# Data Sharing:
# Why is it still so hard?

John Muschelli

Assistant Scientist

Department of Biostatistics

Johns Hopkins Bloomberg School of Public Health

# Data Types

Raw

- GT3X
- CWA
- Bin files

Raw-ish

- HDF5
- CSV - gzipped vs. not

Some need applications/licenses - ideal to have open source solutions to open

# Data Types

Interpolated/Filtered/Aggregated

- What epoch?
    - 1 second?
    - 5 second?


- Measure to use?
    - $r = \sqrt{(x^2 + y^2 + z^2)}$
    - Mean of x, y, z?
    - Variance?

Can't invert these - can't go back!

# RAW

Interpolated/Filtered/Aggregated

- What epoch?
    - 1 second?
    - 5 second?

- Measure to use?
    - $r = \sqrt{(x^2 + y^2 + z^2)}$
    - Mean of x, y, z?
    - Variance?

Can't invert these - can't go back!

# Journals: Guidelines

- Data Availability Upon Request

- Have it "somewhere"
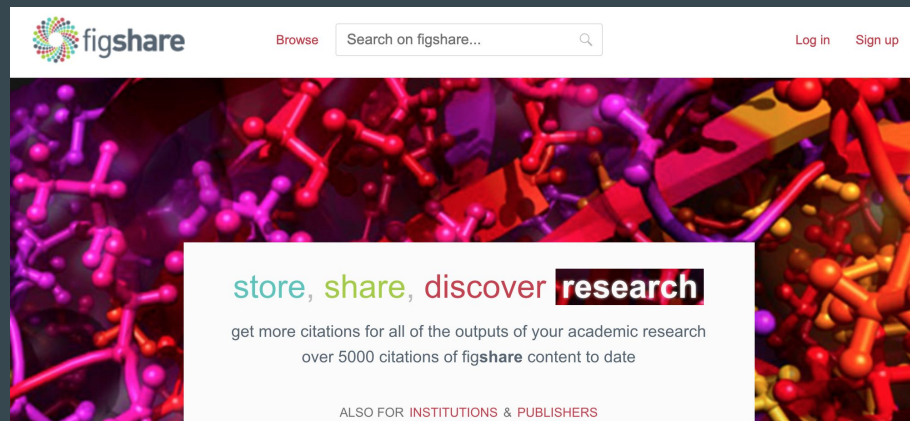
OR 

# But where?

Example of biologically-related repos

- Physionet

General data

- Figshare
- Zenodo
- Synapse
- Datacite

What are your other favorites?

# Get a Digital Object Identifier (DOI)

Make sure it's

- FAIR - Findable, Accessible, Interoperable, and Reusable

But our commodity is paper/citations (it's true, so needs):

- Citable




If you're good at something, never do it for free.

# Funding Agencies

- No central repository
- Other fields (e.g. genomics) have one that you **\*need\*** to put in
    - dbGaP - https://www.ncbi.nlm.nih.gov/gap/

One of the main issues starting one: funding
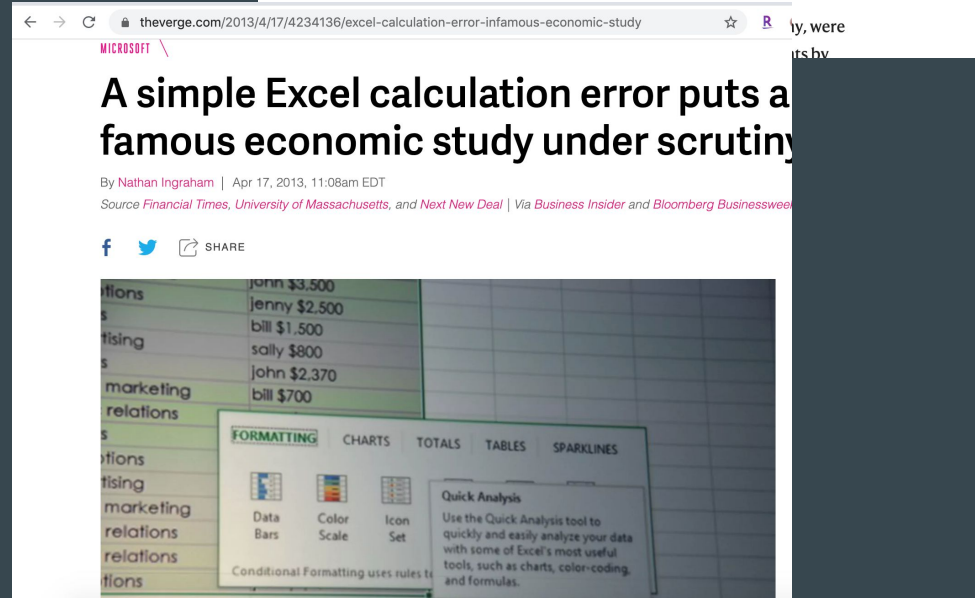
# Downsides: Why people don't share

Issue with non-replicability

- People see if your results replicate
- High-profile examples of errors
  - https://www.theverge.com/2013/4/17/4234136/excel-calculation-error-infamous-economic-study
  - https://www.nature.com/articles/nj7396-137a
- It's a Legitimate concern, but
  - Want to do good science or not be wrong?

# Downsides: Why people don't share

Many excuses

- "I wanted to do that analysis", but many times →
- They'll scoop me

PRIMARY DATA COLLECTORS
SHOULD GET PRIORITY

- But the embargo time can't be ∞
- 5 years?
- What if publicly-funded?



NEVER GONNA HAPPEN

# Why people don't share: I'm not ready yet



Ready? No one is ever ready, my boy. But some do what they plan to do and some never will. The difference between the two is that the first group understand that they need to start somewhere, so they do so. Straight away.

Chris Murray

quotefancy

https://quotefancy.com/quote/1655560/Chris-Murray-Ready-No-one-is-ever-ready-my-boy-But-some-do-what-they-plan-to-do-and-some

# Upside: Why should you still share

Replication

- Other researchers are helping identifying issues in your dataou with your data

More citations

- Hopefully good ones or people finding new things
- Consortia - UK Biobank

Pushing science

# Upside: Why should you still share

Data Aggregation

- You have data X, I have data Y
- Can show better out-of-sample validity,
  or …. you can have MORE POWER:





WE'RE MORE POWERFUL AS A GROUP.
THE PATH hulu

# Consent

Make sure your study documents have this in the language

- Talk to your local IRB
- Ask others
- Ask the repository
- Ask the funders

# De-identification/Anonymization

LOADS of open questions - what is "identifiable"?

- HIPAA stuff obviously, but what about

Gait?

GPS?

Data also hides in things sometimes (e.g. header of gt3x)
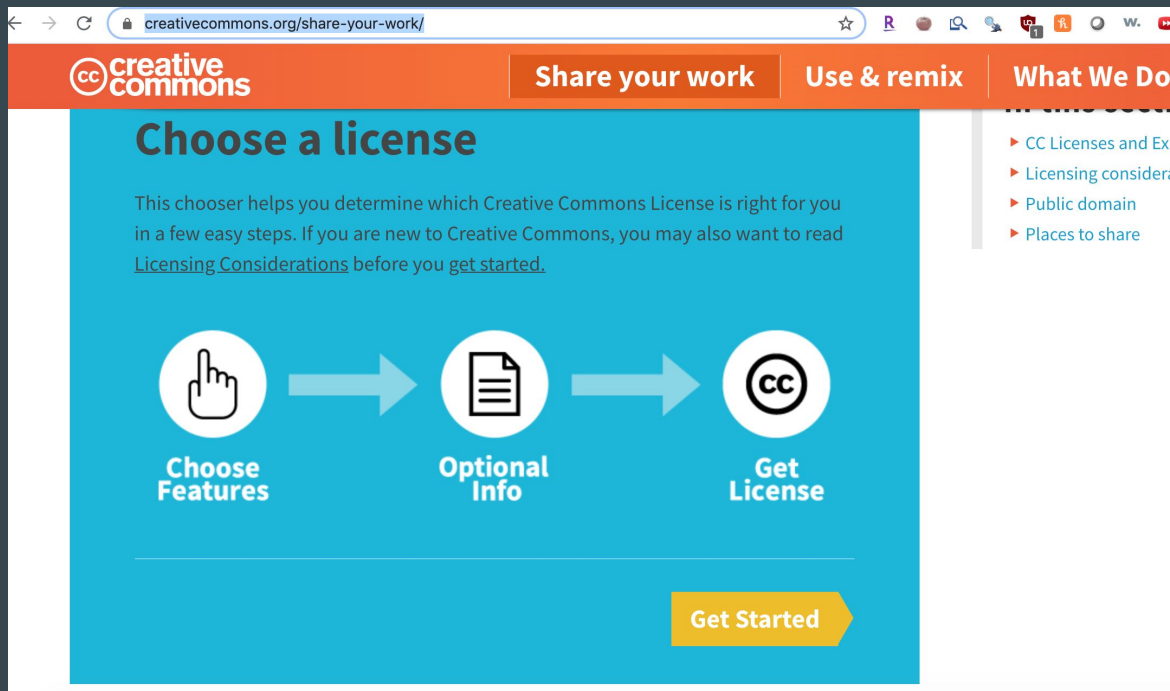
# Licenses:

https://creativecommons.org/share-your-work/

- Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0)

- MIT
- Apache
- BSD
- GPL

# Open Data

Upper limb activity of twenty myoelectric prosthesis users and twenty healthy anatomically intact adults.

- https://springernature.figshare.com/collections/Upper_limb_activity_of_twenty_myoelectric_prosthesis_users_and_twenty_healthy_anatomically_intact_adults_/4457855

Newcastle polysomnography and accelerometer data

- https://zenodo.org/record/1160410

# Open Data: Why I care

Having open data allows us to make tutorials/workflows for:

1. Our Students
2. Students in our courses
3. Software packages
4. Collaborators
5. The rest of science

(But maybe gives away "secret sauce")


TEACH ME.

# If you're releasing data, see if



- Can you read it open source
  - gt3x files (R: AGread, read.gt3x, Python: ActiGraph-ActiWave-Analysis)
  - Most other file types: GGIR (e.g. Axtivity CWA)
  - Everything can read CSV
  - Mostly everything can read HDF5 - can do ALL of these
- Data Structure/naming convention: BIDS
  - https://www.openmhealth.org/schemas/omh_acceleration/

# Gifs/Images

- https://imgur.com/gallery/Cr2L8s1
- https://painepublishing.com/measurementadvisor/wp-content/uploads/sites/4/2018/03/too-much-data.jpg
- https://media1.giphy.com/media/G5X63GrrLjjVK/giphy.gif?cid=ecf05e4780a131d84268014d0b5604c95823e8fc843a62e2&rid=giphy.gif
- https://media.giphy.com/media/9LZpYKd2CmV6ICxs17/giphy.gif
- https://media.giphy.com/media/3oz8xRD1irfRHMnVD2/giphy.gif
- https://media.giphy.com/media/3orifdO6eKr9YBdOBq/giphy.gif
- https://media.giphy.com/media/l41Yr7JNU9U7upjNu/giphy.gif
- https://media.giphy.com/media/FqAwoNjVneJxK/giphy.gif
- https://media.giphy.com/media/3oKIP5F6l0hfcdOn6M/giphy.gif
- https://media.giphy.com/media/26AHPxxnSw1L9T1rW/giphy.gif
- https://media.giphy.com/media/WUCeN3kP3X3JZu52vg/giphy.gif
- 
-

# Thanks!