# R-eproducible workflows

**1-day workshop**
**Morning lecture**



Brendan Palmer,

Statistics & Data Analysis Unit,

Clinical Research Facility - Cork

@B_A_Palmer

# Plan for the day

**Morning: Reproducible research through R/RStudio**

10 am: Lecture - Project orientated workflows

10.30 am: Coffee break – Discussion

10.45 am: 2 × 30 minute tutorials

      - Joined up thinking when writing R code

      - R-projects as means to organise your research

11.45 pm: Lecture - Introduction to the tidyverse

**12.30 pm: Lunch break**

**Afternoon: A crash course in the tidyverse**

1.15 pm: 3 × 45 minute hands-on tidyverse tutorials including;

    - Differences between the tidyverse and base R code

    - Example scripts and problem sheets exploring R packages

    - Useful add on packages
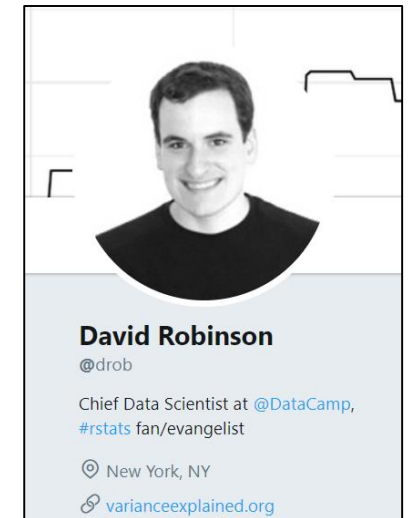
**3.30 pm: Closing remarks, questions**

# Disciamer 1



**Hadley Wickham** ✔
@hadleywickham

R, data, visualisation.

⊙ Houston, TX
🔗 hadley.nz

**HADLEY WICKHAM**    TEACHING  CODE  PERSONAL

I also teach in person workshops from time-to-time; see the RStudio workshops page for more details.

## CODE

Most of my work is in the form of open source R code, which you can find on my github. You can roughly divide my work into three categories: tools for data science, tools for data import, and software engineering tools.

**DATA SCIENCE**
- ggplot2 for visualising data.
- dplyr for manipulating data.
- tidyr for tidying data.
- stringr for working with strings.
- lubridate for working with date/times.

**DATA IMPORT**
- readr for reading .csv and fwf files.
- readxl for reading .xls and .xlsx files.
- haven for SAS, SPSS, and Stata files.
- httr for talking to web APIs.
- rvest for scraping websites.
- xml2 for importing XML files.

**SOFTWARE ENGINEERING**
- devtools for general package development.
- roxygen2 for in-line documentation.
- testthat for unit testing

**Jenny Bryan**
@JennyBryan

Software engineer @rstudio, humane
#rstats, adjunct prof @UBC where I
created @STAT545, part of @ropensci

**STAT 545**    Home    FAQ    Syllabus    Topics    People

# Data wrangling, exploration, and analysis with R

## UBC STAT 545A and 547M

Learn how to

- explore, groom, visualize, and analyze data
- make all of that reproducible, reusable, and shareable
- using R

**David Robinson**
@drob

Chief Data Scientist at @DataCamp,
#rstats fan/evangelist

⊙ New York, NY
🔗 varianceexplained.org

**VARIANCE EXPLAINED**    ABOUT ME    POSTS    LEARN R    TEXT MINING IN R    INTRODUCTION TO EMPIRICAL BAYES

**David Robinson**

*Chief Data Scientist at
DataCamp, works in R and
Python.*

☐ Email
☐ Twitter
○ Github
✏ Stack Overflow

This is the homepage and blog of David Robinson, Chief Data Scientist at DataCamp. For more about me, see here.
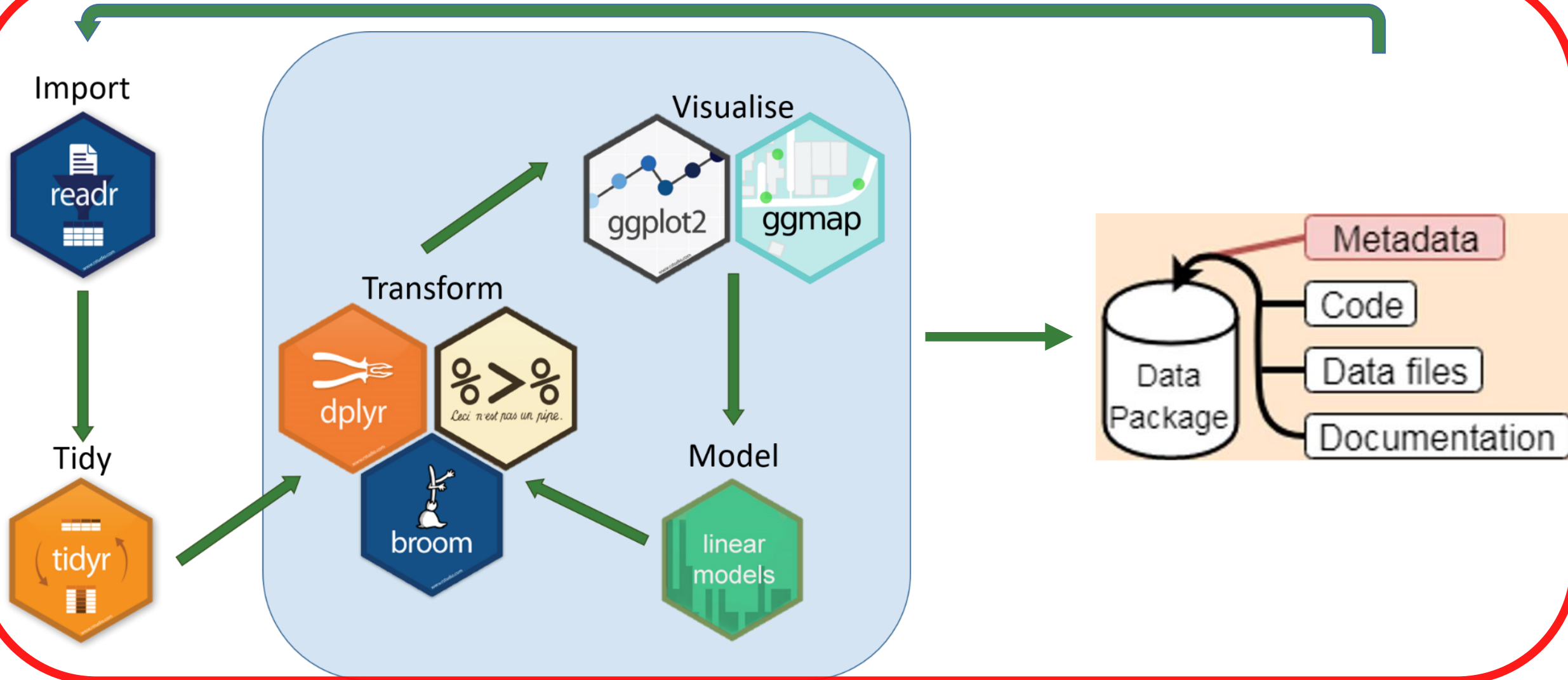
**Recent Posts**

Exploring college major and income: a live data analysis in R    *October 16, 2018*
A live screencast of an exploratory data analysis from the Tidy Tuesday series. This one explores college major and income data from 538.

Who wrote the anti-Trump New York Times op-ed? Using tidytext to find document similarity    *September 06, 2018*
An analysis of an anonymous op-ed in the New York Times, using document similarity metrics to match it to Twitter accounts.

Scientific debt    *May 10, 2018*
Introducing an analogy to 'technical debt' for data scientists.
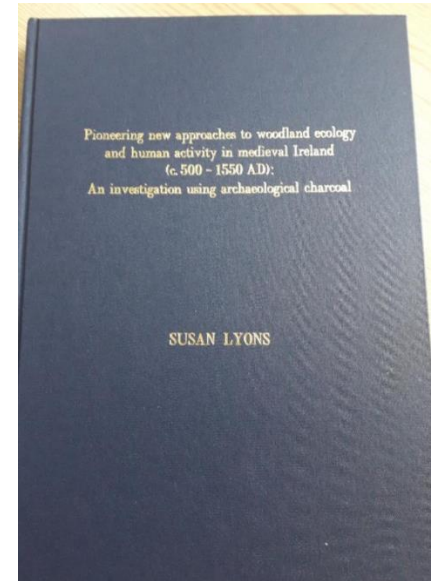
# Putting the pieces together

- Data analysis in a tidyverse nutshell

# How is research presented?

**Papers**

**Books**

**Theses**

**Talks**

**Posters**

# But what does it look like under the bonnet?

# The explosion of data



Growth of DNA Sequencing



Congratulations to Dr Katie Bouman!
This is the woman who created the algorithm to crunch the 5 petabytes of data from 500 kg of hard drives from 8 radio telescopes to make the first image of the #EHTBlackHole #BlackHole

Stevens et. al., 2015, Plos Biology

# Large quantities of data ≠ high quality of science

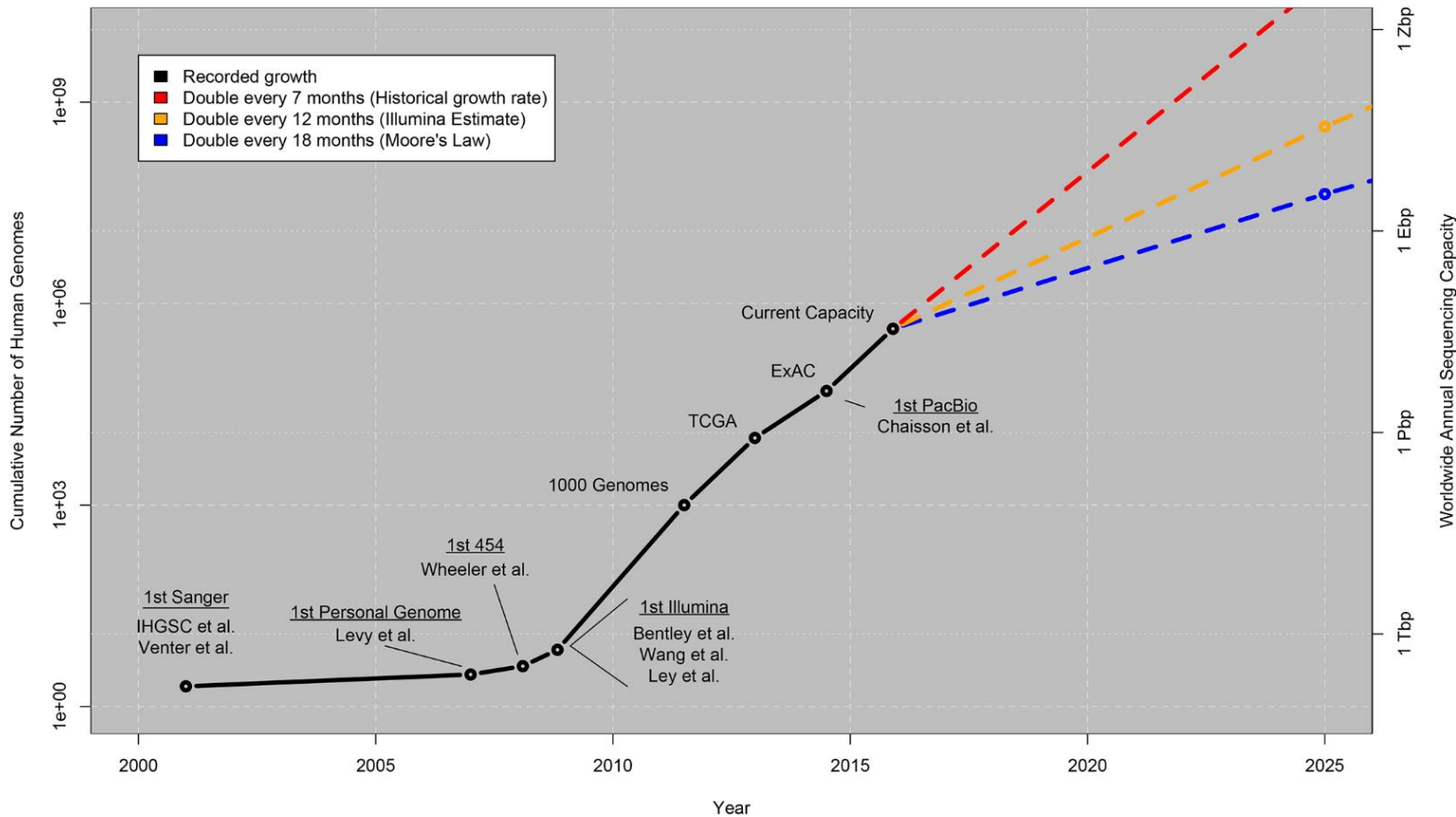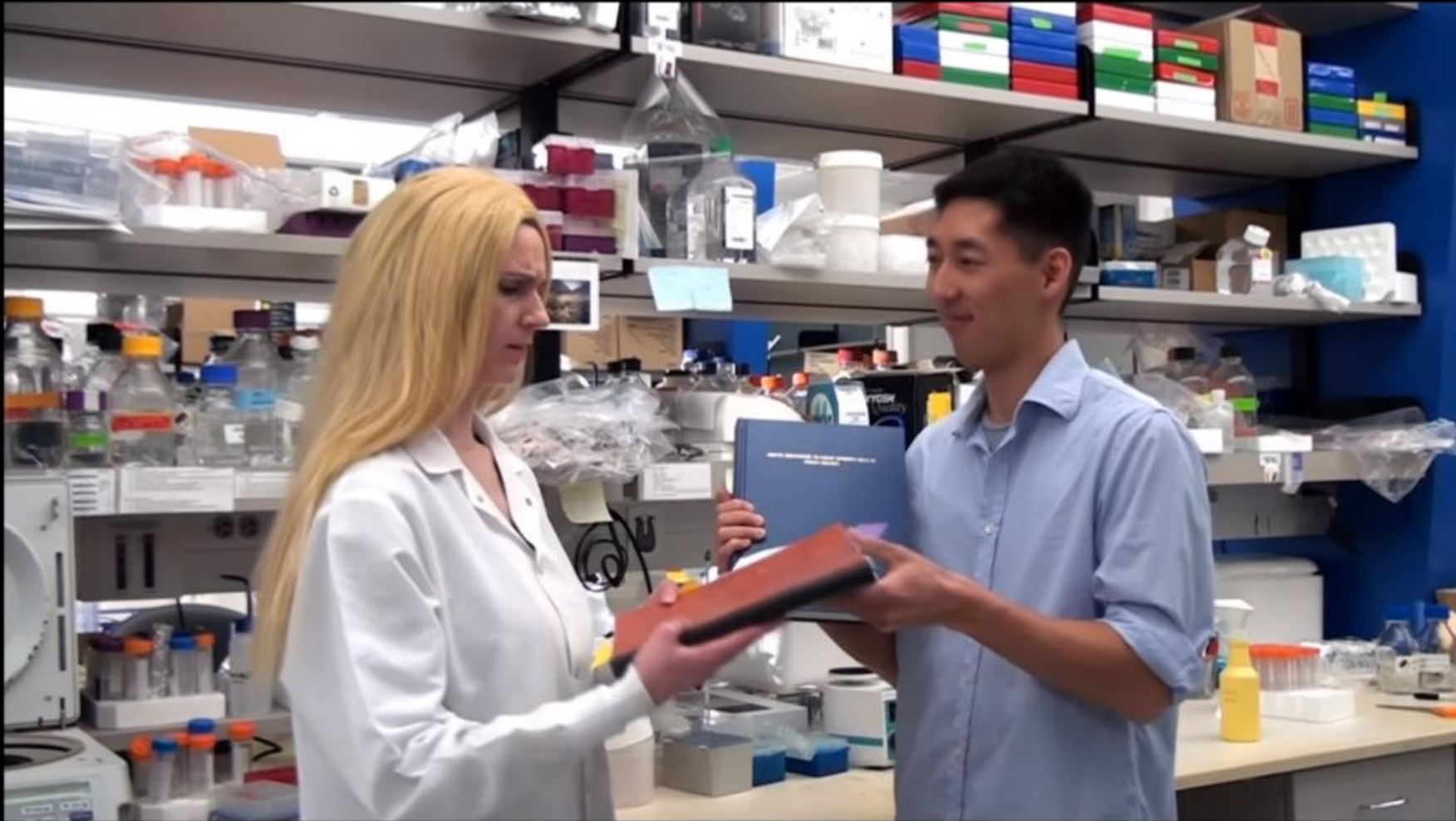## The association between adolescent well-being and digital technology use

Amy Orben [1]* and Andrew K. Przybylski [1,2]

The widespread use of digital technologies by young people has spurred speculation that their regular use negatively impacts psychological well-being. Current empirical evidence supporting this idea is largely based on secondary analyses of large-scale social datasets. Though these datasets provide a valuable resource for highly powered investigations, their many variables and observations are often explored with an analytical flexibility that marks small effects as statistically significant, thereby leading to potential false positives and conflicting results. Here we address these methodological challenges by applying specification curve analysis (SCA) across three large-scale social datasets (total $n = 355,358$) to rigorously examine correlational evidence for the effects of digital technology on adolescents. The association we find between digital technology use and adolescent well-being is negative but small, explaining at most 0.4% of the variation in well-being. Taking the broader context of the data into account suggests that these effects are too small to warrant policy change.

Orben & Przybylski (2019), Nature Human Behaviour

You were defending, one foot out the door

# This session

- **Project structure**

  - **Naming conventions**

    - **Scripted workflows**

      - R Markdown

        - **Reproducible research**

THIS PERSON IS likely to be YOU BTW!!

Jorge Cham | www.phdcomics.com

# Still haven't found what I'm looking for

- Help your future-self

# R-projects

# Define a generic project structure

- STEP 1: Give your research projects a shared structure

# Work from the raw data ALWAYS!!

**Tom Webb** @tomjwebb · 16 Jan 2015
If you could tell a new PhD student one thing to help make their data more useful/shareable, what would it be?

💬 27    🔁 11    ♡ 7    ✉

**Dr Gavin Simpson**
@ucfagls

Follow

Replying to @tomjwebb

@tomjwebb don't, not even with a barge pole, not for one second, touch or otherwise edit the raw data files. Do any manipulations in script

7:15 AM - 16 Jan 2015

# Give your files informative names

- STEP 1: Give your research projects a shared structure

# Everything in its right place

- STEP 2: Make you file names machine readable, human readable and work with default ordering

## NO

Name

- All unique 4a amino acid Sequences (B-N).fas
- All unique 4a amino acid Sequences (B-N).meg
- All_AA_haplotypes.meg
- All_AA_haplotypes_with_clonal_sequences.meg
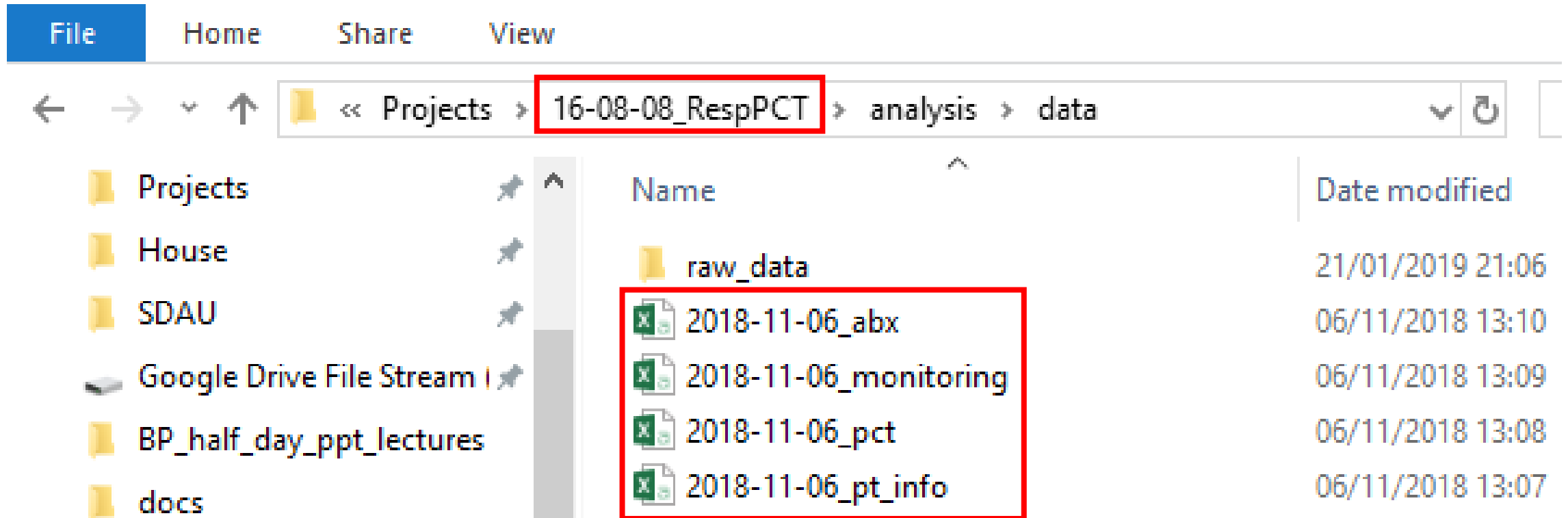- BS100_AA_with_clones
- BS100_AA_with_clones.nwk
- BS1000_AA_pyro&clones
- BS1000_AA_pyro&clones.nwk
- BS1000_AA_pyro_only
- BS1000_AA_pyro_only.nwk
- BS1000_Unique_Clonal_AA
- BS1000_Unique_Clonal_AA.nwk
- BS1000_Unique_Pyro_AA
- BS1000_Unique_Pyro_AA.nwk
- pic

## Yes

← → ∨ ↑ « Projects › 16-08-08_RespPCT › analysis › scripts

- Projects 📌
- House 📌
- SDAU 📌
- Google Drive File Stream 📌
- BP_half_day_ppt_lectures
- docs
- My Drive
- Screenshots

Name

- ℝ 01_clean_data
- ℝ 02_plots
- ℝ 03_tables
- ℝ 04_stats_analysis
- ℝ 05_post_hoc_stats
- ℝ functions
- ℝ randomization
- ℝ tables

# Outline a file naming convention

**Machine readable:**
- Inherent order
- Avoid spaces
- Avoid punctuation
- Remove case-sensitivity


**Human readable:**
- Contains info on content
- Avoid spaces
- Avoid punctuation
- Remove case sensitivity


**Metadata:**
Separate with underscores ("_")
- Avoid punctuation
- Remove case-sensitivity

```
01_marshal-data.r

02_pre-dea-filtering.r

03_dea-with-limma-voom.r

04_explore-dea-results.r

90_limma-model-term-name-fiasco.r

helper01_load-counts.r

helper02_load-exp-des.r

helper03_load-focus-statinf.r

helper04_extract-and-tidy.r
```

# Outline a file naming convention

**Chronological order:**

```
2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H01.csv
2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H02.csv
2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H03.csv
2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_platefile.csv
2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A01.csv
2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A02.csv
2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A03.csv
2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A04.csv
```

**Logical order:**

```
01_marshal-data.r
02_pre-dea-filtering.r
03_dea-with-limma-voom.r
04_explore-dea-results.r
90_limma-model-term-name-fiasco.r
helper01_load-counts.r
helper02_load-exp-des.r
helper03_load-focus-statinf.r
helper04_extract-and-tidy.r
```
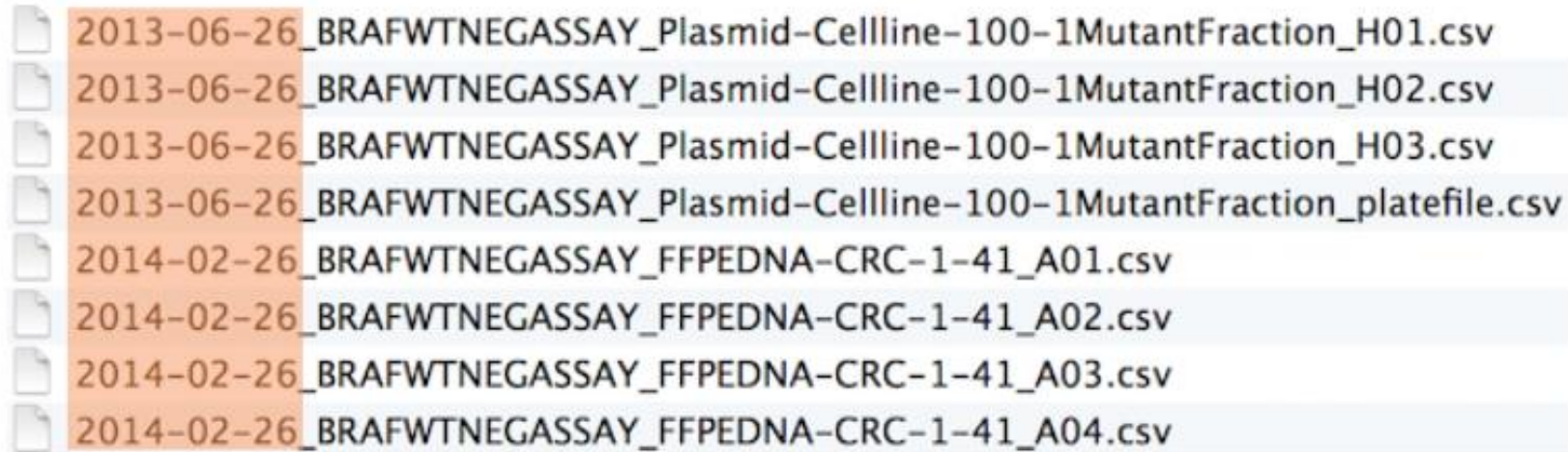
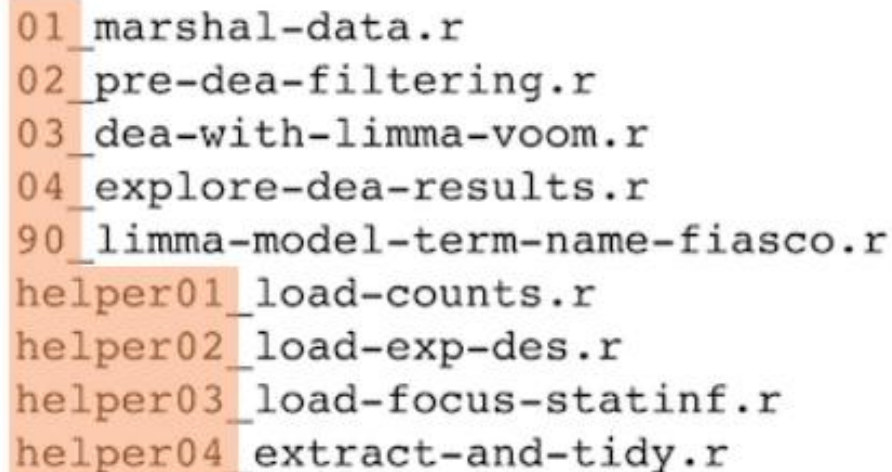# Joined up thinking

- The R scripts you generate should be human readable
    - Annotate the code
    - Break up the scripts into dedicated tasks
    - Interlink with other within project scripts

```r
# Script: 04_stats_analysis.R

# Data ----
# Four tibbles will be returned from scripts/01_clean_data.R
# 1. abx => details of the antibiotic consumption by type
# 2. monitoring => patient condition over time. Also WCC, CRP
# 3. pct => PCT values from the PCT arm of the trial
# 4. pt_info => general patient information

# Load the cleaned data sets
source("scripts/01_clean_data.R")

#Load the necessary add-on packages
library(knitr)
library(broom)
library(survminer)
```

# R Markdown

- R Markdown combines the code you wrote, the output produced and you own comments

- You can view it as a digital lab notebook, where you are both recording what you're doing, and what you were thinking while you were thinking it!

- R Markdown outputs can take many forms
        - Word documents, PDFs, slideshows etc.

# R Markdown

YAML header

```
---
title: "Diamond sizes"
date: 2016-08-25
output: html_document
---
```

Chunks of code

````
```{r setup, include = FALSE}
library(ggplot2)
library(dplyr)
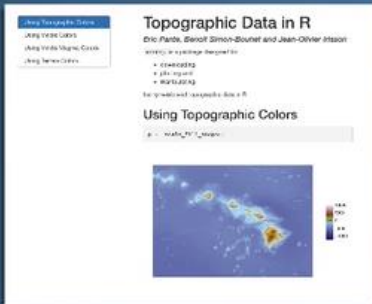smaller <- diamonds %>%
filter(carat <= 2.5)
```
````

Plain text with integrated outputs from R

```
We have data about `r nrow(diamonds)`
diamonds. Only
`r nrow(diamonds) - nrow(smaller)` are
larger than
2.5 carats. The distribution of the
remainder is shown below:
```
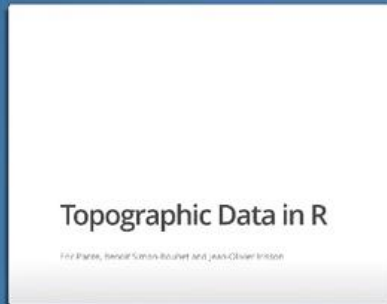
Chunks of code

````
```{r, echo = FALSE}
smaller %>%
ggplot(aes(carat)) +
geom_freqpoly(binwidth = 0.01)
```
````

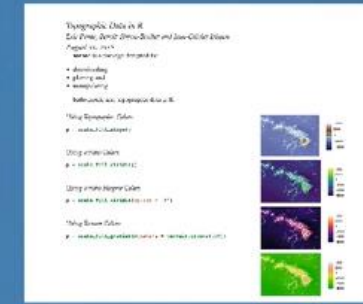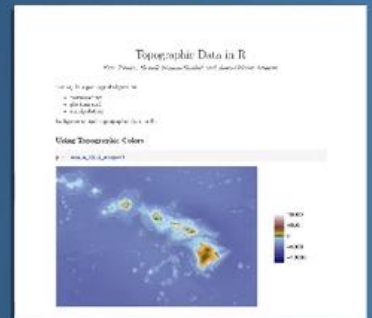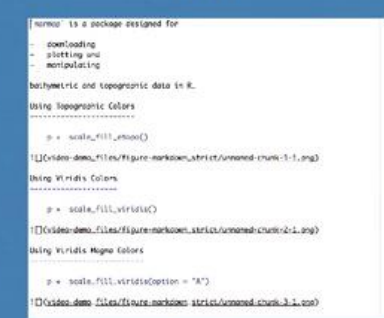# What has R Markdown ever done for us?



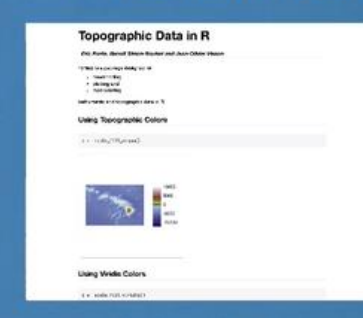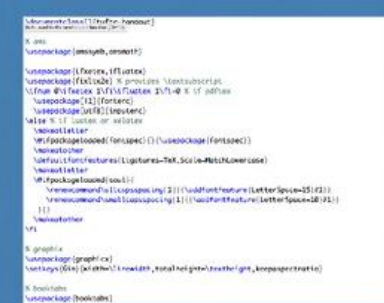html · ioslides · reveal.js · rtf · tufte handout · book · pdf · dashboard · slidy · markdown · package vignette · website · Word · notebook · beamer · latex · custom template · shiny app