# How *NOT* to share your data:
## Avoiding data horror stories

Rosie Higman

Office of Scholarly Communication

8th March 2017

UNIVERSITY OF CAMBRIDGE

1. Where?
2. What?
3. File formats
4. Formatting your spreadsheet*
5. Document and describe your data!

\* Based on Avoiding data disasters course by Mark Dunning, CRUK-CI
http://bioinformatics-core-shared-training.github.io//avoid-data-disaster/

- Every discipline is different

- These are general principles

- Application will vary according to your research

# Where NOT to share your data
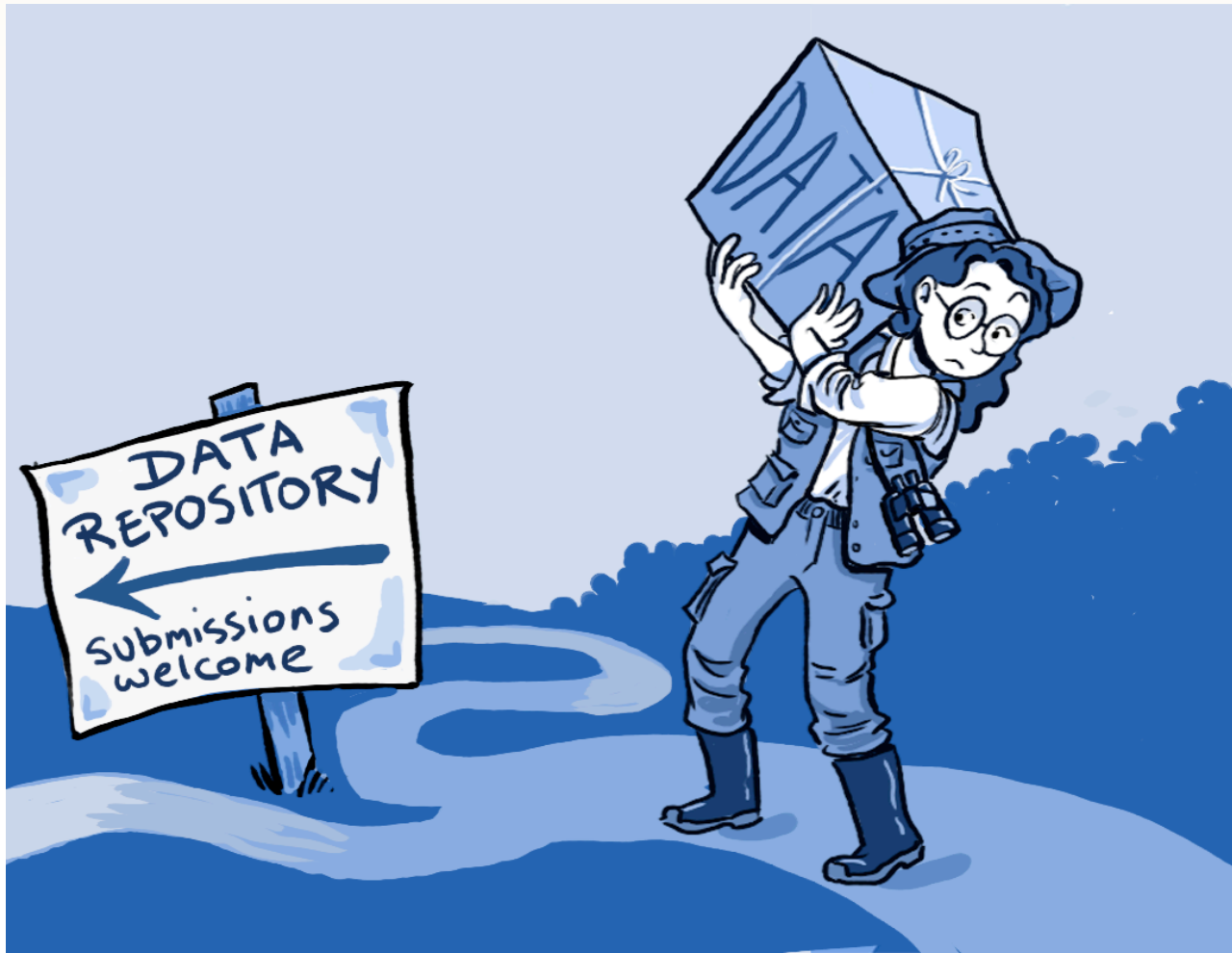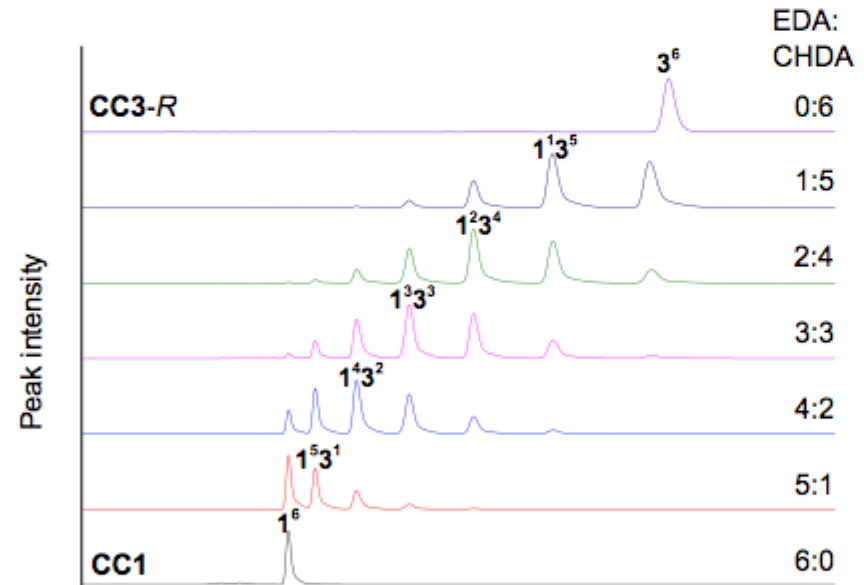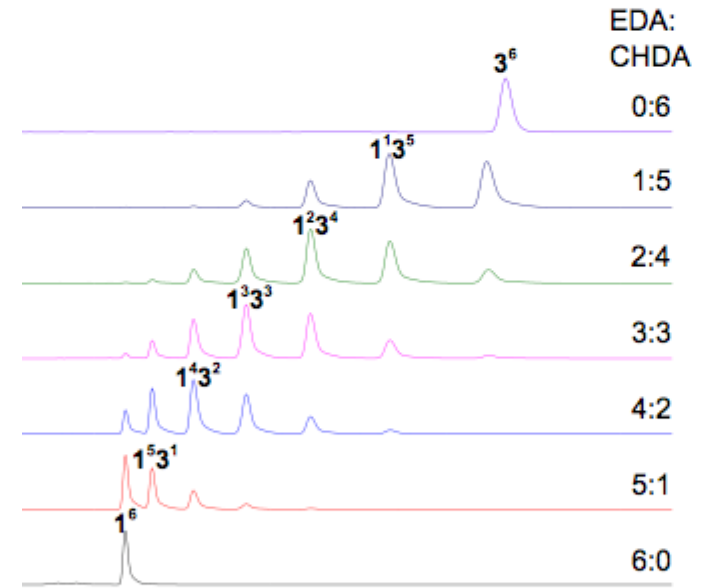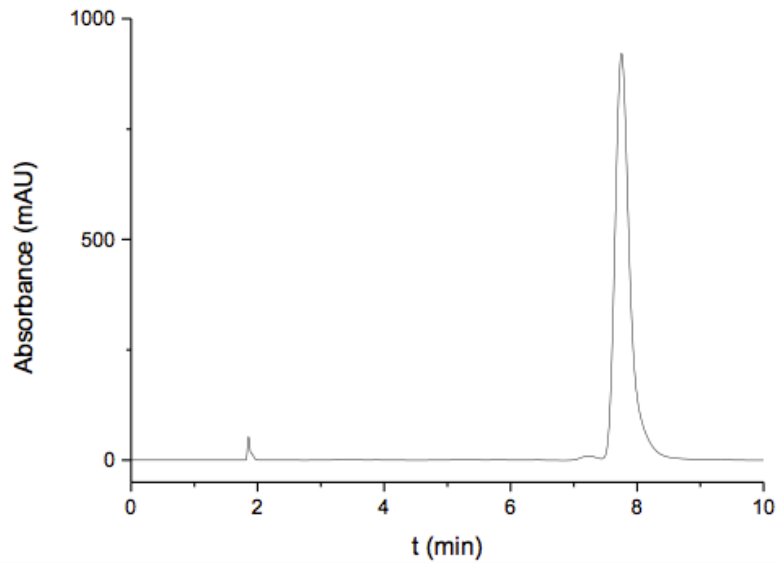
# Where SHOULD you share your data?

Roche DG, Lanfear R, Binning SA, Haff TM, Schwanz LE, et al. (2014) - Roche DG, Lanfear R, Binning SA, Haff TM, Schwanz LE, et al. (2014) Troubleshooting Public Data Archiving: Suggestions to Increase Participation. PLoS Biol 12(1): e1001779. doi:10.1371/journal.pbio.1001779, CC BY 4.0,
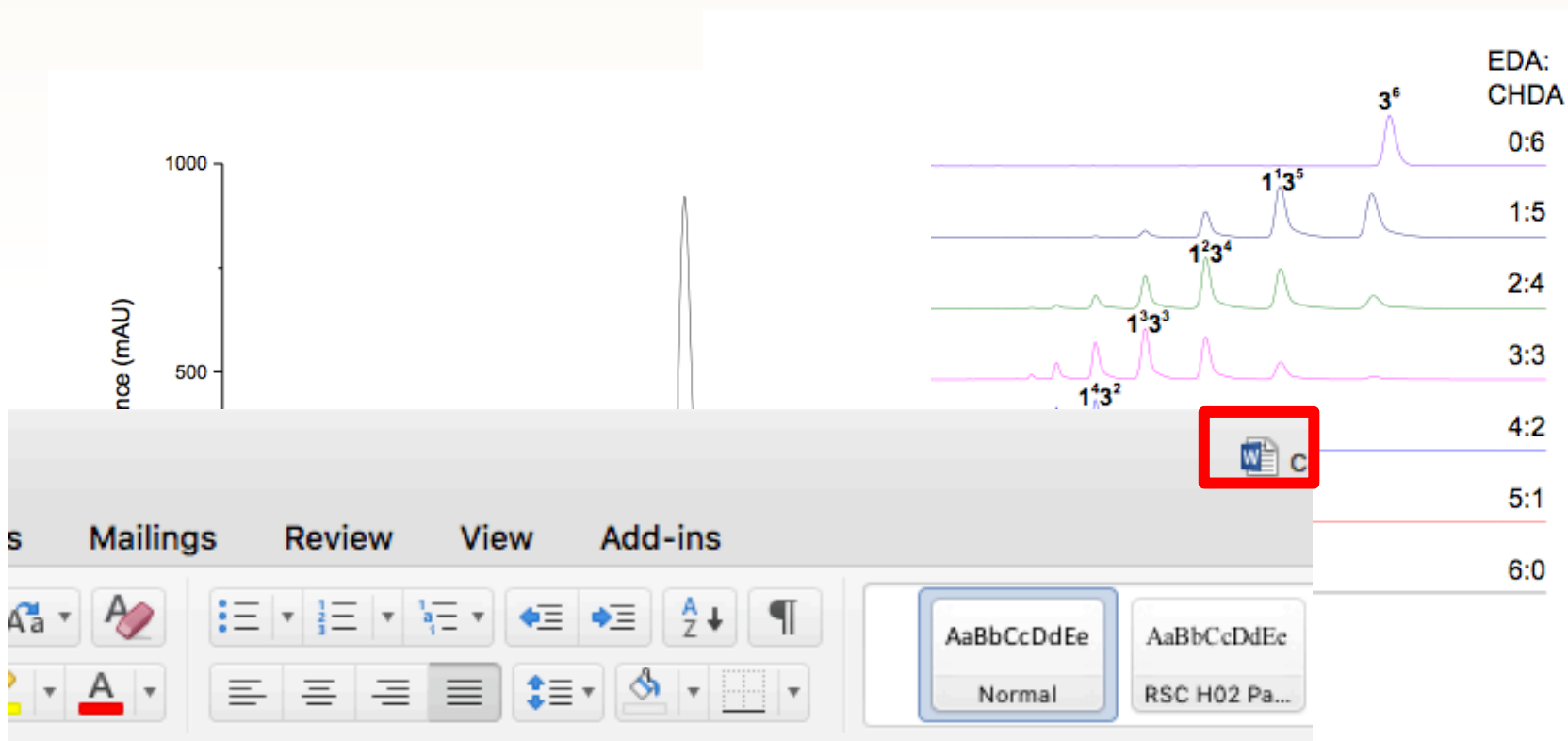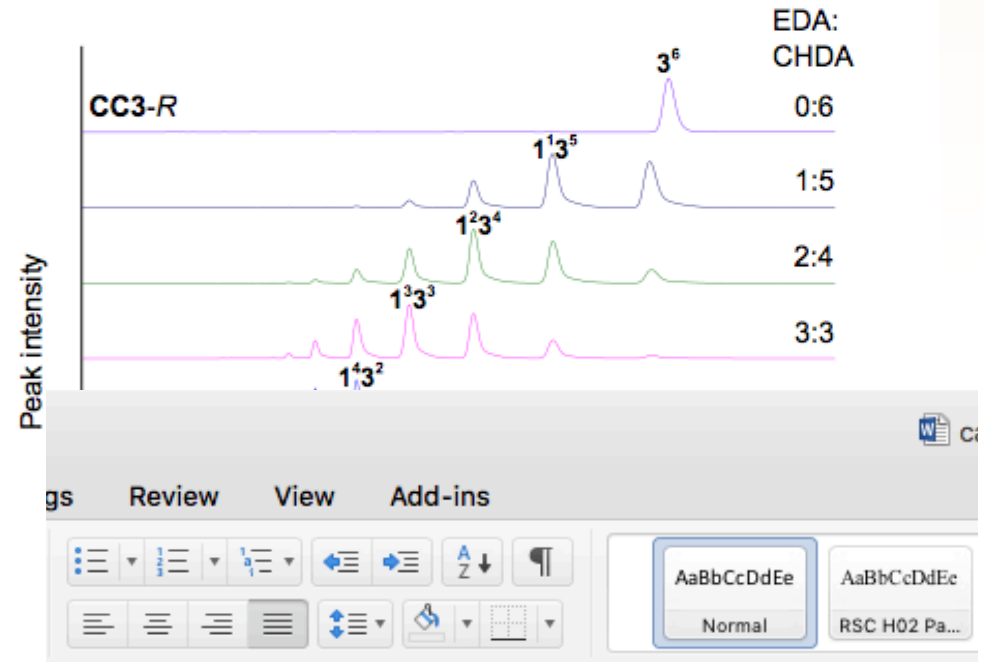
OSC

# What data should you include?

| Entry | Flow rate A (mL/min) | Flow rate B (mL/min) | Total flow rate (mL/min) | Reactor volume (mL) | Residence time (min) | Reactor temperature (°C) | Peak area, CC3-R (% a/a) |
|-------|------|------|------|------|------|------|------|
| 1 | 0.6 | 0.4 | 1 | 10 | 10 | 40 | 50.2 |
| 2 | 0.6 | 0.4 | 1 | 10 | 10 | 60 | 64.6 |

More than your figures!

Data and code necessary to recreate your results

# Powerpoint is for presentations NOT data!

**osc**

Instead:

Original image files
Appropriate formats
Annotations embedded in
separate PDF/csv/txt file
(README file)

**Textual data** = XML, TXT, HTML, PDF/A (Archival PDF)

**Tabular data (spreadsheets)** = CSV
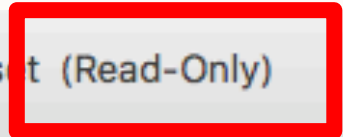
**Databases** = XML, CSV

**Images** = TIFF, PNG, JPEG*

**Audio** = FLAC, WAV, MP3

**Think! Preservation vs access/re-use**

Graphs in separate sheet

No highlighting

No colours

No formulas*

*In your raw data

No blank cells

1 piece of data per cell

Keep units out of cells

Use data validation

| | A | B | C |
|---|---|---|---|
| 1 | Age | Weight | Test score |
| 2 | 25 | 65kg | 93 |
| 3 | 36y4m | 62.4 | |
| 4 | | 70 | 40 |
| 5 | 47 | 82000g | 31 |
| 6 | 33 | 77 | 49.7 |
| 7 | 28.4 | | 89 |
| 8 | | | |
| 9 | | | |

# Gene name errors are widespread in the scientific literature

Mark Ziemann, Yotam Eren and Assam El-Osta ✉

## Abstract

The spreadsheet software Microsoft Excel, when used with default settings, is known to convert gene names to dates and floating-point numbers. A programmatic scan of leading genomics journals reveals that approximately one-fifth of papers with supplementary Excel gene lists contain erroneous gene name conversions.

### Keywords

Microsoft Excel – Gene symbol – Supplementary data

The problem of Excel software (Microsoft Corp., Redmond, WA, USA) inadvertently converting gene symbols to dates and floating-point numbers was originally described in 2004 [1]. For example, gene symbols such as *SEPT2*

## Description

Data supporting publication '███████████████████████
██████████████'

## Software

excel, matlab

# How you SHOULD describe your data

## Citation

Thwaites, A., Nimmo-Smith, I., Wieser, E., Soltan, A., & Marslen-Wilson, W. D. *Measurement datasets 1-3.01 for the "Kymata Atlas"* [dataset]. https://doi.org/10.17863/CAM.1660

## Description

The electromagnetic measurements of the human cortex used in the creation of the Kymata Atlas (datasets 1-3.01). The recordings were made at the MRC Cognition and Brain Sciences Unit, using an Elekta Neuromag MEG (306 ch.) and an EasyCap EEG (70 ch.). Electromagnetic brain signals are recorded from participants as they experience passive, naturalistic, stimuli. The participants involved are asked to watch a movie and/or listen to the radio (without any further tasks asked of them) and the recordings are made during this period. Data is anonomised and averaged over participants. Due to millisecond differences in stimulus delivery, there are two sets of recordings, one syncronised to the sound stimulus, and one to the visual stimulus.

OSC

```
README.txt
==========

Data supporting the conference paper: Wireless sensor monitoring of Paddington Station Box Corner

URI:              doi:10.1680/tfitsi.61279.209 (paper)
                  https://www.repository.cam.ac.uk/handle/1810/254928 (dataset)


This data consists of displacement and inclination sensor data from an excavation at a construction
and transmitted using a wireless sensor network. Accompanying this data is a location of each of the
has been used to generate the figures presented in the paper "Wireless sensor monitoring of Paddingt


Archive structure:
------------------

    paddington-wsn-data.zip
      README.txt (this file)
      Data/
        paddington-2014-02-17-1.csv
        paddington-2014-02-17-2.csv
        paddington-2014-02-17-4.csv
        ...
      Figures/
      Fig2/
        Fig2a/
        Fig2b/
      Fig6/
      Fig7/
      Fig8/
      Fig9/
        Fig9a/
        Fig9b/
        Fig9c/


The dataset consists of all the data in the Data directory. Subsets of this data, together with with
.xlsx files and Origin Project .opj files used to generate the figures for the paper, are also prese
do not form part of the dataset as such, but are included as an example of how the dataset can be us
```

Choose a repository.
Choose open file formats.
Choose sharing more than your figures.
Choose a tidy spreadsheet.
Choose to describe your data.
Choose decent documentation so your research is reproducible.

## CHOOSE DATA SHARING

info@data.cam.ac.uk   @CamOpenData    www.data.cam.ac.uk