# R: A Hitchhikers Guide to Reproducible Research

## - Everything in it's right place

Brendan Palmer,
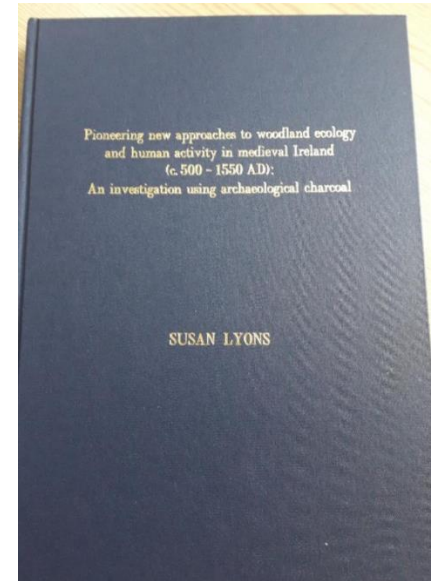
Clinical Research Facility - Cork &

School of Public Health

@B_A_Palmer

CRF-C
HRB Clinical Research Facility Cork
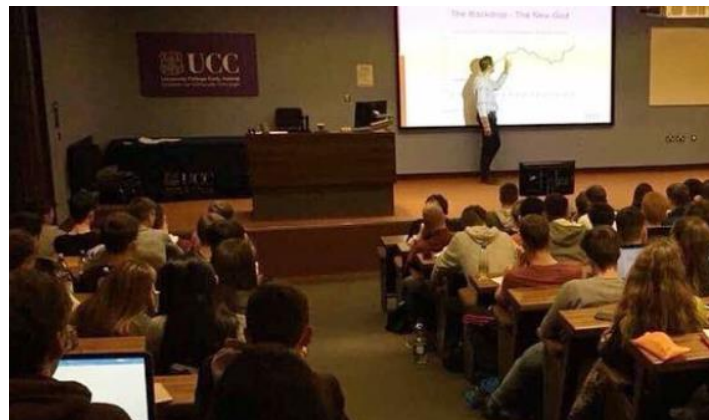
UCC
University College Cork, Ireland
Coláiste na hOllscoile Corcaigh

School of
Public Health

# How is research presented?

**Theses**



**Papers**



**Books**



**Posters**



**Talks**

# But what does it really look like?

You were defending, one foot out the door

I got your project and its problems galore

I hate my life,

THIS PERSON IS likely to be YOU BTW!!

Jorge Cham | www.phdcomics.com

# Reproducibility comes in many forms

Lab Handbook

*Candice Morey*

*2019-05-16*

## Chapter 1   Joining a lab

> No man is an island,
>
> Entire of itself,
>
> Every man is a piece of the continent,
>
> A part of the main.
>
> — John Donne

Although you will be familiar with the names of a handful of scientific heroes, science does not actually advance from the rapid insights of rare geniuses. Scientific knowledge accumulates through the consistent, often painstaking, efforts of groups of people. Across the world, webs of laboratories focusing on related topics work toward the common goal of understanding human memory and communicating how best to use our emerging knowledge to improve peoples' lives. You have opened this manual because you have joined one such group. The purpose of this manual is to help you understand your role in this endeavor and how to contribute in a way that makes your work maximally useful both to your local colleagues, your international colleagues, and the public, both now and branching into the future.

# 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

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

- Help your future-self

# R-projects

**Create Project**

**New Directory**
Start a project in a brand new working directory

**Existing Directory**
Associate a project with an existing working directory

**Version Control**
Checkout a project from a version control repository

# Step 1: Define a generic project structure

# Step 2: Give your files informative names

« example_project > data

Name

📁 raw_data

📊 2019-05-02_clean_who_tb_data

📄 README

📄 who_tb_dictionary

# Step 3: 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

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

# Step 4: 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
```

Data carpentry – File naming

# Step 4: 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
```
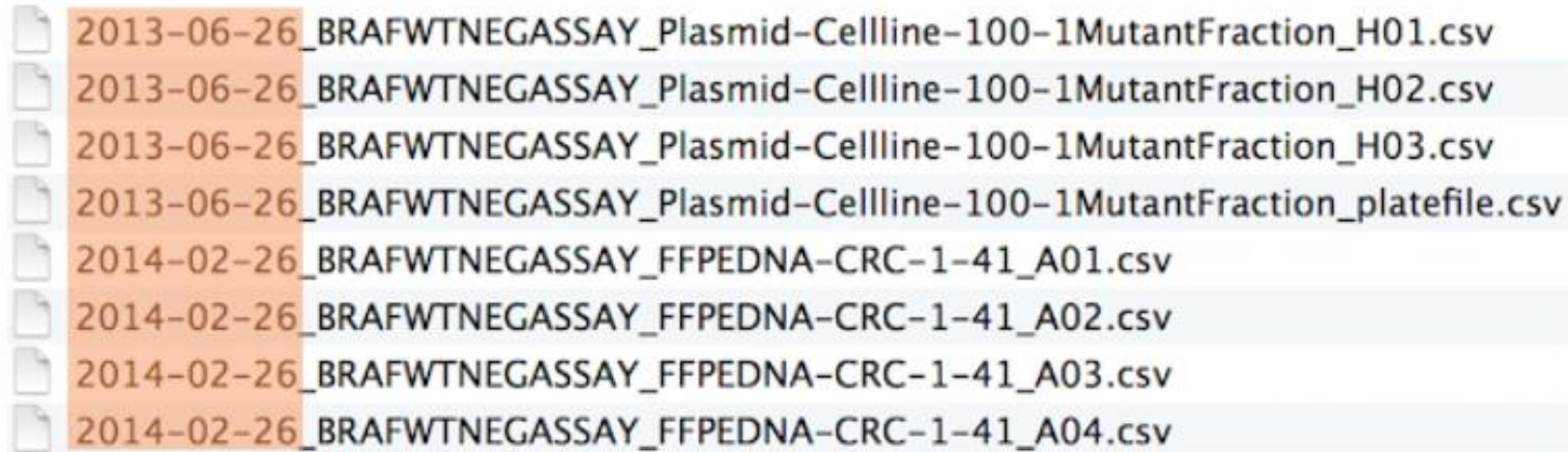
Data carpentry - File naming

# Step 5: 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)
```

# Step 6: R-projects



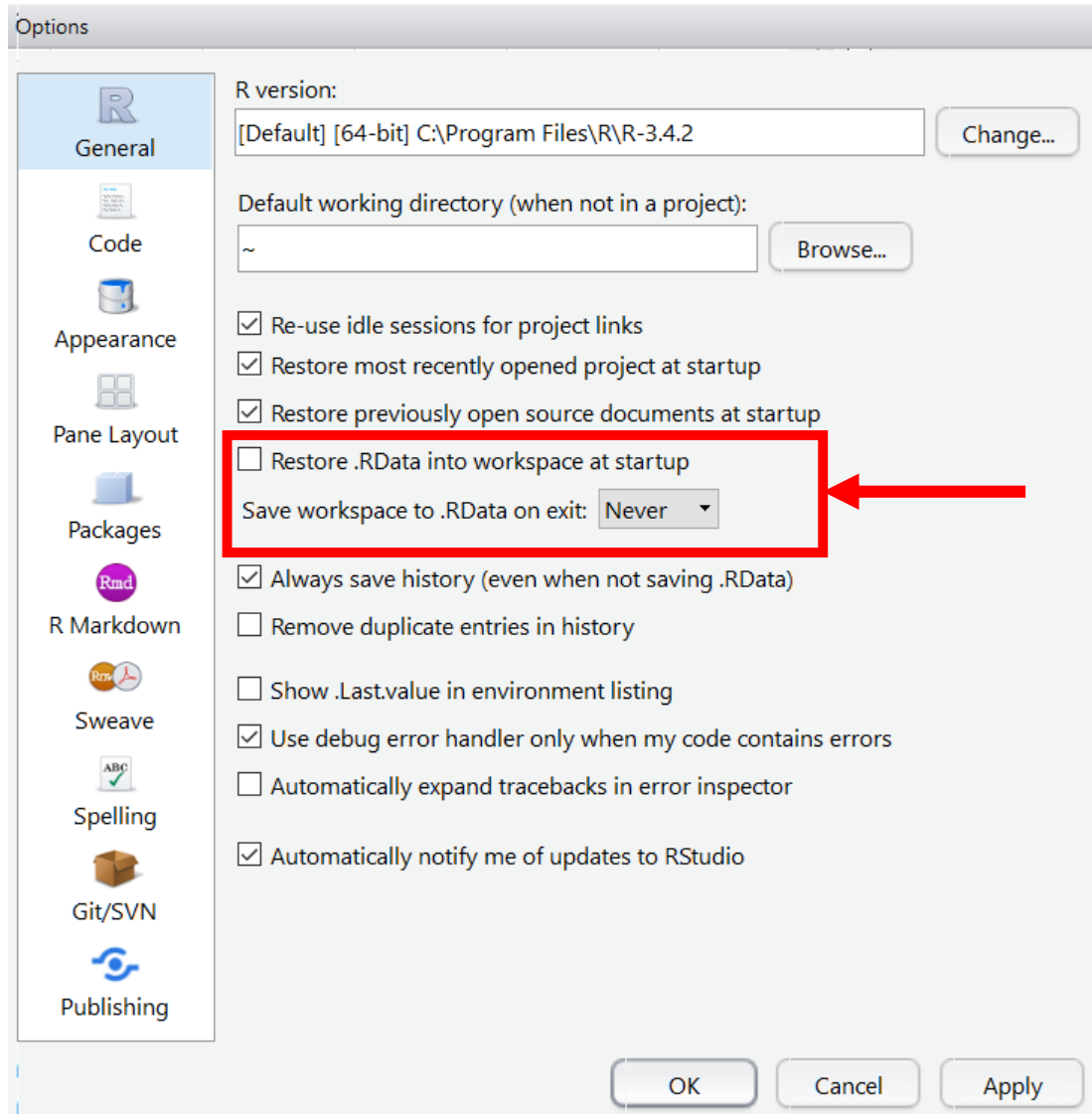R-A_Hitchhikers_Guide_to_Reproducible_Research  >  Day_1  >  example_project

Name

📁 data
📁 docs
📁 figures
📁 scripts
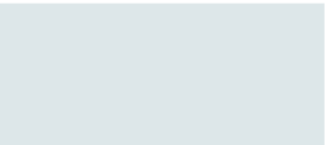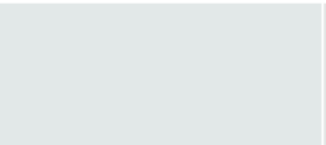📁 tables
Ⓡ all_together_now
🅡 example_project  ⟵

- Switch to the R-project file…
  Day_1/example_project/example_project.Rproj
- Open the scripts 01_eg_clean_data.R, 02_eg_figures.R and
  03_eg_analysis.R

# Other points to note



- You might consider your environment as "real"

- If you continue to use R, it is better for you to consider your R scripts as "real", as these should recreate the environment

- You may suffer short term pain

- This will prevent long term agony

# Is too much choice good or bad?

| | | | | |
|---|---|---|---|---|
| Blue Horizon SW 6497 | Sky High SW 6504 | Snowdrop SW 6511 | Ski Slope SW 6518 | Rarified Air SW 6525 |
| Byte Blue SW 6498 | Atmospheric SW 6505 | Balmy SW 6512 | Hinting Blue SW 6519 | Icelandic SW 6526 |
| Stream SW 6499 | Vast Sky SW 6506 | Take Five SW 6513 | Honest Blue SW 6520 | Blissful Blue SW 6527 |
| Open Seas SW 6500 | Resolute Blue SW 6507 | Respite SW 6514 | Notable Hue SW 6521 | Cosmos SW 6528 |
| Manitou Blue SW 6501 | Secure Blue SW 6508 | Leisure Blue SW 6515 | Sporty Blue SW 6522 | Scanda SW 6529 |
| Loch Blue SW 6502 | Georgian Bay SW 6509 | Down Pour SW 6516 | Denim SW 6523 | Revel Blue SW 6530 |
| Bosporus SW 6503 | Loyal Blue SW 6510 | Regatta SW 6517 | Cammodore SW 6524 | Indigo SW 6531 |

# Inconsistent function names, inconsistent syntax

- R is a very versatile language
    - Sometimes it can be too versatile
    - Do you want to use…
                            row.names or rownames
                            rowSums or rowsum
                            Sys.time, system.time

    - Should it be written as…
                            newobject or new.Object
                            x = 5 or x <- 5
                            mapping=aes(x,y) or mapping = aes(x, y)

# Variable selection

```
summary(starwars$name)

summary(starwars$"name")

summary(starwars["name"])

summary(starwars[ , "name"])

summary(starwars[1])

summary(starwars[ , 1])

summary(starwars[[1]])
```

- Open the script 04_too_much.choice.R

# Motivation to move on from poorly written code



```r
21  sites1<-as.list(unique(RL6.7$Var1))
22  sites2<-as.list(unique(RL6.7$Var2))
23
24  sites<-as.data.frame(t(merge(sites1,sites2)))
25  colnames(sites)[1]<-"Position"
26
27  for(i in 1:nrow(sites)){
28  ans<-(sites$Position[i]<=65)
29  sites$E1[i]<-ans
30  }
31
32  # Start building network
33  RL6.7_topology<-subset(RL6.7[2:3])
34  g2<-graph.data.frame(RL6.7_topology,vertices=sites,directed=FALSE)
35  V(g2)$color<-ifelse(V(g2)$E1==TRUE,"white","grey")
36  V(g2)$color<-ifelse(V(g2)$E1==TRUE,"white","grey")
37  plot(g2,vertex.label.color="black",vertex.size=20,edge.color="black",edge.width=1.5)
38
```

Lack of annotation
Poor naming conventions
Poor readability
Spacing absent

Cluttered environment
Intermediate objects

- Open the script 05_bad_habits.R

# Writing clearer code

- Annotation

- Object names
    - should use only lowercase letters, numbers, and "_"

- Spacing
    - Put a space before and after =
    - Put a space after a ,
    - Operators should be surrounded by spaces e.g. ==, <-, +

- For a more complete list visit
    - http://style.tidyverse.org/syntax.html

- Open the script 06_good_habits.R