# UCLA
## Presentations

**Title**

How and why do scientists reuse others' data to produce new knowledge?

**Permalink**

https://escholarship.org/uc/item/6rx42812

**Authors**

Borgman, Christine
Pasquetto, Irene V.

**Publication Date**

2018-09-15

**Copyright Information**

# How and why do scientists reuse others' data to produce new knowledge? Background, Foreground, and Beyond

## Christine L. Borgman, PhD

Distinguished Research Professor
Center for Knowledge Infrastructures
University of California, Los Angeles
http://christineborgman.info
https://knowledgeinfrastructures.gseis.ucla.edu
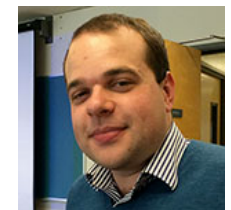@scitechprof

## Irene V. Pasquetto, PhD

Research Affiliate
UCLA Center for Knowledge Infrastructures

Fringe Event
Cochrane Colloquium
Edinburgh, 15 September 2018



Christine Borgman    Bernie Boscoe    Peter Darch

Milena Golshan    Irene Pasquetto    Michael Scroggins

Cheryl Thompson    Morgan Wofford

**Cochrane**
Trusted evidence.
Informed decisions.
Better health.

**UCLA** Center for Knowledge Infrastructures

# Data sharing policies

- Research Councils of the UK

- European Union

- U.S. Federal research policy

- Australian Research Council

- Individual countries, funding agencies, journals, universities

# Data Stewardship: The Ideal



**F**indable    **A**ccessible    **I**nteroperable    **R**eusable

Wilkinson, et al. (2016). The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data*, *3*, http://dx.doi.org/10.1038/sdata.2016.18
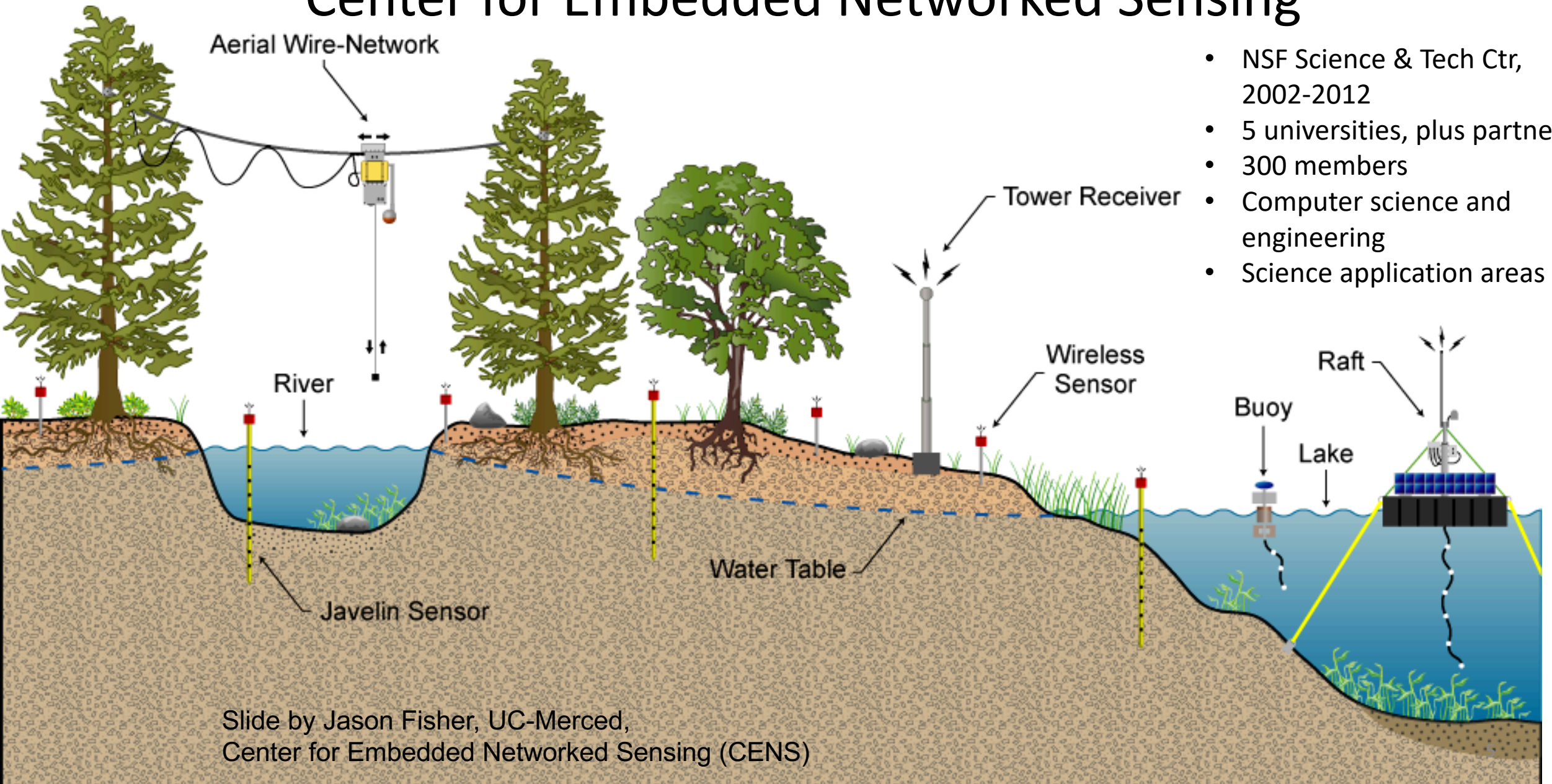
# What is "data reuse"?



Pasquetto, I. V., Randles, B. M., & Borgman, C. L. (2017). **On the Reuse of Scientific Data**. *Data Science Journal*, *16*. https://doi.org/10.5334/dsj-2017-008

# Center for Embedded Networked Sensing



Aerial Wire-Network

Tower Receiver

Wireless Sensor

Raft

Buoy

Lake

River

Water Table

Javelin Sensor

- NSF Science & Tech Ctr, 2002-2012
- 5 universities, plus partners
- 300 members
- Computer science and engineering
- Science application areas

Slide by Jason Fisher, UC-Merced,
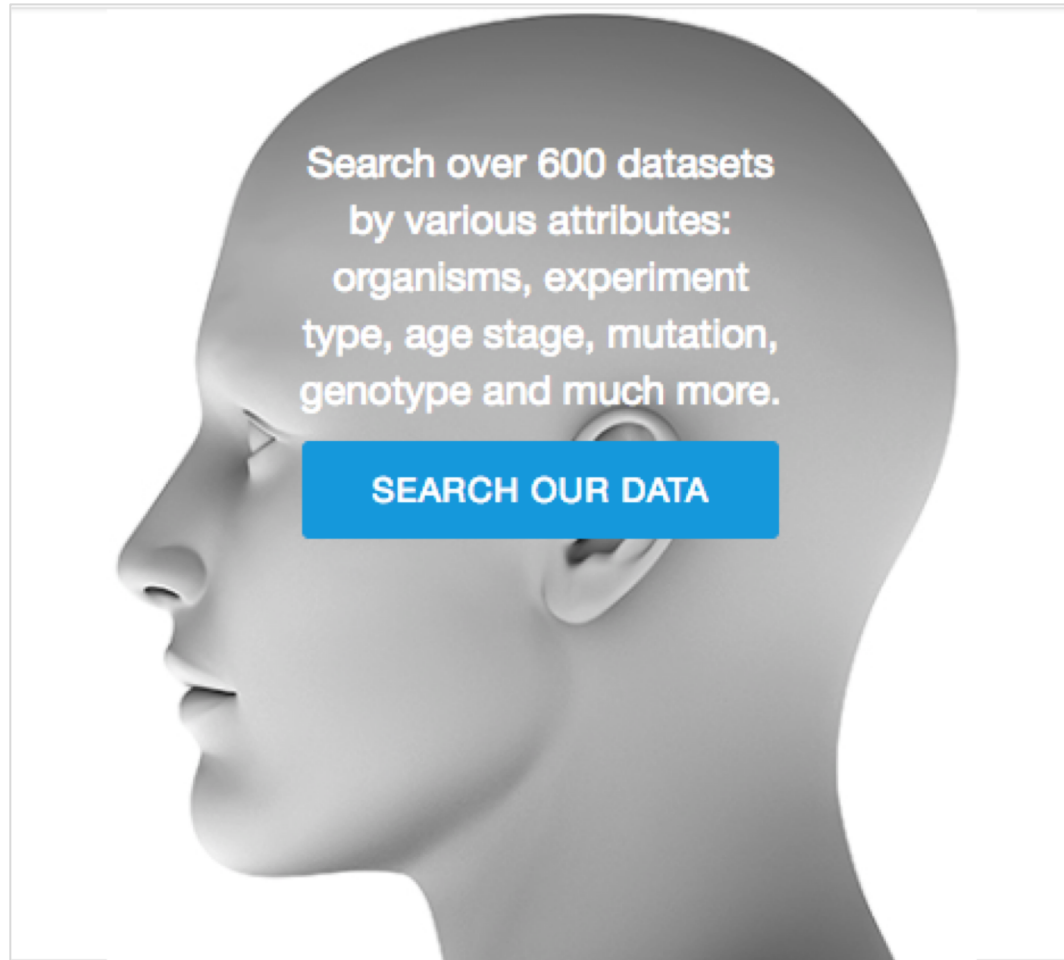Center for Embedded Networked Sensing (CENS)

# Background and Foreground Reuses of data at CENS



Images: CKI and NSF archives

# The DataFace Consortium for Data Sharing



Search over 600 datasets by various attributes: organisms, experiment type, age stage, mutation, genotype and much more.

SEARCH OUR DATA

**GOAL:**

Collect and release high-throughput "hypothesis free" biomedical data related to the genetics of facial formation and development in humans and animals.
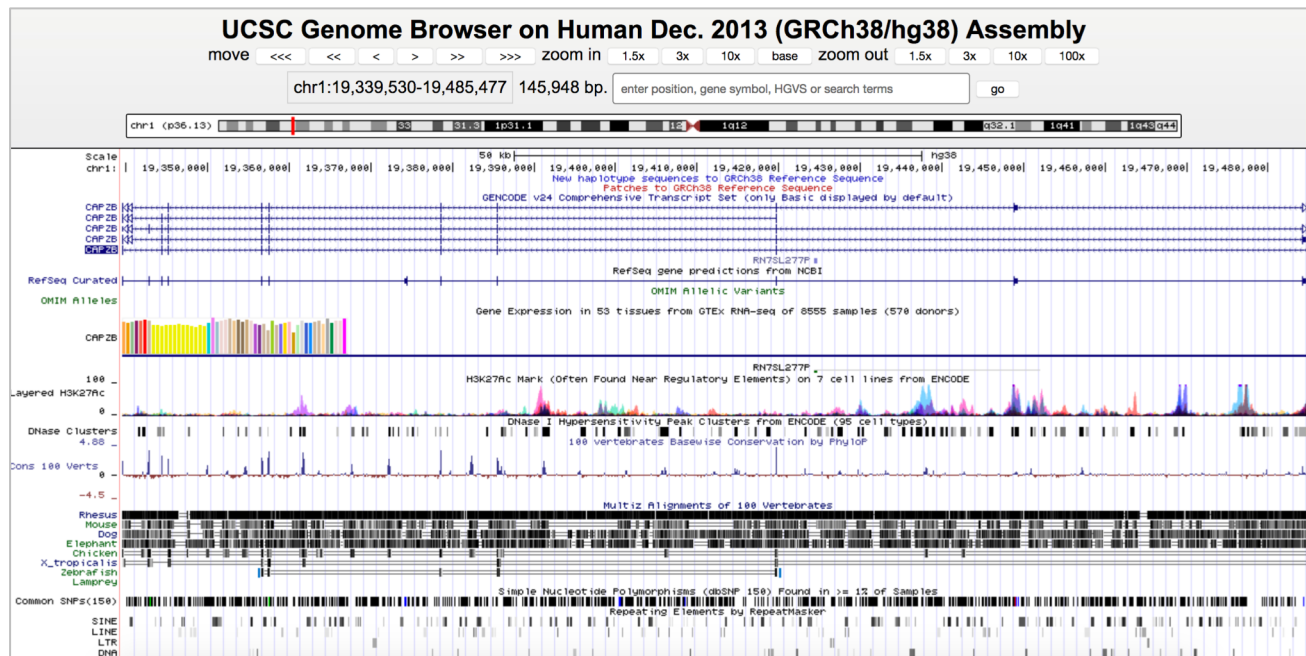
**DATA TYPES:**

Whole genome sequences, gene expression data from ChiP-seq, RNA-seq, and microarrays, genotypes and phenotypes from GWAS studies, etc.
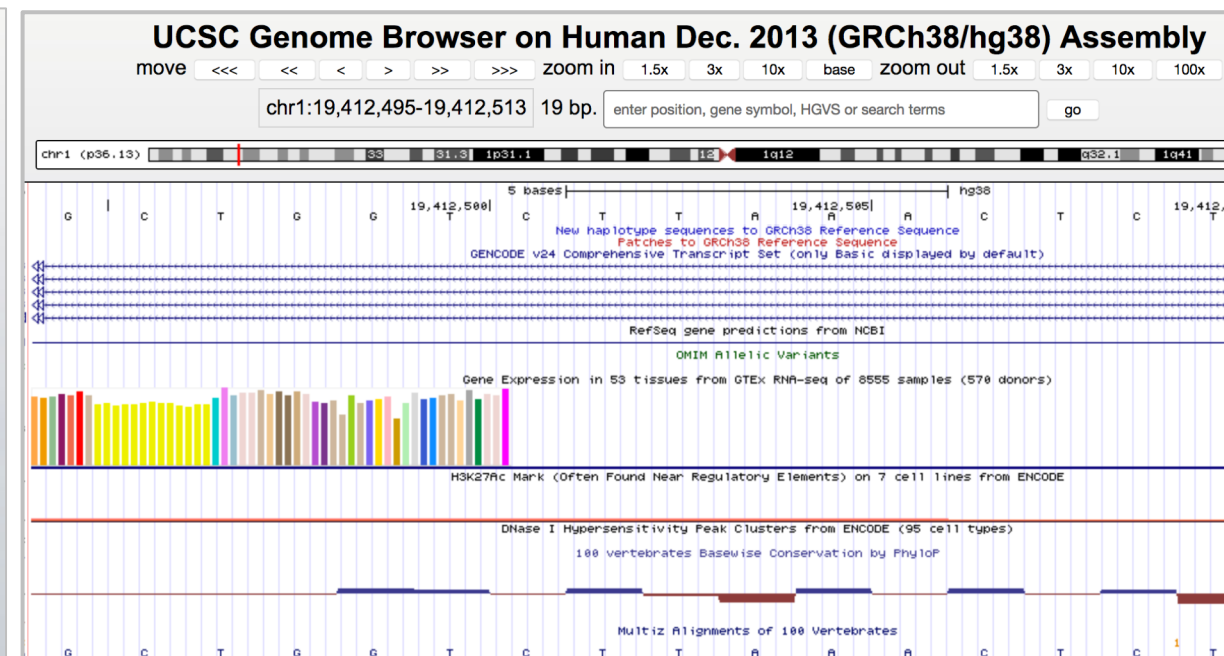
**DOMAINS:**

Developmental biologists, evolutionary experts, human geneticists, computational biologists, surgeons, etc.

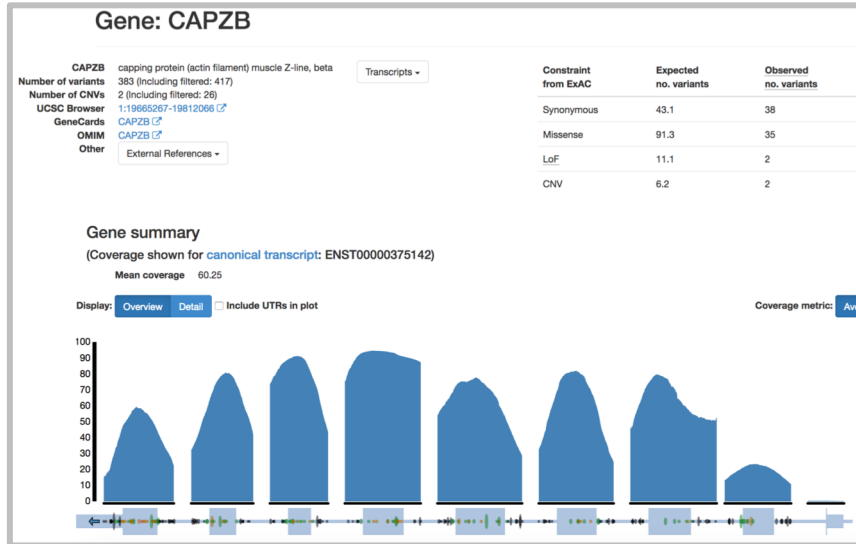# Background Reuse at DataFace: Comparison, control, verification. (I)



UCSC Genome Browser – Search example
(CAPZB gene)

UCSC Genome Browser – Zoom IN

Irene Pasquetto, DataFace Study, 2018

8

# Background Reuse at DataFace: Comparison, control, verification. (II)

# Foreground Reuse at DataFace: Data Analysis



Aligner software pairs "reads" using reference assemble genome

Data processing tool summarizes BAM information to compute likelihood of each possible genome

In-house script takes the ratio of mutant and allele frequencies to find the highest peak

R studio calculates elative frequency and generate plotting graphs

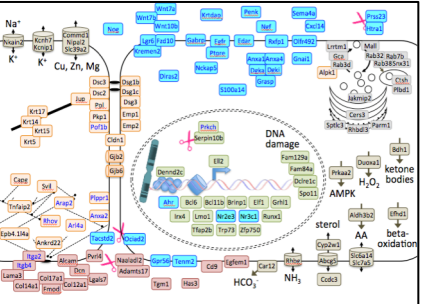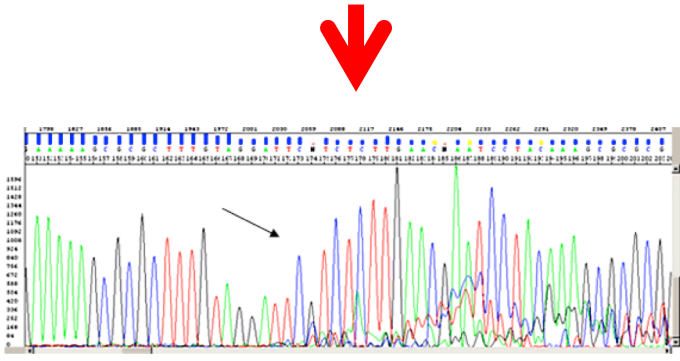Annotation tool predicts consequences of gene function

Variants are annotated by gene names, variant impact, and type of variant

"RAW" DATA

Pipeline

RESULTS

*Having access to the contact information of those who collected the data increases rates of foreground reuse.*

# The "Data Creator Advantage"

- Creator has most current annotations about the dataset
- Creator has most specialized knowledge of relevant literature
- Creator may have software pipelines locally customized for the dataset



Image source: https://www.siteminder.com

12

|  | **BACKGROUND Reuse of Data** | **FOREGROUND Reuse of data** |
|---|---|---|
| **Goal of reuse** | "Ground truthing:" calibrate, compare, confirm | Analysis: identify patterns, correlations, causal relationships |
| **Example of reuse** | Instrument calibration, sequence annotation, review summary-level data | Meta-analyses, novel statistical analyses |
| **Frequency of reuse** | Frequent, routine practice | Rare, emergent practice |

Pasquetto, I. V., Borgman, C. L., & Wofford, M. F. (in review). The Who, What, When, and Why of Reusing Data in Scientific Practice. *Harvard Data Science Review*.

|  | **BACKGROUND Reuse of Data** | **FOREGROUND Reuse of data** |
|---|---|---|
| **Goal of reuse** | "Ground truthing:" calibrate, compar~~e~~ confirm | Analyses: identify patterns, co~~rrelatio~~ns, causal r~~elation~~ |
| **Example of reuse** | Instrume~~nt~~ seque~~ncing~~ r~~~~ ~~prim~~ary-level da~~~~ | Me~~~~ |
| **Frequency of reuse** | Freq~~ue~~nt - routine practice | Rare - emergent practice |

*INDEPENDENT REUSE OF DATA*

*COLLABORATIVE REUSE WITH DATA CREATORS*

Pasquetto, I. V., Borgman, C. L., & Wofford, M. F. (in review). The Who, What, When, and Why of Reusing Data in Scientific Practice. *Harvard Data Science Review*.

# Questions: Trusted Evidence?

- When to reuse open data independently?
- When to collaborate with data creators?
- What information is needed, when, to trust evidence?

**Cochrane**

Trusted evidence.
Informed decisions.
Better health.

# Questions: Informed decisions?

- What do you need to know about the data to inform decisions?

- When are data sufficient for decision making?

- When is further information about about data needed?

- How should data sharing and reuse be governed?

# Questions: Better health?

- Where should community invest in data sharing and reuse?
- How should data resources be governed?
- Who should be responsible for sustaining access to health data?
- What are reasonable licensing agreements?
- What are appropriate funding models for data resources?

**Cochrane**

Trusted evidence.
Informed decisions.
Better health.

# Acknowledgements



UCLA Center for Knowledge Infrastructures

Christine Borgman

Bernie Boscoe

Peter Darch

Milena Golshan

Irene Pasquetto

Michael Scroggins

Cheryl Thompson

Morgan Wofford