#### Data Driven Science at NIH: A Conversation

Philip E. Bourne, PhD, FACMI Associate Director for Data Science

National Institutes of Health

**Federal Demonstration Partnership** May 11, 2015, Washington DC





### What Drives Our Strategic Thinking?

- Be Prepared Responding to take advantage of the opportunities offered by a major disruption in the biomedical research enterprise arising through digitization and exponential growth
- Accelerating discovery during this time of disruptive development



 Continually catalyzing a cultural shift towards a more analytical enterprise while managing expectations



#### Let Me Give You 4 Examples of What Drives Us ...





#### 1. We are at a Point of Deception ...



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







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





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

## **Example - Photography**



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



NIH

#### 2. Democratization Will Follow The Story of Meredith



Congress Unplug WATCH FULL VIDEO More from this conference: Sage More videos from this partner: Sage

#### How DREAM Challenge Recognition Can Help

**3rd Sage B Commons** April 20-21, 201 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**

## 3. Disruption Can Occur





#### Raise standards for preclinical cancer research

C. Glenn Begley and Lee M. Ellis propose how methods, publications and incentives must change if patients are to benefit.

Efforts over the past decade to characterize the genetic alterations in human cancers have led to a better understanding of molecular drivers of this

trials in oncology have the highest failure rate compared with other therapeutic areas. Given the high unnet need in oncology, it is understandable that barriers to clinical is understandable that barriers to clinical

# 47/53 "landmark" publications could not be replicated

IBegley, Ellis Nature, 483, 2012]

#### Must try harder

Too many sloppy mistakes are creeping into scientific papers. Lab heads must lead the data - and at themselves.

## **Error prone**

Biologists must realize the pitfalls of work or massive amounts of data.

#### If a job is worth doing, it is worth doing twice

Researchers and funding agencies need to put a premium on ensuring that results are reproducible, argues Jonathan F. Russell.

The case for open computer programs

Six red flags for suspect work

C. Glenn Begley explains how to recognize the preclinical papers in which the data won't stand up.

Know when your numbers are significant

#### [Carole Goble]

#### 4. Demonetization, Democratization?



View larger 🔀 (PDF - 163KB)



"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**

Vision: Build a broad research program to encourage creative approaches to precision medicine, test them rigorously, and, ultimately, use them to build the evidence base needed to guide clinical practice.

Near Term: apply the tenets of precision medicine to a major health threat – cancer

Longer Term: generate the knowledge base necessary to move precision medicine into virtually all areas of health and disease





#### **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





#### Jensen et al 2014 Nat Comm 5:4022

Office of Biomedical 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.

## **Overall Goals by 2020**

- Enable major scientific discovery through the BD2K initiative
  - Establish and provide evidence of a more sustainable,
    efficient and productive data science ecosystem
    both internal and external to NIH
- Establish and provide evidence of a well-trained and diverse workforce able to use and develop biomedical data science tools and methods
- Build upon NIH's leadership and reputation in data science





# The BD2K Program is Central to the Mission





NIH Big Data to Knowledge Initiative

NIIH National Institutes of Health

I believe the future of research into health and wellbeing is going to be tied very much to our ability to sustain. trust, integrate, analyze/discover disseminate/ visualize and comprehend digital data. Philip E. Bourne, Ph.D Associate Director or Data Science

#### **Elements of The Digital Enterprise**



#### **Elements of The Digital Enterprise** Policies Communities Virtuous Research Cycle Intersection: • Sustainability Efficiency • Collaboration • Training Infrastructure

NIH

#### Consider an example...







na	nature International weekly journal of science						
Home	News & Com	ment Researc	h Careers & Jobs	Current Issue	Archive	Audio & Video	For A
Researc	h > Letters	Article					
ARTICLE PREVIEW							
			view full access op	tions <b>&gt;</b>			
NATURI	:   LETTER					< ⊠	8

Common genetic variants influence human subcortical brain structures

 Big Data: The study involved MRI images & GWAS data from over 30,000 people

ENIGM

- Collaboration: Data came from many different sights affiliated with the ENIGMA consortium
- Methods: To homogenize data from different sites, the group designed standardized protocols for image analysis, quality assessment, genetic imputation, and association

- Found five novel genetic variants
- Results provided insight into the variability of brain development, and may be applied to study of neuropsychiatric dysfunction



- Community Enigma, BD2K
- Policy
  - Improved consent methods
  - Cloud accessibility for human subjects data
  - Trusted partners
  - Data sharing



NIH

- Infrastructure
  - Standards, compute resources, software



## **Communities: Thus Far**

- 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



![](_page_21_Picture_12.jpeg)

![](_page_22_Picture_0.jpeg)

## **Communities: 2015 Activities**

- New FOAs with outreach to new communities – math, stats, comp science etc.
- Work with e.g GA4GH, RDA, FORCE11, NDS ....
- IDEAS lab with NSF
- Competition with international funders
- Software carpentry, hackathons, Pi Day

![](_page_22_Picture_7.jpeg)

## **Communities: Questions?**

Societies of the modern age?

How to enable these groups?

How to marry the funding of individuals with the funding of communities?

![](_page_23_Picture_4.jpeg)

## **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

![](_page_24_Picture_7.jpeg)

![](_page_24_Picture_8.jpeg)

![](_page_24_Picture_9.jpeg)

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

![](_page_25_Picture_8.jpeg)

dbGaP in the cloud (done!)

![](_page_25_Picture_10.jpeg)

![](_page_26_Figure_0.jpeg)

#### **The Commons**

![](_page_27_Figure_1.jpeg)

Vivien Bonazzi George Komatsoulis

#### **The Commons: Compute Platforms**

![](_page_28_Figure_1.jpeg)

- The Broad
- Bionimbus

#### The Commons: Business Model

![](_page_29_Figure_1.jpeg)

![](_page_29_Picture_2.jpeg)

![](_page_29_Picture_3.jpeg)

#### [George Komatsoulis]

#### Infrastructure: Standards

- 2013 Workshop on Frameworks for Community-Based Standards
- August 2014 Input on Information Resources for Data-Related Standards Widely Used in Biomedical Science – 30 responses
- Feb 2015 Workshop Community-based Data and Metadata Standards

![](_page_30_Picture_4.jpeg)

Internal CDE Registry project

#### **Elements of The Digital Enterprise**

![](_page_31_Figure_1.jpeg)

#### **Elements of The Digital Enterprise**

![](_page_32_Figure_1.jpeg)

#### **Sustainability 101**

![](_page_33_Figure_1.jpeg)

![](_page_33_Figure_2.jpeg)

Growth of Biological Databases

NIH

## **Workforce Training**

![](_page_34_Picture_1.jpeg)

![](_page_34_Picture_2.jpeg)

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

![](_page_35_Picture_11.jpeg)

![](_page_35_Picture_12.jpeg)

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

## – Woodrow Wilson

![](_page_36_Picture_2.jpeg)

![](_page_36_Picture_3.jpeg)

#### Associate Di

Scientific Data Council

Echefon

#### **Data Science**

External Advisory Board

Cc

ion

#### Programmatic Theme

Deliverable

Exa

## The Biomedical Research Digital Enterprise

![](_page_38_Picture_0.jpeg)

#### philip.bourne@nih.gov

# Turning Discovery Into Health

![](_page_38_Picture_3.jpeg)

![](_page_38_Picture_4.jpeg)

![](_page_38_Picture_5.jpeg)

![](_page_38_Picture_6.jpeg)