# Advancing Knowledge Through Secondary Data Use
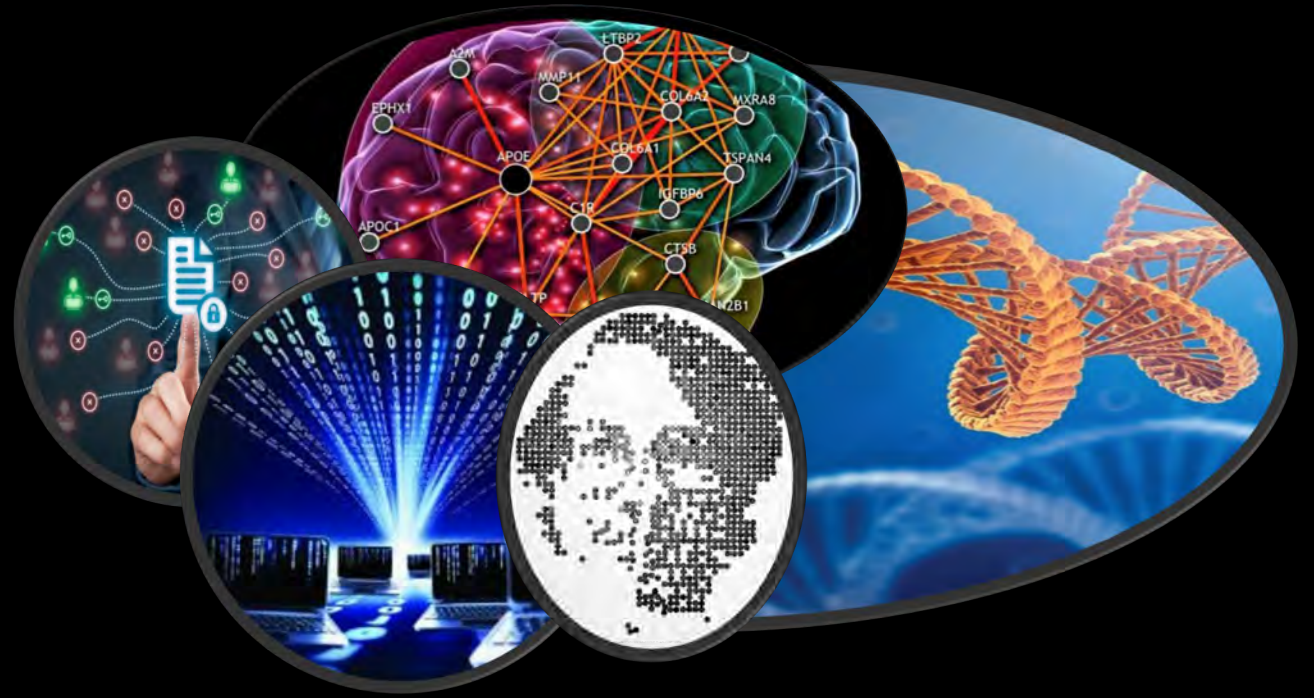
*NIDDK Data Management & Sharing Webinar Series*
*03 March 2023*

**Vivian OTA WANG, Ph.D., CGC, FACMG**
*Office of Data Science Strategy*
*DPCPSI/Office of the Director*
*National Institutes of Health, DHHS*

- **The Drivers**
*Human Rights and Open Science*

- **The Data and Data Science**
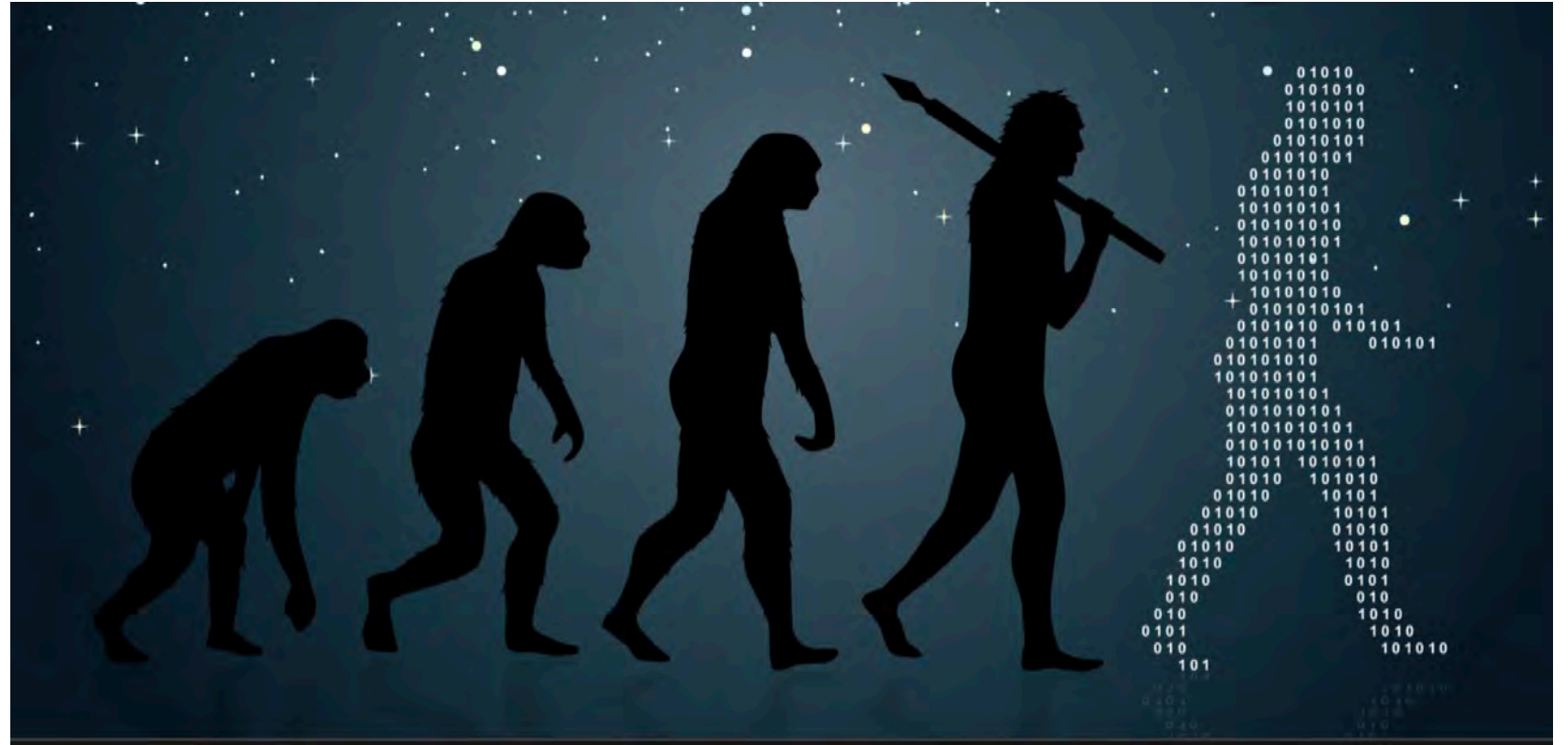*Volumes of Varied and Complex Data*

- **The Challenges**
*Economic, Ethical, Legal, and Social Implications*

- **Next Steps**
*You*

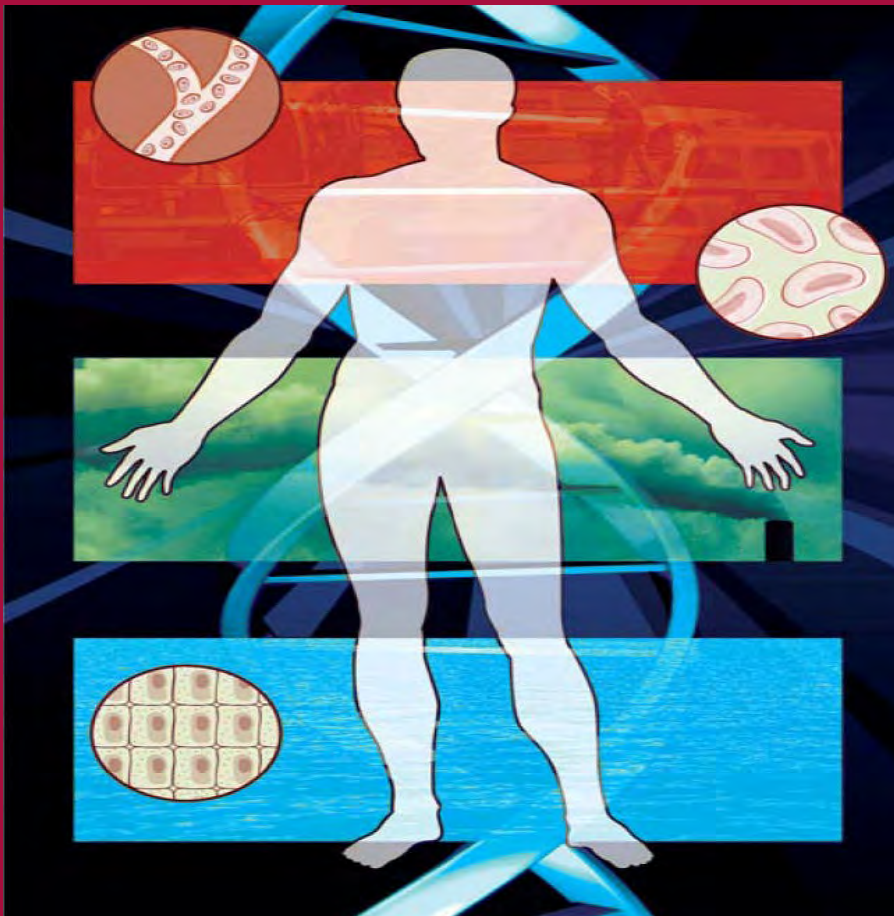THE DRIVERS

# Top Ten Leading Causes of Death *(2020-2022)*

| | Category | Total deaths (Jan.-Sept. 2022) | Total deaths (2021) | Total deaths (2020) |
|---|---|---|---|---|
| 1 | Heart disease | 572,336 | 767,937 | 764,512 |
| 2 | Cancer | 454,176 | 604,358 | 599,607 |
| 3 | COVID-19 | 234,434 | 475,059 | 343,566 |
| 4 | Accidents | 170,166 | 226,987 | 203,033 |
| 5 | Stroke | 123,215 | 162,769 | 159,248 |
| 6 | Chronic respiratory | 107,559 | 141,906 | 152,051 |
| 7 | Alzheimer | 87,866 | 119,442 | 134,271 |
| 8 | Diabetes | 74,716 | 103,197 | 101,355 |
| 9 | Other respiratory | 50,635 | 66,381 | 66,053 |
| 10 | Renal failure | 42,596 | 53,057 | 51,221 |

Notes: For 2022, the total death sum for each category is for January 1 - September 30, 2022, except deaths from accidents and suicides are from January - September 2021. Chronic respiratory is chronic lower respiratory disease.

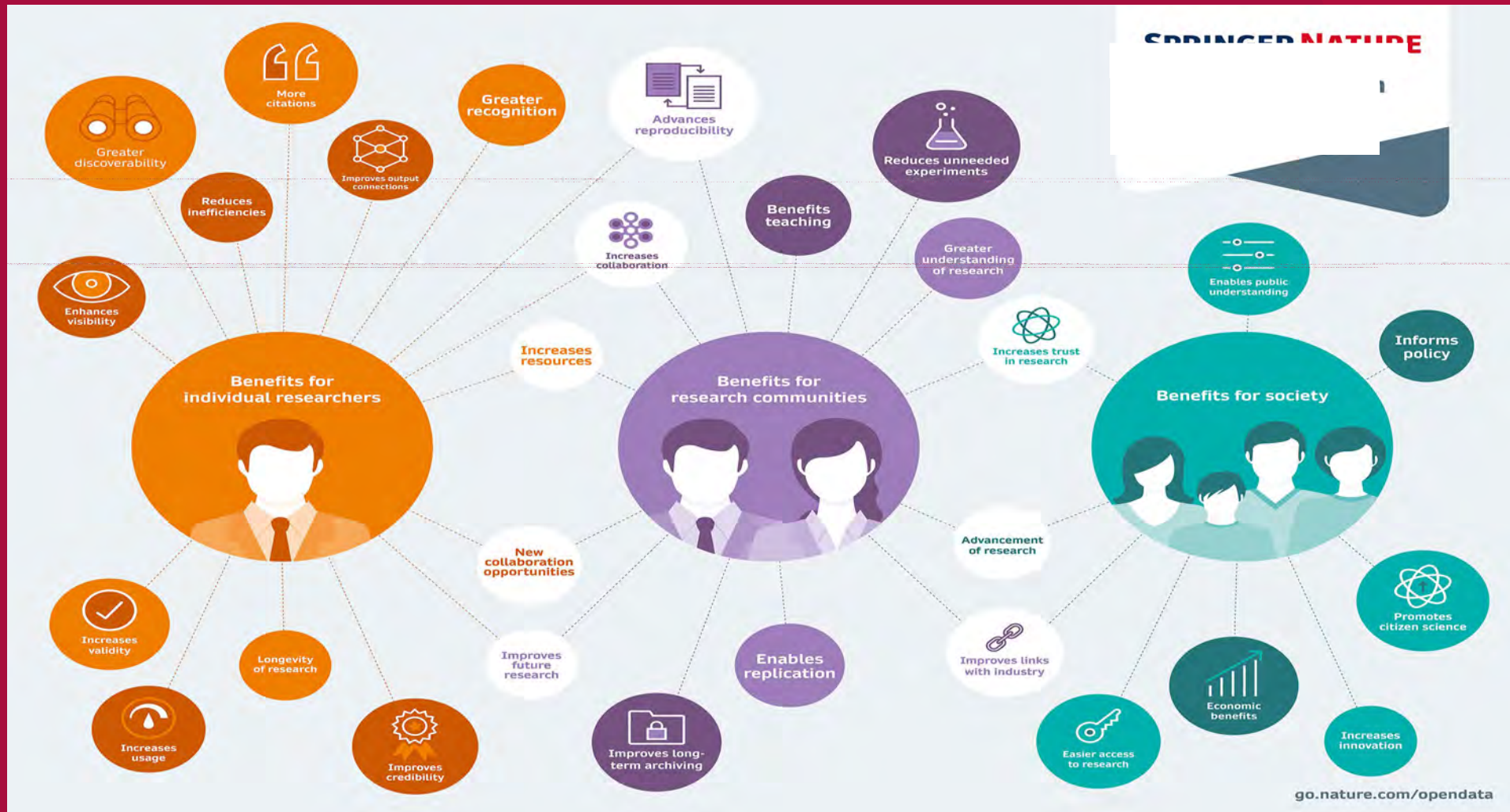https://www.healthsystemtracker.org/brief/covid-19-leading-cause-of-death-ranking/

# The People Factor: Top Ten Modifiable Behaviors Contributing to Mortality



- **Tobacco use**
- **Diet**
- **Physical activity**
- **Alcohol misuse**
- **Microbial agents**
- **Toxic agents**
- **Firearms**
- **Sexual behavior**
- **Motor vehicle accidents**
- **Substance abuse**

# Human Rights and the Democratization of Knowledge

# Open Science



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.

## *DATA SHARING AND INNOVATION*

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

# Open Science

- Facilitates innovation of research tools and methods

- Increases statistical power

- Improves research quality through validation and replication

# Open Science

- Facilitates innovation of research tools and methods

- Increases statistical power

- Improves research quality through validation and replication

# Open Science

- Facilitates innovation of research tools and methods

- Increases statistical power

- Improves research quality through validation and replication



47/53 "landmark" publications could not be replicated

# Open Science and Data Sharing

## Precision Health



Courtesy of P. Kuhn (USC)

- Accounts for social and cultural complexity influencing underlying biology

- Requires
  - Biological understanding
  - Inclusion of social, cultural, and psychological factors
  - Scientific methods advancements
  - Instrumentation advancements
  - Technology advancements
  - Data management and computation advancements

- -omic, imaging, clinical, laboratory, etc data

- Can *change* disease classifications and treatments

# Open Science and Data Sharing

- Increases scientific value by exploring, combining, and analyzing data from multiple sources

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

# Open Science and Data Sharing

- Increases scientific value by exploring, combining, and analyzing data from multiple sources

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

**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.

(Adapted from Glass, TA & McAttee, MJ. (2006) *Social Science & Medicine, 62, 1650-1671).*

# Data: Variety, Volume, Velocity, and Veracity

## Real-time, Real-world Data Capacities

- Proteomics
- Metabolomics
- Microscopy
- Imaging
- Electronic Medical Records
- Mobile Devices
- Psychological/behavioral/self-report
- Other technologies

**Observational Phenotype**

**Genomic**    **Other 'omic'**

**Location**    **Behavior**

**Imaging**    **Clinical/EMR**

**Exposures**

**Mobile Sensors**

# Data: Variety, **Volume**, Velocity, and Veracity
## *Health Knowledge Doubling Time*



Doubling Time of Health Knowledge

| | |
|---|---|
| 1950 | 50 years |
| 1980 | |
| 2010 | |
| 2020 | 73 days |

days

(Densen, P. *Trans Am Clin Climatol Assoc* (2011)

# Data: Variety, Volume, Velocity, and Veracity
## *Personal Health Data*



https://databricks.com/blog/2022/03/09/introducing-lakehouse-for-healthcare-and-life-sciences.html

# Data: Variety, Volume, Velocity, and Veracity
## *Improvements to Store and Process Data*



**A Advances in Storage Capacity**

Size    Storage Capacity

**IBM 305 RAMAC (1956)**
Storage: 5–10 MB
Size: Housed in a room ~30x50 ft
Weight: 20,000 lb (10 tons)
Price: $3,200/month (equivalent to $30,900 in 2021)

**Cray-2 Supercomputer (1985)**
Storage: ~32 GB
Size: 48x66 in
Weight: 5500 lb (2494 kg)
Price: $30 million

**iPhone XS (2018)**
Storage: 512 GB
Size: 5.8x3.05 in
Weight: 0.3 lb (0.136 kg)
Price: $900

Time

**B Advances in Speed**

**Figure 1. Improvements over 50 Years in the Ability of Computers to Store and Process Data.**
Panel A shows advances in data storage, in terms of both physical size and cost per unit of storage. RAMAC denotes random access method of accounting and control. Panel B shows advances in the speed of computing. Each dot represents an individual machine type and the approximate year of its introduction. These improvements in storage and speed have allowed machine learning to progress from a dream to reality. Data in both panels are estimates from many types of system architecture and are derived from multiple public sources.

# Data: Variety, Volume, Velocity, and Veracity

- Trustworthy
- Accurate
- Reliable

*[Information, information uses],* and "definitions belong to the definers, not the defined.

- Toni Morrison

*Beloved*

SOME
SCIENCE
*DATA SCIENCE*

# Download Computation Resources and Data
## *(circa ~pre-2014)*



Public Data

Network Download

Local Data

Local Storage & Compute Resources

Publicly Available Software

UNIVERSITY

Locally Developed Software

- Only large institutions had ability to utilize data
- Storage/data protection cost~$2M/yr
- Data download at 10 Gb/second  (*23 days*)
- Increase rate of data generation

# Co-Locate Computation Resources and Data
## *(circa ~2014-2016)*



**Computational Capacity**

Core Data (TCGA)

Application Programming Interface (API) Secure Data Access

User Data

- Access large data sets without downloading data
- Bring tools and pipelines to data
- Combine own data and analyze with existing data
- Workspace to save, share, and analyze data

# Co-Locate Computation Resources and Data
## *(circa 2017 – now)*



- Democratize data sharing and access
- Cost-effective Scalable Computational Capacity

# Data Reuse:
# Isn't Only a Data and Technology Challenge



The Open Data Iceberg

F.A.I.R[1]

C.A.R.E[2]

CARE
Collective Benefit · Authority to Control · Responsibility · Ethics

The Open Data Iceberg

Technology → The Technical Challenge

Processes & Organisation
- The Ecosystem Challenge
- The Funding Challenge
- The Support Challenge

People
- The Skills Challenge
- The Incentives Challenge
- The Mindset Challenge

Ethics
The Ethical Challenge
Inclusion
Respect
Equity

Developed from: Deetjen, U., E. T. Meyer and R. Schroeder (2015).

# Data Reuse:
## Isn't Only a Data and Technology Challenge
### *Difficulties Accessing, Analyzing, and Integrating Data*

Data Reuse:
Isn't Only a Data and Technology Challenge
*Difficulties Accessing, Analyzing, and Integrating Data*

Technology

Social

Scalable & Secure Environments

Data Harmonization & Organization

Data Sharing & Collaboration

Data Analysis Fluency

NIH

# Data Reuse:
# Isn't Only a Data and Technology Challenge

## *Communication*

# Data Reuse:
# Isn't Only a Data and Technology Challenge

*Multidisciplinary Teams with Diverse Expertise and Resources*

**Biology/Social/Psychology Researcher**
- select a data subset based on clinical, molecular, clinical characteristics
- explore all data for a specific pathways or models
- compare one cohort to another
- upload a small private dataset to analyze in conjunction with existing dataset



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

# Data Reuse:
# Isn't Only a Data and Technology Challenge

*Multidisciplinary Teams with Diverse Expertise and Resources*



Biology/Social/Psychology Researcher

**Computational Scientist**
- interactive data exploration
- use R or Python to perform custom analyses
- develop new tools
- share new tools
- publish new tools (including interactive)
- develop/customize pipelines

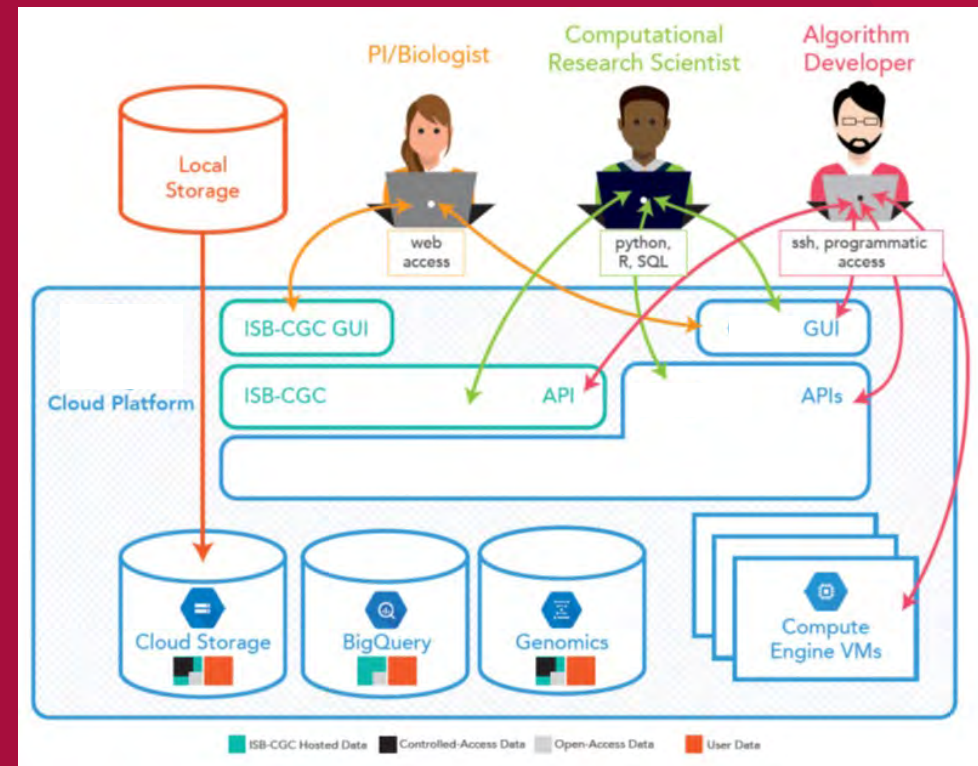(Courtesy: Adapted from A. Kerlavage, CBIIT-NCI-NIH)

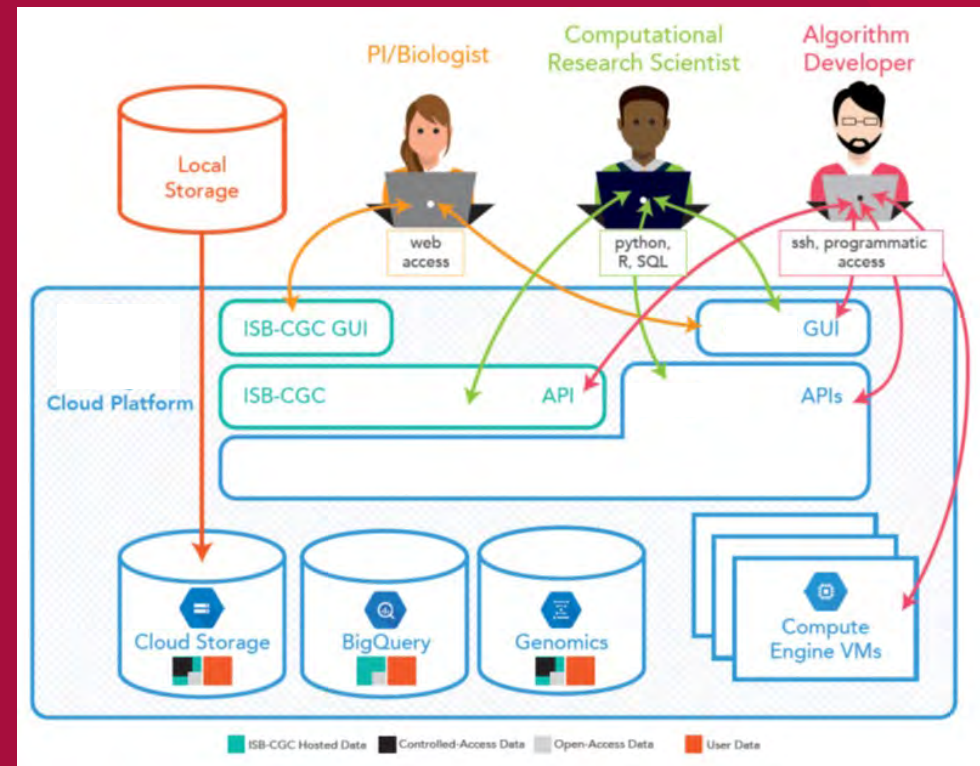# Data Reuse:
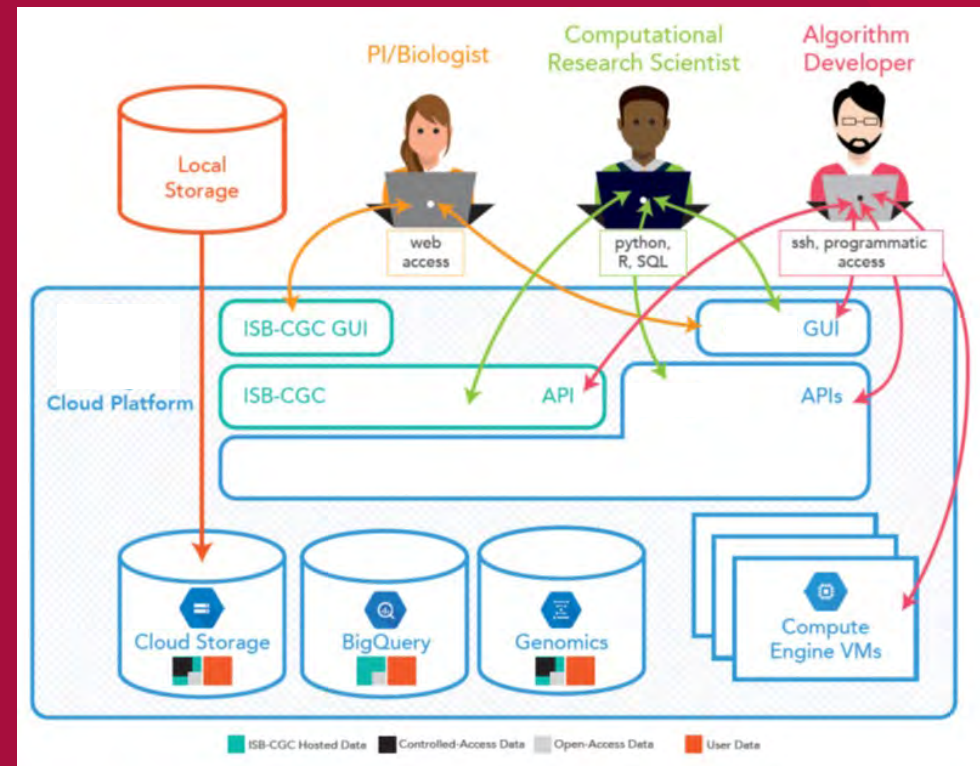# Isn't Only a Data and Technology Challenge

*Multidisciplinary Teams with Diverse Expertise and Resources*

**Biology/Social/Psychology Researcher**

**Computational Scientist**

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



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

INNOVATION AND REVOLUTION

transforming material | transforming energy | transforming information

Progress

Computing information (knowledge & algorithms)

Communicating & storing information

Combustion power

Electric power

Steam power

Water power

Stone tools

Bronze tools

Iron tools

2,000,000bc | 3,300bc | 1,200bc | ... | 1780 | 1848 | 1895 | 1940 | 1973 | 2008

Based on Hilbert (2020). Digital technology and social change: The digital transformation of society from a historical perspective. Dialogues in Clinical Neuroscience, 22(2), 189–194. https://doi.org/10.31887/DCNS.2020.22.2/mhilbert

# The Data Revolution



- Research is now a *data intensive enterprises*

- *Rapidly changing* scale
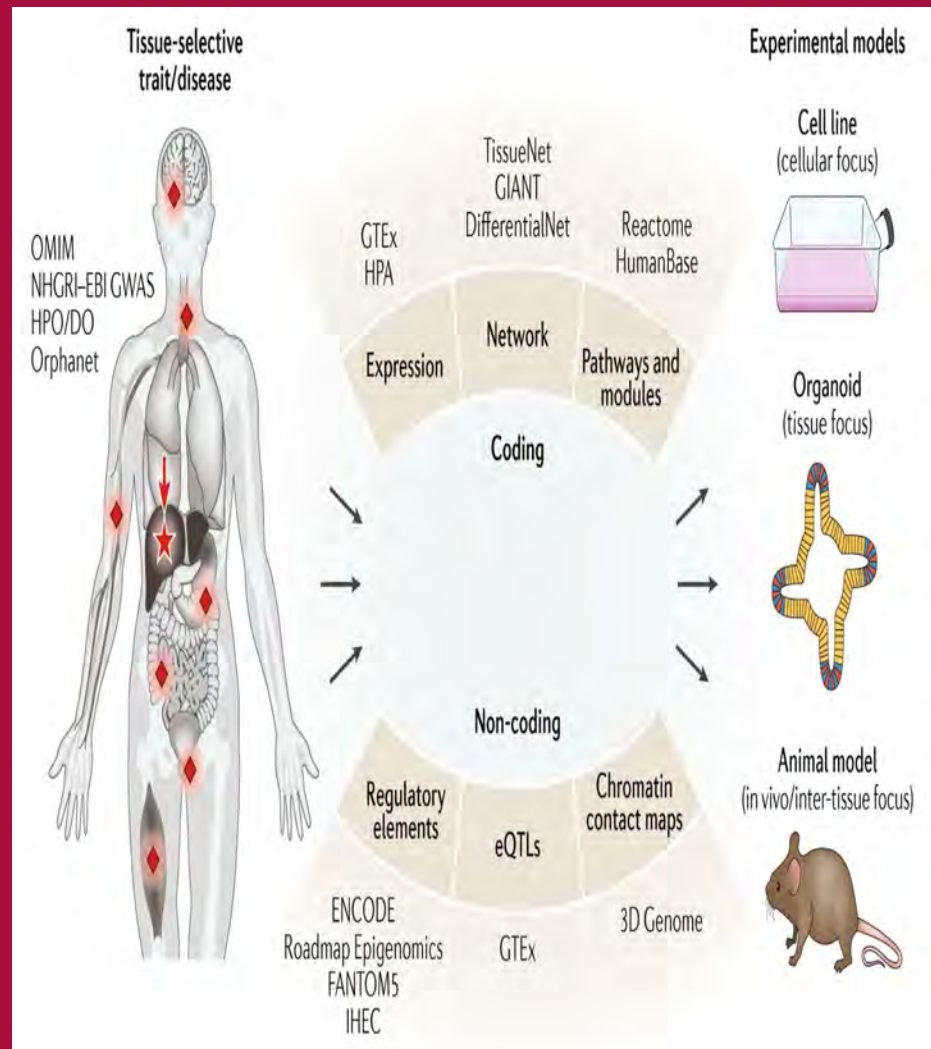
- Technology, *data computing* & *information technology (IT)* are *pervasive* in the *lab, clinics, and homes*

# The Data Revolution



Hekselman, I. and Yeger-Lotem, E. (2020). *Nature Reviews Genetics,* (21) 137-150.
https://doi.org/10.1038/s41576-019-0200-9

- *Changing functional roles of genes across tissues*

- *Relationships among diseases*

- *Relationships of behavior and health*

# The Data Revolution



- *Changing functional roles of genes across tissues*

- *Relationships among diseases*

- *Relationships of behavior and health*

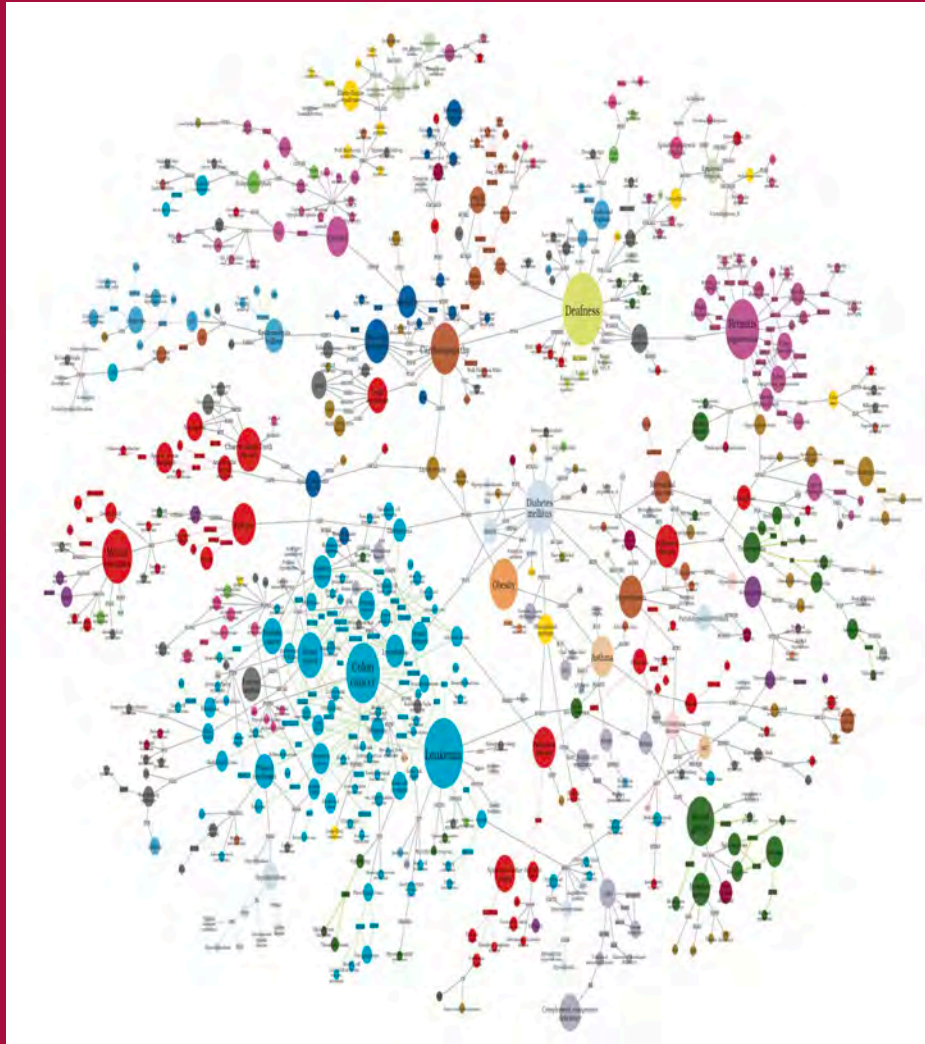# The Data Revolution



- *Changing functional roles of genes across tissues*

- *Relationships among diseases*

- *Relationships of behavior and health*

# The Data Revolution



Christakis NA, Fowler JH. N Engl J Med 2007;357:370-379
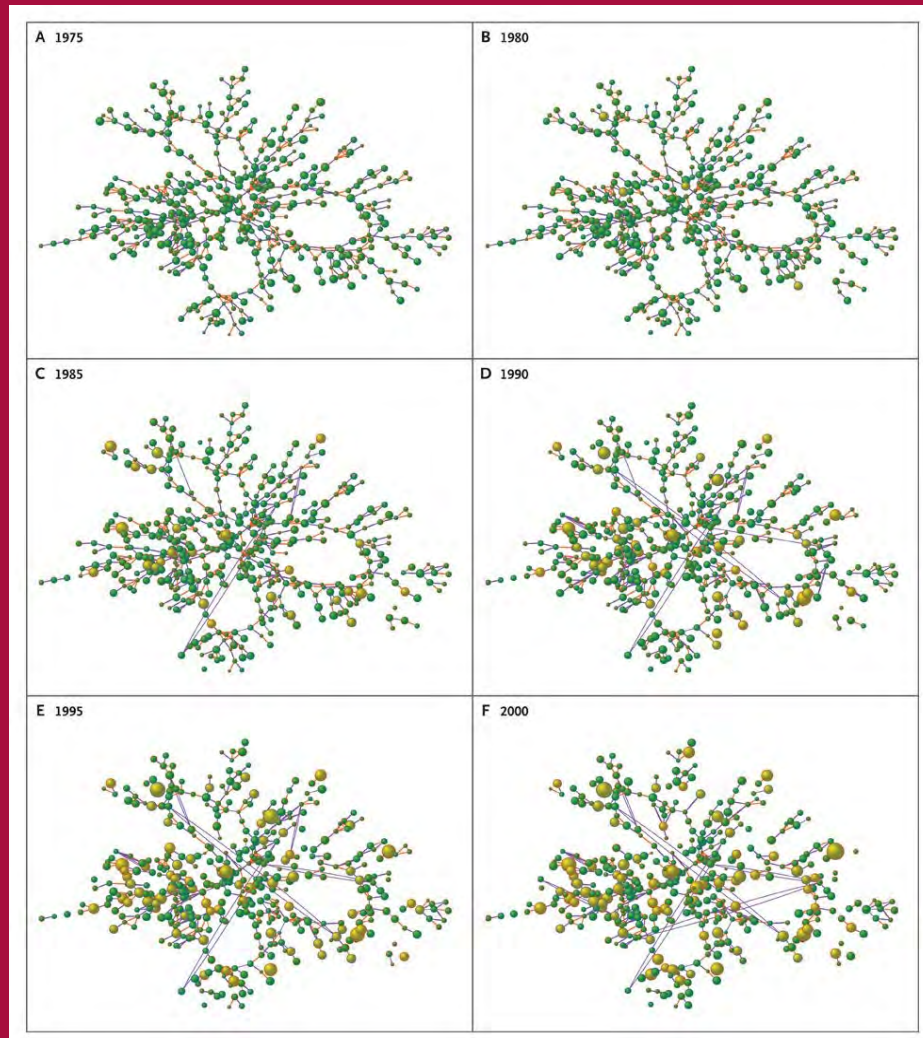
- *Changing functional roles of genes across tissues*

- *Relationships among diseases*

- *Relationships of behavior and health*

# The Data Revolution



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

THE BIG DATA REVOLUTION

# The Paradigm Shift



50TH ANNIVERSARY EDITION

THE STRUCTURE OF SCIENTIFIC

REVOLUTIONS

THOMAS S. KUHN

WITH AN INTRODUCTORY ESSAY BY IAN HACKING

# The Paradigm Shift

The Paradigm Shift

Hypothesis Confirmation TO Hypothesis Generation

The Paradigm Shift



The World's **Cheapest Car** | 23 Hot **Summer Gadgets**

**Get Ready for the Google Phone**

WIRED

JUL 2008

*How do we find* DARK MATTER?

*Do we need* THEORY?

*Is Earth safe from* ASTEROIDS?

*Where will* TERRORISTS *strike next?*

*How can you win any* LAWSUIT?

*Who will be the next* PRESIDENT?

*How will we grow enough* FOOD?

*How can we protect* BUSINESS *from risk?*

## The End of Science

The quest for knowledge used to begin with grand theories.
Now it begins with massive amounts of data. Welcome to the Petabyte Age.

# Ethical, Economic, Legal, Social Implications (EELSI)

- Nature vs Nurture
- Privacy and Confidentiality
- Research and other People Protections
- Informed consent
- Risk Assessment/Decision-making
- Benefits and Harms
- Predictive/Prognostic Screening/Testing
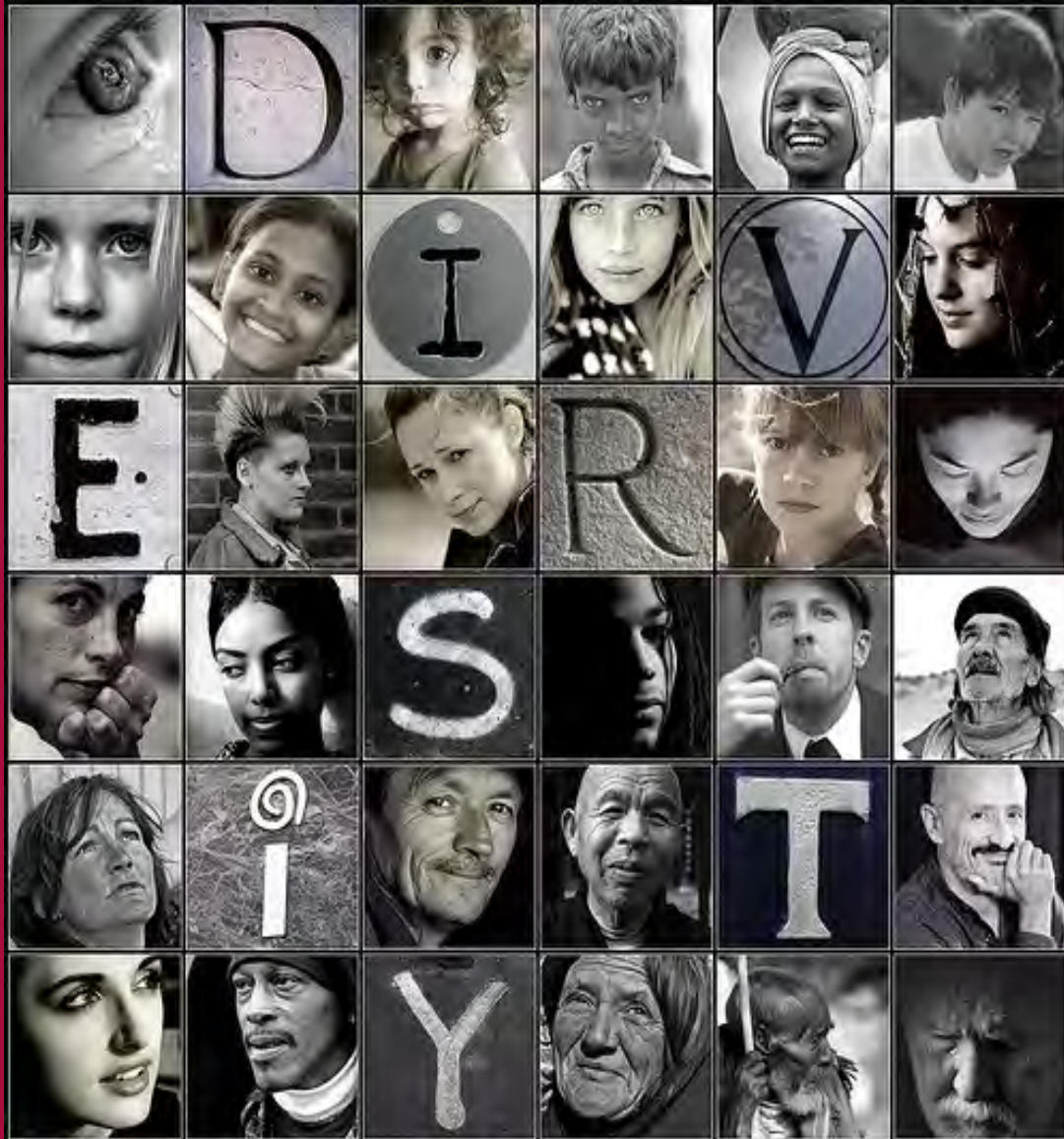- Intellectual Property
- Dual Uses *(Forensics/Surveillance)*

# Equity and Disparity Issues

- Data and Information are not Neutral
    - Stigma:  *People/Groups/ Communities/Phenotypes*
- Inclusion: *Basic/Applied/Clinical Trial Research*
    - Diversity and Workforce *Issues*
        - Citizen Science *and Community Engagement*
            - Inclusion, Equity, *and The Haves and Have Nots*

*Data*

If not designed
to address equity,
research (and data)
will perpetuate
disparities
and injustices

*- me*