DTU

*Leon Eyrich Jessen*

$$P_{RG} = \frac{AP+Sp-1}{Se+Sp-1}$$

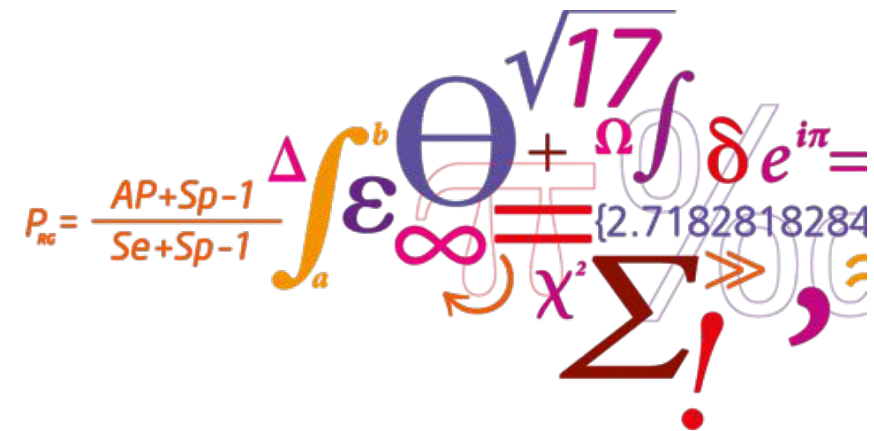**Immunoinformatics and Machine Learning**
Department of Bio and Health Informatics

- *Eur J Heart Fail: "Our findings do not support a major role for fish intake in the prevention of heart failure"*

- Eur J Clin Nutr: *"Moderate consumption of fatty fish … were associated with a lower rate of first HF hospitalization or death"*

1. Eur J Heart Fail. 2009 Oct;11(10):922-8. doi: 10.1093/eurjhf/hfp126
2. Eur J Clin Nutr. 2010 Jun;64(6):587-94. doi: 10.1038/ejcn.2010.50. Epub 2010 Mar 24

- *Eur J Heart Fail: "Our findings do not support a major role for fish intake in the prevention of heart failure"*

- *Eur J Clin Nutr: "Moderate consumption of fatty fish … were associated with a lower rate of first HF hospitalization or death"*

**FAKE NEWS?**

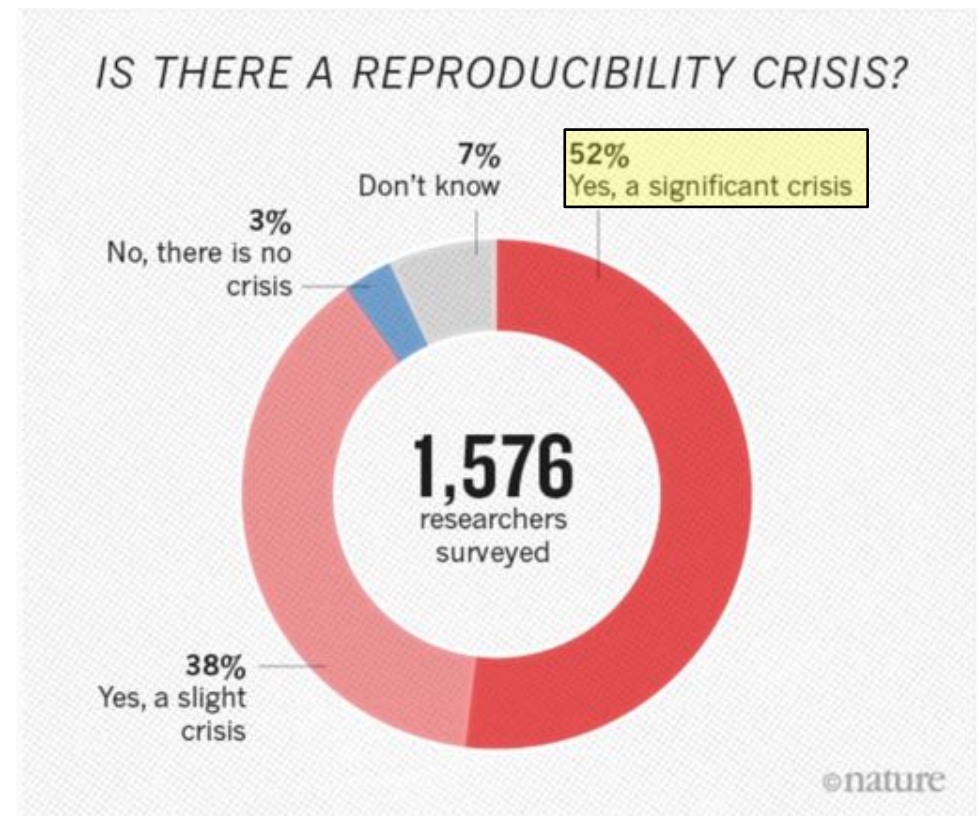1. Eur J Heart Fail. 2009 Oct;11(10):922-8. doi: 10.1093/eurjhf/hfp126
2. Eur J Clin Nutr. 2010 Jun;64(6):587-94. doi: 10.1038/ejcn.2010.50. Epub 2010 Mar 24

# Plenum

- *Take 1 minute to discuss with your neighbour:*

  - *Possible <u>consequences</u> of contradictory research results?*

# Nature | News Feature (May 2016)

- 1,500 scientists lift the lid on reproducibility

- Based on questionnaire

- 52% "Yes, a significant crisis"

- "More than 70% of researchers have tried and failed to reproduce another scientist's experiments"

- "More than 50% have failed to reproduce their own experiments"
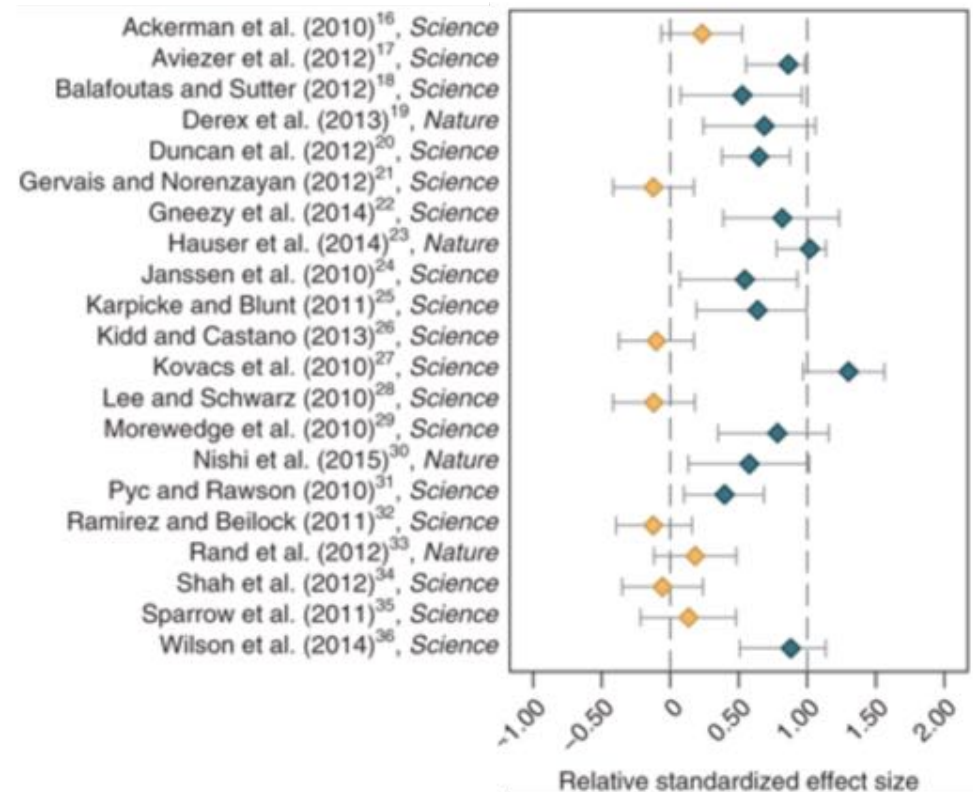


IS THERE A REPRODUCIBILITY CRISIS?

7% Don't know

52% Yes, a significant crisis

3% No, there is no crisis

1,576 researchers surveyed

38% Yes, a slight crisis

©nature

*Baker, M., Nature 533, 452–454 (26 May 2016) doi:10.1038/533452a*

# Nature | Human Behaviour (Aug. 2018)

- Evaluating the replicability of social science experiments in _Nature_ and _Science_ between 2010 and 2015

- Based on re-analysis of published results

- "There is a significant effect in the same direction as in the original study for 13 out of 21 replications" (62%)

- I.e. 38% of the studies failed to replicate

◆ Replicated    ◆ Not replicated



Ackerman et al. (2010)[16], _Science_
Aviezer et al. (2012)[17], _Science_
Balafoutas and Sutter (2012)[18], _Science_
Derex et al. (2013)[19], _Nature_
Duncan et al. (2012)[20], _Science_
Gervais and Norenzayan (2012)[21], _Science_
Gneezy et al. (2014)[22], _Science_
Hauser et al. (2014)[23], _Nature_
Janssen et al. (2010)[24], _Science_
Karpicke and Blunt (2011)[25], _Science_
Kidd and Castano (2013)[26], _Science_
Kovacs et al. (2010)[27], _Science_
Lee and Schwarz (2010)[28], _Science_
Morewedge et al. (2010)[29], _Science_
Nishi et al. (2015)[30], _Nature_
Pyc and Rawson (2010)[31], _Science_
Ramirez and Beilock (2011)[32], _Science_
Rand et al. (2012)[33], _Nature_
Shah et al. (2012)[34], _Science_
Sparrow et al. (2011)[35], _Science_
Wilson et al. (2014)[36], _Science_

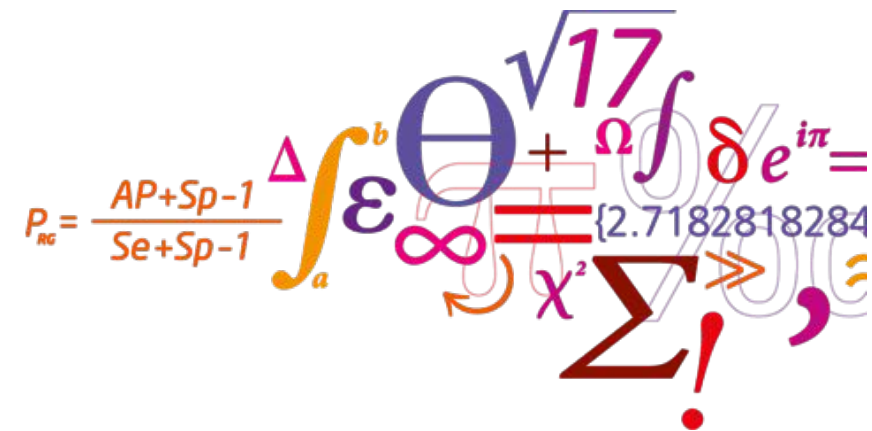Relative standardized effect size

# Plenum

- *Take 1 minute to discuss with your neighbour:*

  – *Possible <u>reasons</u> for irreproducible research results?*

# Reasons for lack of reproducibility

- Many!

- Many of which cannot be controlled

- Focus for this talk is on one aspect we can control

- Namely
  - Reproducible data analysis workflow

# Motivation - Why bother?

- You are obliged to make reproducible research, it's the cornerstone of what we do

- During a review process, you might be asked to redo and/or expand your analysis

- If you revisit an unorganised project, 2 years later, you will loose *a lot* of time redoing everything

- Even worse, if you revisit someone else's project 2 years later…

- So, spend time to save time

- …and once again, we *really have to be able to redo an analysis* as effortlessly as possible and it's your responsibility, not your supervisors

- IMHO: Every detail of the analysis in a paper should be open source (100% transparency)
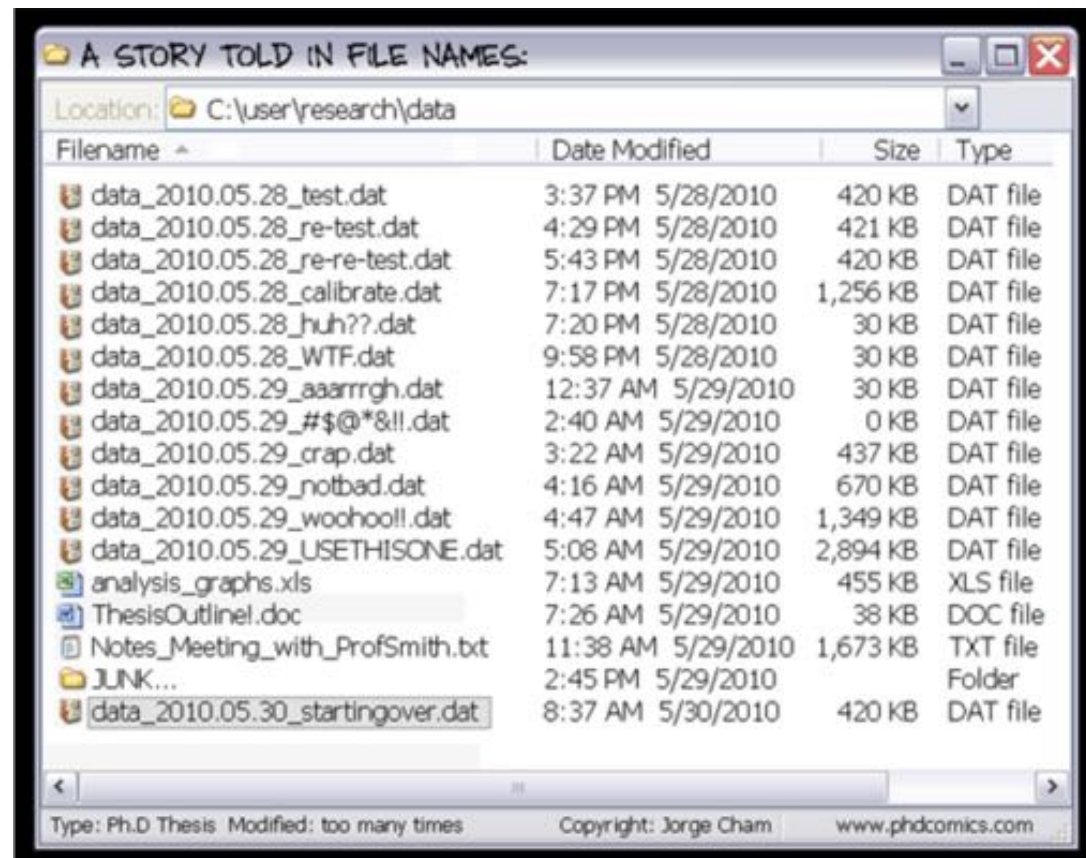
# Learning Objectives

- After this session, you should be able to:

  1. *Define what a reproducible data analysis workflow is*

  2. *List the elements of a reproducible data analysis workflow*

  3. *Explain the meaning and purpose of each of the elements in a reproducible data analysis workflow*

  4. *When presented with a pre-made workflow, determine if it constitutes a reproducible data analysis workflow*
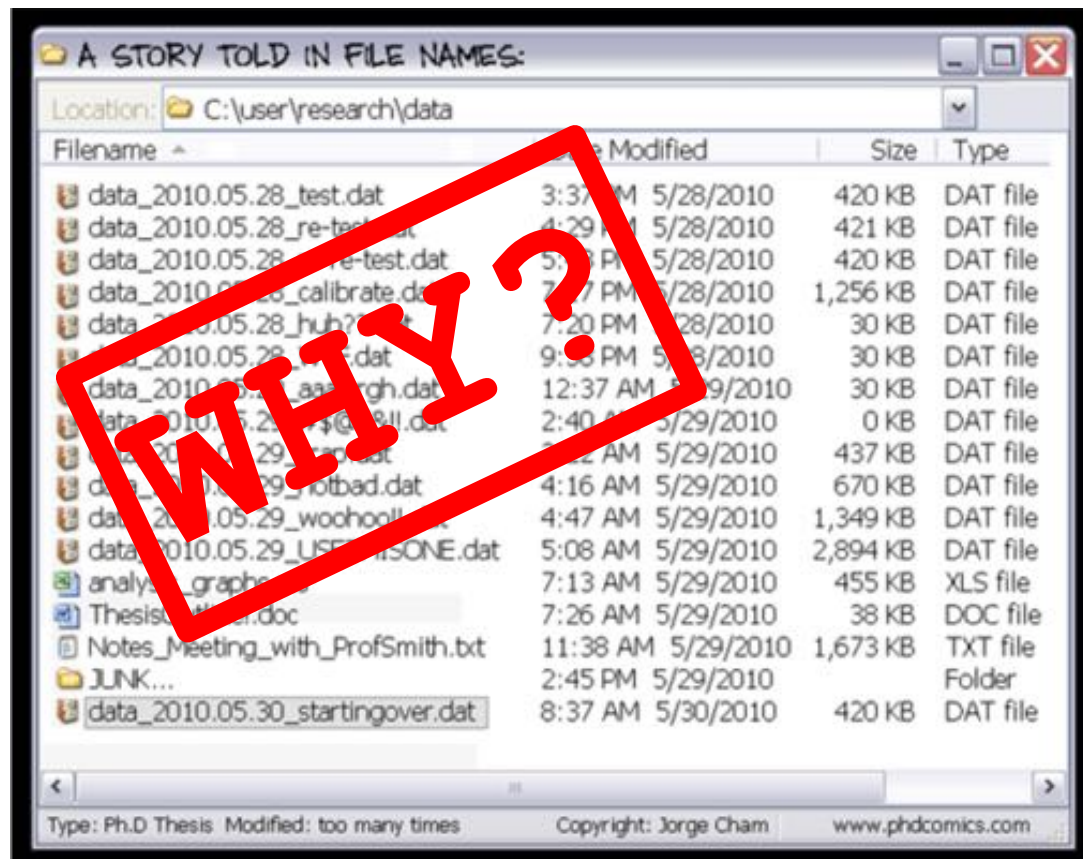
# Definition

- *A reproducible data analysis workflow is when you can go from the <u>raw data</u> to recreating all the <u>figures, tables and numbers</u> in your paper automatically and consistently*

# I've seen several data dirs like this

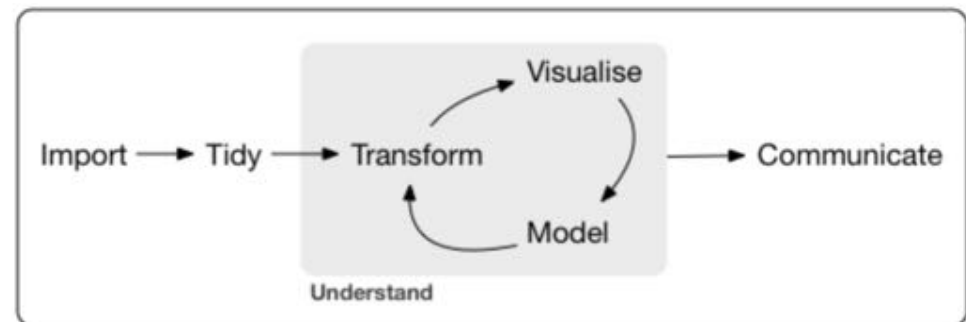# Most research deals with data, so...

# Something Essential You (most likely) Were Not Taught

- *"In practice, the principles behind organizing and documenting computational experiments are often learned on the fly"*

OPEN ACCESS Freely available online 2009                    PLoS COMPUTATIONAL BIOLOGY

**Education**

# A Quick Guide to Organizing Computational Biology Projects

**William Stafford Noble**[1,2*]

1 Department of Genome Sciences, School of Medicine, University of Washington, Seattle, Washington, United States of America, 2 Department of Computer Science and Engineering, University of Washington, Seattle, Washington, United States of America

# The Elements of Data Analysis

- Import

- Tidy

- Transform

- Visualise

- Model

- Communicate

**DTU Bioinformatics, Technical University of Denmark**

December 3rd 2018

# The Elements of Data Analysis

- Import
  - Import and combine your raw data from (potentially) multiple sources
- Tidy
  - Clean variables (e.g. missing data), setup observations as rows and variables as columns
- Transform
  - Compute new variables (augment the data), reduce data to desired focus
- Visualise
  - Explore and understand your data by seeing it, generate questions
- Model
  - Answer initial and generated questions, extract value, gain insight
- Communicate
  - Condense and communicate gained insight via essential, well defined and focused plots

*http://r4ds.had.co.nz/introduction.html*

# The Elements of Data Analysis

- Import
  - Import and combine your raw data from (potentially) multiple sources

- **Tidy**
  - **Clean variables (e.g. missing data), setup observations as rows and variables as columns**
- Transform
  - Compute new variables (augment the data), reduce data to desired focus
- Visualise
  - Explore and understand your data by seeing it, generate questions
- Model
  - Answer initial and generated questions, extract value, gain insight
- Communicate
  - Condense and communicate gained insight via essential, well defined and focused plots

*http://r4ds.had.co.nz/introduction.html*

# The Elements of Data Analysis

- Import
  - Import and combine your raw data from (potentially) multiple sources
- Tidy
  - Clean variables (e.g. missing data), setup observations as rows and variables as columns
- Transform
  - Compute new variables (augment the data), reduce data to desired focus
- Visualise
  - Explore and understand your data by seeing it, generate questions
- Model
  - Answer initial and generated questions, extract value, gain insight
- Communicate
  - Condense and communicate gained insight via essential, well defined and focused plots

*http://r4ds.had.co.nz/introduction.html*

# The Elements of Data Analysis

• Import
   – Import and combine your raw data from (potentially) multiple sources
• Tidy
   – Clean variables (e.g. missing data), setup observations as rows and variables as columns
• Transform
   – Compute new variables (augment the data), reduce data to desired focus

• Visualise
   – Explore and understand your data by seeing it, generate questions
• Model
   – Answer initial and generated questions, extract value, gain insight
• Communicate
   – Condense and communicate gained insight via essential, well defined and focused plots
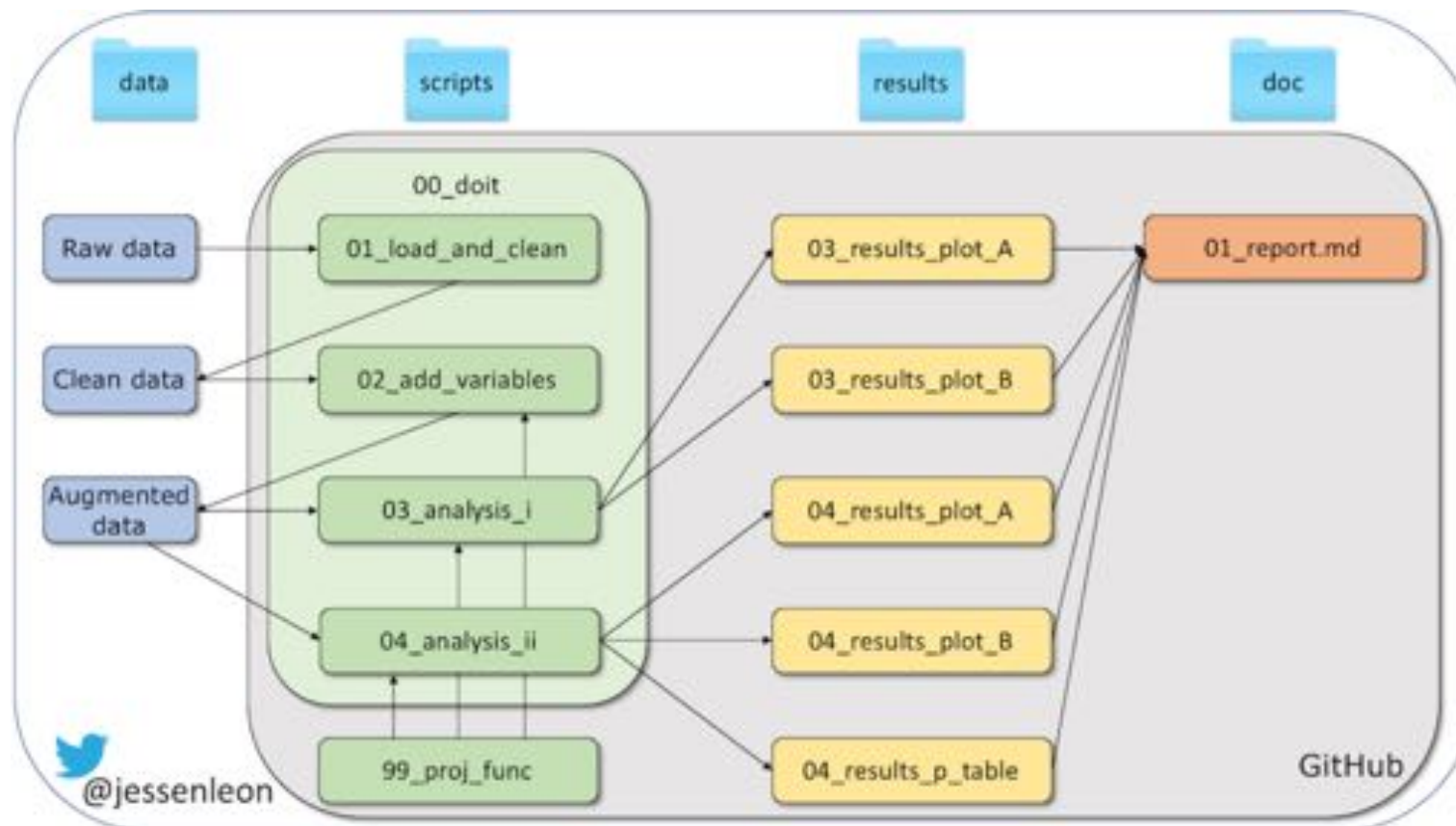
*http://r4ds.had.co.nz/introduction.html*

**DTU Bioinformatics, Technical University of Denmark**
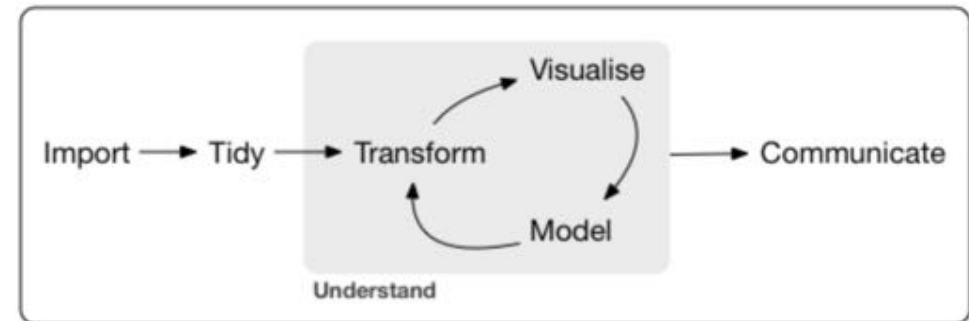
December 3rd 2018

# The Elements of Data Analysis

- Import
  - Import and combine your raw data from (potentially) multiple sources
- Tidy
  - Clean variables (e.g. missing data), setup observations as rows and variables as columns
- Transform
  - Compute new variables (augment the data), reduce data to desired focus
- Visualise
  - Explore and understand your data by seeing it, generate questions
- Model
  - Answer initial and generated questions, extract value, gain insight
- Communicate
  - Condense and communicate gained insight via essential, well defined and focused plots
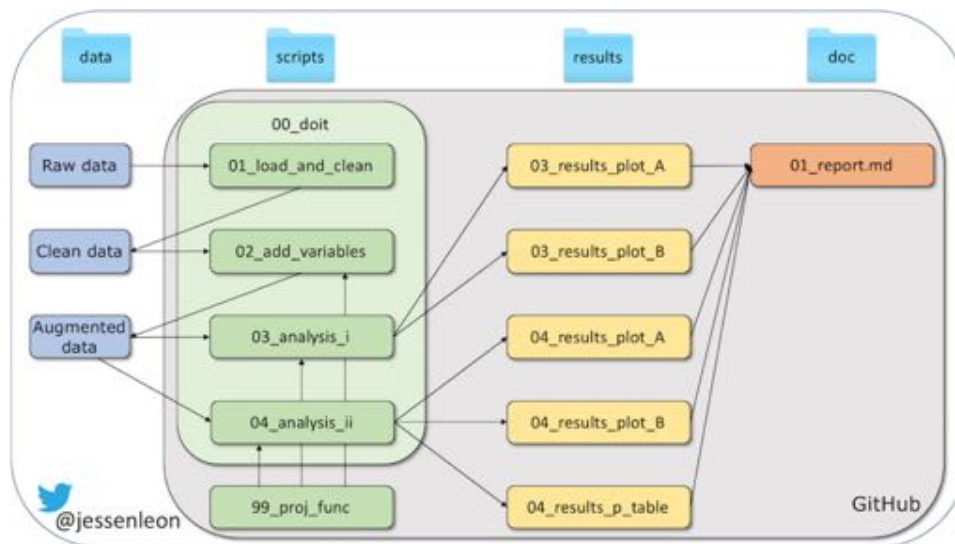
*http://r4ds.had.co.nz/introduction.html*

# The Elements of Data Analysis

- Import
  - Import and combine your raw data from (potentially) multiple sources
- Tidy
  - Clean variables (e.g. missing data), setup observations as rows and variables as columns
- Transform
  - Compute new variables (augment the data), reduce data to desired focus
- Visualise
  - Explore and understand your data by seeing it, generate questions
- Model
  - Answer initial and generated questions, extract value, gain insight
- Communicate
  - Condense and communicate gained insight via essential, well defined and focused plots
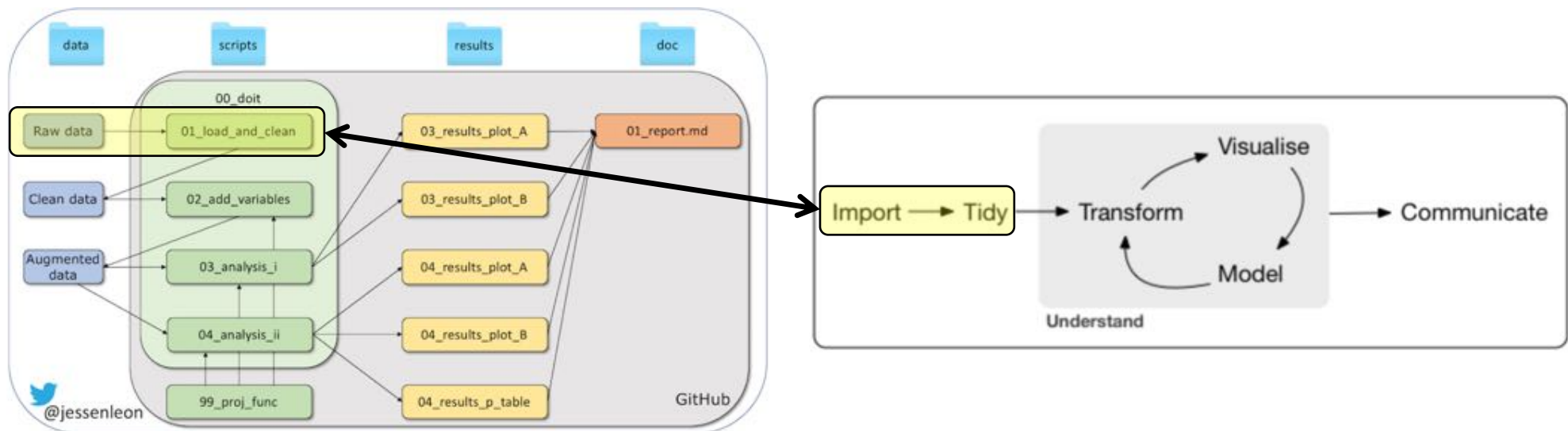
*http://r4ds.had.co.nz/introduction.html*
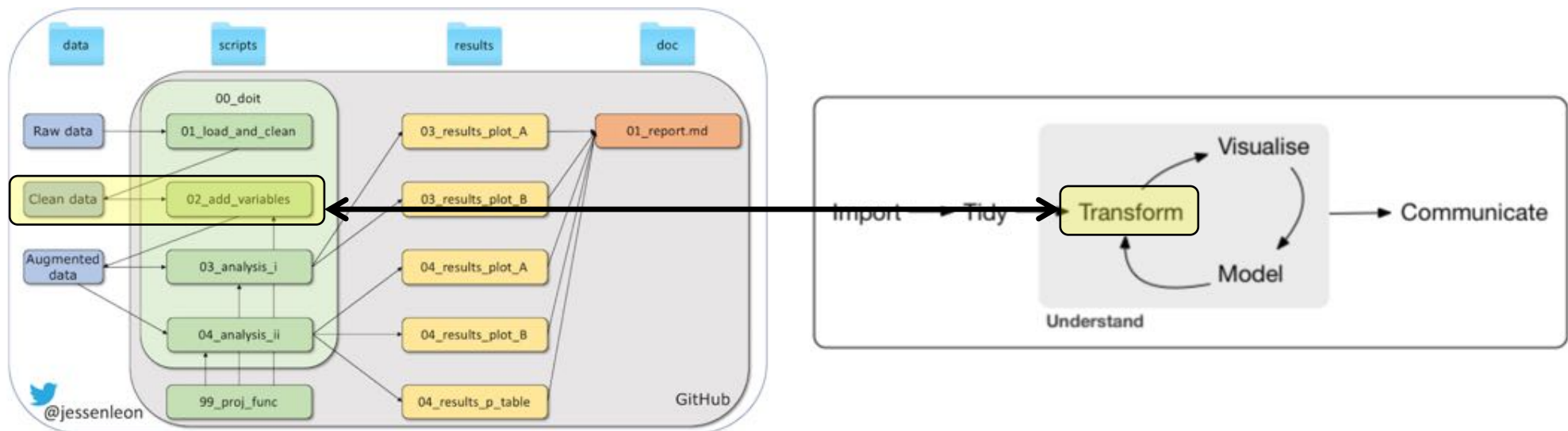
# Organising Your Data Analysis Project



**DTU Bioinformatics, Technical University of Denmark**                    December 3rd 2018

# Combining Project Organisation with DA elements



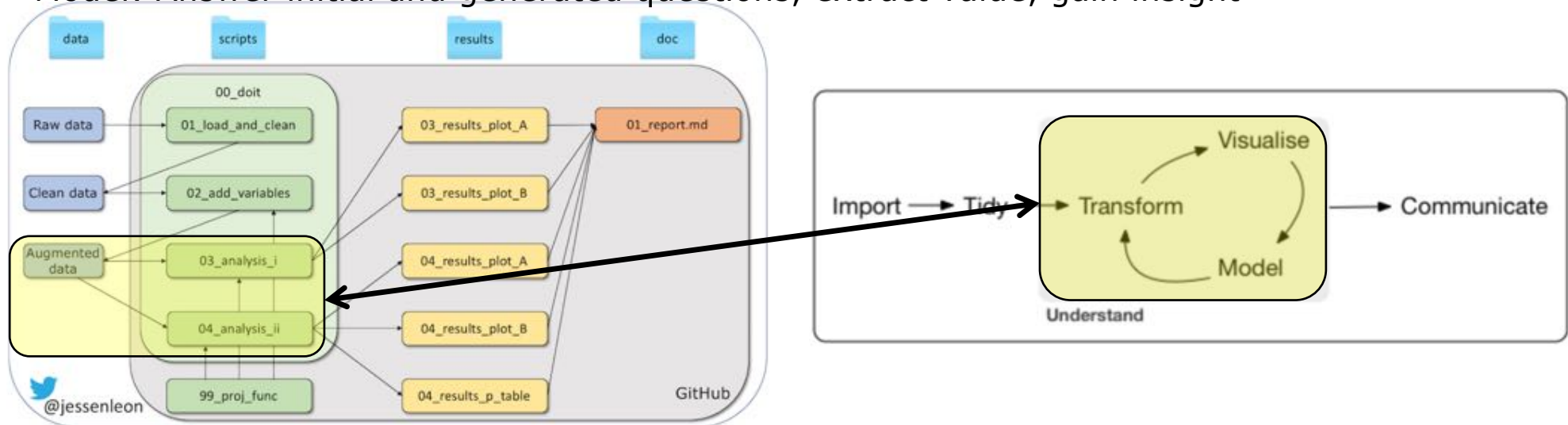**DTU Bioinformatics, Technical University of Denmark**

- Import: Import and combine your raw data from (potentially) multiple sources
- Tidy: Clean variables (eg. missing data), setup observations as rows and variables as columns



**DTU Bioinformatics, Technical University of Denmark**                                                    December 3rd 2018
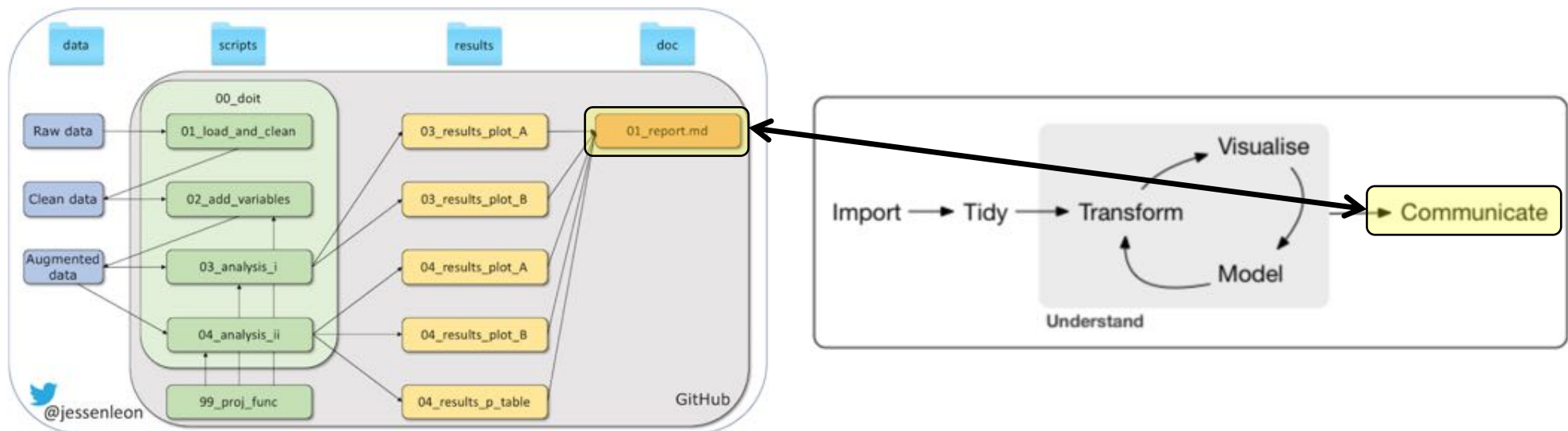
- Transform: Compute new variables (augment the data), reduce data to desired focus

- Transform: Compute new variables (augment the data), reduce data to desired focus
- Visualise: Explore and understand your data by seeing it, generate questions
- Model: Answer initial and generated questions, extract value, gain insight
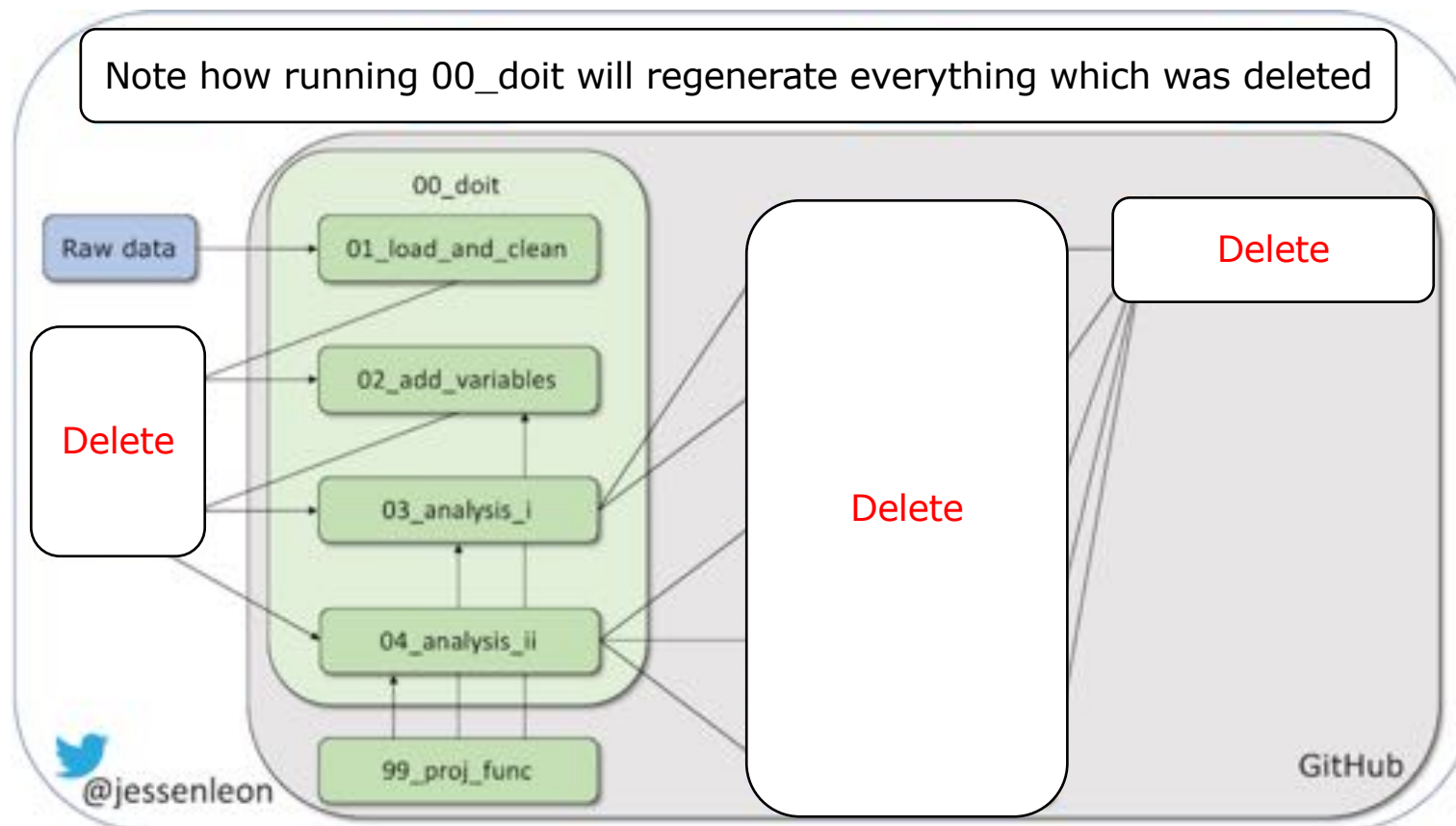
- Communicate: Condense and communicate gained insight via essential, well defined and focused plots



**DTU Bioinformatics, Technical University of Denmark**                                December 3rd 2018

- Now, your paper is published, what then to do with all the files you generated?

# Cleaning Up Your Data Analysis Project



Note how running 00_doit will regenerate everything which was deleted

# Plenum

- *Take 1 minute to discuss with your neighbour*

  - *What are possible factors which may influence the reproducibility of a workflow, which we have not touched upon?*

# Is this a reproducible workflow?

• When looking at a workflow, ask your self:

– Can the entire workflow be run without manual intervention?

– Is the workflow start data static or dynamic?

– What are the dependencies of the workflow?

**DTU Bioinformatics, Technical University of Denmark** December 3rd 2018

# Summary - Learning Objectives Revisited

- Define what a reproducible data analysis workflow is
  - *"A reproducible data analysis workflow is when you can go from the raw data to recreating all the figures, tables and numbers in your paper automatically and consistently"*

- List the elements of a reproducible data analysis workflow
  - "Import → Tidy → Transform → Visualise → Model → Communicate"

- Explain the meaning and purpose of each of the elements in a reproducible data analysis workflow

- When presented with a pre-made workflow, determine if it constitutes a reproducible data analysis workflow
  - Run without manual intervention? Workflow start data static or dynamic? Workflow dependencies?

**DTU Bioinformatics, Technical University of Denmark**                                            December 3rd 2018
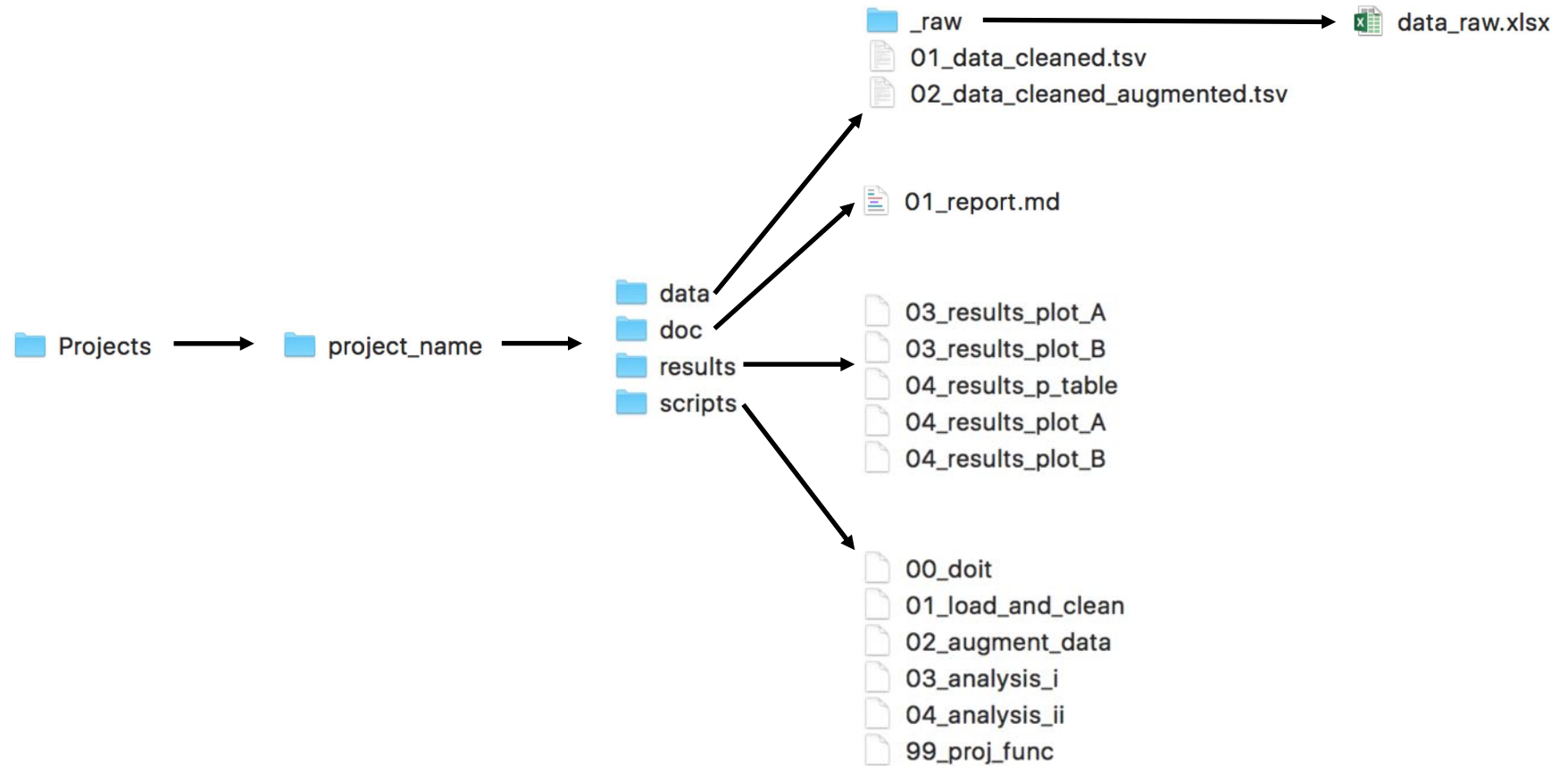
# Reasons for lack of reproducibility

- Many!

- Many of which cannot be controlled

- Focus for this talk is on one aspect we can control

- Namely
  - Reproducible data analysis workflow

- Hopefully, now you have an idea about how you can control your data analysis workflow

Think about readability of your code. Every project you work on is fundamentally collaborative. Even if you are not working with any other person, you are always working with future you and you really do not want to be in a situation where future you has no idea what past you was thinking, because past you will not respond to any emails!
- Hadley Wickham

# Summary and open discussion



**DTU Bioinformatics, Technical University of Denmark**                                    December 3rd 2018

# Redo analysis across projects



**DTU Bioinformatics, Technical University of Denmark** December 3rd 2018