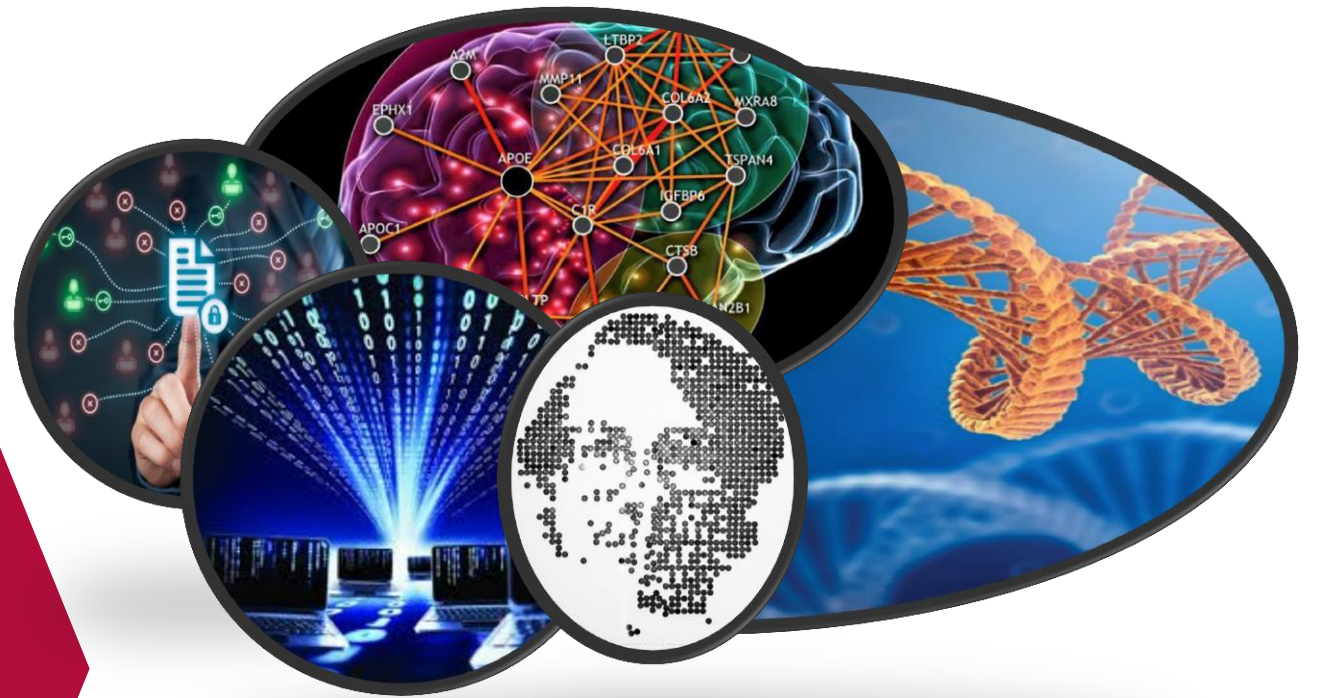# YOU and Big Data in Data Wonderland: It's NOT a Data Jabberwocky

*Vivian OTA WANG, Ph.D., CGC, FACMG*

*Deputy Director, Office of Data Sharing*
*Center for Biomedical Informatics &*
*Information Technology (CBIIT)*
*National Cancer Institute-NIH*

*Professional Development Workshop and Mock Review*
*03 June 2019*

NIH NATIONAL CANCER INSTITUTE

# Agenda

## Data Wonderland

*Big Data*
*Human Rights*
*Open Science*

## Jabberwocky

*The Policies*
*Office of Data Sharing*

## Science

*The Data*
*The People*
*The Ethics*

NIH NATIONAL CANCER INSTITUTE

# Data Wonderland

Big Data
Open Science
Human Rights

Science

VOLUME 361, ISSUE 6400
26th July 2018

Big Data in Psychology

Special issue of Psychological Methods

Vol. 21, No. 4, December 2016

**Item #:** 2272104

**ISBN:** 978-1-4338-9024-6

**Format:** Hard copy

**Other Format:** PDF

NIH National Cancer Institute

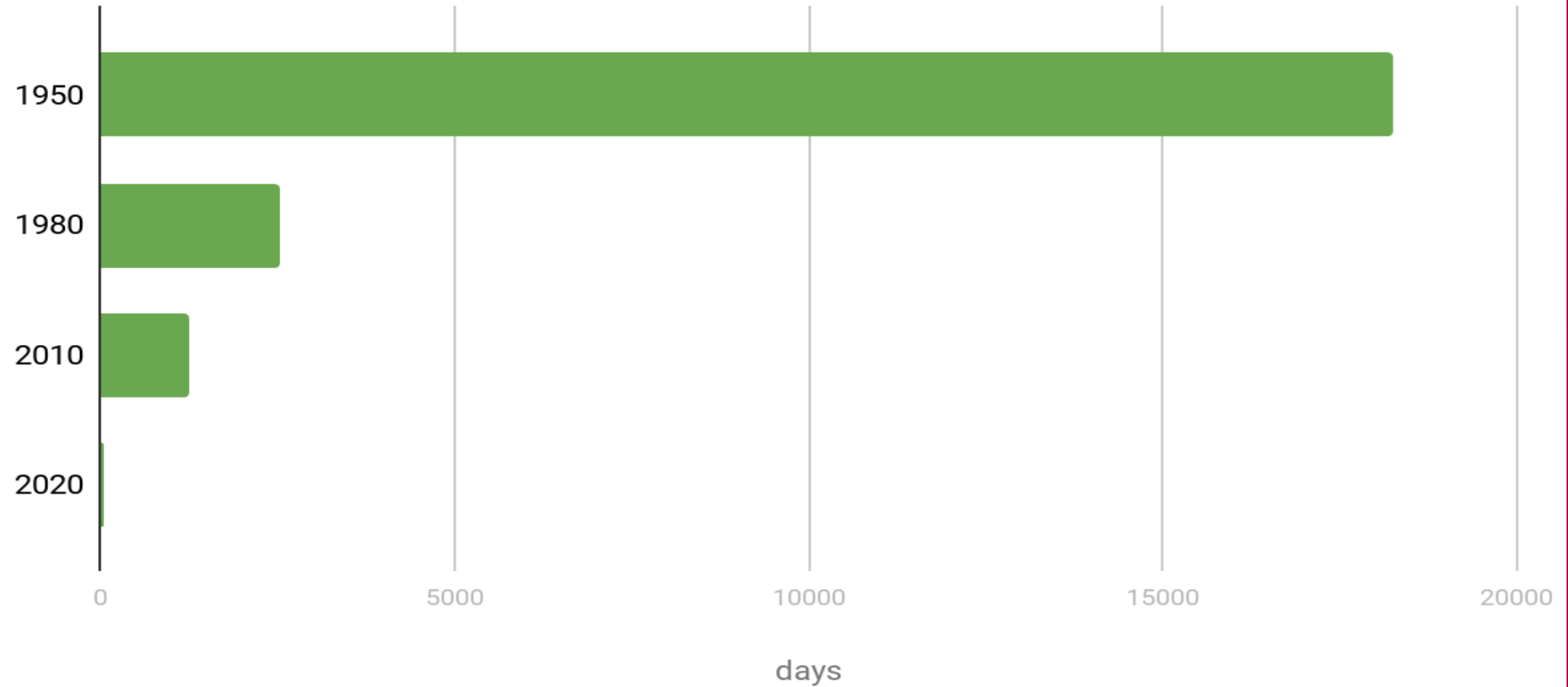# Big Data: Human Rights and the Democratization of Knowledge

# Big Data



- Large amounts of data and data types
  - *Mobile devices, tracking systems, RFID, sensor networks, social networks, Internet searches, automated record keeping, video archives, e-commerce*
- Secondary analyses of primary and derived data
- Identify trends
- Improve research quality

# Big Data



Doubling Time of Health Knowledge

# Open Science



The National Academies of
SCIENCES · ENGINEERING · MEDICINE

CONSENSUS STUDY REPORT

OPEN SCIENCE BY DESIGN

Realizing a Vision for 21st Century Research

## *DATA SHARING AND INNOVATION*

- Open access
  - Accessible research & data to all levels of society *(e.g., amateurs, citizen scientists, and professionals)*
- Open data
- Open sources

National Academies of Sciences, Engineering, and Medicine. 2018. *Open Science by Design: Realizing a Vision for 21st Century Research*. Washington, DC: The National Academies Press. https://doi.org/10.17226/25116.

NIH NATIONAL CANCER INSTITUTE

# Open Science and Data Sharing

Facilitates innovation of research tools and methods



NIH | NATIONAL CANCER INSTITUTE

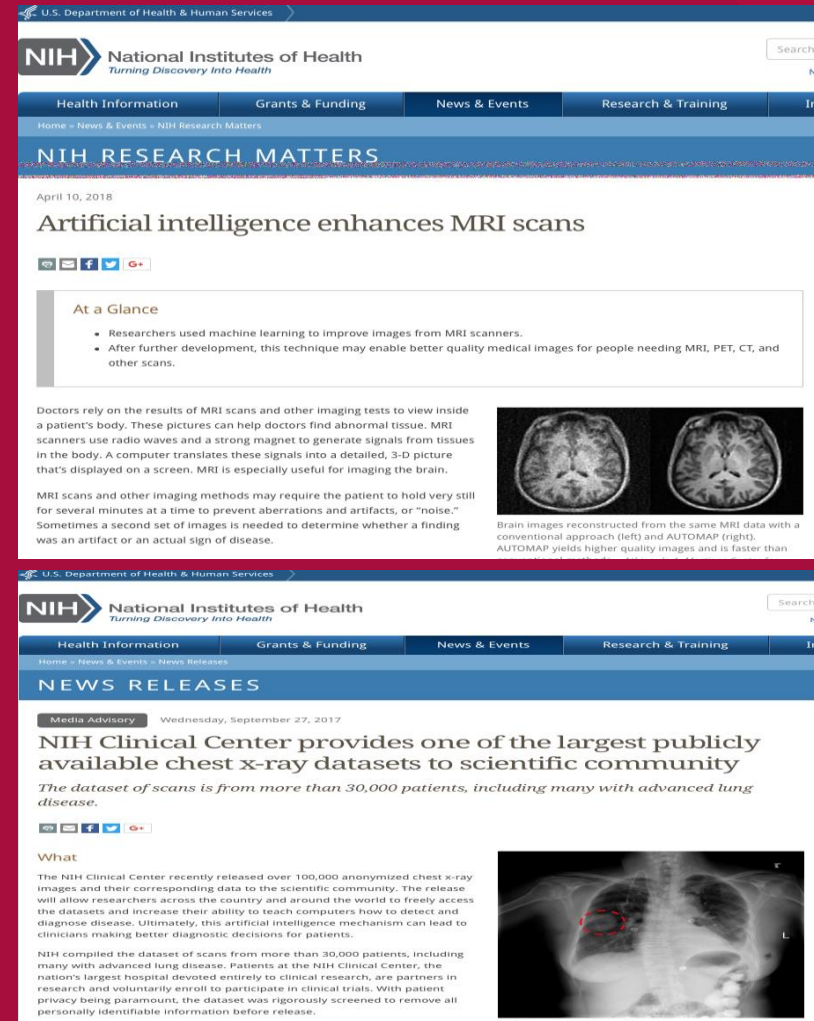Courtesy: (Adapted from Paltoo, Harris, Luetkemeier, OSP-Office of Science Policy – NIH)

# Open Science and Data Sharing

Facilitates innovation of research tools and methods

- Increases statistical power

NIH NATIONAL CANCER INSTITUTE

# Open Science and Data Sharing

- Facilitates innovation of research tools and methods

- Increases statistical power

- Improves research quality through validation and replication

# Open Science and Data Sharing

Increases scientific value and analyses by enabling data from multiple studies to be combined and explored

# Open Science and Data Sharing

- Increases scientific value and analyses by enabling data from multiple studies to be combined and explored

- Increases scale of studies, # publications, and types of scientists from a broader range of disciplines
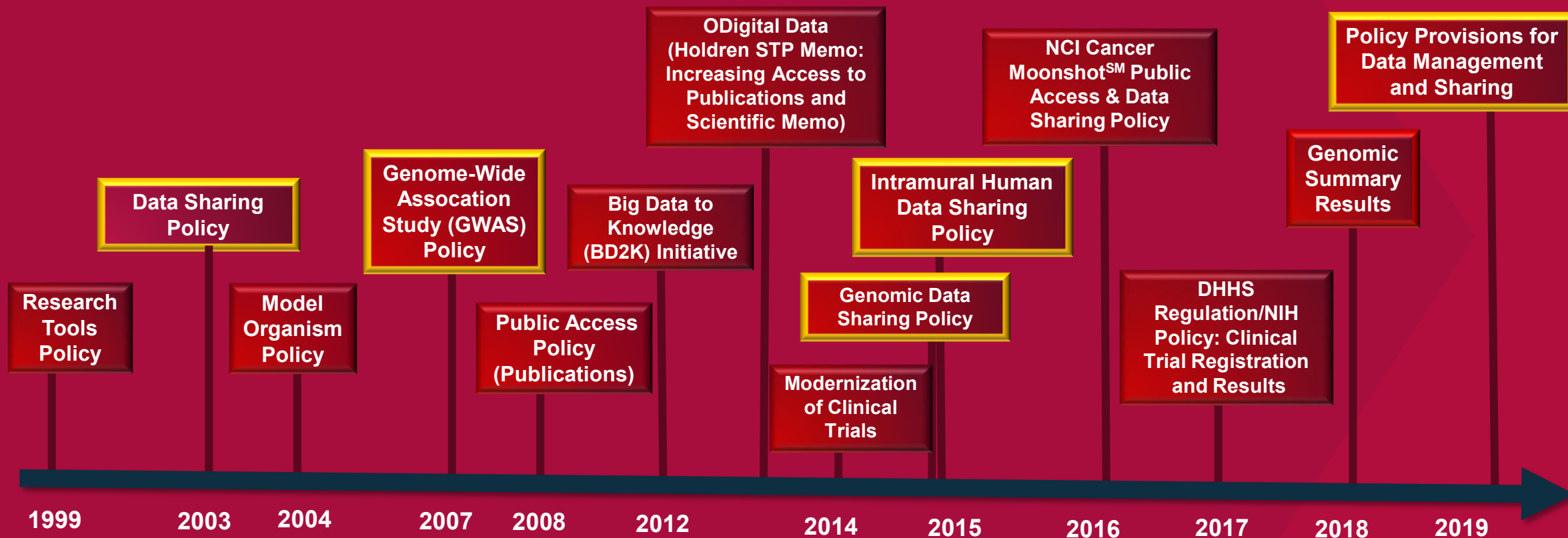
NIH NATIONAL CANCER INSTITUTE

# Open Science and Data Sharing



- Biology and Medicine are now *data intensive enterprises*

- *Rapidly changing* scale

- Technology, *data computing* and *information technology (IT)* are *pervasive* in the *lab, clinic, and home*
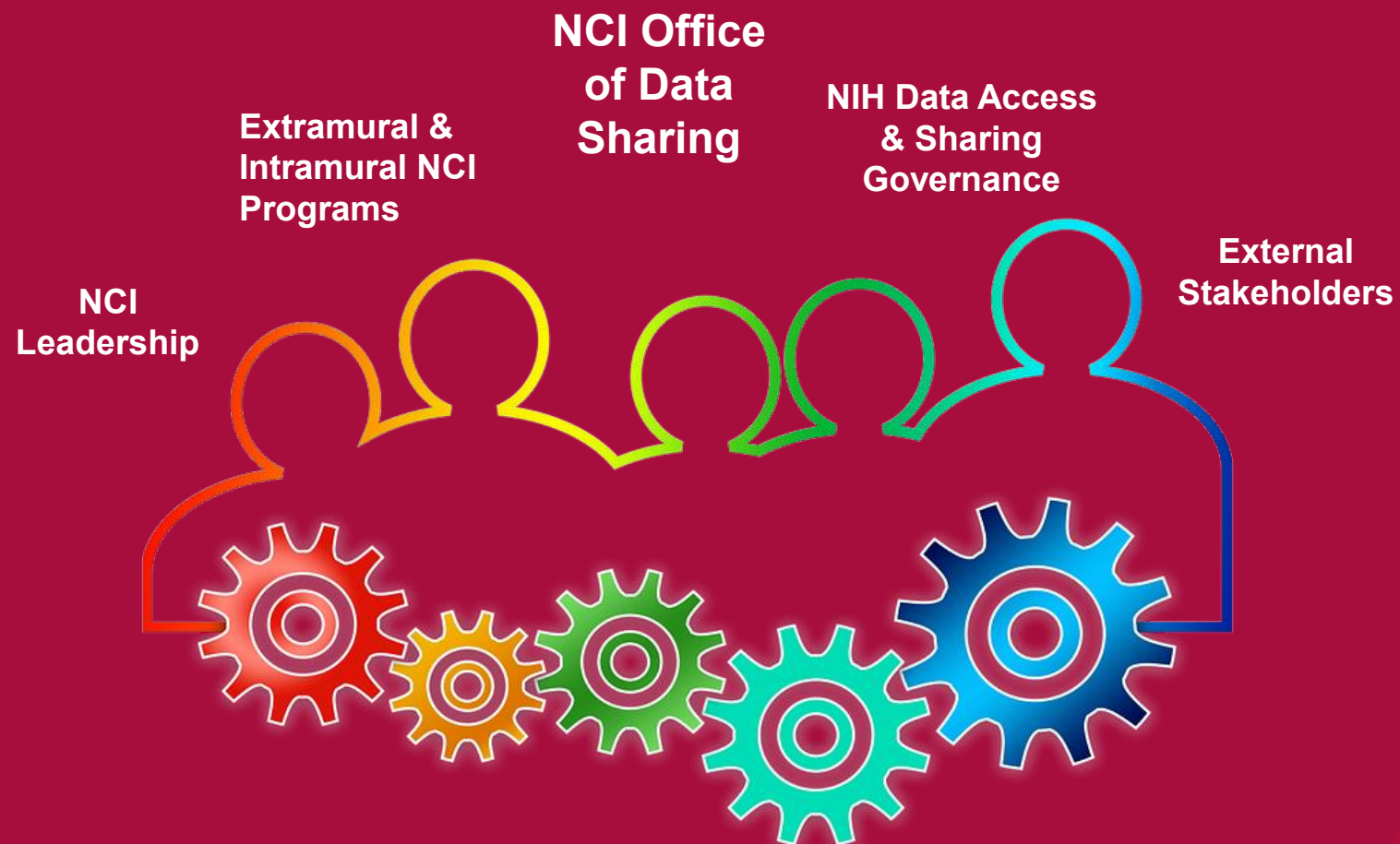
# The Jabberwocky

Policies
NCI Office of Data Sharing

# NIH Data Sharing Policies



- **Research Tools Policy** — 1999
- **Data Sharing Policy** — 2003
- **Model Organism Policy** — 2004
- **Genome-Wide Association Study (GWAS) Policy** — 2007
- **Public Access Policy (Publications)** — 2008
- **Big Data to Knowledge (BD2K) Initiative** — 2012
- **ODigital Data (Holdren STP Memo: Increasing Access to Publications and Scientific Memo)** — 2012
- **Modernization of Clinical Trials** — 2014
- **Genomic Data Sharing Policy** — 2015
- **Intramural Human Data Sharing Policy** — 2015
- **NCI Cancer Moonshot[SM] Public Access & Data Sharing Policy** — 2016
- **DHHS Regulation/NIH Policy: Clinical Trial Registration and Results** — 2017
- **Genomic Summary Results** — 2018
- **Policy Provisions for Data Management and Sharing** — 2019

NIH NATIONAL CANCER INSTITUTE

(Courtesy: (Adapted from D.Paltoo and L.L. Rodriguez)

(Courtesy: (Adapted from J Guidry Auvil, NCI – NIH)

# NCI Office of Data Sharing
## *nciofficeofdatasharing@nih.gov*

Provide *leadership* and *guidance* to enhance data sharing for NCI and the cancer research community.

Guide NCI approach to *implementation and interpretation* of NIH and NCI data management and sharing policies.

*Coordinate registration, submission, and access procedures* for NCI datasets/repositories.

Advise on considerations for *ethical and minority and health disparity issues* related to data access and sharing for the cancer community.

Encourage *participation* in major data sharing initiatives.

*Create data sharing resources* to inform and guide the cancer communities.

NIH NATIONAL CANCER INSTITUTE

(Courtesy: (Adapted from J Guidry Auvil, NCI – NIH)

# Total Number of Actual and Projected NCI dbGaP Data Access Requests By Year

**# of DATA ACCESS REQUESTS**

40000
35000
30000
25000
20000
15000
10000
5000
0

NCI DACs

221
896
1949
2806
4873
6030
13110
17615
19381
23558
30653
35314

2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018

**YEAR**

NIH NATIONAL CANCER INSTITUTE

(Courtesy: Adapted from J Guidry Auvil, NCI – NIH)

# Total Number of Actual and Projected NCI dbGaP Data Access Requests By Year



(Courtesy: Adapted from J Guidry Auvil, NCI – NIH)

# Percentage of Data Access Requests By NIH I/Cs 01/01218 –05/31/2019

# NCI Office of Data Sharing*



**Pre-centralization to Phase I**
(6/25/17 - 6/22/18)
- 24.2
- 53.4

**Phase I to Phase II**
(6/25/18 – 8/17/18)
- 5.1
- 26.4

**Post Phase III**
(8/18/18 – 5/31/19)
- 0.3
- 90.7

**Phase I**
- Move TCGA from NHGRI to NCI
- Completed June 25, 2018

**Phase II**
- Combine iNCI +TCGA DACa
- Completed June 25, 2018

**Phase III**
- Incorporate with eNCI DAC
- Completed August 20, 2018
- KidsFirst DAC is separate

Time (Days)
Data Access Requests (x100)

* NCI DAC receives ~30% of DARs to dbGaP; results include significant efficiency in DAR processing times and eliminating 1000+ DAR backlog

# The Science

## The Data
## The People
## The Ethics

Unstructured Info → Machine Learning → Interface

Machine learning technology analyzes millions of unstructured sources in real-time and...

selects and synthesizes that knowledge so you can...

see trends easily & quickly.

https://www.cbinsights.com/research/team-blog/data-network-effects/

# Data: Variety, Volume, Velocity, and Veracity

- More scientific data domains are emerging with capacities to capture real time health information
  - Proteomics
  - Metabolomics
  - Microscopy
  - Medical imaging
  - Other various technologies



National Library of Medicine
Twenty Seven Years of Growth:
NCBI Data and User Services

# Human Disease Networks (2015)



- Identify disease gene-phenotype associations at higher cellular and organismal levels

# Precision Medicine



Courtesy of P. Kuhn (USC)

- Learning system that accounts for complexity of underlying biology

- Requires
  - Deep biological understanding
  - Advances in scientific methods
  - Advances in instrumentation
  - Advances in technology
  - Advances in data management and computation

- Can *change* disease classifications

- Genomic, imaging, clinical, and laboratory data

# Standard Model of Computational Analyses
## *(circa 2014)*

# NCI Cloud Resources

- Democratize access to NCI-generated data
Create cost-effective scalable computational capacity



- Access large data sets without downloading data
- Bring tools and pipelines to the data
- Bring and combine own data and analyze with existing data
- Workspace to save and share data and results of analyses



**Data**
- Access and analyze 11,000 TCGA samples without having to download data
- Upload your own data for analysis

**Compute**
- Perform large scale analysis using the elastic compute power of commercial cloud platforms

**Security**
- dbGaP-authorized users can access controlled TCGA data
- Systems meet strict Federal security guidelines

# National Cancer Data Ecosystem

(Courtesy: Adapted from A. Kerlavage, NCI – NIH)

# The People

# Data Access and Sharing:
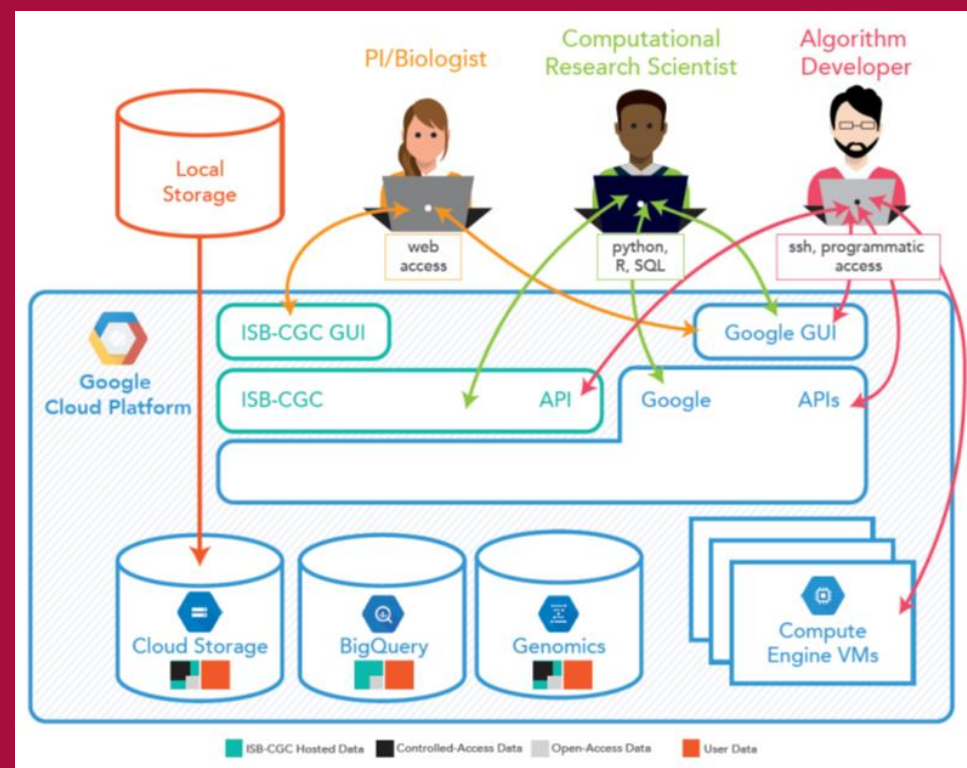## Isn't Only a Technology Challenge
### *Multidisciplinary Teams with Diverse Expertise and Resources*



**Biology/Social/Psychology Researcher**
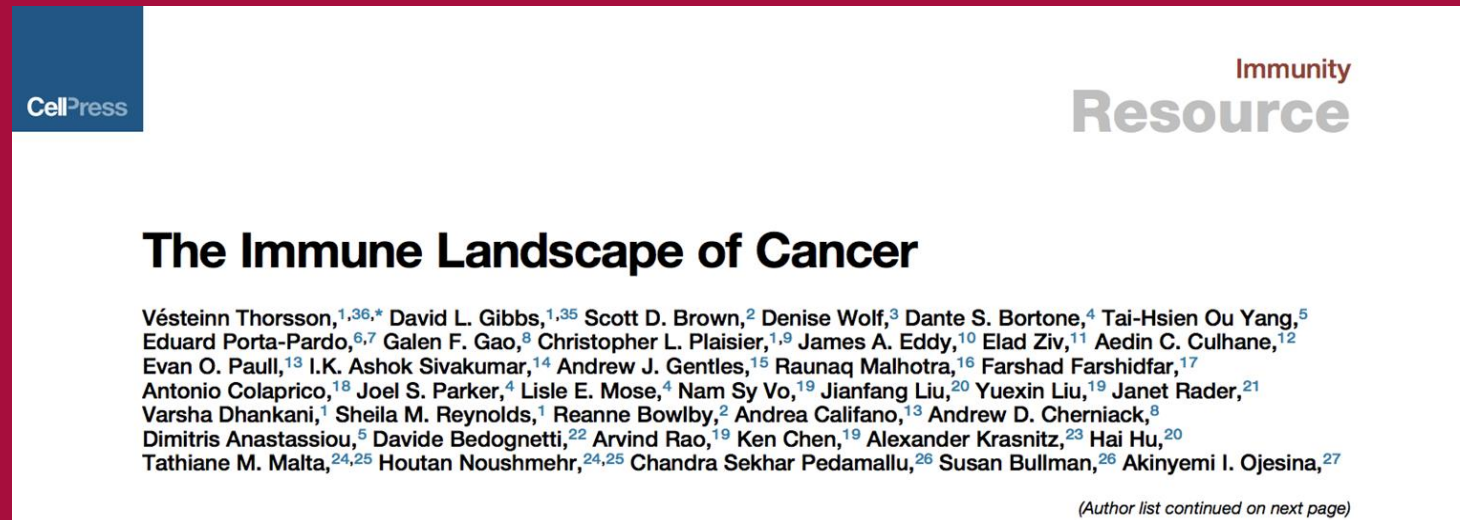
**Computational Scientist**

**Algorithm Developer**
- test new algorithm on hundreds or thousands of BAM or FASTQ files
- run novel image segmentation method across whole-slide images

NIH NATIONAL CANCER INSTITUTE

Wilkinson, M. D. *et al.* (2016). The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* 3:160018 doi: 10.1038/sdata.2016.18

# *d*atasharing **E**thical, **E**conomic, **L**egal, **S**ocial **I**mplications (dEELSI)

- Data and Information are not Neutral

# *d*atasharing Ethical, Economic, Legal, Social Implications (dEELSI)

- Data and Information are not Neutral
  - Identity, Phenotypes, and Bias
  - Identifiability and Privacy
  - Bias –Machine Learning and Artificial Intelligence
  - Incidental Findings and Return of Results
  - Informed Consent and Broad Data Uses
  - Governance, and Trustworthiness

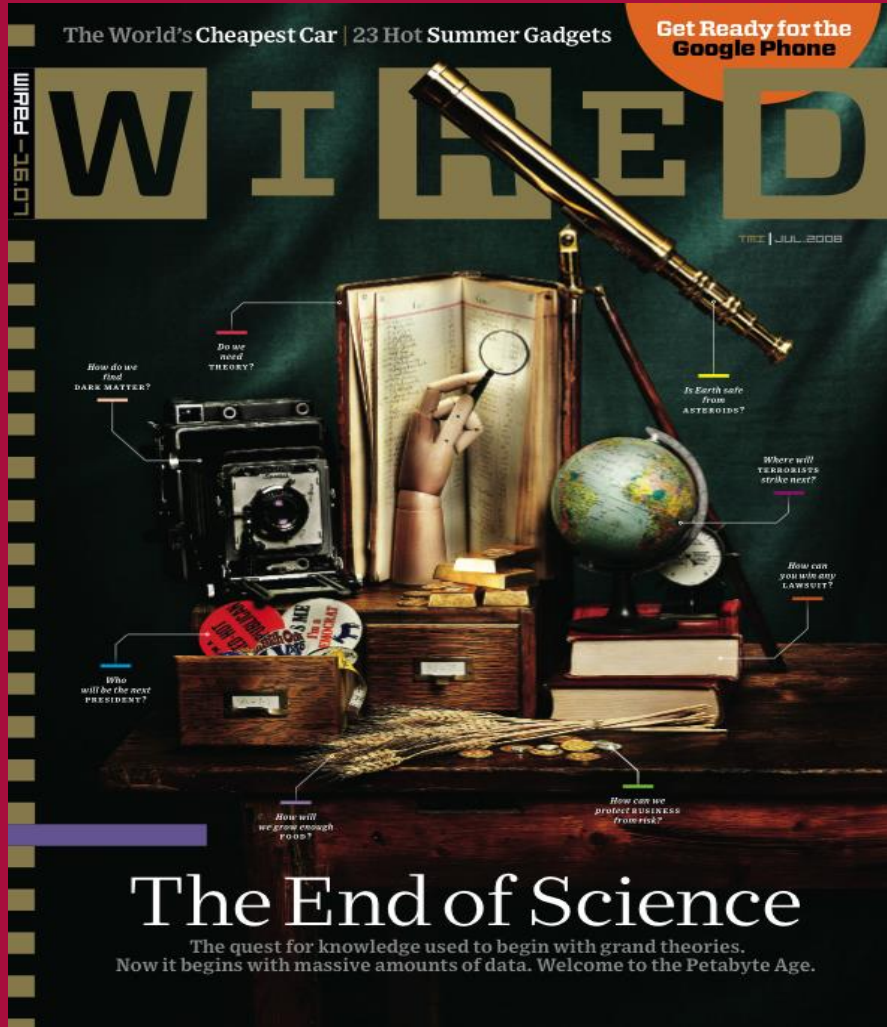# Minority and Health Disparity Issues

- Data and Information are not Neutral
  - Stigma: People/groups/communities/diagnoses/phenotypes
  - Inclusion: Data collected in basic/applied/clinical trial research
  - Diversity and Workforce issues
- Citizen Science and *community and patient engagement*
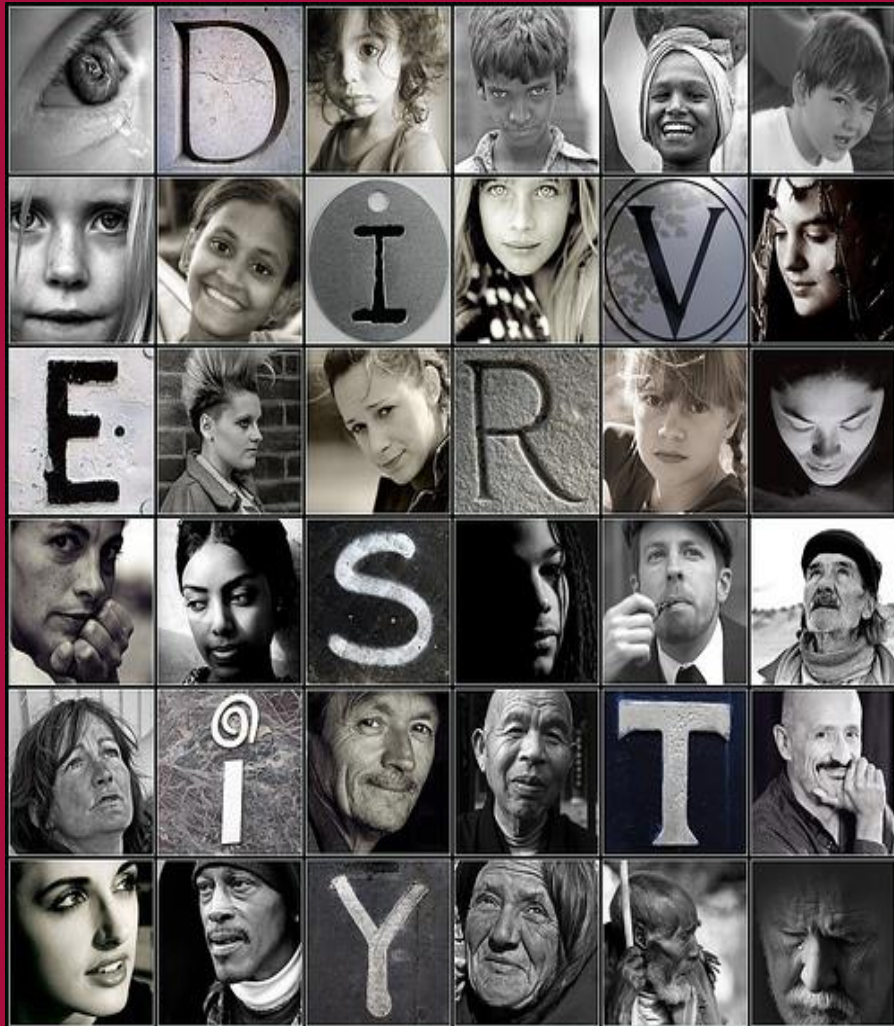- Inclusion, Equity, and Data Access – The Haves and Have Nots

# Challenges

## Data are difficult to

- *collect*
- *store*
- *delete*
- *search*
- *share*

- *visualize*
- *curate*
- *process*
- *analyze*

*with current available databases and tools*

# Challenges



Genomics is failing on diversity

An analysis by **Alice B. Popejoy** and **Stephanie M. Fullerton** indicates that some populations are still being left behind on the road to precision medicine.

A 2009 analysis revealed that 96% of participants in genome-wide association studies (GWAS) were of European descent[1]. Such studies scan the genomes of thousands of people to find variants associated with disease traits. The finding prompted warnings that a much broader range of populations should be investigated[2] to avoid genomic medicine being of benefit merely to "a privileged few".

Seven years on, we've updated that analysis. Our findings indicate that the proportion of individuals included in GWAS who are not of European descent has increased to nearly 20%. Much of this rise, however, is a result of more studies being done in Asia on populations of Asian ancestry. The degree to which people of African and Latin American ancestry, Hispanic people and indigenous peoples are represented in GWAS has barely shifted.
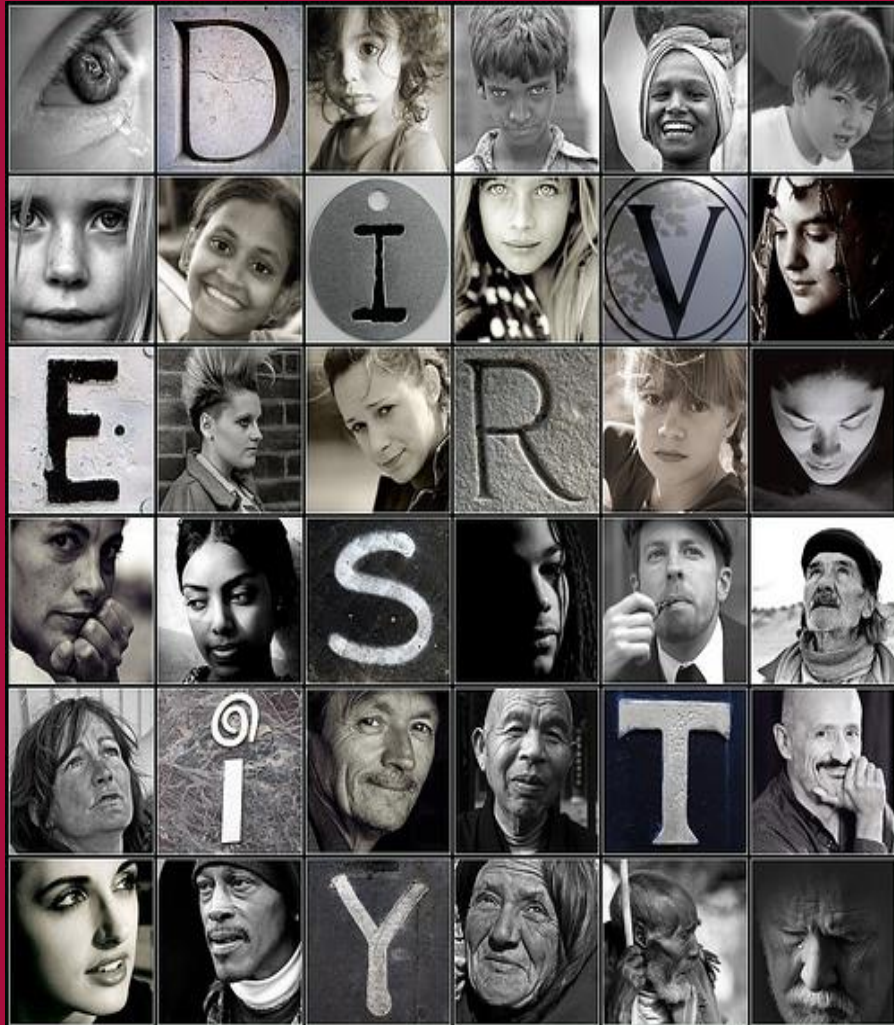
Thus, more than 20 years after the

US National Institutes of Health (NIH) mandated the inclusion of diverse participants in the biomedical research it funds, GWAS funded by the NIH and other sources are continuing to miss a vast portion of the world's genetic variation.

Over the past decade, GWAS have been the preferred tool for discovering the genetic factors involved in common diseases. Tens of thousands of significant associations between genetic variants and biological traits have ▸

# Challenges



How to have integrity
in [data
and data sharing]
in a world
that does not
affirm everyone's
humanity

*- Adapted from Thomas A. Parham*

# Challenges



*Science Isn't Broken*

"…it's just…a lot harder than we give it credit for…"

NIH NATIONAL CANCER INSTITUTE

# Challenges



https://www.economist.com/leaders/2013/10/21/how-science-goes-wrong

## *Science Isn't Broken*

"…it's just…a lot harder than we give it credit for…"

"Now it needs to change itself…"

# The Challenge



PARADIGM
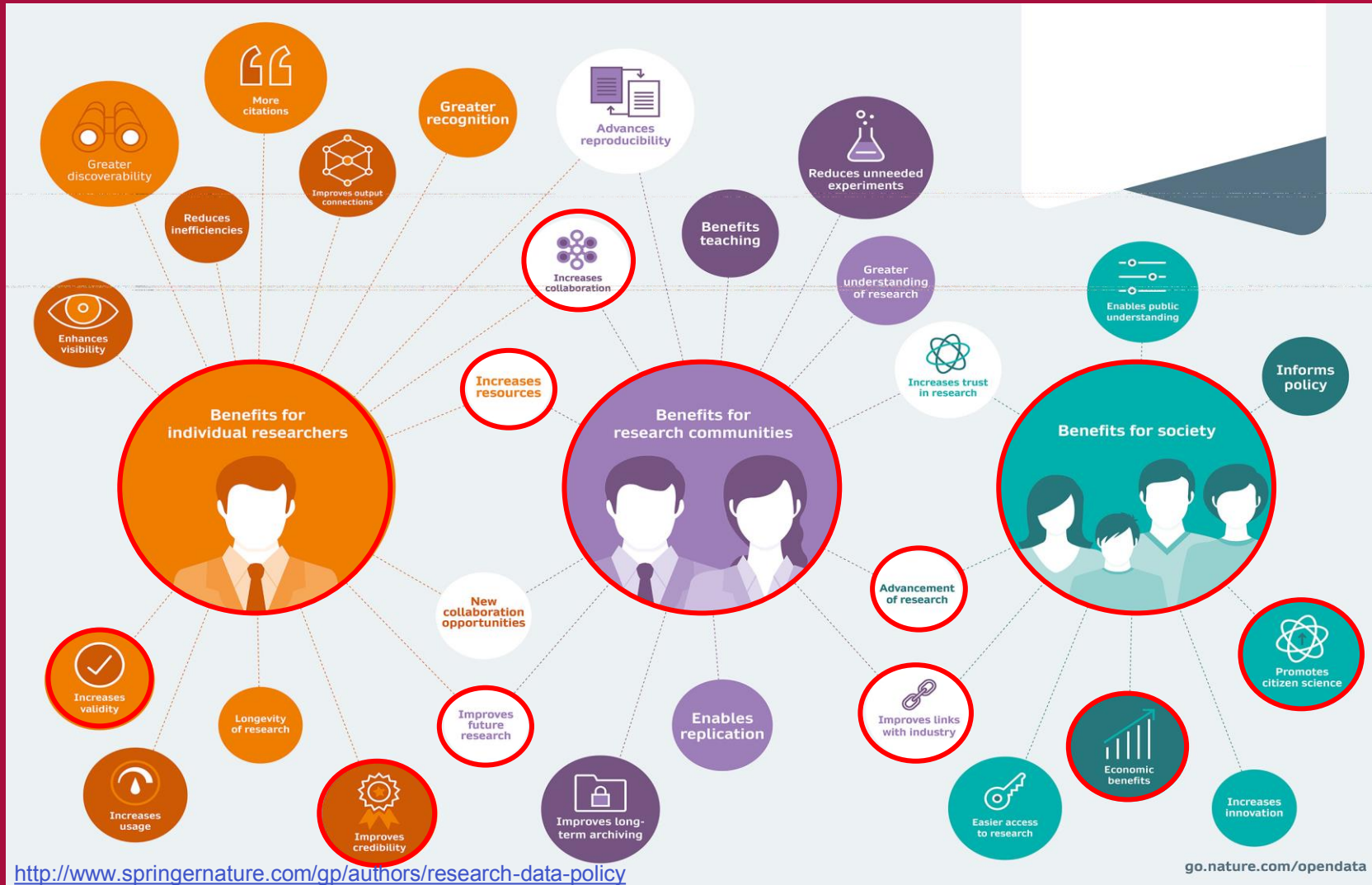SHIFT

Hypothesis
Confirmation
TO
Hypothesis
Generation



50TH ANNIVERSARY EDITION

THE STRUCTURE OF SCIENTIFIC

REVOLUTIONS

THOMAS S. KUHN
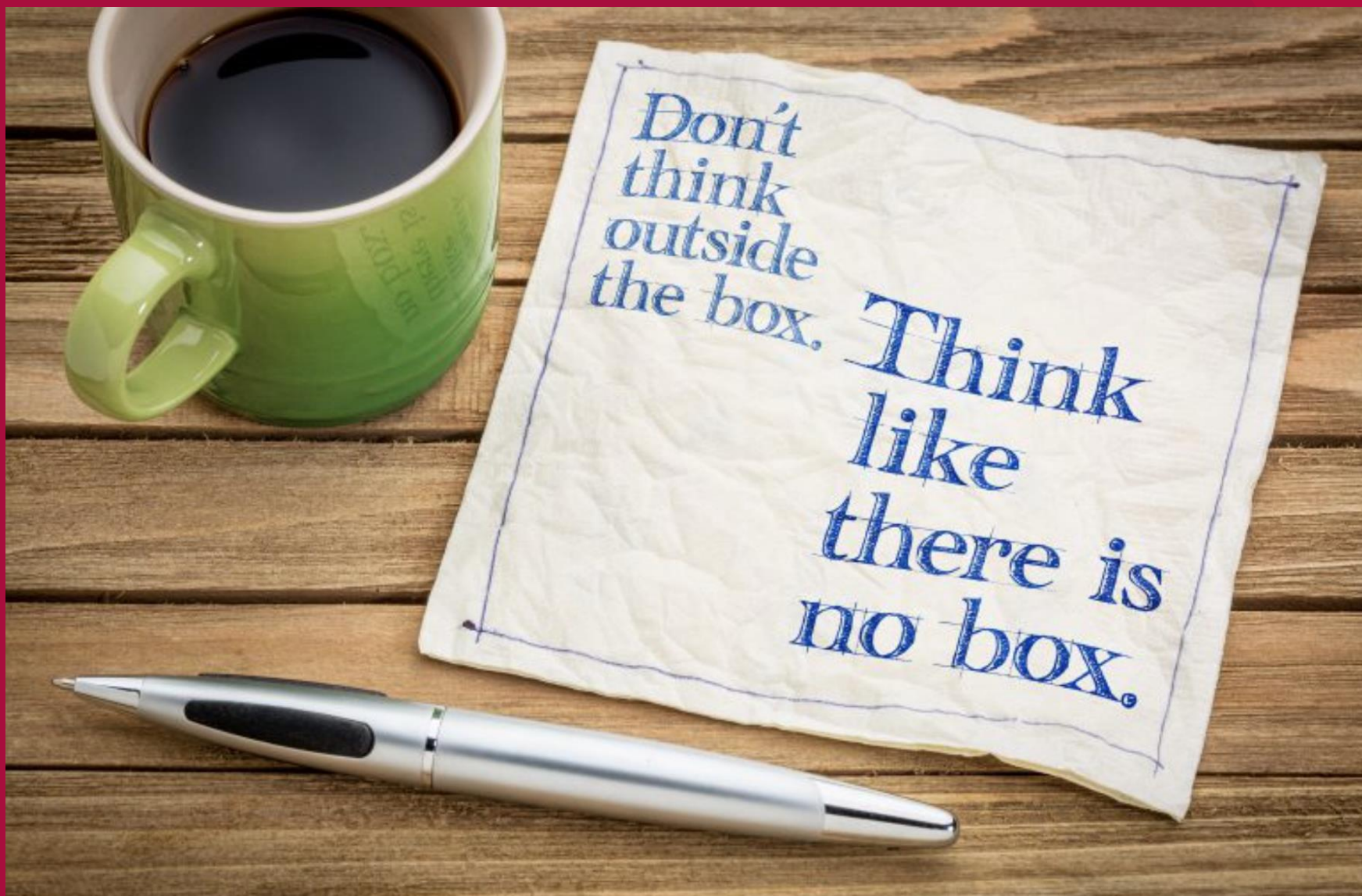
WITH AN INTRODUCTORY ESSAY BY IAN HACKING

# The Challenge



http://www.springernature.com/gp/authors/research-data-policy

Glass and MCAtee (2006).

# Setting You and your Data Free

- Catalyze new collaborations
- Increase confidence in results
- Larger projects (scope & sample size)
- Local vs global (>external validity)
- Generate greater recognition
- Credit
  - Digital Object Identifier (DOI) enables independent discoverable citability for researcher credit
  - Data tracking on the impact of research
  - Journal data availability statements
  - Open data badges



MENU ⌄   **nature**
International journal of science

**CAREER FEATURE** · 13 MAY 2019

## Data sharing and how it can benefit your scientific career

*Open science can lead to greater collaboration, increased confidence in findings and goodwill between researchers.*

**Gabriel Popkin**

Find a new job

Data sharing can be complex for scientists to navigate, but the rewards are often career-enhancing. Credit: Hero Images/Getty

**NIH) NATIONAL CANCER INSTITUTE**

Don't think outside the box. Think like there is no box.

NATIONAL CANCER INSTITUTE

**https://datascience.cancer.gov/data-sharing**

NATIONAL CANCER INSTITUTE

DEPARTMENT OF HEALTH & HUMAN SERVICES · USA

NIH

NATIONAL CANCER INSTITUTE

NIH NATIONAL CANCER INSTITUTE

www.cancer.gov          www.cancer.gov/espanol