# Developing a New Culture: NIH Policies in Data Management, Sharing & Access

JAIME M. GUIDRY AUVIL, PH.D.
Office of Data Sharing, NCI, NIH
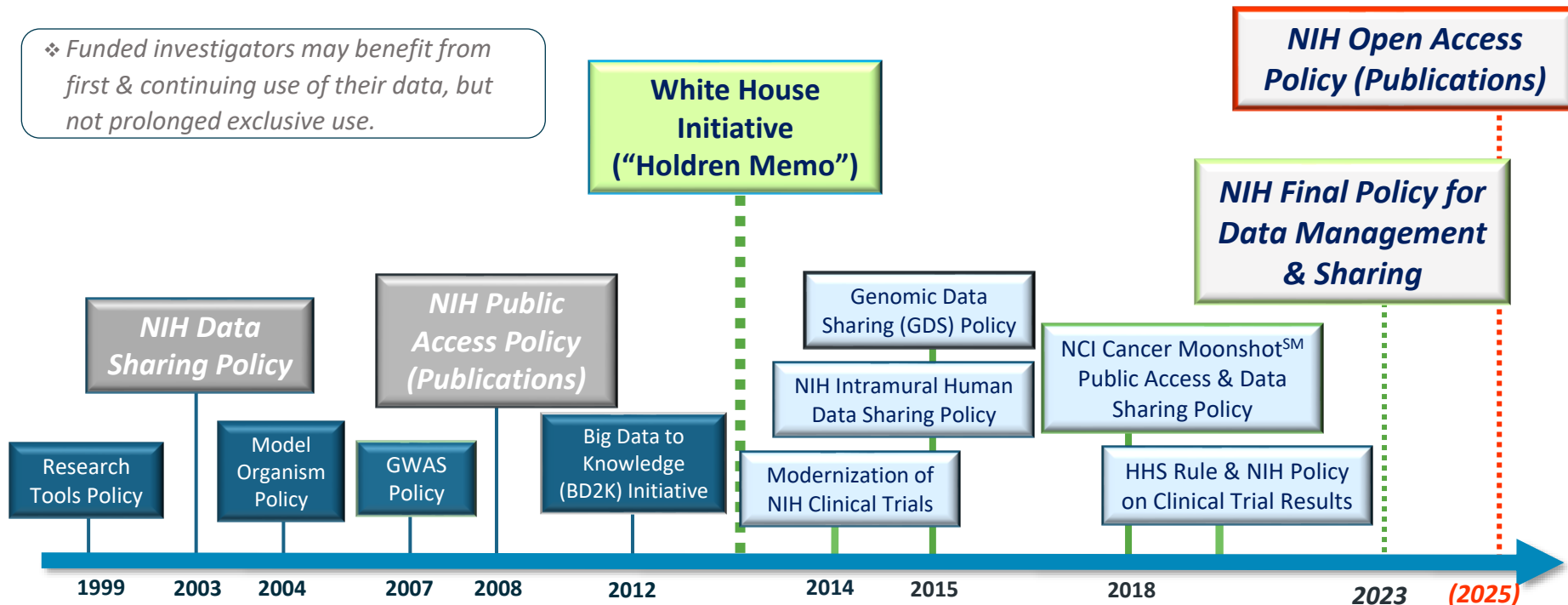
January 18, 2023

# What is "Open Science"?

A "movement" to make **scientific research** (including ***publications, DATA, physical samples, and software***) and dissemination **accessible** to all levels of society, amateur or professional.

- Open science is transparent and accessible knowledge that is shared and developed through **collaborative networks**.

- It encompasses practices such as:
  - publishing open research & campaigning for **open access**,
  - encouraging scientists to practice open-notebook science (such as **openly sharing data and code**),
  - broader dissemination and engagement in science, and
  - generally making it easier to publish, access and communicate scientific knowledge.

- Usage of the term varies substantially across disciplines, with a notable prevalence in the STEM disciplines.

# NIH History of Data Sharing Policies

❖ *Funded investigators may benefit from first & continuing use of their data, but not prolonged exclusive use.*

**NIH Open Access Policy (Publications)**

**White House Initiative ("Holdren Memo")**

**NIH Final Policy for Data Management & Sharing**

*NIH Data Sharing Policy*

*NIH Public Access Policy (Publications)*

Genomic Data Sharing (GDS) Policy

NCI Cancer Moonshot[SM] Public Access & Data Sharing Policy

NIH Intramural Human Data Sharing Policy

Research Tools Policy

Model Organism Policy

GWAS Policy

Big Data to Knowledge (BD2K) Initiative

Modernization of NIH Clinical Trials

HHS Rule & NIH Policy on Clinical Trial Results

| 1999 | 2003 | 2004 | 2007 | 2008 | 2012 | 2014 | 2015 | 2018 | *2023* | *(2025)* |

*Investigators must share any information necessary to understand, develop or reproduce published research (raw data, statistical methods, tools, source code)*

# Key Messages for Final NIH DMS Policy

Promote **open science**, stimulate new **discovery**, enable **rigor** & **reproducibility**, and provide **transparency**

***Driving A Cultural Shift*** through planning for consistent, collaborative & impactful data management and sharing as a critical part of all research

NIH is taking a "***learning approach***" (i.e., phased and iterative implementation in the years to come).

Over time, thoughtful DMS plans will inform clear guidance on the highest value data types beyond genomics (repository, timelines, etc.)

NIH is providing resources, training, and guidance, to the extent possible, for initial rollout & will continue to develop these over time

# Final NIH Policy for Data Management & Sharing (DMS)

**Goal: Share scientific data supported by the NIH broadly and immediately**

## Scope

- **All NIH-supported research that generates <u>scientific data</u>**
  - Intramural & Extramural research
  - Grants, Contracts & OTAs
  - Human & non-human data; no budget threshold
- Submitted on/after **Jan. 25, 2023**
- Complementary to current data sharing policies (does not replace)

## Expectations

- ✓ **Investigators must submit <u>plans to manage & share data</u> (research applications)**
- PIs should comply with app[...] DMS Plans, and share dat[...]
  - In set repositories, to extent[...] possible **(FAIR Principles)**
  - At publication or end of awa[...] whichever sooner

## Public Access

- Final manuscripts to PubMed Central w/i 12 months of acceptance

***Public Access policies to be updated asap (by Dec. 31, 2025) → agencies make <u>publications</u> & supporting <u>data</u> from federally funded research publicly accessible with <u>no embargo</u> on their <u>free & broad release</u>***

# DMS Policy Applies to All *Scientific Data*

> *"The recorded factual material commonly accepted in the scientific community as of sufficient quality to validate and replicate research findings, regardless of whether the data are used to support scholarly publications."*



## *Not Scientific Data Under the DMS Policy:*

- *Data not necessary for, or of sufficient quality, to validate & replicate research findings*
- *Laboratory notebooks*
- *Preliminary analyses*
- *Completed case report forms*

- *Drafts of scientific papers,*
- *Plans for future research*
- *Peer reviews*
- *Communications with colleagues, or*
- *Physical objects, (e.g., laboratory specimens)*

# Responsibilities and Expectations: Investigators



Researchers to prospectively plan for how scientific data will be managed and shared through submission of a DMS Plan that considers any potential restrictions or limitations (no separate GDS Plans)



Researchers to maintain alignment with the approved DMS Plans by the NIH ICO (in the awards) and with the evolution of Plans during project period (through RPPR)

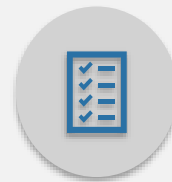# NIH DMS: Supplemental Guidance

## Recommended Elements of a NIH DMS Plan

Data Types

Tools for Access/ Manipulation

Data Access/Reuse Considerations

Data Standards

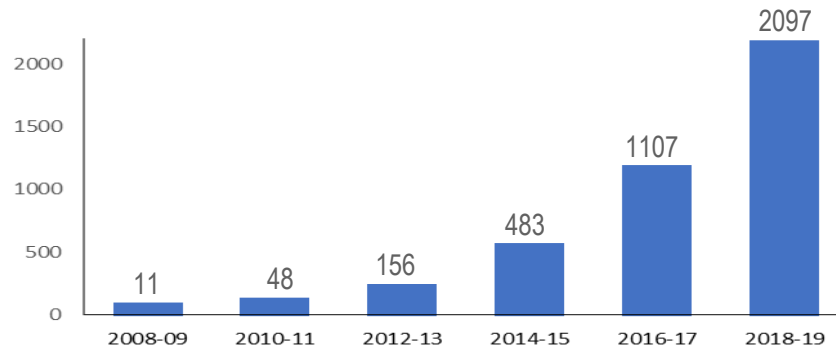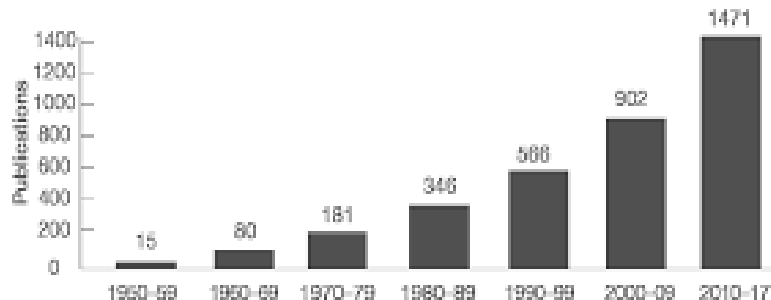Sharing Oversight Methods

Data Access & Timelines

Guidance for Selecting a Repository for Data

Guidance for Allowable Costs of Data Management & Sharing

# Driving Science through Publications, Data & Collaboration





| | Framingham Heart Study | The Cancer Genome Atlas |
|---|---|---|
| *Study Length* | 70 years | 12 years |
| *Cases Studied* | 15,144 | 11,429 |
| *Publications* | **3,698 (~38,000 PMC)** | **3,747 (~62,000 PMC)** |
| *Controlled-access Data* | Consortia; HMB (+IRB/MDS, 2K=NPU) | Collaborative Teams & Public Use of Data; GRU |
| *Authorized Users* | 715 | 3,335 |
| *Open Data Use & Availability Timing* | Little Open Data; mostly available with publication | Some Open Data; All data immediately available to community |

# The Cancer Moonshot: Success in Mission-Driven Science

**Cancer Moonshot℠:**
Accelerate discovery, increase collaboration, and expand data sharing

**In the Cancer Moonshot's first 4 years** (2017–2021):
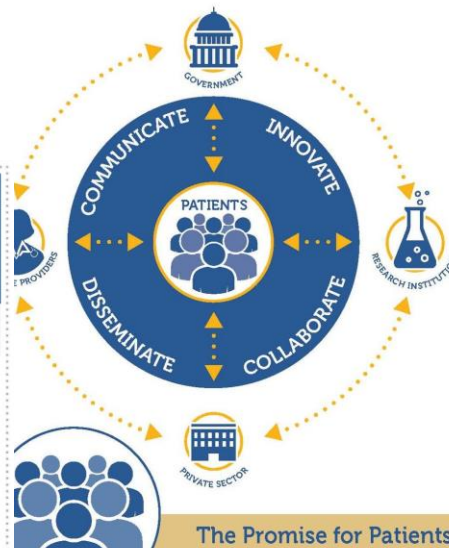
**>2,000** Publications

**49** Clinical Trials

**>30** Patent Filings

CANCER MOONSHOT

**INITIATIVES 2017–2022**

OVER **70** CONSORTIUMS OR PROGRAMS

OVER **240** RESEARCH PROJECTS

CANCER MOONSHOT

COMMUNICATE • INNOVATE • COLLABORATE • DISSEMINATE
PATIENTS
GOVERNMENT • RESEARCH INSTITUTIONS • PRIVATE SECTOR

**MISSION**
Dramatically accelerate efforts to prevent, diagnose, and treat cancer—to achieve a decade's worth of progress in 5 years

**WHY NOW**
New scientific understanding and vast amounts of rich data just waiting to be transformed into solutions

Immense science and technological capabilities positioning us for a quantum leap

A shared national commitment to harness the intellectual creativity and innovation of the American people

**The Promise for Patients**

New and improved treatment options

More sensitive screening measures

Improved use of effective prevention strategies

Better information for making medical decisions

Increased tools for community care providers

New ways to track and share health information

Together, we can end cancer as we know it.

To learn more, please visit WH.gov/cancermoonshot

**\*\*Take Home Message: purposeful, broad, early access to data leads to much faster and impactful outcomes**

# NIH Resources to Support Investigators

https://sharing.nih.gov

➢ **A Central Data Sharing Website**: A one-stop-shop for information on NIH data sharing policies & related resources

- *NIH Policy Notices & Supplemental Guidance*
- *Relevant Forms & Templates*
- *Frequently Asked Questions (FAQs)*
- *Sample DMS Plans*
- *Policy Decision Tool*
- *Links to Individual NIH Institutes & Centers for program-specific references*

➢ **Other Resources and Tools** [e.g., DMPTool (https://dmptool.org/), Repositories]

➢ **NIH Training**: News Events

➢ **A Central Mailbox:** Help answer DMS questions (Sharing@nih.mail.gov).

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

# Additional Slides

# Benefits of Broad Data Sharing

## Collaborator Sharing

- Between investigator to investigator (e.g., sharing upon publication and request to the author)

## Consortium Sharing

- Within large collaborative groups (e.g., sharing between investigators within a consortium/ network)

## Broad Sharing

- Ensures fair and equitable access and secondary use of data by the wider research community (e.g., NIH's Genomic Data Sharing Policy)
- Has the most impact on *driving scientific innovation and discovery and ensuring replication of results*
- Broad sharing ≠ Open access data

# Framingham Study: Success in Data Collection Over Time



**BY THE NUMBERS:** Uncovering the Mysteries of the Heart

By American Heart Association News

**70** — Years the Framingham Heart Study continues to break new ground on cardiovascular disease

**3** — Generations who have participated in the study

**15,447** — Participants over the past 70 years

**1960** — Year the study pinned cigarette smoking as a risk factor for heart disease

**5,209** — Initial volunteer participants
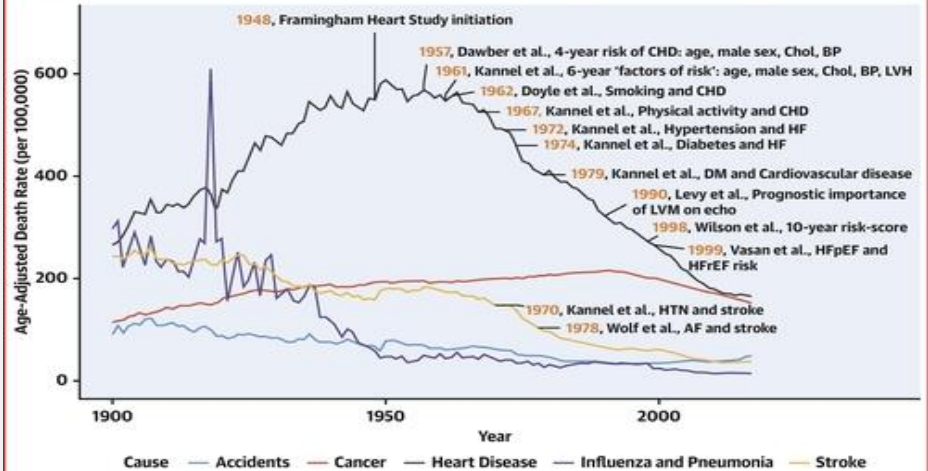
**3,698** — Published journal articles based on Framingham Heart Study data
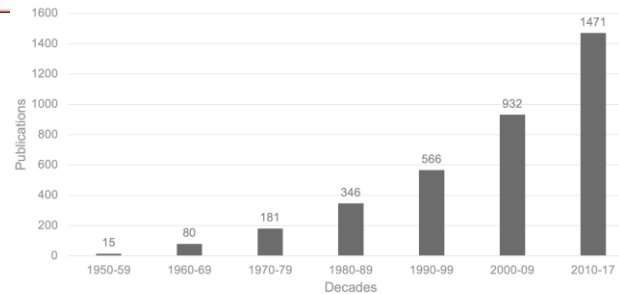
**802** — Participants who have donated or registered to donate their brain for further study

Sources: Framingham Heart Study, Boston University
Published Oct. 10, 2018

**CENTRAL ILLUSTRATION:** Age-Adjusted Death Rates for the Leading Causes of Death in the United States and the Framingham Heart Study

1948, Framingham Heart Study initiation
1957, Dawber et al., 4-year risk of CHD: age, male sex, Chol, BP
1961, Kannel et al., 6-year "factors of risk": age, male sex, Chol, BP, LVH
1962, Doyle et al., Smoking and CHD
1967, Kannel et al., Physical activity and CHD
1972, Kannel et al., Hypertension and HF
1974, Kannel et al., Diabetes and HF
1979, Kannel et al., DM and Cardiovascular disease
1990, Levy et al., Prognostic importance of LVM on echo
1998, Wilson et al., 10-year risk-score
1999, Vasan et al., HFpEF and HFrEF risk
1970, Kannel et al., HTN and stroke
1978, Wolf et al., AF and stroke

Cause — Accidents — Cancer — Heart Disease — Influenza and Pneumonia — Stroke

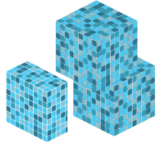Andersson, C. et al. J Am Coll Cardiol. 2021;77(21):2680-92.

Total articles published through November 2017 = 3,561

# The Cancer Genome Atlas: Success in Open Team Science

## TCGA BY THE NUMBERS

TCGA produced over

**2.5**
PETABYTES
of data

To put this into perspective, **1 petabyte** of data is equal to

**212,000**
DVDs

TCGA data describes

**33**
DIFFERENT
TUMOR TYPES

...including

**10**
RARE
CANCERS

...based on paired tumor and normal tissue sets collected from

**11,000**
PATIENTS

...using

**7**
DIFFERENT
DATA TYPES

## TCGA RESULTS & FINDINGS

| | | | |
|---|---|---|---|
| MOLECULAR BASIS OF CANCER | Improved our understanding of the genomic underpinnings of cancer | For example, a TCGA study found the basal-like subtype of breast cancer to be similar to the serous subtype of ovarian cancer on a molecular level, suggesting that despite arising from different tissues in the body, these subtypes may share a common path of development and respond to similar therapeutic strategies. |
| TUMOR SUBTYPES | Revolutionized how cancer is classified | TCGA revolutionized how cancer is classified by identifying tumor subtypes with distinct sets of genomic alterations.* |
| THERAPEUTIC TARGETS | Identified genomic characteristics of tumors that can be targeted with currently available therapies or used to help with drug development | TCGA's identification of targetable genomic alterations in lung squamous cell carcinoma led to NCI's Lung-MAP Trial, which will treat patients based on the specific genomic changes in their tumor. |

## THE TEAM

**20**
COLLABORATING
INSTITUTIONS
across the United States
and Canada

*TCGA's analysis of stomach cancer revealed that it is not a single disease, but a disease composed of four subtypes, including a new subtype characterized by infection with Epstein-Barr virus.
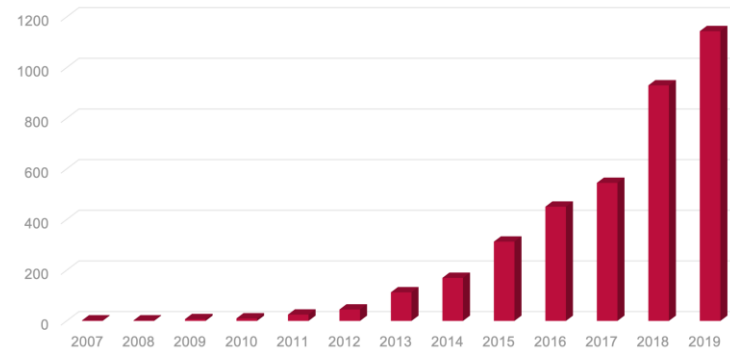
## WHAT'S NEXT?

The Genomic Data Commons (GDC) houses TCGA and other NCI-generated data sets for scientists to access from anywhere. The GDC also has many expanded capabilities that will allow researchers to answer more clinically relevant questions with increased ease.
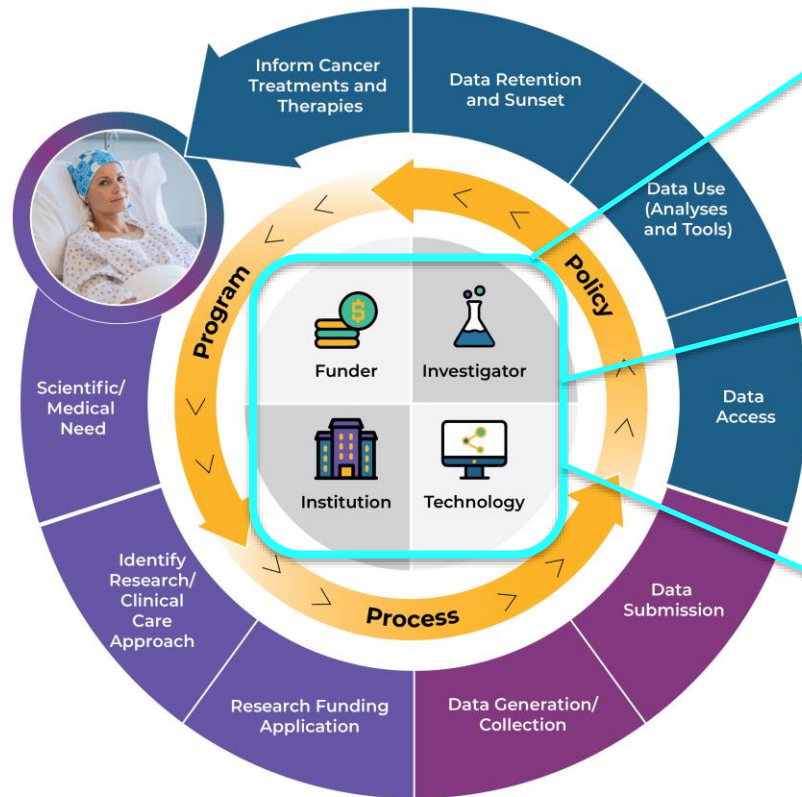
www.cancer.gov/ccg

## Number of Publications Using TCGA Data

# Scientific Data Lifecycle: Keys to Impactful Discovery



**Critical Questions to Answer**

Programs that define therapeutic needs and essential scientific gaps to be filled using structured datasets.
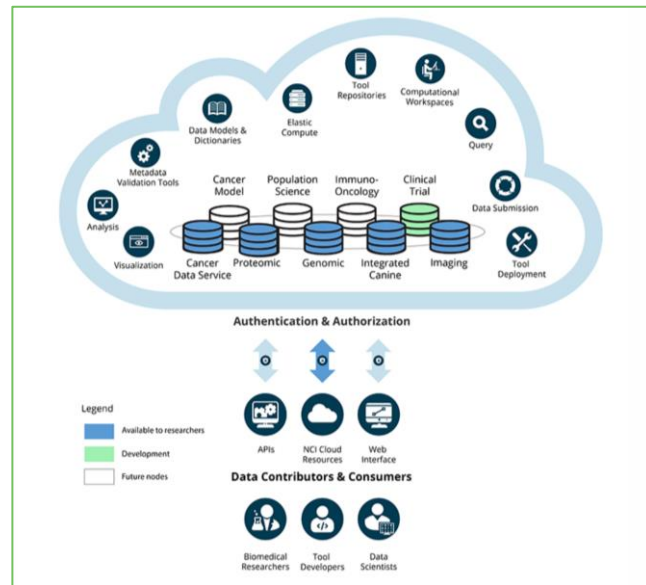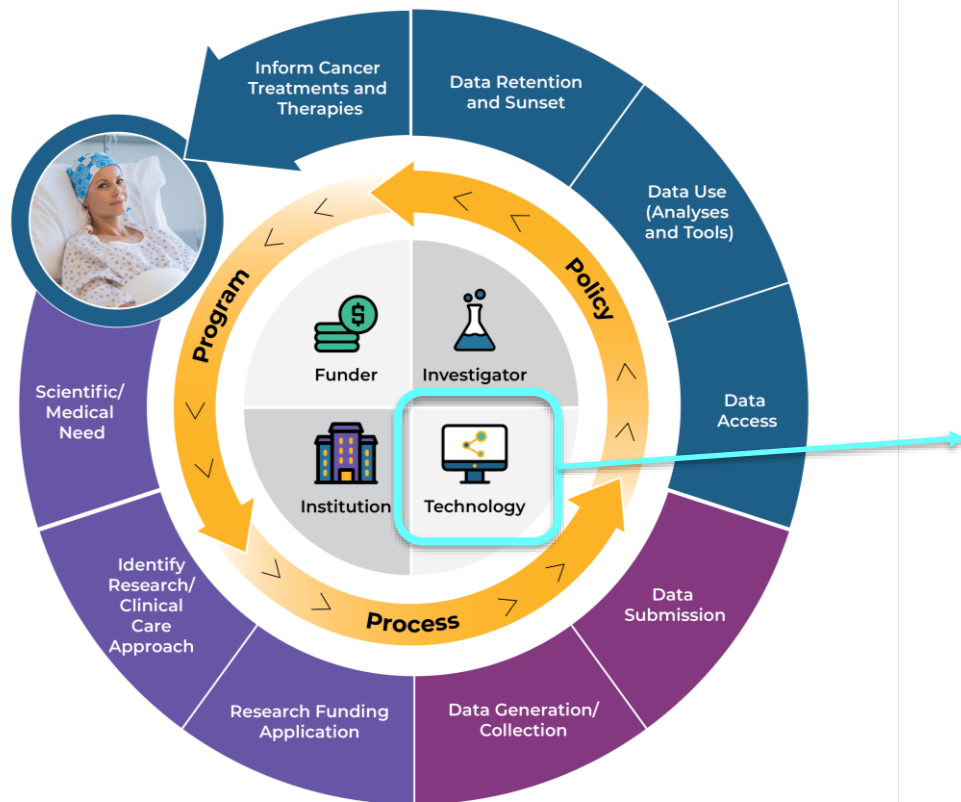
**Policies to Promote Broad Use**

Implementation of aggressive data management, sharing and access policies that ensure rapid, free and immediate access to all types of data.

**Infrastructure to Support FAIR Principles**

Technology platforms and tools that employ standards to make data findable, accessible, interoperable and reusable.

# Sharing Data Openly through a Cancer Data Ecosystem



Cancer Research Data Ecosystem

# National Data Ecosystem: Integrating Cancer Research

# NIH Guidance for the New DMS Policy

**Recommended DMS Plan Elements**

**Selecting Data Repositories**

**Allowable Costs of Data Management & Sharing**

# How to "Plan" for NIH Data Management & Sharing

## Data Collection

**_What data types will be generated in your research project?_**

_~Align w/ Specific Aims in application (e.g. omics, clinical, imaging, flow cytometry)_

## Data Management

**_How will those data be managed?_**

_~Outline plan to maintain data throughout project (e.g. lab, Institute data core, DCC, hard drive/ cloud)_

## Data Sharing

**_Which data filesare most useful for broad uses in research community?_**

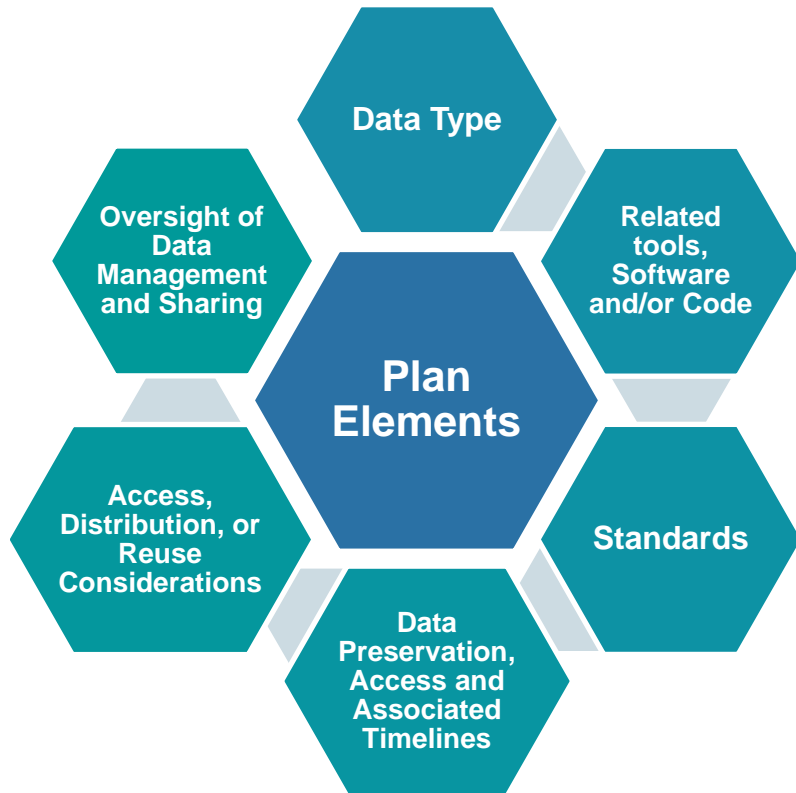_~Determine data needed for new research projects to avoid duplication or supplement experiments_

## Plan Writing

**_Are all DMS Plan elements addressed (data, repositories, costs)?_**

_~Submit all costs in Budget Justification (e.g. data formatting, curation)_

**_*Not all managed data need to be shared for secondary use_**

# Six Elements of A DMS Plan



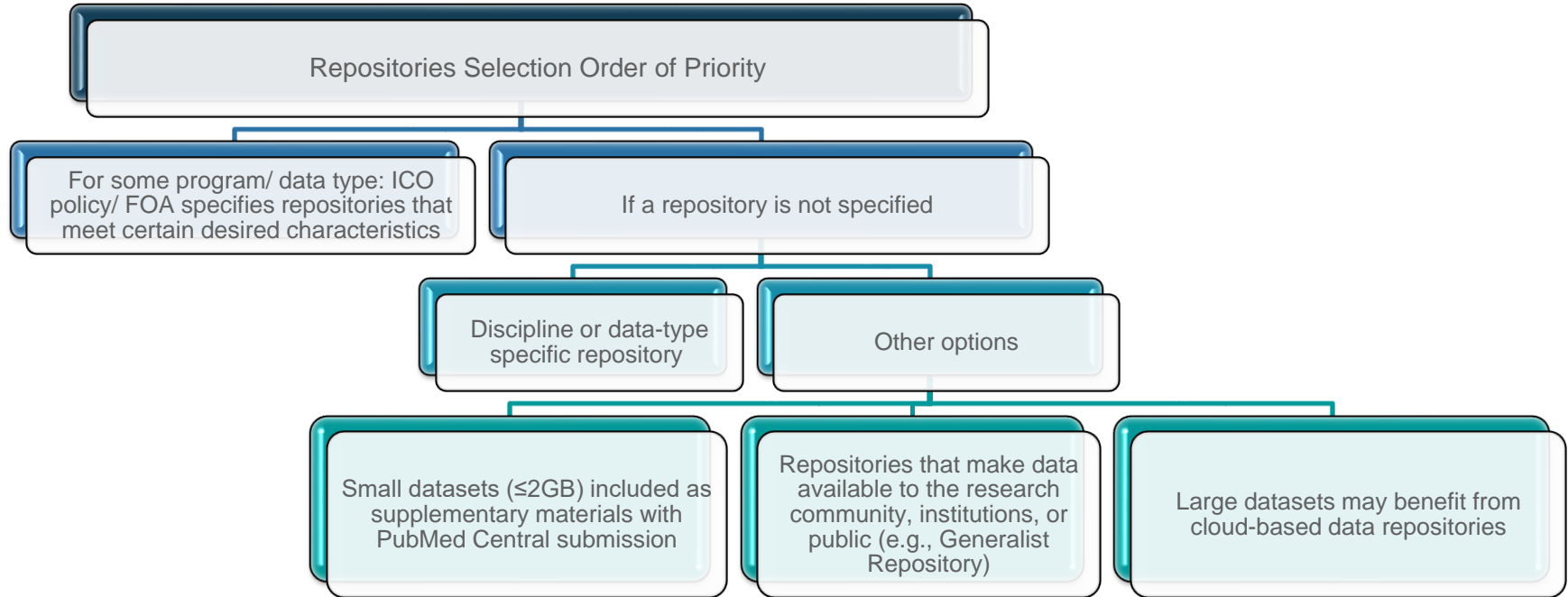See Writing a Data Management & Sharing Plan for details

- **Data type**
  - Identifying data to be preserved and shared and listing metadata
- **Standards**
  - Standards to be applied to scientific data and metadata (i.e., data formats, data dictionaries, data identifiers, definitions, unique identifiers, and other data documentation)
- **Data preservation, access, timelines**
  - Repository to be used, persistent unique identifier, and when/ how long data will be available
- **Access, distribution, reuse considerations**
  - Description of factors for data access, distribution, or reuse
- **Related tools, software, code**
  - Tools and software needed to access and manipulate data
- **Oversight of data management and sharing**
  - Indicates how compliance with the DMS plan will be monitored and managed

# Selecting A Data Repository

*Promoting the use of established data repositories to improve FAIRness of data*

```
Repositories Selection Order of Priority
```

```
For some program/ data type: ICO policy/ FOA specifies repositories that meet certain desired characteristics
```

```
If a repository is not specified
```

```
Discipline or data-type specific repository
```

```
Other options
```

```
Small datasets (≤2GB) included as supplementary materials with PubMed Central submission
```

```
Repositories that make data available to the research community, institutions, or public (e.g., Generalist Repository)
```

```
Large datasets may benefit from cloud-based data repositories
```

***What is not a Repository? Supplemental material in a journal or Lab website***

# NIH Resources to Guide Repository Selection

## NIH-Supported Repositories

- Filterable list of 75+ NIH Repositories that are open to taking in data of various types & formats

- PubMed holds up to 2Gb of data (as supplemental information to manuscripts)

| Institute or Center | Repository Name | Repository Description |
|---|---|---|
| NCI | | Keyword Filter |
| NCI | Cancer Nanotechnology Laboratory (caNanoLab) | caNanoLab is a data sharing portal designed to facil... sharing in the biomedical nanotechnology research expedite and validate the use of nanotechnology in caNanoLab provides support for the annotation of n characterizations resulting from physico-chemical, i assays and the sharing of these characterizations an |

## Other Repository Resources

- Generalist repositories (9 listed by NIH as part of GREI)

- Nature's Data Repository Guidance

- Registry of Research Data Repositories

re3data.org
REGISTRY OF RESEARCH DATA REPOSITORIES

See **Repositories for Sharing Scientific Data**

# Data Management and Sharing Costs

## ALLOWABLE COSTS

- Curating data/developing supporting documentation

- Preserving/sharing data through repositories

- Local data management considerations

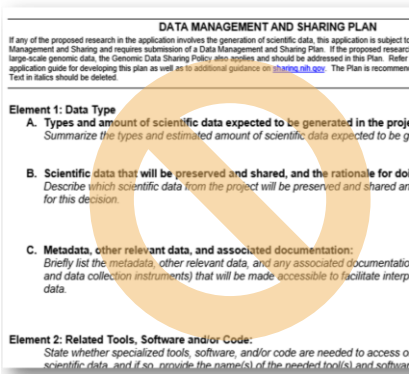- *IMPORTANT*: Must be incurred during the performance period

## UNALLOWABLE COSTS

- Infrastructure costs typically included in indirect costs

- Costs associated with the routine conduct of research (e.g., costs of gaining access to research data)

- *Additional Example*: Data storage costs in NIH/ NCI-supported data repositories

See **Budgeting for Data Management & Sharing** for details

# Responsibilities and Expectations: Peer Reviewers



- Peer Review Will NOT assess DMS Plans*

  - DMS Plan assessment is an administrative issue, reviewers assess only scientific merit of applications

- Peer Review Will NOT see full Plan or be providing comments on Plans*

  - DMS Plan will be excluded from the application image shown to peer reviewers

- Peer Review will NOT factor DMS Plans into their impact scores*

*Exceptions for FOAs where data sharing is core to the science, e.g., FOAs for data dissemination centers)*

# What Should a DMS Plan Look Like (OER)?

- ✓ Plans should be no more than 2 pages in length

- ✓ Applications subject to both DMS and GDS Policies submit a single Plan

- ✓ Optional format page available

- ✓ Exploring structured form in future *(NCI & NICHD have each developed forms; possible pilot)*

**DATA MANAGEMENT AND SHARING PLAN**

If any of the proposed research in the application involves the generation of scientific data, this application is subject to the NIH Policy for Data Management and Sharing and requires submission of a Data Management and Sharing Plan. If the proposed research in the application will generate large-scale genomic data, the Genomic Data Sharing Policy also applies and should be addressed in this Plan. Refer to the detailed instructions in the application guide for developing this plan as well as to additional guidance on sharing.nih.gov. The Plan is recommended not to exceed two pages. Text in italics should be deleted.

**Element 1: Data Type**
  A. **Types and amount of scientific data expected to be generated in the project:**
     *Summarize the types and estimated amount of scientific data expected to be generated in the project,*

  B. **Scientific data that will be preserved and shared, and the rationale for doing so:**
     *Describe which scientific data from the project will be preserved and shared and provide the rationale for this decision.*

  C. **Metadata, other relevant data, and associated documentation:**
     *Briefly list the metadata, other relevant data, and any associated documentation (e.g., study protocols and data collection instruments) that will be made accessible to facilitate interpretation of the scientific data.*

**Element 2: Related Tools, Software and/or Code:**
     *State whether specialized tools, software, and/or code are needed to access or manipulate shared scientific data, and if so, provide the name(s) of the needed tool(s) and software and specify how they*

*DMS Plan format page will be added to list of Format Pages and incorporated into FORMS-H application instructions by Fall 2022*

# A Decision Support Tool to Guide Plan Assessment

- <u>Optional</u> tool to guide POs through DMS Plan assessment

- Stepwise list of considerations for each element

- Additional ICO or program-specific guidance may also apply

**Coming Soon**

| ELEMENT | PROPOSED ASSESSMENT CRITERIA |
|---|---|
| Element 1, Data Type | • ☐ Y ☐ N - Does the Plan identify the types and estimated amounts of scientific data to be generated and/or used in the research (e.g., 256-channel EEG data and fMRI images from ~50 research participants)? <br><br> • ☐ Y ☐ N - Does the Plan identify which scientific data will be preserved and shared with the scientific community? <br><br> • ☐ Y ☐ N - Has a rationale been provided justifying the decision of which scientific data will be preserved and shared, based on ethical, legal, and technical factors? <br><br> • ☐ Y ☐ N - Does the Plan list the metadata, other relevant data, and any associated documentation (e.g., study protocols and data collection instruments) that will be made accessible to facilitate interpretation of the scientific data? <br><br> **For research subject to the GDS policy:** <br><br> • ☐ Y ☐ N – Does the Plan describe whether the proposed research involves human data, non-human data, or both. <br><br> • ☐ Y ☐ N – Does the Plan list the type(s) and amounts of genomic data that will be shared (e.g., sequence, transcriptomic, epigenomic, and/or gene expression data) and whether it is individual-level data, aggregate-level data, or both)? <br><br> **If the answer is "no" for ANY of the questions above, more information is needed from the applicant to accurately assess the Plan** |

# Harmonization of DMS and Genomic Data Sharing Plans

| Investigators | Program Staff |
|---|---|
| • One Plan: Genomic data sharing considerations to be addressed in DMS Plans using the expected DMS Plan elements.<br>• Timeline: At the time of application or Just-in-Time.<br>• *For human data subject to GDS*: Applicants should complete the DMS Plan anticipating sharing according to the criteria in the <u>Institutional Certification (IC).</u> | ▪ IC or provisional certification must be submitted and accepted before issuing the award. Investigators should state in the DMS Plan what data can be shared and how if IC criteria cannot be met.<br>▪ Plans will be reviewed by Program Staff. Peer reviewers will no longer comment on the Plans but rather on the reasonableness of the DMS budget.<br>▪ Submission timelines for non-human and human data subject to GDS will remain unchanged.<br>▪ Compliance and enforcement terms for awards subject to GDS will be handled in accordance with that under DMS (i.e., "end of the performance period" is the *latest* opportunity to submit data). |

**Notice Number:** <u>NOT-OD-22-198</u>

# How Do the DMS and Genomics Data Policies Compare?

| | **2015 NIH GDS Policy** | **2023 NIH DMS (New Policy)** |
|---|---|---|
| **Effective Dates** | ▪ **New applications**: received on/after January 25, 2015 | ▪ **New applications & Competitive renewals** (Type 2, previously funded awards): For submission on/after January 25, 2023 |
| **Scope** | ▪ <u>**Large-scale genomics data**</u> & some smaller studies *(programmatic priority, rare diseases)*<br>▪ **Thresholds** & **assay types** defined | ▪ All NIH-supported research that generates <u>**scientific data**</u><br>▪ Not applicable to other activities (e.g. training, infrastructure development)<br>▪ No detailed expectations defined outside of GDS, CT or programs |
| | ▪ Intramural & Extramural research mechanism (grants, cooperative agreements, contracts, other transactions)<br>▪ Human and non-human data; no budget threshold | |
| **NIH Expects (for PIs)** | ▪ Submit & release data through NIH repositories<br>▪ Use guidelines for standard formats, levels of data | ▪ **Prospectively plan** how to preserve & share scientific data<br>▪ Outline in **Data Management & Sharing (DMS) Plan** (w/ applications), and<br>▪ Report data sharing progress in annual submission of **RPPR** |
| | *<u>**Only one plan will be submitted**</u>*; **DMS plans** will **include genomics** elements along with other data types (i.e. no separate GDS plan submitted after January 25, 2023) | |
| **Timelines** | ▪ Lower-level primary data: released by <u>**9 months**</u> (full, QC dataset generation)<br>▪ Summary analyses at publication | ▪ Shared scientific data: **as soon as possible** → <u>by time of associated publication, or end of performance period</u> (whichever comes first) |