

# Reproducible Research

---

Marco Chiapello

July 4, 2016

Center for Proteomics  
University of Cambridge  
*mc983@cam.ac.uk*

# Overview

Introduction

Reproducible research Reasons

Reproducible research Rules

Reproducible research Tools

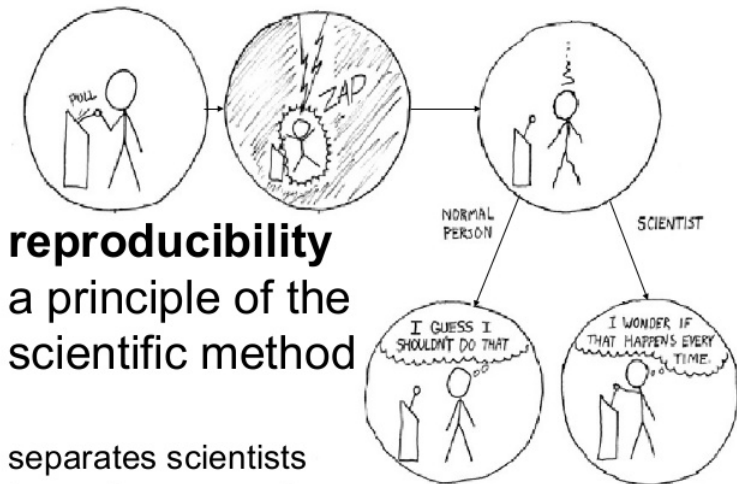
Conclusion

# Introduction

---

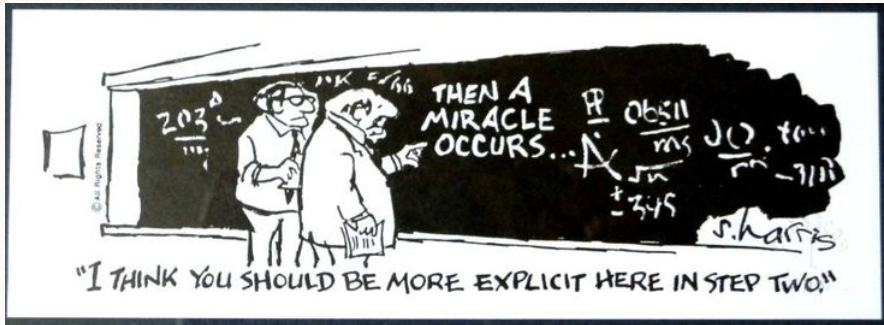
**Replication** is the ultimate standard by which scientific claims are judged [2]  
The fact that an analysis is reproducible does not guarantee the quality, correctness,  
or validity of the published results.

# What reproducible research is



<http://xkcd.com/242/>

# What reproducible research is



- This is exactly how it seems when you try to figure out how authors got from a large and complex data set to a dense paper with lots of busy figures.

Without access to the **data and the analysis code**, a miracle occurred.

- And there should be NO MIRACLES IN SCIENCE. [1]

$$DATA + ANALYSIS \rightarrow RESULTS$$

---

Common practice of writing statistical reports:

- We import a dataset into Excel
- Run a procedure to get all results
- Copy and paste selected pieces into a typesetting program
- Add a few descriptions
- Finish a report

# What reproducible research is

There are obvious dangers and disadvantages in this process:

1. It is **error-prone** due to too much manual work;
2. It requires lots of human effort to do **tedious jobs**;
3. The workflow is barely recordable, therefore it is **difficult to reproduce**;
4. A **tiny change** of the data source in the future will require the author(s) to go through the same procedure again;
5. The analysis and writing are separate, so close attention has to be paid to the **synchronization of the two parts**.

### What is Reproducible Research?

THE ABILITY TO REPRODUCE SOMEONE ELSE  
RESULTS

---

What do you need?

- Analytic data
- Analytic code
- **Documentation for data and code**



# What reproducible research is

## REPRODUCIBLE VS REPLICABLE

		DATA	
		Same	Different
CODE	Same	Reproducible	Replicable
	Different	Robust	Generalisable

Ref: <https://github.com/KirstieJane/ReproducibleResearch>

# What reproducible research is

## Reproducibility/reproduce

A study is reproducible if there is a specific set of computational functions/analyses (usually specified in terms of code) that **exactly reproduce all of the numbers in a published paper from raw data**.

## Replication/replicate

A study is only replicable if you perform the exact same experiment (at least) twice, collect data in the same way both times, perform the same data analysis, and **arrive at the same conclusions**.

---

**Reproducibility** is, to some extent, a technical challenge, while replication gives the results scientific validity. **Replicability** requires new samples and new data<sup>1</sup>, which introduces new variability, and additional risks of errors.

---

Ref: <https://github.com/lgatto/TeachingMaterial/tree/master/open-rr-bioinfo-best-practice>

<sup>1</sup>in particular biological replicates

# Reproducible research Reasons

---

How does working reproducibly help to achieve more as a scientist [1]

## REPRODUCIBILITY [1]

### Idealist:

1. It is the foundation of science!
2. The world would be a better place if everyone worked transparently and reproducibly!

## REPRODUCIBILITY [1]

### Idealist:

1. It is the foundation of science!
2. The world would be a better place if everyone worked transparently and reproducibly!

### Realist:

1. It helps to avoid disaster
  - You need to record in detail how you got there
  - Work reproducibly early on will save you time later

## REPRODUCIBILITY [1]

### Idealist:

1. It is the foundation of science!
2. The world would be a better place if everyone worked transparently and reproducibly!

### Realist:

1. It helps to avoid disaster
  - You need to record in detail how you got there
  - Work reproducibly early on will save you time later
2. It makes it easier to write papers
  - To have very transparent data and code, it costs just few minutes to spot a mistake (if any)

## REPRODUCIBILITY [1]

### Idealist:

1. It is the foundation of science!
2. The world would be a better place if everyone worked transparently and reproducibly!

### Realist:

1. It helps to avoid disaster
  - You need to record in detail how you got there
  - Work reproducibly early on will save you time later
2. It makes it easier to write papers
  - To have very transparent data and code, it costs just few minutes to spot a mistake (if any)
3. It helps reviewers see it your way
  - Made the data and well-documented code easily accessible to the reviewers

## REPRODUCIBILITY [1]

### Idealist:

1. It is the foundation of science!
2. The world would be a better place if everyone worked transparently and reproducibly!

### Realist:

1. It helps to avoid disaster
  - You need to record in detail how you got there
  - Work reproducibly early on will save you time later
2. It makes it easier to write papers
  - To have very transparent data and code, it costs just few minutes to spot a mistake (if any)
3. It helps reviewers see it your way
  - Made the data and well-documented code easily accessible to the reviewers
4. It enables continuity of your work
  - How can you ensure the continuity of work in your lab if progress is not documented reproducibly?
  - No proof of reproducibility, no result!



## REPRODUCIBILITY [1]

### Idealist:

1. It is the foundation of science!
2. The world would be a better place if everyone worked transparently and reproducibly!

### Realist:

1. It helps to avoid disaster
  - You need to record in detail how you got there
  - Work reproducibly early on will save you time later
2. It makes it easier to write papers
  - To have very transparent data and code, it costs just few minutes to spot a mistake (if any)
3. It helps reviewers see it your way
  - Made the data and well-documented code easily accessible to the reviewers
4. It enables continuity of your work
  - How can you ensure the continuity of work in your lab if progress is not documented reproducibly?
  - No proof of reproducibility, no result!
5. It helps to build your reputation
  - To build a reputation for being an honest and careful researcher

# Reproducible research Rules

---

– based on Sandve et al., 2013 [3]

## Rule 1

FOR EVERY RESULT, KEEP TRACK OF HOW IT WAS  
PRODUCED

# Rule 1

## FOR EVERY RESULT, KEEP TRACK OF HOW IT WAS PRODUCED

- The **full sequence** of pre- and post-processing steps are often critical in order to reach the achieved result
- **Every detail** that may influence the execution of the step **should be recorded**
- Include the name and version of the program, as well as the exact parameters and inputs

*As a minimum, you should at least record sufficient details on programs, parameters, and manual procedures to allow yourself, in a year or so, to approximately reproduce the results*

## Rule 2

AVOID MANUAL DATA MANIPULATION STEPS

## Rule 2

### AVOID MANUAL DATA MANIPULATION STEPS

- Manual procedures are not only inefficient and error-prone, they are also difficult to reproduce
- Manual modification of files can usually be replaced by the use of standard UNIX commands or scripts
- Manual tweaking of data files to attain format compatibility should be replaced by format converters that can be reenacted and included into executable workflows
- Manual operations like the use of **copy and paste** between documents should also be avoided

*If manual operations cannot be avoided, you should as a minimum note down which data files were modified or moved, and for what purpose*

## Rule 3

ARCHIVE THE EXACT VERSIONS OF ALL EXTERNAL  
PROGRAMS USED

### ARCHIVE THE EXACT VERSIONS OF ALL EXTERNAL PROGRAMS USED

- In order to exactly reproduce a given result, it may be necessary to use programs in the **exact versions used originally**
- It is not always trivial to get hold of a program in anything but the current version

*As a minimum, you should note the exact names and versions of the main programs you use*



VERSION CONTROL ALL CUSTOM SCRIPTS

### VERSION CONTROL ALL CUSTOM SCRIPTS

- **Only that exact state of the script may be able to produce that exact output**, even given the same input data and parameters
- The standard solution to track evolution of code is to use a version control system
  - A version control system is a repository of files with monitored access.  
*Every change made to the source is tracked, along with who made the change, why they made it*

*As a minimum, you should archive copies of your scripts from time to time*

## Rule 5

RECORD ALL INTERMEDIATE RESULTS, WHEN POSSIBLE IN  
STANDARDIZED FORMATS

## Rule 5

### RECORD ALL INTERMEDIATE RESULTS, WHEN POSSIBLE IN STANDARDIZED FORMATS

- In principle, as long as the **full process** used to produce a given result is tracked, all **intermediate data can also be regenerated**
- In practice, having easily **accessible intermediate results** may be of great value
- When the full process is not readily executable, it allows parts of the process to be rerun
- It **allows critical examination** of the full process behind a result

*As a minimum, archive any intermediate result files that are produced when running an analysis*

## Rule 6

FOR ANALYSES THAT INCLUDE RANDOMNESS, NOTE  
UNDERLYING RANDOM SEEDS

## Rule 6

### FOR ANALYSES THAT INCLUDE RANDOMNESS, NOTE UNDERLYING RANDOM SEEDS

- Many analyses and predictions include some element of randomness, meaning the same program will typically give **slightly different results** every time it is executed
- Given the **same initial seed**, all random numbers used in an analysis will be equal, thus giving identical results every time it is run

*As a minimum, you should note which analysis steps involve randomness, so that a certain level of discrepancy can be anticipated when reproducing the results*

ALWAYS STORE RAW DATA

### ALWAYS STORE RAW DATA

- ALWAYS store in a safe place the raw data
- NEVER touch or modify the raw data



GENERATE HIERARCHICAL ANALYSIS OUTPUT, ALLOWING  
LAYERS OF INCREASING DETAIL TO BE INSPECTED

### GENERATE HIERARCHICAL ANALYSIS OUTPUT, ALLOWING LAYERS OF INCREASING DETAIL TO BE INSPECTED

- The final results that make it to an article, be it plots or tables, often represent highly summarized data
- In order to validate and fully understand the main result, it is often useful to inspect the detailed **values underlying the summaries**
- When working with summarized results, you should as a minimum at least once generate, inspect, and validate the detailed values underlying the summaries

## Rule 9

CONNECT TEXTUAL STATEMENTS TO UNDERLYING RESULTS

## Rule 9

### CONNECT TEXTUAL STATEMENTS TO UNDERLYING RESULTS

- The results of analyses and their corresponding textual interpretations are clearly interconnected but often **lie in different places**
- Results usually live on a personal computer, while interpretations live in text documents
- To allow efficient retrieval of details behind textual statements, we suggest that **statements are connected to underlying results** already from the time the statements are initially formulated
- **Integrate reproducible analyses directly into textual documents**

PROVIDE PUBLIC ACCESS TO SCRIPTS, RUNS, AND RESULTS

### PROVIDE PUBLIC ACCESS TO SCRIPTS, RUNS, AND RESULTS

- All input data, scripts, versions, parameters, and intermediate **results should be made publicly and easily accessible**
- Making reproducibility of your work by peers a realistic possibility sends a **strong signal of quality, trustworthiness, and transparency**

# Reproducible research Tools

---

*Let us change our traditional attitude to the construction of programs:  
Instead of imagining that our main task is to instruct a computer what to  
do, let us concentrate rather on explaining to humans what we want the  
computer to do.*

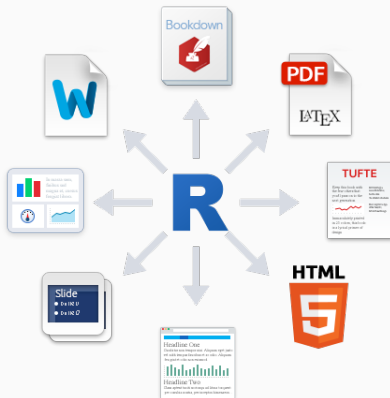
*– Donald E. Knuth Literate Programming, 1984*

**Literate programming** is a methodology that combines a programming language with a documentation language

- Write program code to do computing
- Write narratives to explain what is being done by the program code



## RMarkdown



Ref: <http://rmarkdown.rstudio.com/index.html>

## RMarkdown

```
---  
title: "Untitled"  
author: "Marco Chiapello"  
date: "10 June 2016"  
output: html_document  
---
```

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
```{r}  
summary(cars)  
```
```

You can also embed plots, for example:

```
```{r, echo=FALSE}  
plot(cars)  
```
```

Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.

## FOLDER ORGANIZATION

+– README

+– codeBook

+== rawdata

+== rscript

+== analysis

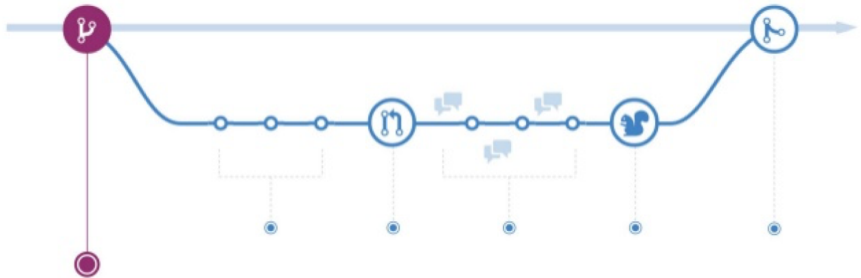
+== docs

+== manuscript

+== tmp

# VERSION CONTROL

## GitHub Flow - Create a Branch



Learning to use these tools will require **commitment** and a **massive investment of your time and energy**.

A priori it is not clear why the benefits of working reproducibly outweigh its costs.

**Does reproducibility sound like extra work?**

It can be, particularly when one is first trying to do it, that is, to break one's own previous nonreproducible habits

## Conclusion

---

# MY ADVICE IS:

Learn the tools of reproducibility as quickly as possible and use them in every project.

## References

---

- [1] Markowetz, F. (2016). Five selfish reasons to work reproducibly. *Genome biology*, pages 1–4.
- [2] Peng, R. D. (2011). Reproducible research in computational science. *Science (New York, NY)*, 334(6060):1226–1227.
- [3] Sandve, G. K., Nekrutenko, A., Taylor, J., and Hovig, E. (2013). Ten Simple Rules for Reproducible Computational Research. *PLoS Computational Biology*, 9(10):e1003285–4.