# High Fidelity:

**Connecting information for Better Research Reproducibility**

Terrie Wheeler, AMLS
Director
Samuel J. Wood Library
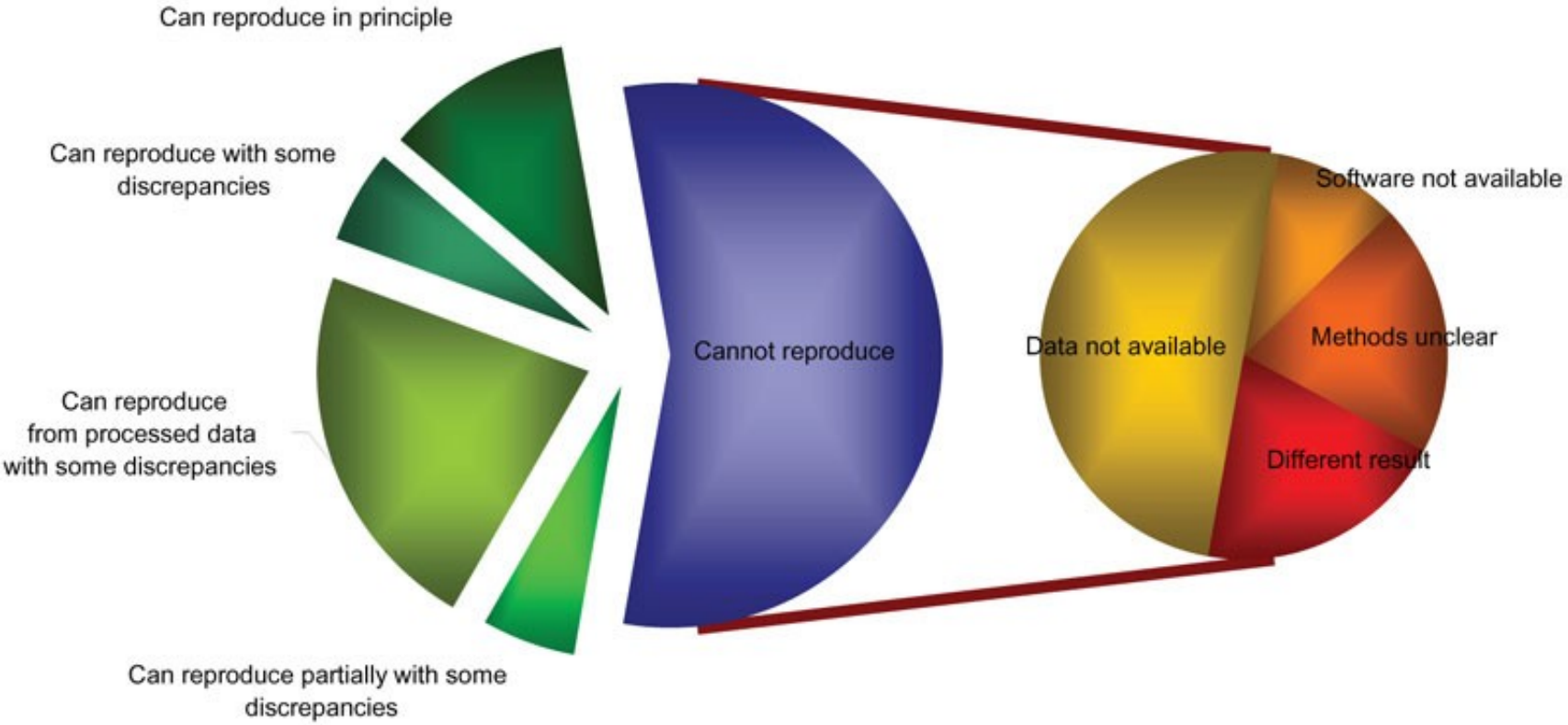
Peter Oxley, PhD
Associate Director for Research Services
Samuel J. Wood Library

# There are three fires burning
## that motivate our actions towards better reproducibility



1. "Reproducibility Crisis"

2. Data Retention Mandates

3. Allegations of Misconduct

# Attempting to reproduce research proves to be difficult, when you can't get the data...



Can reproduce in principle

Can reproduce with some discrepancies

Can reproduce from processed data with some discrepancies

Can reproduce partially with some discrepancies

Cannot reproduce

Data not available

Software not available

Methods unclear

Different result

Ioannidis *et al.* (2009) Nat. Gen. 41:149

"OMB Circular A–110 states that the retention period is th[e]
from the date the final financial report is submitted."

"NSF states in its General Grant Conditions that records must be retained
for three years after the submission of all required reports"

"in the case of research misconduct involving NIH fund[s]
records must be retained for six years
after the final resolution date of the case."

"retain research data pertinent to patented inventions
for the life of the patent"

"about 2% of scientists admitted to have fabricated, falsified or mod
 data or results at least once"
"Up to one third admitted a variety of other  questionable practices"

Fanelli (2009) PLOS One 0005738

"3.8% of published papers contained problematic figures, with at lea
 exhibiting features suggestive of deliberate manipulation"

Bik *et al.* (2016) mBio 00809–16



Brainard and You (2018) Science 00809–16

# There are three fires burning
that motivate our actions towards better reproducibility



1. "Reproducibility Crisis"

2. Data Retention Mandates

3. Allegations of Misconduct

Maximizing data value and ethical research conduct

We are not automatically making research reproducible,
but providing the infrastructure that helps make research re

**Reporting Standards**  **Principles of reproducibility**

PRISMA  FAIR

MIAME  Literate Programming (Knuth 1992)

MIQE  Compendiums (Gentleman and Lang, 2003)

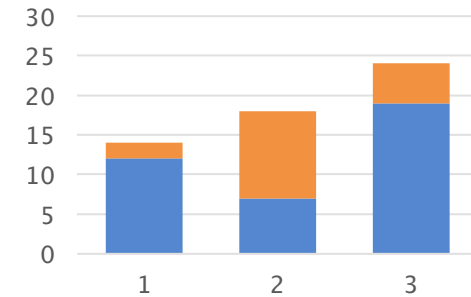**Data Publication standards** **Communities of practice**

Open Science  The Turing Way

# How can we help researchers capture
# the data and workflows that lead to publishable results?



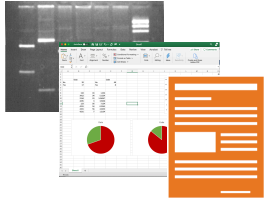**DATA**  **WORKFLOWS**  **RESULTS**

Hurdle #1: Dealing with confidential data

Hurdle #2: Researchers are very busy

Hurdle #3: Data and workflows are very diverse

Hurdle #4: Maintaining a high quality solution

# Piece #1: Electronic Lab Notebooks capture (small) data and while remaining flexible to researcher needs

 Direct storage of images, data files, analysis files, workflows

 File versioning and immutable timestamps

 Integration with Jupyter Notebooks (on roadmap)

 Shareable and transferable

# Piece #2: A file management system for tagging, tracking, a then archiving files on the institutional storage systems

Data associated with a project/publication assigned a unique ta
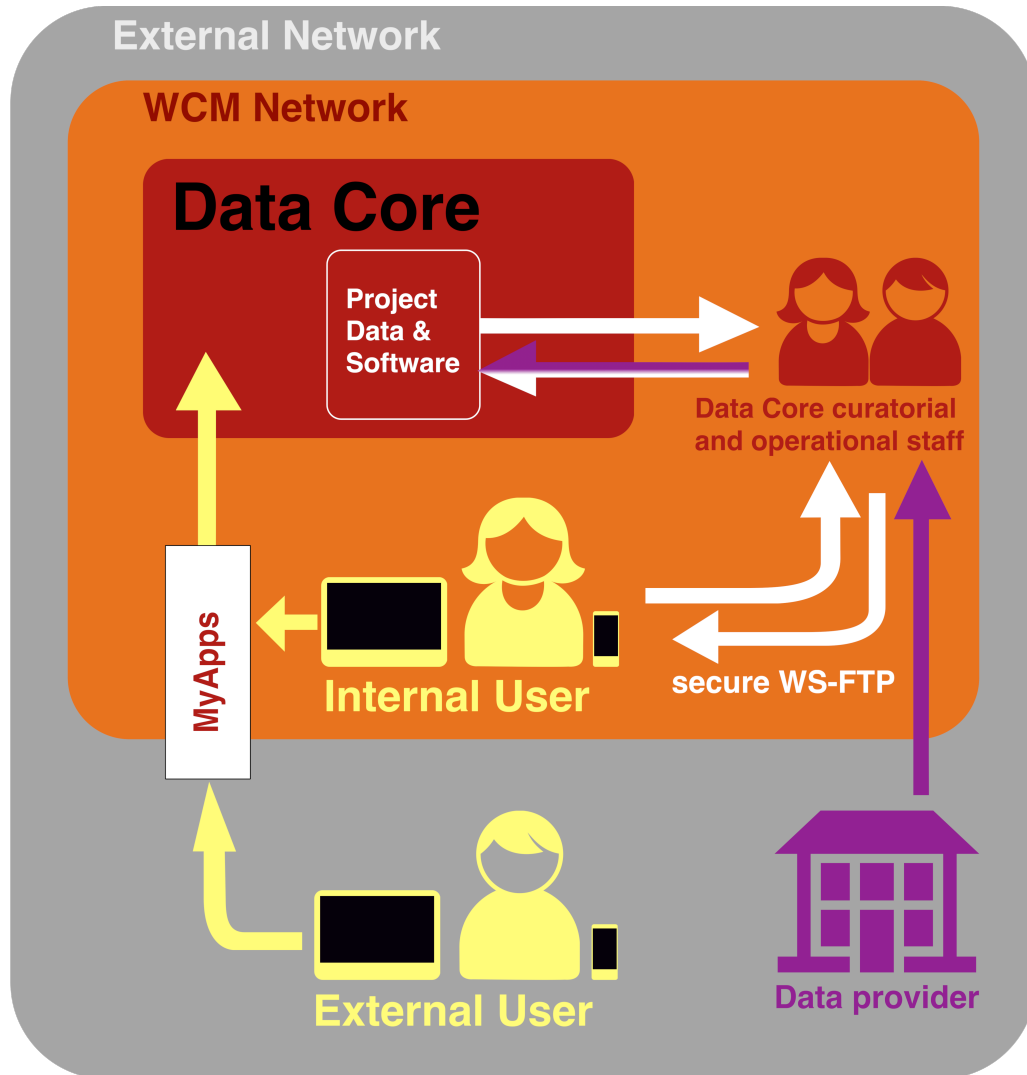
File identity managed by hashing

File location tracked within institutional storage

Actionable scripts for marking project complete, and archiving

# Piece #3: Secure file access management and computation the institutional Data Core
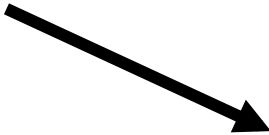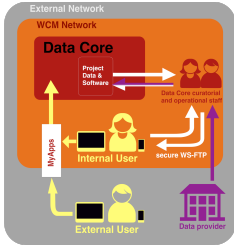


Secured, collaborative, flexible

Project governance and monitoring

Data curation for import and export

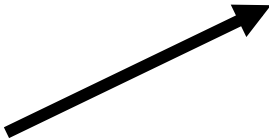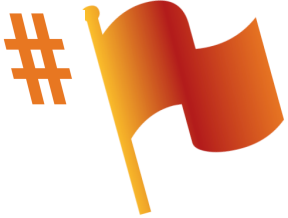# Holding the pieces together: an institutional Data Catalog that connects data, workflows, governance and access condi



**COVID-19 Data Lake**

Data publisher: Architecture for Research Computing in Health (ARCH)

A data lake to support the study of patients tested for SARS-CoV-2. The set includes standard electronic health record (EHR) data in addition to novel collection of SARS-CoV-2 data via REDCap. Available EHR data includes: patient demographics, patient visits, medications, prescriptions, transfers, labs, orders/procedures, diagnoses, observations, vitals and computed scores, notes (both outpatient and discharge summary), ventilator usage, transfusions. Available SARS-CoV-2 data, for patients that tested positive for the virus, includes: HIV regimen, pregnancy and gestational age, immunosuppression (binary variable and details), date of first symptoms and details of symptoms, prior utilization (was patient discharged from other hospital or ED before

| | |
|---|---|
| Local contact | Tom Campion (thc2015) |
| Record Period Start | Jan. 1, 2015 |
| Record Period End | None |
| Publication Date | April 13, 2020 |
| Confidentiality Impact Level | WCM ITS 11.03: High Risk (Confidential) |
| Data Creators | Architecture for Research Computing in Health (ARCH) |
| Media Types | |
| Landing URL | https://medcornell.sharepoint.com/sites/covid_arch/idr/SitePages/About.aspx |
| Keywords | ARCH COVID |

Discovery layer with the capacity to connect grants, data, publication

# The WCM Data Catalog was built to manage data governance and access conditions

Scope of authorization

User authorization

Data Controls

Reuse scope

# Three triggers to prompt capture and storage



1. Project/Grant completion

2. Publication

3. Faculty member leaves the institution

Move files to archive
Register project, file tag, ELNs in Data Catalog

# Acknowledgements:

Curt Cole (WCM CIO)

Data Core
Alice Chin
Ana Proper
Frank Ashmun
John Ruffing
Jo Hargitai
Lucy Walle
Michael Bales
Heather Kleinschmidt

Electronic Lab Notebooks
Marie Linvill (LabArchives)
Cindy Chen
Danny Tan
Lucy Walle
Tony DiFazio

Data Catalog
The Data Catalog Collaboration (DDC)
Nicole Contaxis (NYU Langone)
Tom Campion

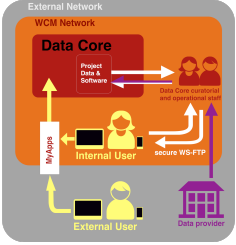File System Management
Cindy Chen
Darrin Stivala (Starfish)
Doug Hughes (Starfish)
Jo Hargitai
Tony DiFazio

Peter Oxley  pro2004@med.cornell.edu
Terrie Wheeler tew2004@med.cornell.edu

External Network
WCM Network
Data Core
Project Data & Software
Data Core curatorial and operational staff
My/Apps
Internal User
secure WS-FTP
External User
Data provider

Data Catalog (Beta)   View Details ▾          pro2004 ▾   Submit New: ▾   Search catalog   Submit

COVID-19 Data Lake                                     Update dataset details

Data publisher: Architecture for Research Computing in Health (ARCH)

A data lake to support the study of patients tested for SARS-CoV-2. The set includes standard electronic health record (EHR) data in addition to novel collection of SARS-CoV-2 data via REDCap. Available EHR data includes: patient demographics, patient visits, medications, prescriptions, transfers, labs, orders/procedures, diagnoses, observations, vitals and computed scores, notes (both outpatient and discharge summary), ventilator usage, transfusions. Available SARS-CoV-2 data, for patients that tested positive for the virus, includes: HIV regimen, pregnancy and gestational age, immunosuppression (binary variable and details), date of first symptoms and details of symptoms, prior utilization (was patient discharged from other hospital or ED before

| | |
|---|---|
| Local contact | Tom Campion (thc2015) |
| Record Period Start | Jan. 1, 2015 |
| Record Period End | None |
| Publication Date | April 13, 2020 |
| Confidentiality Impact Level | WCM ITS 11.03: High Risk (Confidential) |
| Data Creators | Architecture for Research Computing in Health (ARCH) |
| Media Types | |
| Landing URL | https://medcornell.sharepoint.com/sites/covid_arch/idr/SitePages/About.aspx |
| Keywords | ARCH  COVID |