communities    Policies

- Intersection:
  - Sustainability
  - Efficiency
  - Collaboration
  - Training

Infrastructure

# And Lastly Policies …

# Policies: Now & Forthcoming

- Data Sharing
  - Genomic data sharing announced
  - Data sharing plans on all research awards
  - Data sharing plan enforcement
    - Machine readable plan
    - Repository requirements to include grant numbers



http://www.nih.gov/news/health/aug2014/od-27.htm

# Policies - Forthcoming

- Data Citation
  - Goal: legitimize data as a form of scholarship
  - Process:
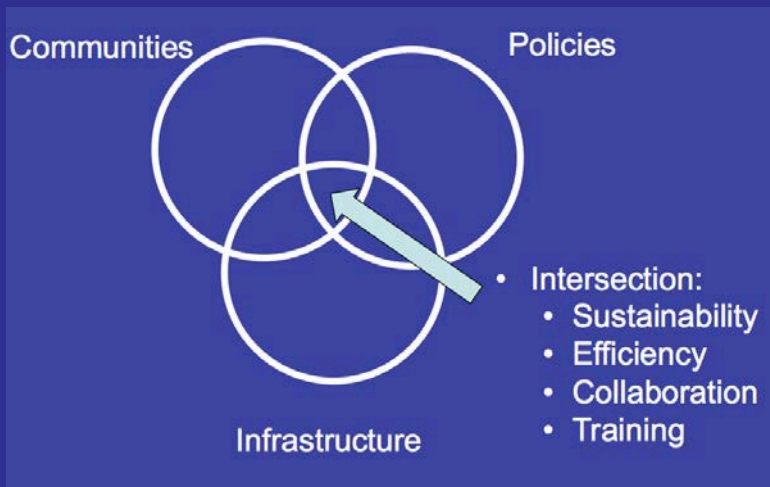    - Machine readable standard for data citation (done)
    - Endorsement of data citation for inclusion in NIH bib sketch, grants, reports, etc.
    - Example formats for human readable data citations
    - Slowly work into NLM/NCBI workflow

- dbGaP in the cloud (done!)

# Issues At The Intersection…

# Workforce Training

**Goal:** *To strengthen the ability of a diverse biomedical workforce to develop and benefit from data science*

**Strengthening a diverse biomedical workforce to utilize data science**

BD2K funding of Short Courses and Open Educational Resources

**Building a diverse workforce in biomedical data science**

BD2K Training programs and Individual Career Awards

**Discovery of Educational Resources**

BD2K Training Coordination Center

**Fostering Collaborations**

BD2K Training Coordination Center, NSF/NIH IDEAs Lab

**Expanding NIH Data Science Workforce Development Center**

Local courses, e.g. Software Carpentry

NIH

# Agenda

- Drivers of change

- Scientific motivators

- The NIH response
    - The Office of Data Science
    - The Big Data to Knowledge (BD2K) initiative

- The 3-legged stool strategy
    - Infrastructure
    - Policies
    - Communities

- Some future activities

- Your thoughts here?

# New BD2K Workshops

- NIH Common Data Elements (CDE) Workshop
  - in conjunction with NCI, NLM, and BMIC
- Academic deans and data science career paths workshop
  - in conjunction with NSF
- Data Science of Citizen Science
  - Possibly in conjunction with the Heart of Data Science (Ping) Center
- National Academies (CATS) workshop on big data inference
  - in conjunction with NSF/CISE
- National Academies (CSTB) workshop on data science curriculum development
  - in conjunction with NSF/CISE

# New BD2K Activities

- **Reference Datasets**
  - Will move important, FAIR, digital resources into the cloud to support increased access and utility.
  - Will release an RFI, to inform details of an FOA.
  - Will follow up and formalize activities started in the Commons Administrative supplements
- **Sustainability of Data Repositories**
  - BD2K sustainability group is doing a financial and portfolio analysis of digital data repositories across NIH, and has been studying current sustainability approaches of such repositories.
  - BD2K is working in conjunction with NIGMS, NHGRI, and BISTI to develop FOA on this topic.
  - Will follow up on initial activities of the interoperability supplements.
- **Software Hardening Resource**
  - BD2K is standing up a group to develop a proposal to ensure useful software can develop from useful academic grade prototypes to more robust, commons-compliant tools.

*FAIR = Findable, Accessible, Interoperable, Reusable*

# Agenda

- Drivers of change

- Scientific motivators

- The NIH response
  - The Office of Data Science
  - The Big Data to Knowledge (BD2K) initiative

- The 3-legged stool strategy

  - Infrastructure

  - Policies

  - Communities

- Diversity

- Your thoughts here?

*I not only use all the brains I have, but all I can borrow.*

**— Woodrow Wilson**

# NIH...

*philip.bourne@nih.gov*

# Turning Discovery Into Health