

BIOMEDICAL DATA MANAGEMENT AT SCALE

the rise of data commons and the quest for FAIRness

A vertical decorative panel on the left side of the slide features a microscopic view of cells and viruses. The background is a gradient of light blue to white. Several spherical cells with small protrusions are visible, along with a larger, more complex structure that resembles a virus or a specialized cell. The overall aesthetic is clean and scientific.

AGENDA

1

Data sharing is difficult

2

Biomedical big data and data commons

3

Metadata and FAIRness

4

The urgency of COVID-19

SHARING BIOMEDICAL DATA IS DIFFICULT



The background of the slide features a dark blue gradient with vertical lines of varying lengths and semi-transparent circles, creating a digital or data-centric aesthetic. In the center, there are faint, light blue silhouettes of server racks or data storage units. Overlaid on this background is a vertical column of binary code (0s and 1s) in a light blue color, with some characters appearing larger and more prominent than others.

Sharing data is difficult

- Data are rarely collected for reuse
- The effort to make data reusable falls the data producer; the benefit accrues to the data consumer
- Incentives are not aligned to support data reuse

Challenges with biomedical data



Privacy and regulatory concerns over data from human subjects



Clinical data generally collected for treatment, not reuse



Biologists tend not to converge on standards



Biomedical Big Data

Cost per Human Genome



NIH National Human Genome Research Institute

genome.gov/sequencingcosts

NATIONAL CANCER INSTITUTE THE CANCER GENOME ATLAS

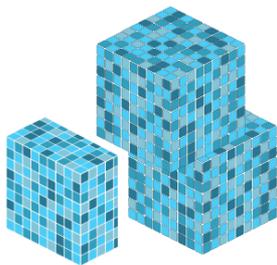
TCGA BY THE NUMBERS

TCGA produced over

2.5

PETABYTES

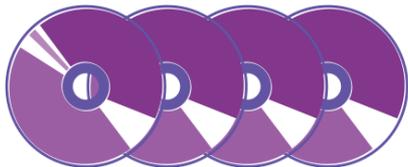
of data



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



2006-2015

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.

2015

portal.gdc.cancer.gov

NIH NATIONAL CANCER INSTITUTE GDC Data Portal

Home Projects Exploration Analysis Repository Quick Search Manage Sets Login Cart GDC Apps

Harmonized Cancer Datasets

Genomic Data Commons Data Portal

Get Started by Exploring:

Projects Exploration Analysis Repository

Q e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2

Data Portal Summary

[Data Release 26.0 - September 08, 2020](#)

PROJECTS 67	PRIMARY SITES 68	CASES 84,375
FILES 590,367	GENES 23,399	MUTATIONS 3,287,299

Cases by Major Primary Site

Primary Site	Cases (approximate)
Adrenal Gland	100
Bile Duct	100
Bladder	100
Bone	100
Bone Marrow	900
Brain	100
Breast	900
Cervix	100
Colorectal	800
Esophagus	100
Eye	100
Head and Neck	100
Kidney	300
Liver	300
Lung	1100
Lymph Nodes	1100
Nervous System	300
Ovary	300
Pancreas	200
Pleura	100
Prostate	100
Skin	200
Soft Tissue	100
Stomach	200
Testis	100
Thymus	100
Thyroid	100
Uterus	200

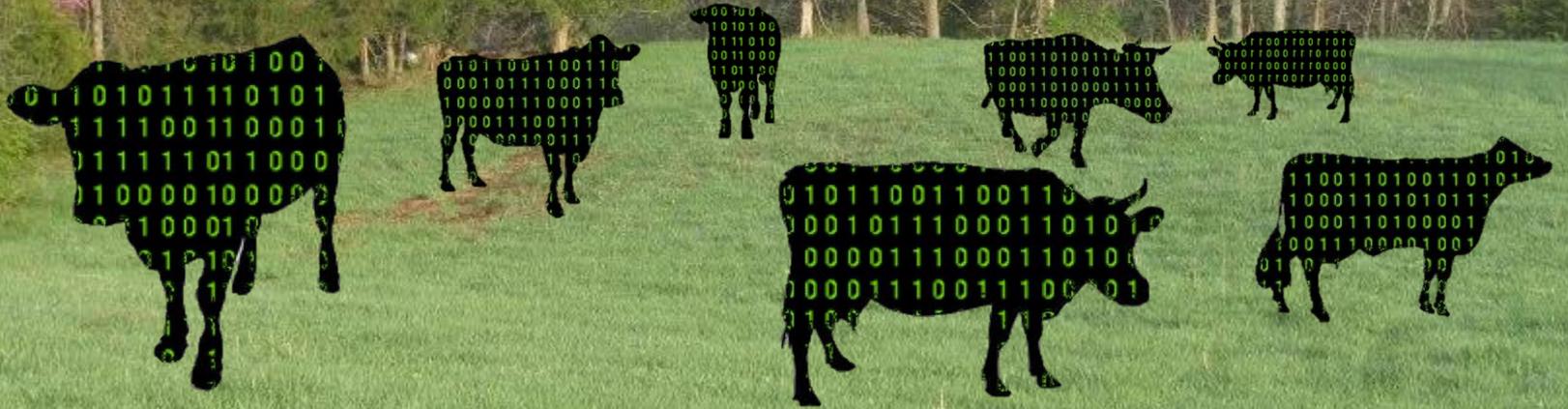
GDC Applications

The GDC Data Portal is a robust data-driven platform that allows cancer researchers and bioinformaticians to search and download cancer data for analysis. The GDC applications include:

- Data Portal
- Website
- API
- Data Transfer Tool
- Documentation
- Data Submission Portal
- Legacy Archive
- Publications

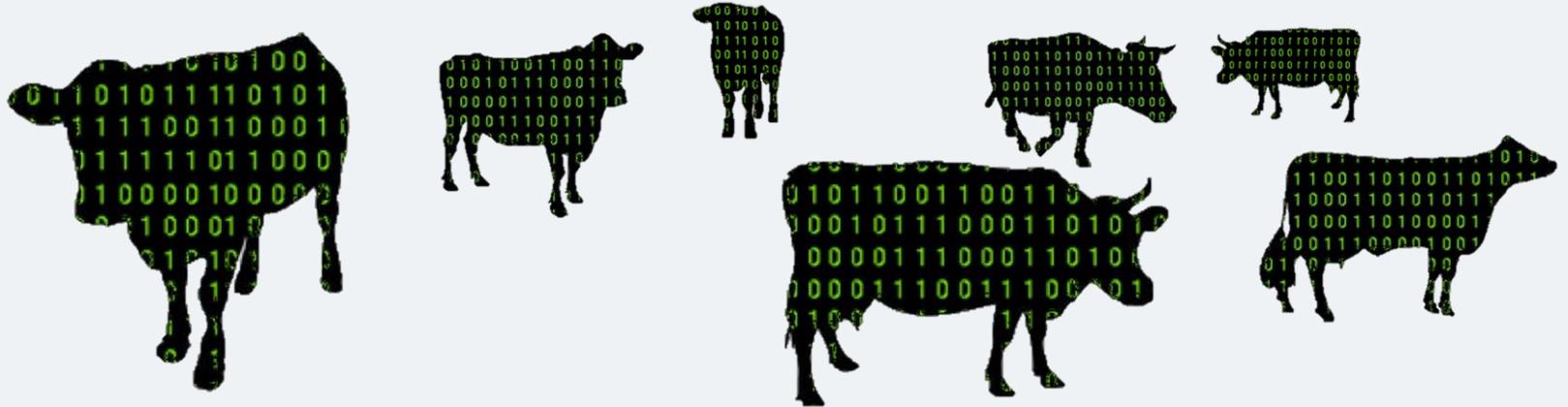
2016

What is a Data Commons?



A Functional Definition

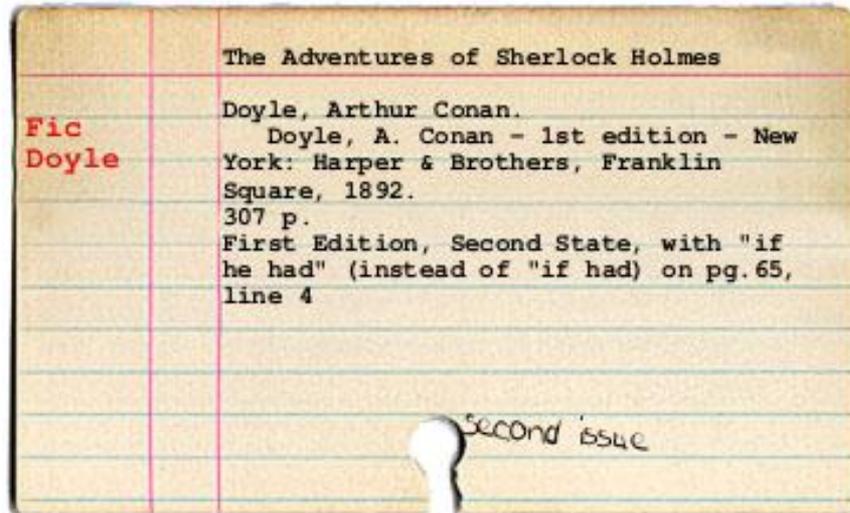
“Data commons are cloud-based software platforms that co-locate: 1) data, 2) computing infrastructure, and 3) commonly used software applications, tools and services to create a resource for managing, analyzing, integrating and sharing data with a community.”



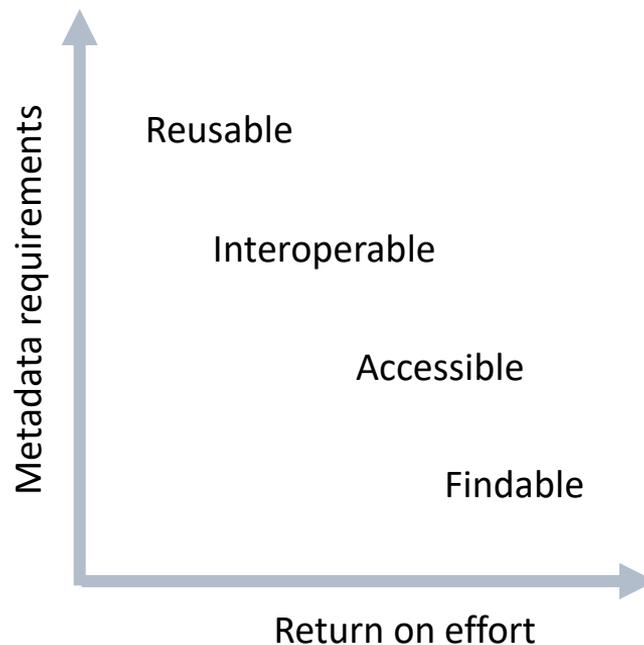


”Our guiding principles for data access, use, and reuse will adhere to the Findable, Accessible, Interoperable, and Reusable (FAIR) guidelines. While the FAIR guidelines are well-accepted, to bring them to practice will require defining and adopting community-based metrics and rubrics so these can be applied to data, and other types of digital objects, hosted within or available through the Data Commons. At the same time, once FAIR metrics and rubrics are defined, these will be used to measure the level of “FAIRness” of repositories, datasets, and other digital objects. Such evaluations will inform and engage both Data Commons users and digital objects producers.”

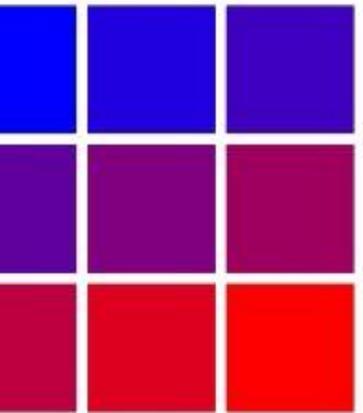
#metadata and FAIRness



Metadata & FAIRness

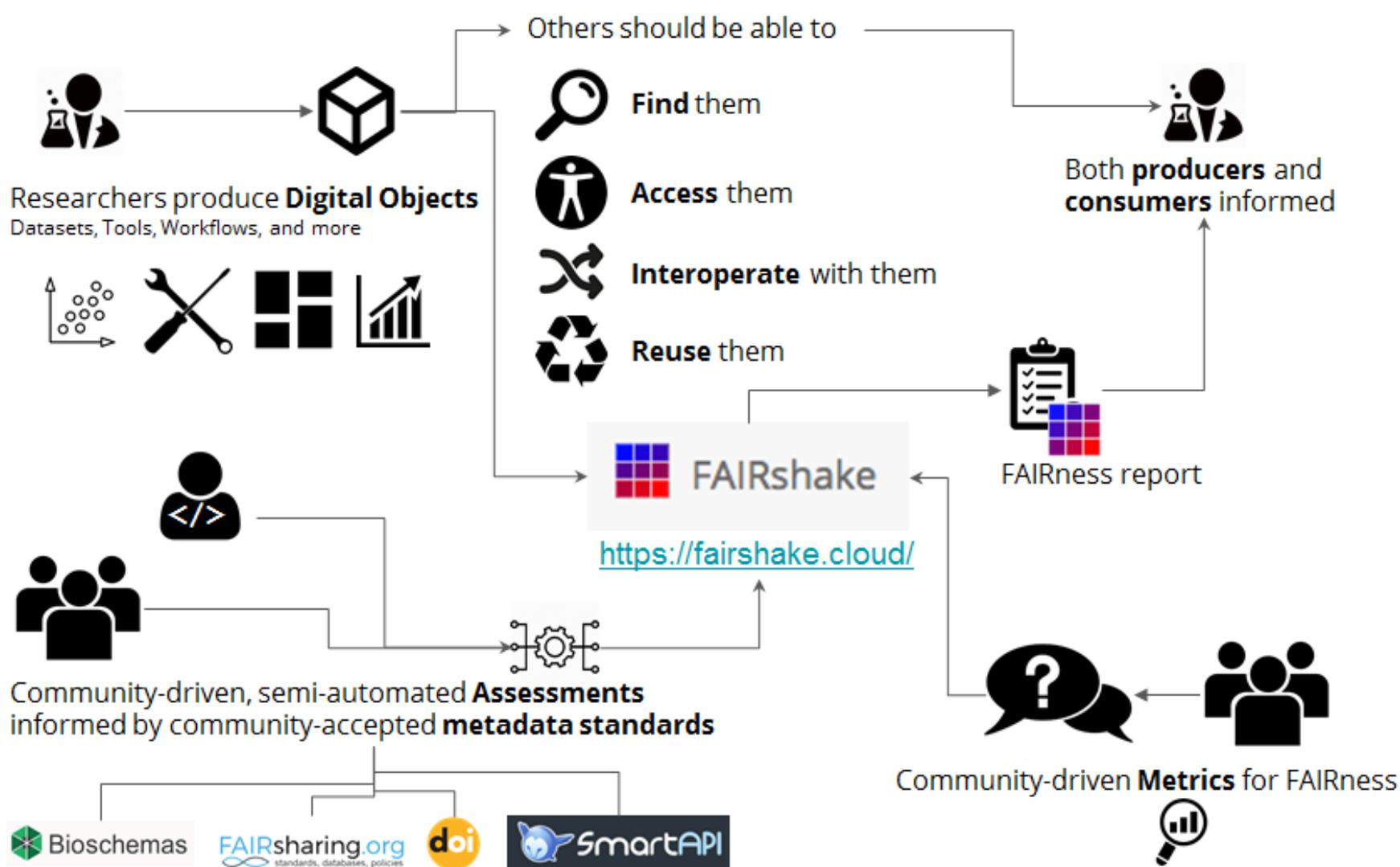


Quantifying FAIRness



FAIRshake

2019





Commentary

FAIRshake: Toolkit to Evaluate the FAIRness of Research Digital Resources

Daniel J.B. Clarke¹, Lily Wang¹, Alex Jones², Megan L. Wojciechowicz¹, Denis Torre¹, Kathleen M. Jagodnik¹, Sherry L. Jenkins¹, Peter McQuilton³, Zachary Flamholz¹, Moshe C. Silverstein¹, Brian M. Schilder¹, Kimberly Robasky⁴, Claris Castillo⁴, Ray Idaszak⁴, Stanley C. Ahalt⁴, Jason Williams⁵, Stephan Schurer⁶, Daniel J. Cooper⁶ ... Avi Ma'ayan¹  

Show more 

<https://doi.org/10.1016/j.cels.2019.09.011>

[Get rights and content](#)

<https://www.alliancegenome.org/>



Score: 78%
Information is provided describing how to cite the dataset.

[View Associations](#)

[View Assessments](#)



FAIRshake



The goal of Common Fund Data Ecosystem CFDE is to federate data from a number of Common Fund Data Coordinating Centers (DCCs) to improve access to data derived from hundreds of studies and samples collected from thousands of human subjects. This project involves a diversity of datatypes has been generated at the genomic, expression, proteomic, metagenomic, and imaging levels, and the DCCs support a tremendous range of scientific discovery efforts.

2020

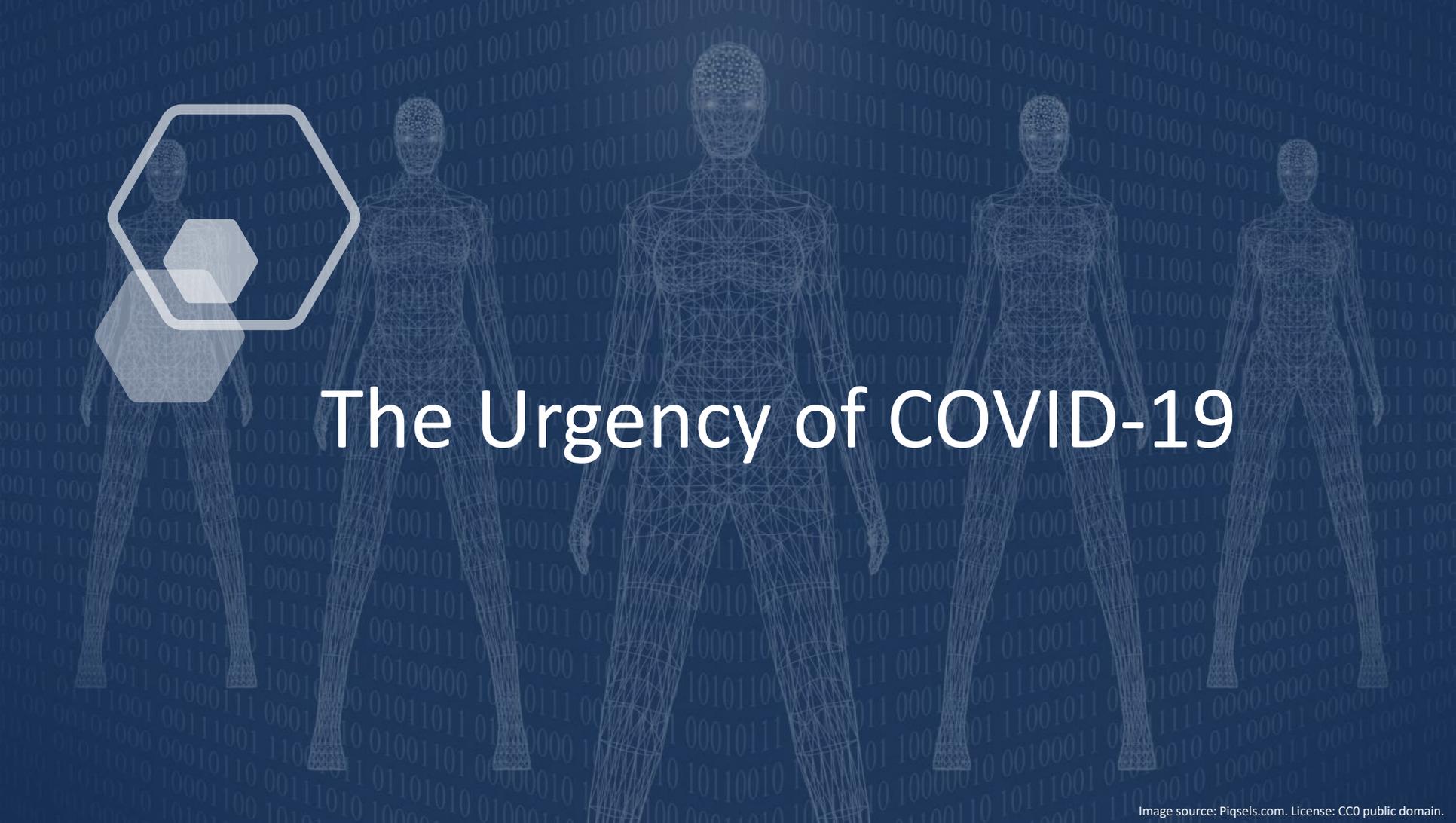


Crosscut Metadata Model “C2M2”

C2M2 Level 0 defines a **minimal valid C2M2 instance**. Data submissions at this level of metadata richness will be the easiest to produce, and will support the simplest available functionality implemented by downstream applications.

C2M2 Level 1 models **basic experimental resources and associations between them**. This level of metadata richness is more difficult to produce than Level 0's flat inventory of digital file assets. As a result, Level 1 metadata offers users more powerful downstream tools than are available for Level 0.

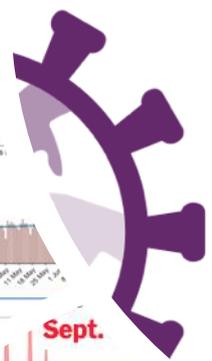
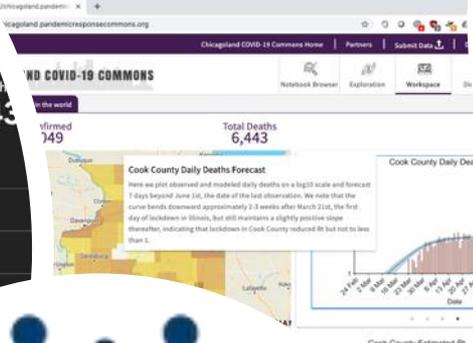
2020



The Urgency of COVID-19



Global Deaths
1,155,433
 225,239 deaths US
 157,134 deaths Brazil
 119,014 deaths India
 88,924 deaths Mexico
 44,800 deaths Spain

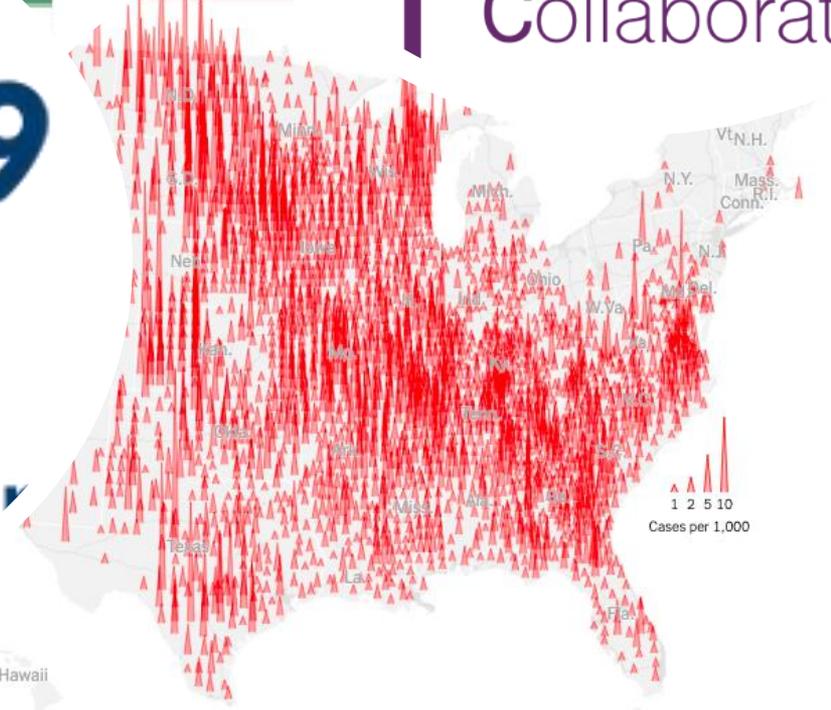


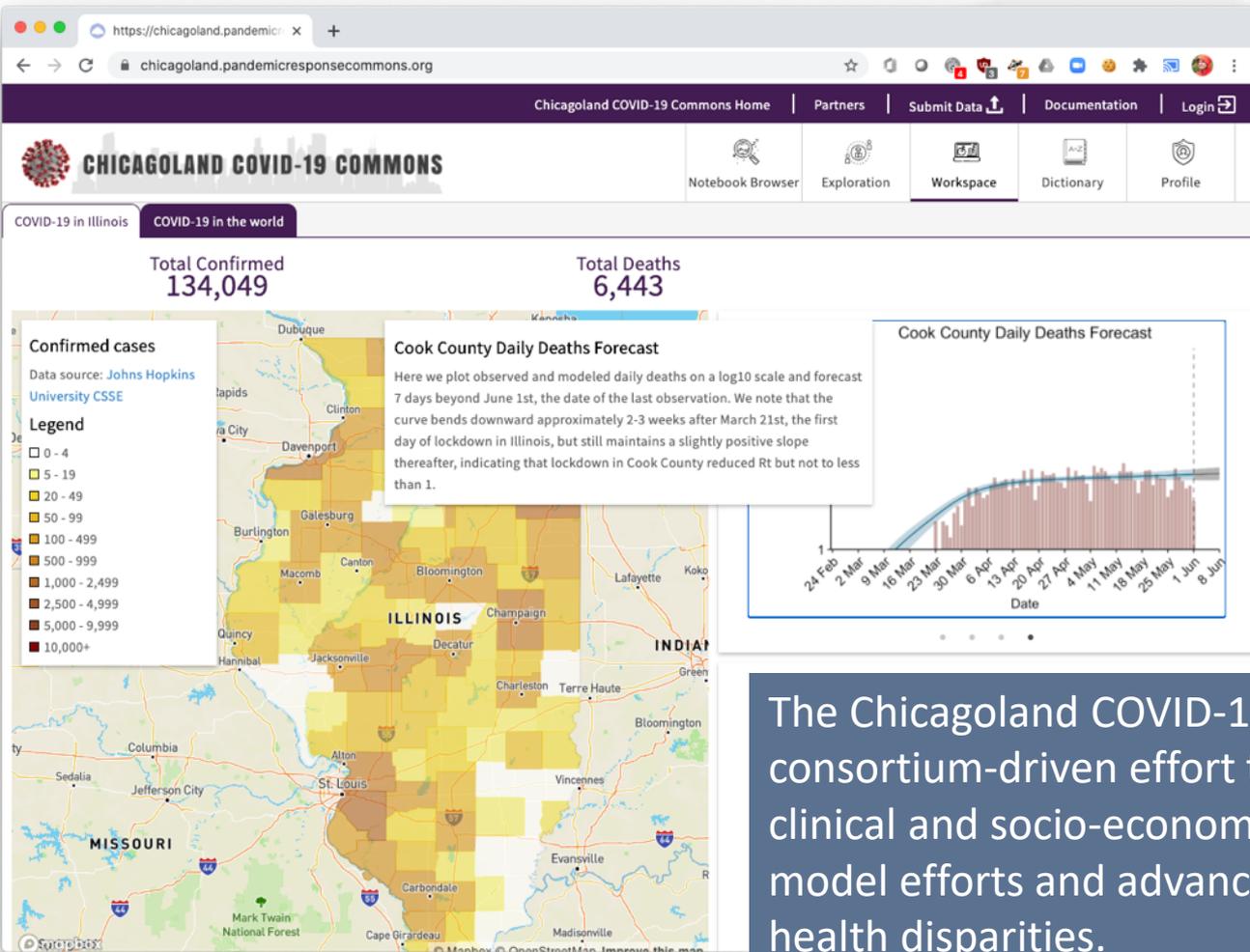
National COVID Cohort Collaborat



COVID-19 Data Portal

COVID-19 & Cancer





The Chicagoland COVID-19 Commons is a consortium-driven effort to collect regional clinical and socio-economic data to drive local model efforts and advance the understanding of health disparities.

Final thoughts

- Data sharing (or enabling data access) is key to advancing biomedical research
- Data sharing is difficult; it requires
 - Active data management
 - Effective data governance
 - Adoption of standards, structured vocabulary
- Most importantly, it requires alignment of incentives



Thank you

<https://chicagoland.pandemicresponsecommons.org/>

Matthew Trunnell
matthew@occ-data.org
@MatthewTrunnell