# Sustaining the 'Big Data' Ecosystem

**Eric Green, M.D., Ph.D.**
**National Human Genome Research Institute**
**National Institutes of Health, USA**

# The Largest Bottleneck In Biomedical Research…
## …Pick Your Favorite Data-Related Metaphor!

# TECHNOLOGY FEATURE

# THE BIG CHALLENGES OF BIG DATA

*As they grapple with increasingly large data sets,*
*biologists and computer scientists uncork new bottlenecks.*

# For My Talk…

- ➢ **Goals:**
    - ▪ **Help to provide context for the workshop**
    - ▪ **Begin to frame the discussion**

- ➢ **Two Overlapping Roles:**
    - ▪ **Represent the U.S. National Institutes of Health (NIH)**
    - ▪ **Director, National Human Genome Research Institute (NHGRI)**

- ➢ **Articulate themes from the 2015 Bourne/Lorsch/Green *Nature* Perspective**

# Challenges Of Sustaining Data Resources



*Nature* (2015)

# Overarching Realities



1. We are victims of our own success

2. Genomics is a 'poster child' for the problem

3. But other data types are quickly becoming similarly problematic

# Myriad Data Types



**Genomic**

**Other 'Omic**

**Imaging**

**Phenotypic**

**Exposure**

**Clinical**

# Overarching Realities



1. We are victims of our own success

2. Genomics is a 'poster child' for the problem

3. But other data types are quickly becoming similarly problematic

4. The 'trends' are particularly troubling

# Budgets Are Mostly Flat (or Worse)... Data Growth Is Anything But Flat!



Note: The 3.7 % Real Annual Growth is based on real compound annual growth between 1971 and 2000. Dollar values are adjusted to 2012 dollars using the Biomedical Research and Development Price Index (BRDPI), http://officeofbudget.od.nih.gov/gbiPriceIndexes.html.
Source: NIH Office of Extramural Research and Office of Budget source data (January 19, 2016)



# Current models for data sustainability do not scale

# Big Data & Data Sustainability



**Analyzing the problem**

**Self-reflection**

# Many Aspects Of The Problem
## 'Slip Between The Cracks'

# Relevant Working Group Reports



ADVISORY COMMITTEE TO THE DIRECTOR

National Institutes of Health

**Data and Informatics Working Group**

Draft Report to
The Advisory Committee to the Director

June 15, 2012

National Institutes of Health
**Advisory Committee to the Director**

National Library of Medicine (NLM) Working Group

FINAL REPORT – JUNE 11, 2015

**MEMBERS:** Eric Green (co-chair), Harlan Krumholz (co-chair), Russ Altman, Howard Bauchner, Deborah Brooks, Doug Fridsma, 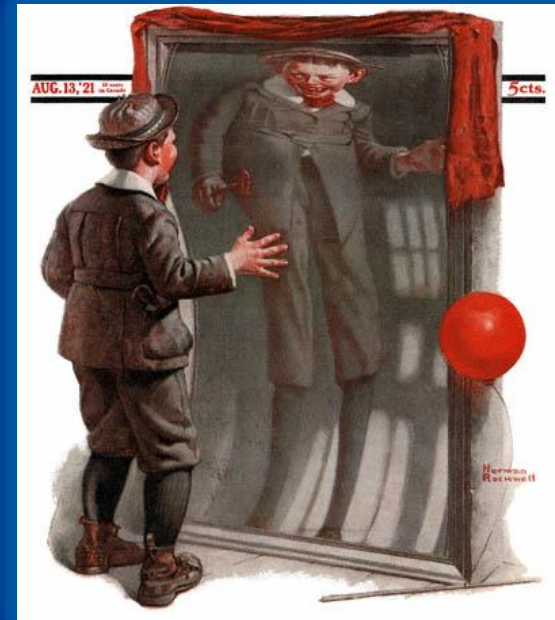Steven Goodman, Eric Horvitz, Trudy MacKay, Alexa McCray, Chris Shaffer, David Van Essen, Joanne Waldstreicher, James Williams, II, Kathy Hudson (ex officio), Lyric Jorgenson (executive secretary) *(titles and affiliations listed in Appendix A)*

EXECUTIVE SUMMARY

The NIH Director charged the National Library of Medicine (NLM) Working Group, hereafter referred to as the Working Group, with articulating a strategic vision for NLM to ensure that NLM remains an international leader in biomedical and health information. Over the course of five months of deliberations, the Working Group reviewed numerous documents and reports pertaining to NLM's mission and activities, consulted with NLM leadership and staff, and solicited public comments and suggestions. The Working Group recognizes that NLM has an important opportunity to play a key leadership role in one of the most exciting periods of biomedical history: data science is increasing rapidly, computational power is expanding at a breathtaking pace, the breadth and depth of digital health data are undergoing unprecedented and accelerating growth, a movement towards more interdisciplinary work and team science continues to gain momentum, a broad commitment to open science is becoming increasingly adopted, and the demand for services to support an ever more engaged and informed public is expanding. To leverage these historic changes, the Working Group, with respect for the outstanding history of NLM and its potential for the future, formulated a series of recommendations to guide the future of NLM:

**acd.od.nih.gov/diwg.htm**        **acd.od.nih.gov/nlm.htm**

# Data Science @ NIH





**Valentina Di Francesco
NHGRI, NIH**

# Model Organism Databases (MODs)



**Valentina Di Francesco
NHGRI, NIH**

# Emerging Themes from NIH Experience

➤ **At a pivotal point:**
   **Risk failing to capitalize on technology advances**
   **Failure to act would be devastating**

➤ **Must implement new models to ensure a sustainable infrastructure for preserving, sharing, analyzing, and integrating data**

➤ **Cultural changes must be part of the solution**

➤ **Long-term commitment is required**

# Challenges Of Sustaining Data Resources

## Sustaining the big-data ecosystem
*Organizing and accessing biomedical big data will require quite different business models,* say **Philip E. Bourne**, **Jon R. Lorsch** *and* **Eric D. Green**.

THE RESEARCH COMMUNITY MUST FIND MORE EFFICIENT MODELS FOR STORING, ORGANIZING AND ACCESSING BIOMEDICAL DATA.

- **Understanding Usage**

- **Fair and Efficient**

- **Business Models**

- **Common Ground**

- **Uniting Funders**

# What Do We Need To Understand?

➢ **What is the relative value of data at different points in time?**

➢ **Where are redundancies?**

➢ **How to better match data curation to usage?**

➢ **How to make more informed decisions about changing data resources (e.g., merging, closing, expanding, etc.)?**

➢ **What is the nature of the major cost drivers?**

# Who Is Responsible For Sustaining Data?

- ➢ **Indexed to data production vs. data usage?**

- ➢ **What component(s) within a funding agency?**

- ➢ **Government vs. other funders vs. industry?**

- ➢ **Individual countries vs. international group(s)?**

**Tackling the *international* dimensions of a problem that the NIH has been intensely working on for ~5 years (…and likely will be addressing for the next ~100 years!)**

# What We Need…

➢ **To recognize that:**

- ▪ **This is a significant and growing problem**

- ▪ **This is <u>not</u> the problem of any one country/agency**

- ▪ **This is <u>not</u> just the funders' problem**

- ▪ **The data ecosystem is complicated and changing**

- ▪ **The answer cannot be 'provide more $$$'**

➢ **A robust 'scoping out' of the problem**

➢ **A willingness to break down traditional boundaries and silos to create new collaborative solutions**

➢ **An openness to new 'new business models'**

➢ **A new framework for <u>international</u> solutions**

➢ **Most important … get to work!**