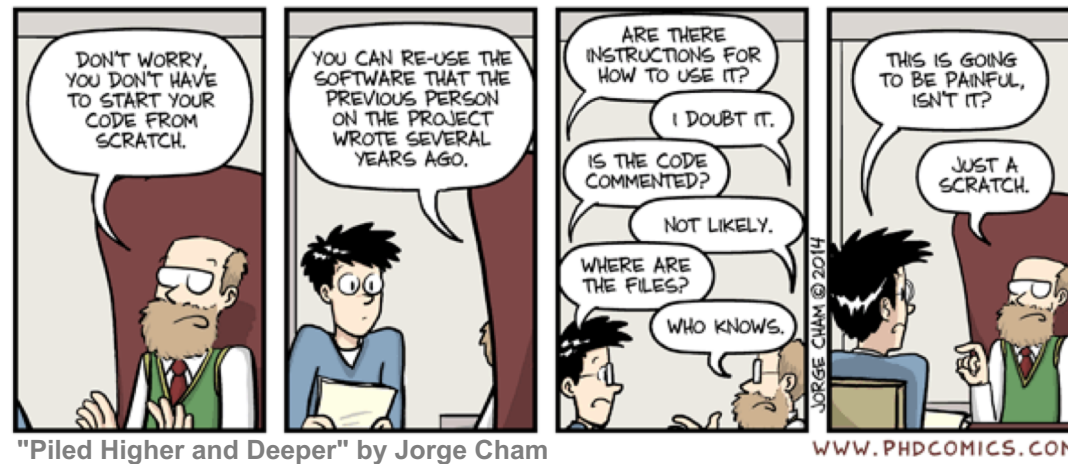


Tools for Reproducible Research

1 ECTS

Göteborg, November 28-29



Teachers



Rasmus



Leif



Viktor

Course content:

- good practices for data analysis and management
- how to use the version control system git to track edits and collaborate on coding
- how to use the package and environment manager Conda
- how to use the workflow manager Snakemake
- how to use R Markdown to generate automated reports
- how to use Jupyter notebooks to document your ongoing analysis
- how to use Docker to distribute containerized computational environments



R Markdown



The screenshot shows a web browser window with the URL `nbis-reproducible-research.readthedocs.io/en/latest/feedback/`. The page title is "Feedback - NBIS Reproducible". The sidebar on the left contains a search bar and a list of links: "Welcome", "About", "The course", "Schedule", "Travel info", "Feedback", "Tutorials", "Introduction to the tutorials", "Conda", "Snakemake", "Git", "Jupyter", "R Markdown", "Docker", and "Take down". A red arrow points to the "Feedback" link. The main content area has a breadcrumb trail "Docs » The course » Feedback" and a heading "Questions, comments, rants". Below the heading is a paragraph: "We will send out a course evaluation form later, which we'd be really happy if you could fill out! That's not the purpose of this message board though. Here you can for example:" followed by a list of bullet points: "Ask if there was something in a lecture that you'd like us to clarify.", "Tell us if you find bugs or inconsistencies in a tutorial.", "Tip us about tools or resources that you think we should know of.", and "Dead links, typos or other issues with the course site." Below the list is another paragraph: "We will keep track of this during the course, so hopefully you will be able to get quick feedback." At the bottom of the form is a "SUBMIT" button and a note: "Never submit passwords through Google Forms."

Introduction to Reproducible Research

Why all the talk about reproducible research?

RESEARCH ARTICLE

Estimating the reproducibility of psychological science

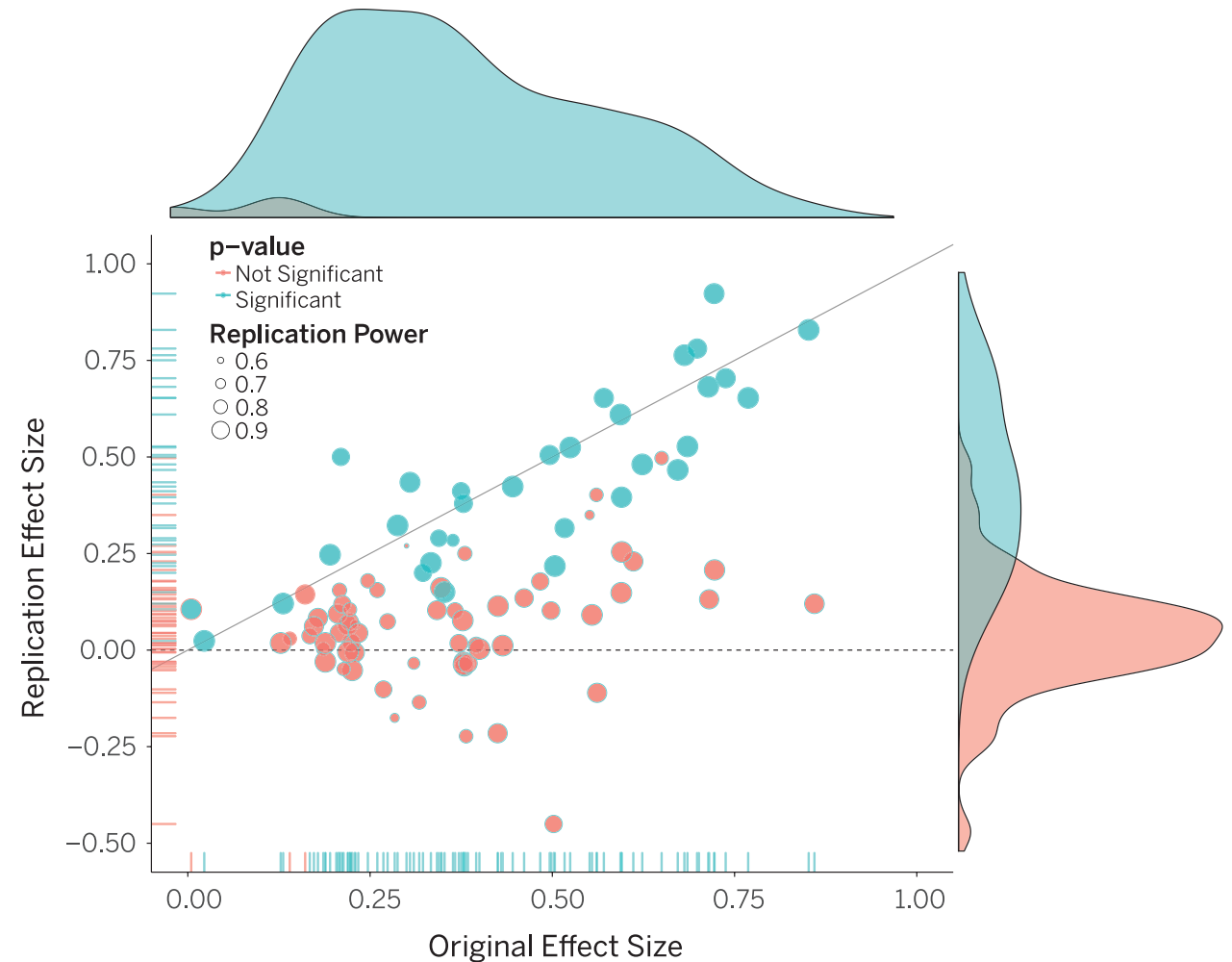
Open Science Collaboration^{*,†}

[†] See all authors and affiliations

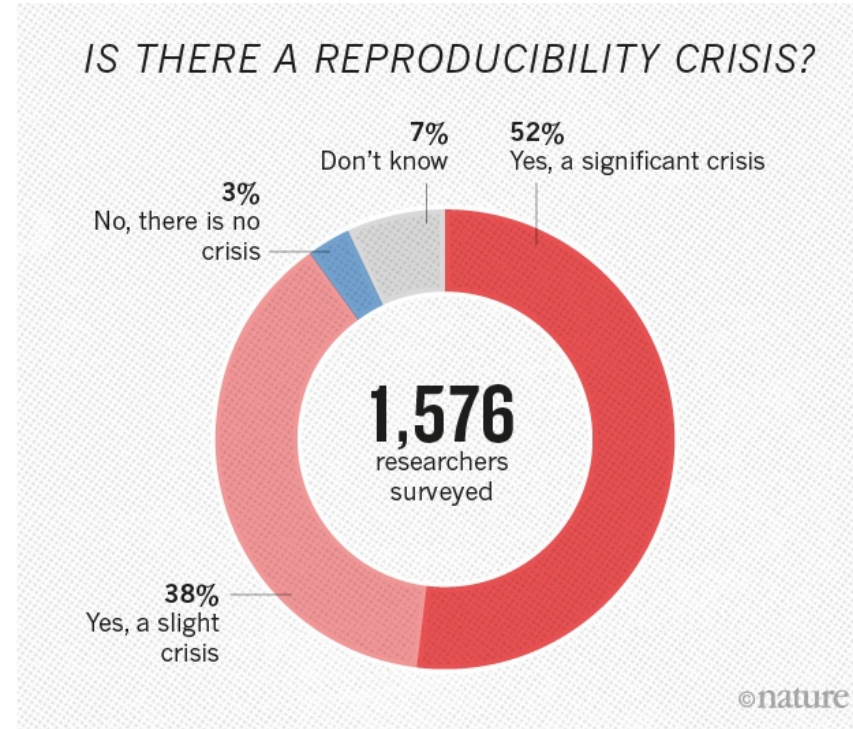
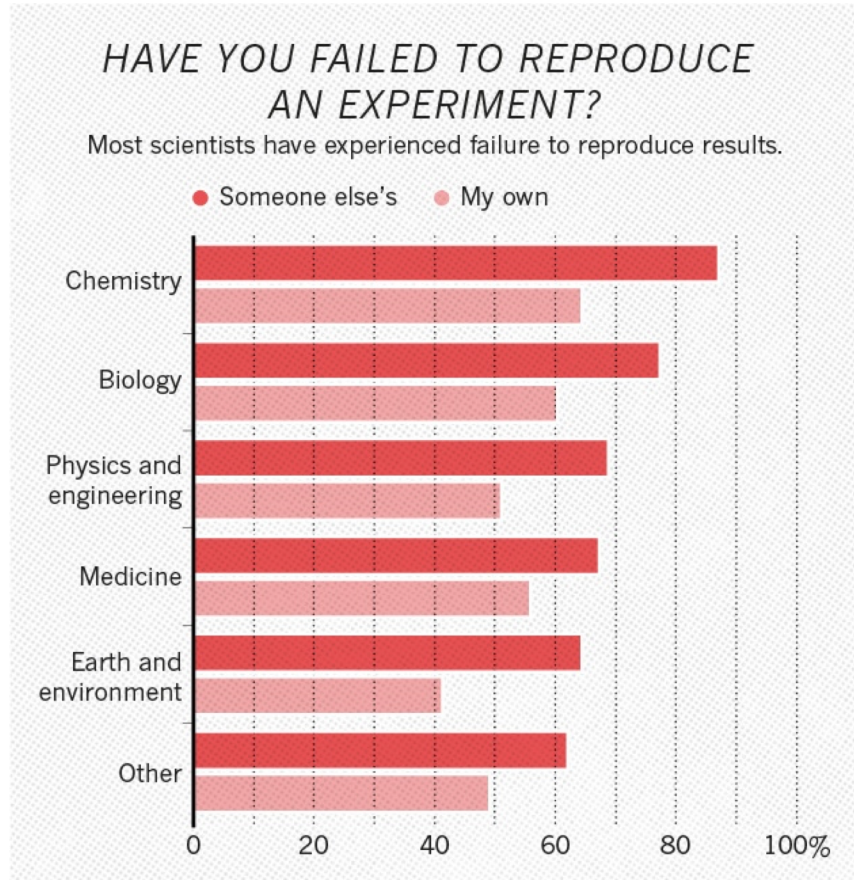
Science 28 Aug 2015:
Vol. 349, Issue 6251, aac4716
DOI: 10.1126/science.aac4716

The *Reproducibility project* set out to replicate 100 experiments published in high-impact psychology journals.

About one-half to two-thirds of the original findings could not be observed in the replication study.



Why all the talk about reproducible research?



A survey in Nature revealed that irreproducible experiments are a problem across all domains of science¹.

Medicine is among the most affected research fields. A study in Nature found that 47 out of 53 medical research papers focused on cancer research were irreproducible².

Common features were failure to show all the data and inappropriate use of statistical tests.

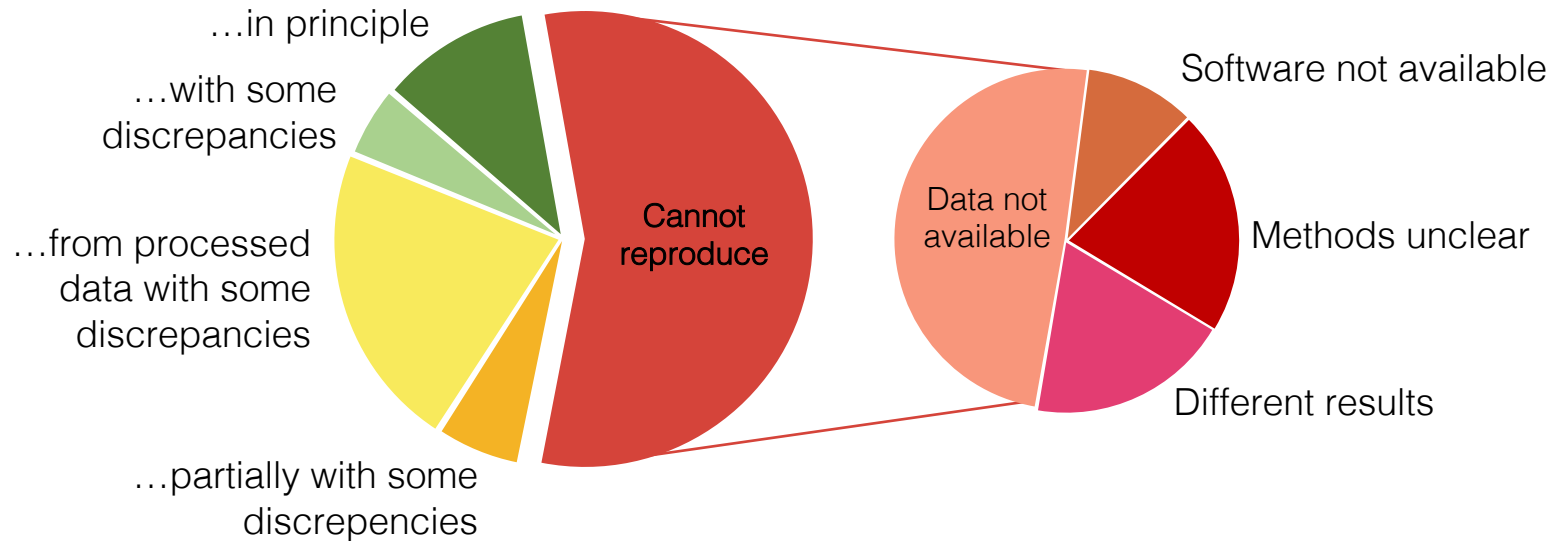
[1] "1,500 scientists lift the lid on reproducibility". Nature. 533: 452–454

[2] Begley, C. G.; Ellis, L. M. (2012). "Drug development: Raise standards for preclinical cancer research". Nature. 483 (7391): 531–533.

Why all the talk about reproducible research?

Replication of data analyses in 18 articles on microarray-based gene expression profiling published in Nature Genetics in 2005–2006:

Can reproduce...



Summary of the efforts to replicate the published analyses.

Adopted from: Ioannidis et al. Repeatability of published microarray gene expression analyses.

Nature Genetics **41** (2009) doi:10.1038/ng.295

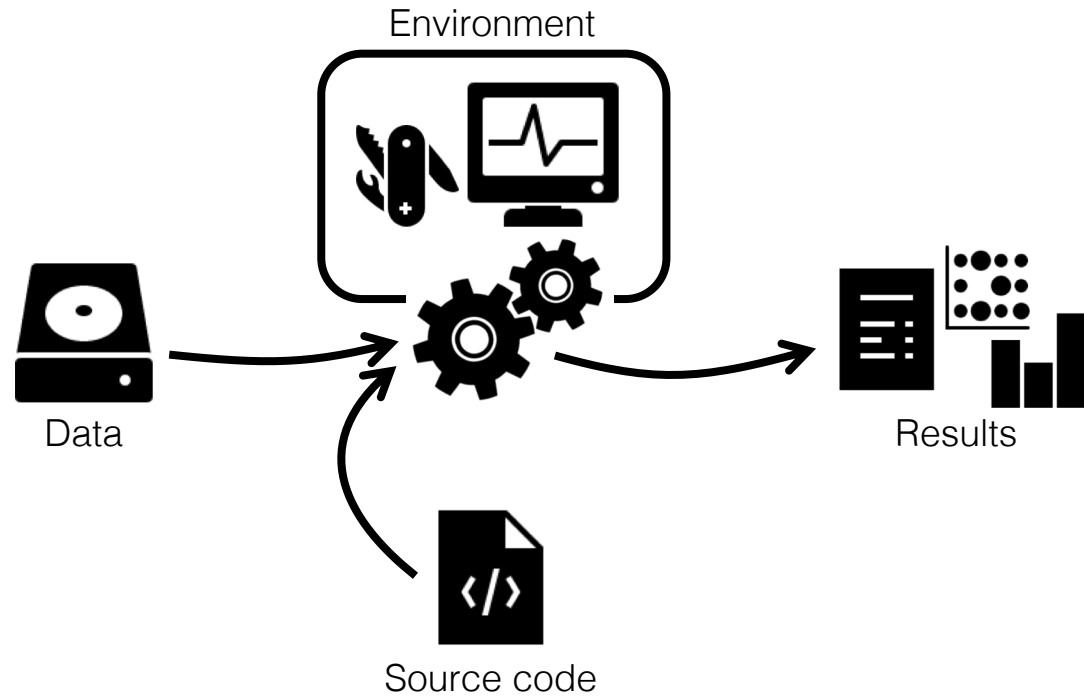
What do we mean with reproducible research?

		Data	
		Same	Different
Code	Same	Reproducible	Replicable
	Different	Robust	Generalisable

"The foundations of knowledge should be constituted by experimentally produced facts, which can be made believable to a scientific community by their reproducibility."

- Robert Boyle, 1627-1691

All parts of a bioinformatics analysis have to be reproducible:



Where does your latest publication fit?



Decent	Getting there...	Well done!
<ul style="list-style-type: none"> • Data available on request. • All meta data required for generating the results available. 	<ul style="list-style-type: none"> • Data deposited in public repositories. • Raw data available in unedited form. • If the raw data needed preprocessing, scripts were used rather than modifying it manually. 	<ul style="list-style-type: none"> • Section in the paper to aid in reproduction. • Used non-proprietary and machine-readable formats, e.g. .csv rather than .xls.
<ul style="list-style-type: none"> • All code for generating results from processed data available on request. 	<ul style="list-style-type: none"> • All code for generating results from raw data is available. • The code is publically available with timestamps/tags. 	<ul style="list-style-type: none"> • All code for generating results from <i>publically available</i> raw data is available. • Code is documented and contains instructions for reproducing results. • Seeds were used and documented for heuristic methods.
<ul style="list-style-type: none"> • Key programs used are mentioned in the methods section. 	<ul style="list-style-type: none"> • List of all programs used, and their respective versions, available. 	<ul style="list-style-type: none"> • Instructions for reproducing the environment publically available.

What's in it for me?

One year in submission loop and reviewer comments are finally back...

Didn't take course

Took course



"It takes some effort to organize your research to be reproducible. We found that although the effort seems to be directed to helping other people stand up on your shoulders, the principal beneficiary is generally the author herself. This is because time turns each one of us into another person, and by making effort to communicate with strangers, we help ourselves to communicate with our future selves."

Schwab et al. Making scientific computations reproducible.
Computing in Science Engineering (2000).

Before project

- Improved structure and organization.
- Forced to think about scope and limitations.

During project

- Easier to rerun analyses and generate results after updating data, tools, parameters, etc.
- Closer interaction between collaborators.
- Much of the manuscript "writes itself".

After project

- Faster resumption of research by others (or your future self), thereby increasing the impact of your work.
- Increased visibility in the scientific community.

Data management

Data (mis)management in practice



Raw data



Data arrives in cumbersome and proprietary format.

Gets converted to format of choice. Original files (and conversion settings) are lost.

Leads a quiet life on the HPC cluster, until the project expires and the data has to be urgently retrieved.

Ends its days on an external hard drive on the researcher's desk.

"Data available upon request".

Meta data



In researcher's lab journal.

Hard-coded in various analysis scripts.

Mailed back and forth between collaborators in ever-changing (but nicely colored) Excel sheets.

Reformatted and included as PDF in the supplementary.

FAIR

Strive to make your data **FAIR** – Findable, Accessible, Interoperable, and Reusable *for both machines and humans*.

Box 2 | The FAIR Guiding Principles

To be Findable:

- F1. (meta)data are assigned a globally unique and persistent identifier
- F2. data are described with rich metadata (defined by R1 below)
- F3. metadata clearly and explicitly include the identifier of the data it describes
- F4. (meta)data are registered or indexed in a searchable resource

To be Accessible:

- A1. (meta)data are retrievable by their identifier using a standardized communications protocol
 - A1.1 the protocol is open, free, and universally implementable
 - A1.2 the protocol allows for an authentication and authorization procedure, where necessary
- A2. metadata are accessible, even when the data are no longer available

To be Interoperable:

- I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- I2. (meta)data use vocabularies that follow FAIR principles
- I3. (meta)data include qualified references to other (meta)data

To be Reusable:

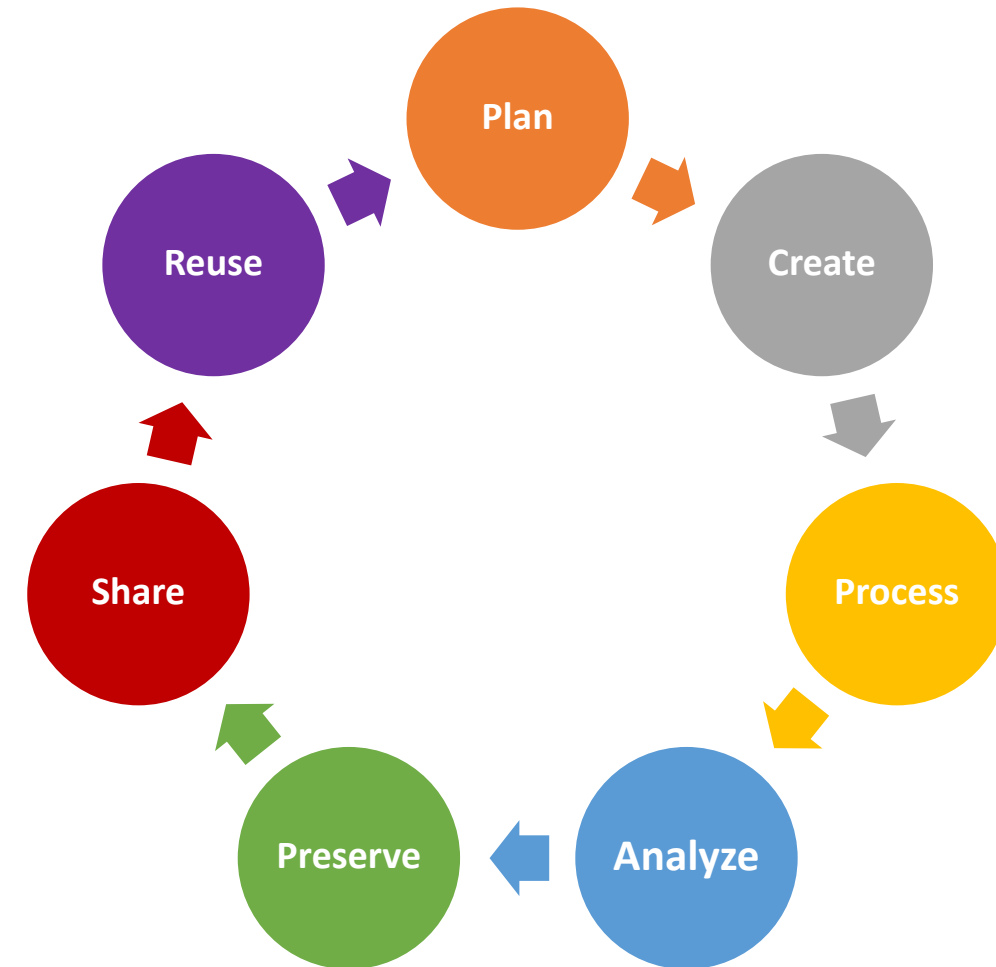
- R1. meta(data) are richly described with a plurality of accurate and relevant attributes
 - R1.1. (meta)data are released with a clear and accessible data usage license
 - R1.2. (meta)data are associated with detailed provenance
 - R1.3. (meta)data meet domain-relevant community standards



Wilkinson, Mark et al. "The FAIR Guiding Principles for scientific data management and stewardship". Scientific Data 3, 160018 (2016) doi:10.1038/sdata.2016.18

Data management plan

- Check the requirements of your funding agency and field of research.
- List the types of data that you expect to produce.
- Decide what data require archiving, and determine how much storage space you will need (short and long term).
- Provide metadata that allows others to understand, cite and reuse your data files.
- Make clear how and when your data can be shared with scientists outside your group.
- If your research involves sensitive data, explain any legal and ethical restrictions on data access and reuse.
- Look for suitable data repositories used by your research community.
- Check what data format and structure the chosen repository might request.



Life cycle for scientific data



Pair up and discuss!

- Does your group have a data management plan in place?
- Do you know "your" repositories and how to submit data to them?

Data acquisition and deposit

- Find the right repository for your data, and strive towards uploading data to its final destination already at the beginning of a project.
- Structure metadata in the format needed by the repository already as the experiments are being performed.
- Stick to non-proprietary and widely used file formats.


Scientific Data (Springer Nature) maintains a list of recommended repositories at www.nature.com/sdata/policies/repositories.

Dedicated repositories:

e.g. SRA, GEO, GenBank, UniProt etc.

Generalist ("long-tail data") repositories:

Research data that doesn't fit in structured data repositories, e.g. Data Dryad, Figshare, Zenodo.

Each dataset can be assigned a Digital Object Identifier (); a persistent identifier used to uniquely identify objects.

- Only 12% of articles from NIH funded research mention data deposited in international repositories
- Estimated 200000+ "invisible" data sets / year

Read et al. (2015) PLoS ONE 10(7) doi:10.1371/journal.pone.0132735



Data acquisition and deposit

- Find the right repository for your data, and strive towards uploading data to its final destination already at the beginning of a project.
- Structure metadata in the format needed by the repository already as the experiments are being performed.
- Stick to non-proprietary and widely used file formats.

	A	B	C	D	E	F	G	H	I	J	K
1	#BLUE headers are required!										
2	#YELLOW columns have a controlled vocabulary										
3	bioproject_accession	sample_name	library_ID	title	library_strategy	library_source	library_selection	library_layout	platform	instrument_model	filetype
4	PRJNA212142	RN4220_empty	RN4220_empty	RN4220_empty; Sta	RNA-Seq	TRANSCRIPTOMIC	cDNA	single	ILLUMINA	Illumina HiSeq 2000	fastq
5											
6											
7											

```
1 ^SAMPLE=RN4220_empty
2 !Sample_title = RN4220_empty
3 !Sample_source_name = S. aureus isolate
4 !Sample_organism = Staphylococcus aureus
5 !Sample_characteristics = strain: RN4220
6 !Sample_characteristics = has_plasmid: False
7 !Sample_characteristics = uMax: 0.2
8 !Sample_molecule = total RNA
9 !Sample_extract_protocol = RNA purified by modified
  Qiagen RNeasy kit Illumina
10 [...]
```




GEO (Gene Expression Omnibus) uses text files in SOFT format.

SRA (Sequence Read Archive) uses a template Excel sheet for metadata.

Data acquisition and deposit

- Find the right repository for your data, and strive towards uploading data to its final destination already at the beginning of a project.
- Structure metadata in the format needed by the repository already as the experiments are being performed.
- Stick to non-proprietary and widely used file formats.

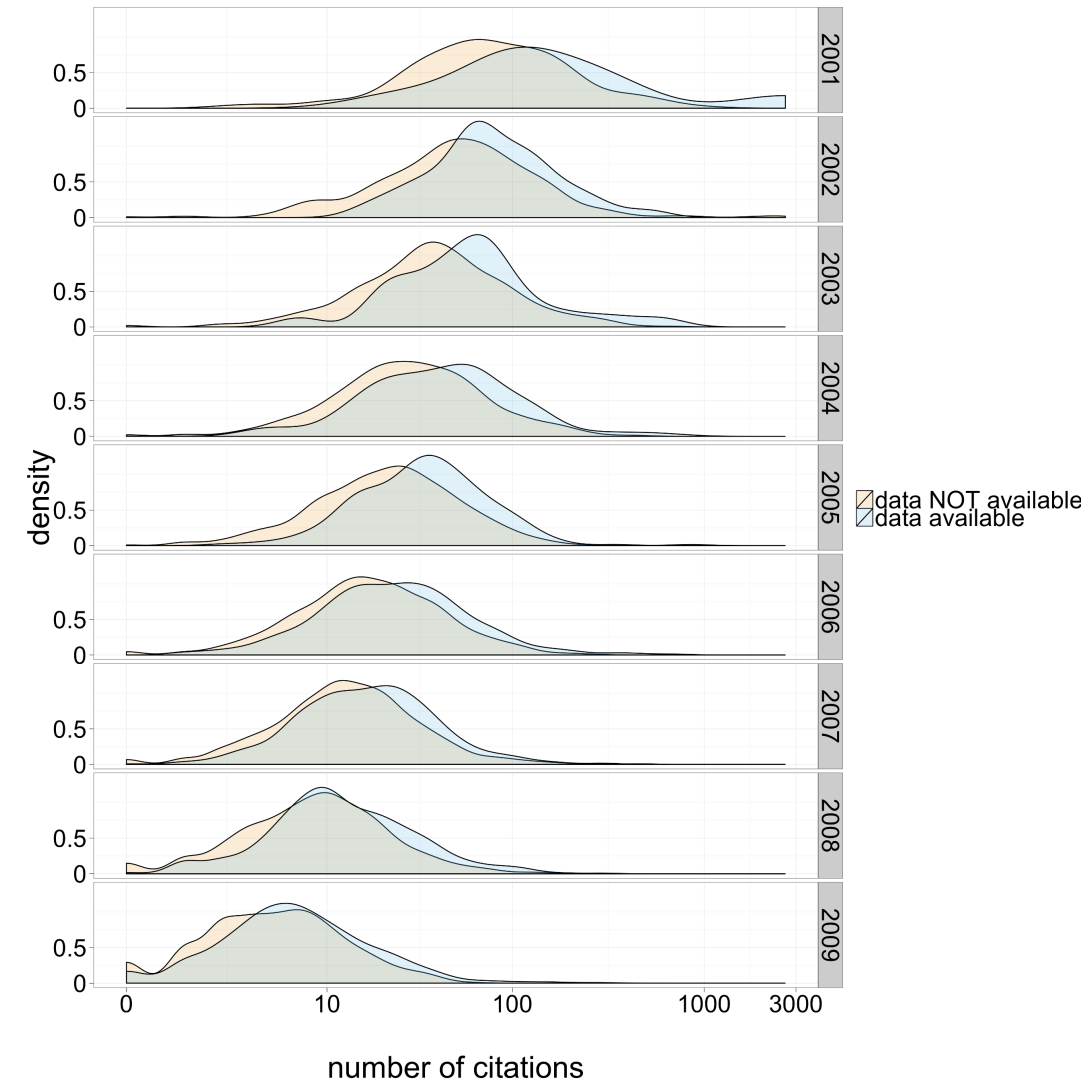
-	+
Binary	Text-based
Proprietary	Open
New kid on the block	Old as the hills
Compressed/encrypted	Uncompressed/unencrypted
Platform dependent	Interoperable
Complex	Simple

			
Raster graphic	wmf, psd	bmp, gif	tiff, png, jpeg
Vector graphic	ai, eps	pdf	svg
Document	doc	docx, tex	odt, utf-8, md
Archive	rar	7z	zip, tar, gz
Tabular data	xls, rds, mat	xlsx, ods	csv

Data sharing

From 10,555 studies with gene expression microarray data:

- Studies that shared data received 9% more citations (after accounting for other covariates).
- Data reuse by other researchers continued for >6 years.
- A very conservative estimate found that 20% of the datasets deposited between 2003 and 2007 had been reused at least once by third parties.



Piwovar and Vision (2013), Data reuse and the open data citation advantage, PeerJ 1:e175, doi:10.7717/peerj.175

Data sharing – Open access

- Democracy and transparency
 - Publicly funded research data should be accessible to all free of charge.
 - Published results and conclusions should be possible to check by others.
- Research
 - Enables others to combine data, address new questions, and develop new analytical methods.
 - Reduce duplication and waste.
- Innovation and utilization outside research
 - Public authorities, companies, and individuals outside academia can make use of the data.
- Citation
 - Citation of data will be a merit for the researcher that produced it.



Data sharing – Ontologies

lauroyl-CoA

dodecanoyl-CoA

C12:0-CoA

lauroyl coenzyme A

coenzyme A, S-dodecanoate

dodecanoyl coenzyme A

C12:0 coenzyme A

dodecanoic acid coenzyme A

lauroylic acid CoA

Dodecanethioic acid, S-ester with coenzyme A

Coenzyme A, S-laurate (7CI,8CI)

12:0, lauroyl-CoA

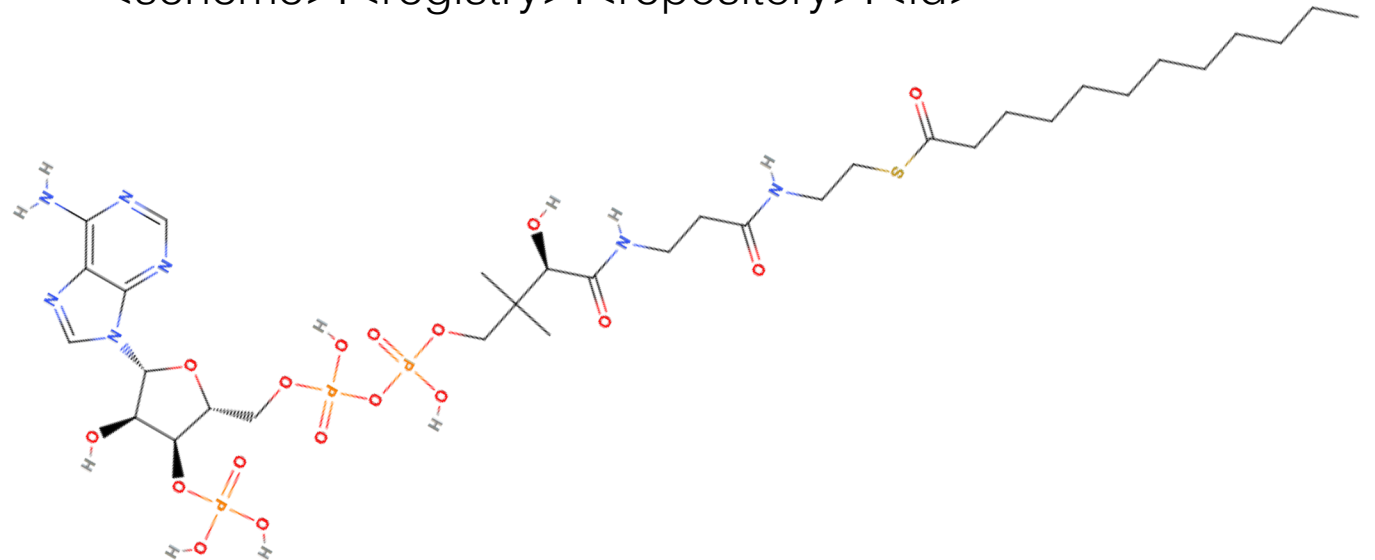
1-undecanecarboxylic acid CoA

vulvic acid CoA

Who am I?

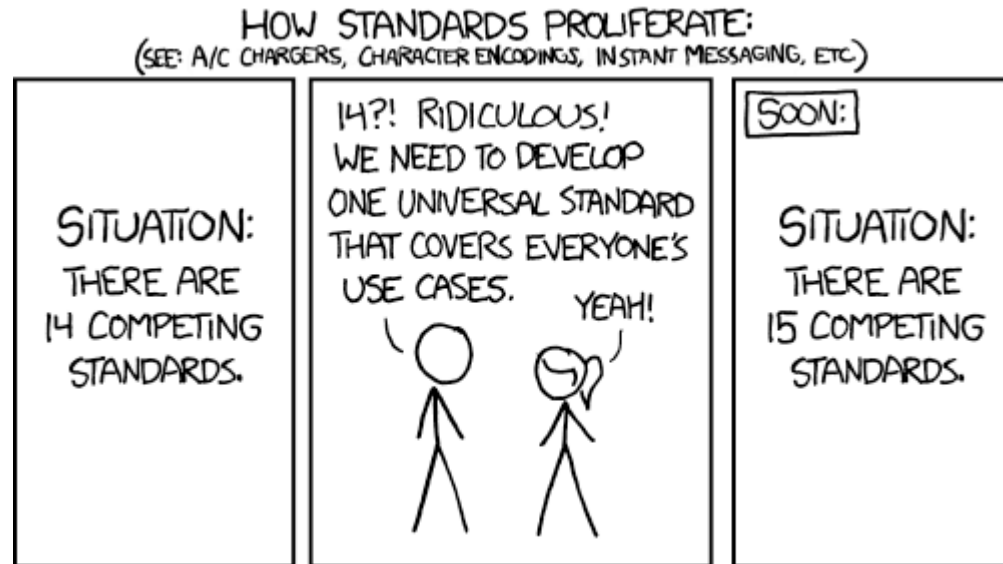
urn.miriam.chebi:15521 of course!

<scheme>.<registry>.<repository>:<id>



3'-phosphoadenosine 5'-(3-((3R)-4-((3-[[2-(dodecanoylsulfanyl)ethyl]amino]-3-oxopropyl)amino]-3-hydroxy-2,2-dimethyl-4-oxobutyl) dihydrogen diphosphate)

Data sharing – Ontologies



FAIRsharing.org
standards, databases, policies

Standards Databases Policies Collections Add/Claim Content Stats Log in or Register

A curated, informative and educational resource on data and metadata *standards*, across all disciplines, inter-related to *databases* and data *policies*.

Find
Recommendations
Standards and/or databases recommended by journal or funder data policies.

Discover
Collections
Standards and/or databases grouped by domain, species or organization.

Learn
Educational
About standards, their use in databases and policies, and how we can help you.

Search FAIRsharing

Advanced Search
Fine grained control over your search.

Search Wizard
Let us guide you to your results.

699 Standards

Terminology Artifact	343
Model/Format	239
Reporting Guideline	117

View all

974 Databases

Life Science	733
Biomedical Science	181
General Purpose	10

View all

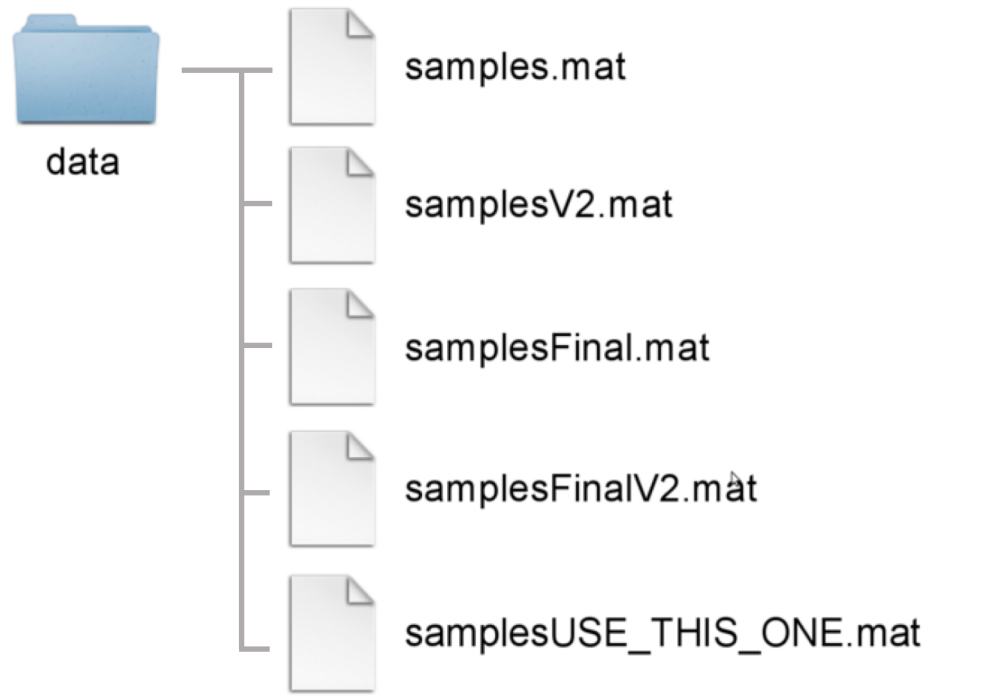
97 Policies

Funder	22
Journal	68
Society	3

View all

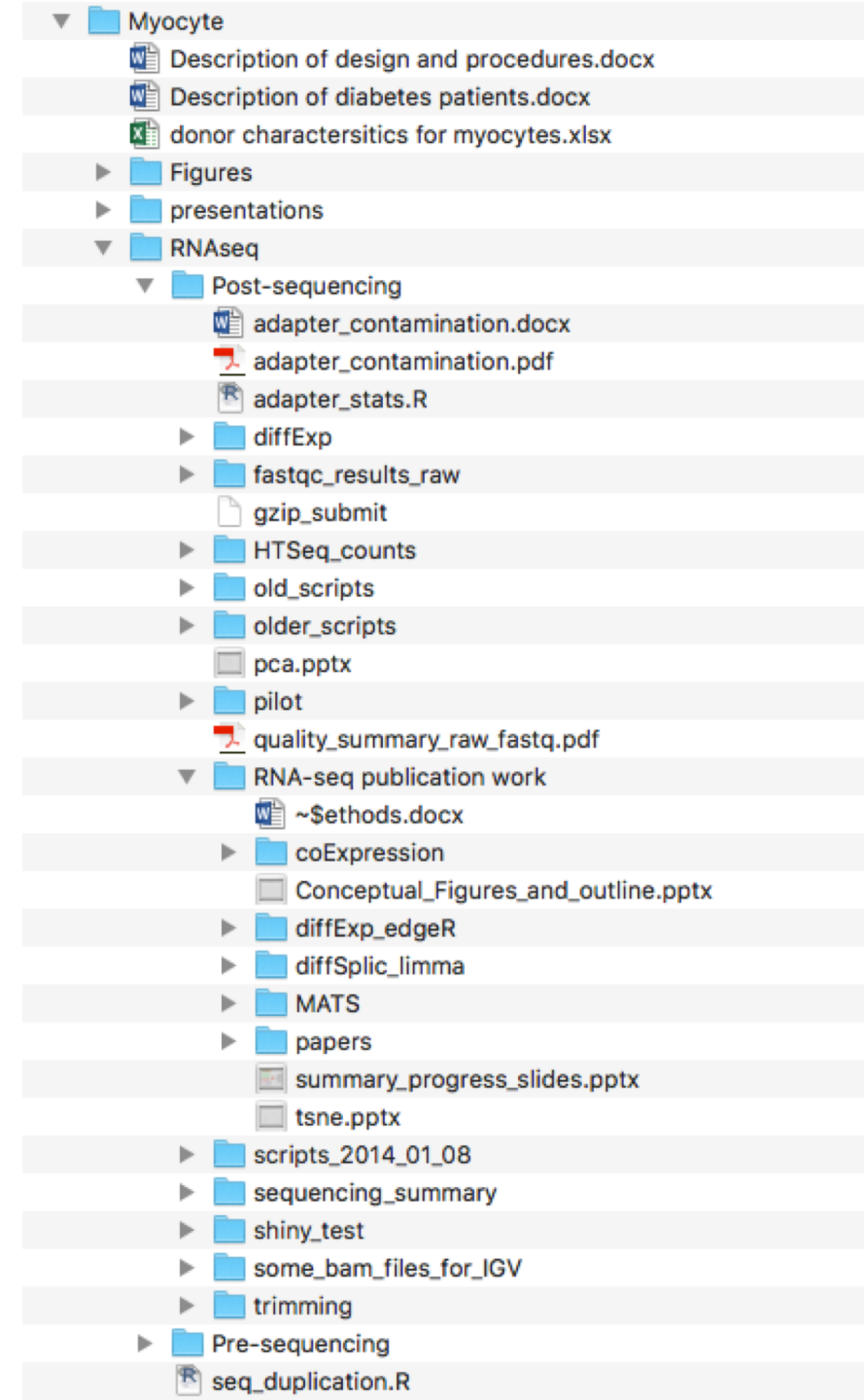
Project organization

The project directory



The first step towards working reproducible: Get organized!

Divide your work into distinct projects and keep all files needed to go from raw data to final results in a dedicated directory with relevant subdirectories.





Pair up and discuss!

- Do you organize your work in distinct projects?
- How do you organize your files in this context?
- Are you happy with the way you work today?

The project directory

project	
- doc/	documentation for the study
- data/	raw and primary data, essentially all input files, never edit!
- raw_external/	
- raw_internal/	
- meta/	
- code/	all code needed to go from input files to final results
- notebooks/	notebooks that document your day-to-day work
- intermediate/	output files from different analysis steps, can be deleted
- scratch/	temporary files that can be safely deleted or lost
- logs/	logs from the different analysis steps
- results/	output from workflows and analyses
- figures/	
- tables/	
- reports/	
- Snakefile	project workflow, carries out analysis contained in code/
- config.yml	configuration of the project workflow
- environment.yml	software dependencies list, used to create a project environment
- Dockerfile	recipe to create a project container

Working in projects



File naming system

Machine readable

- Avoid special characters, e.g.: ~!@#\$%^&*() `; <>?, [] {} ' " |
- Avoid spaces, alternatives:
 - file_name.txt
 - file-name.txt
 - filename.txt
 - FileName.txt

Human readable




- Know the content of a file without opening it, e.g.:
SRR1234.hg19.sorted.trimmed.bam

Control file ordering




- Use dates if appropriate
- Use 01, 02, rather than 1, 2



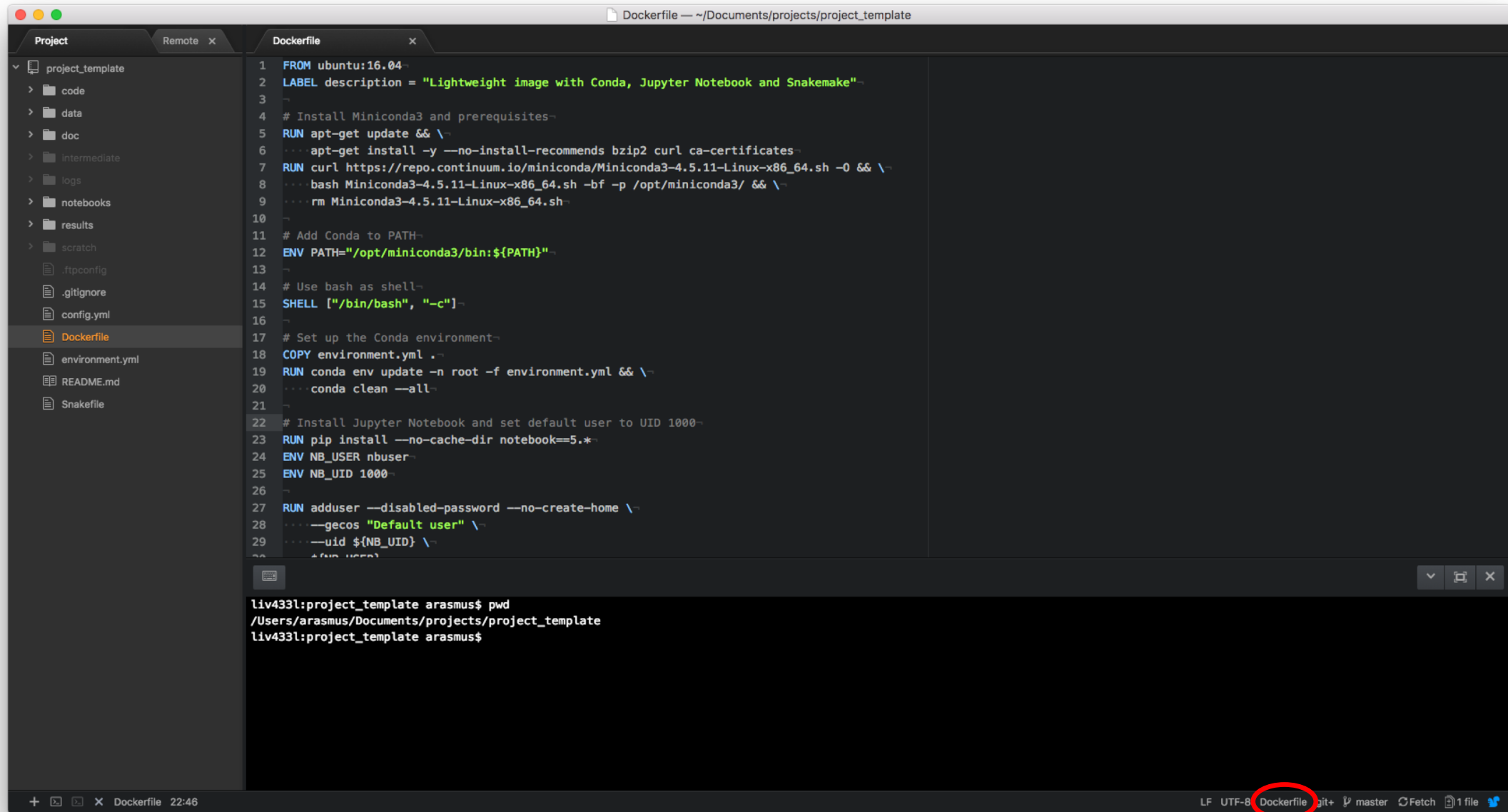
Bad examples:

 reproducible%20research.pptx
 suppl fig 10.png
 Supplementary Figure 9.png

Good examples:

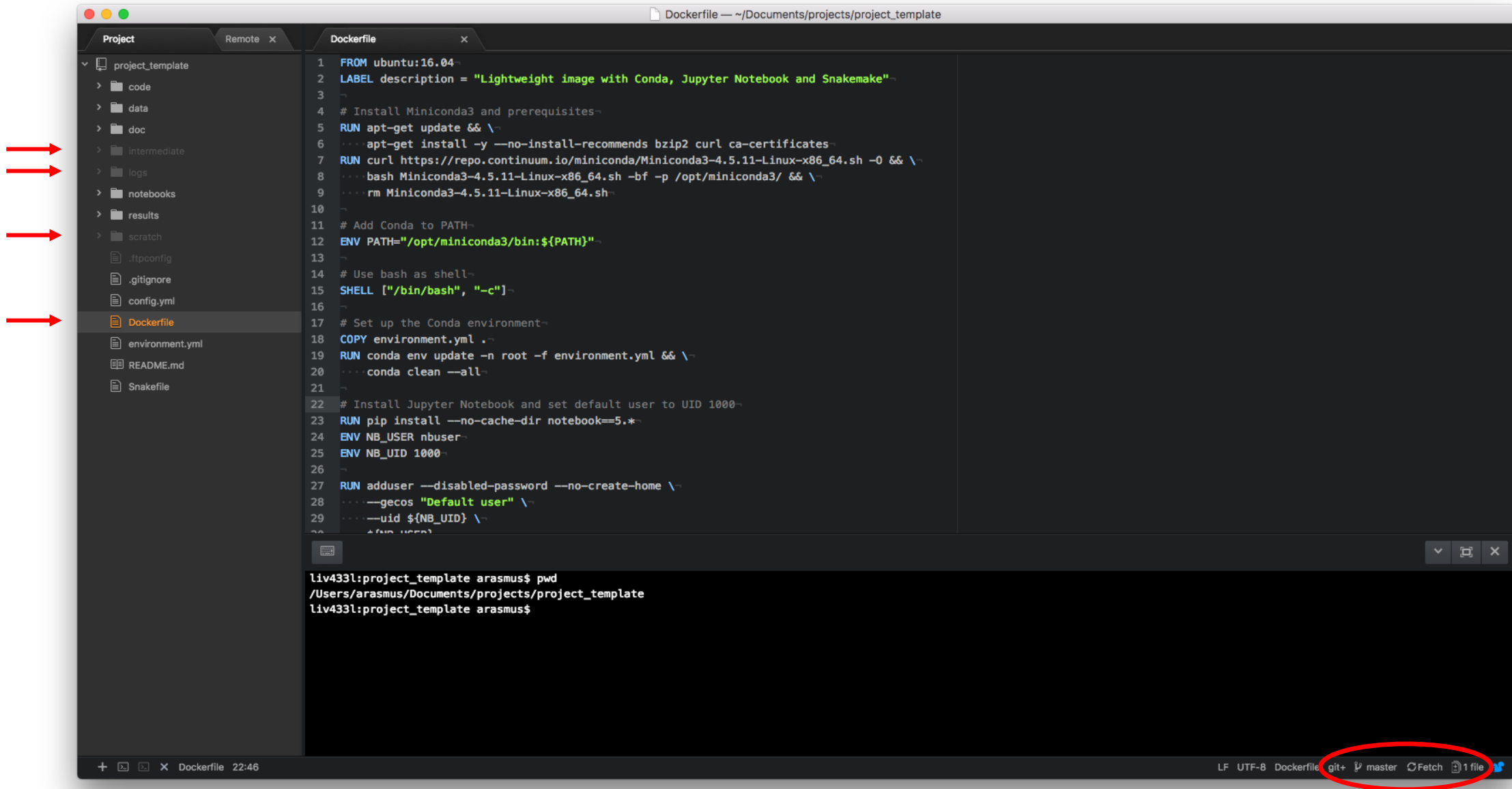
 2018-11-28_Gothenburg_Reproducible_research.pptx
 suppl_fig_09_barplot_alignment_stats.png
 suppl_fig_10_samples_PCA_count_data.png

A project in Atom



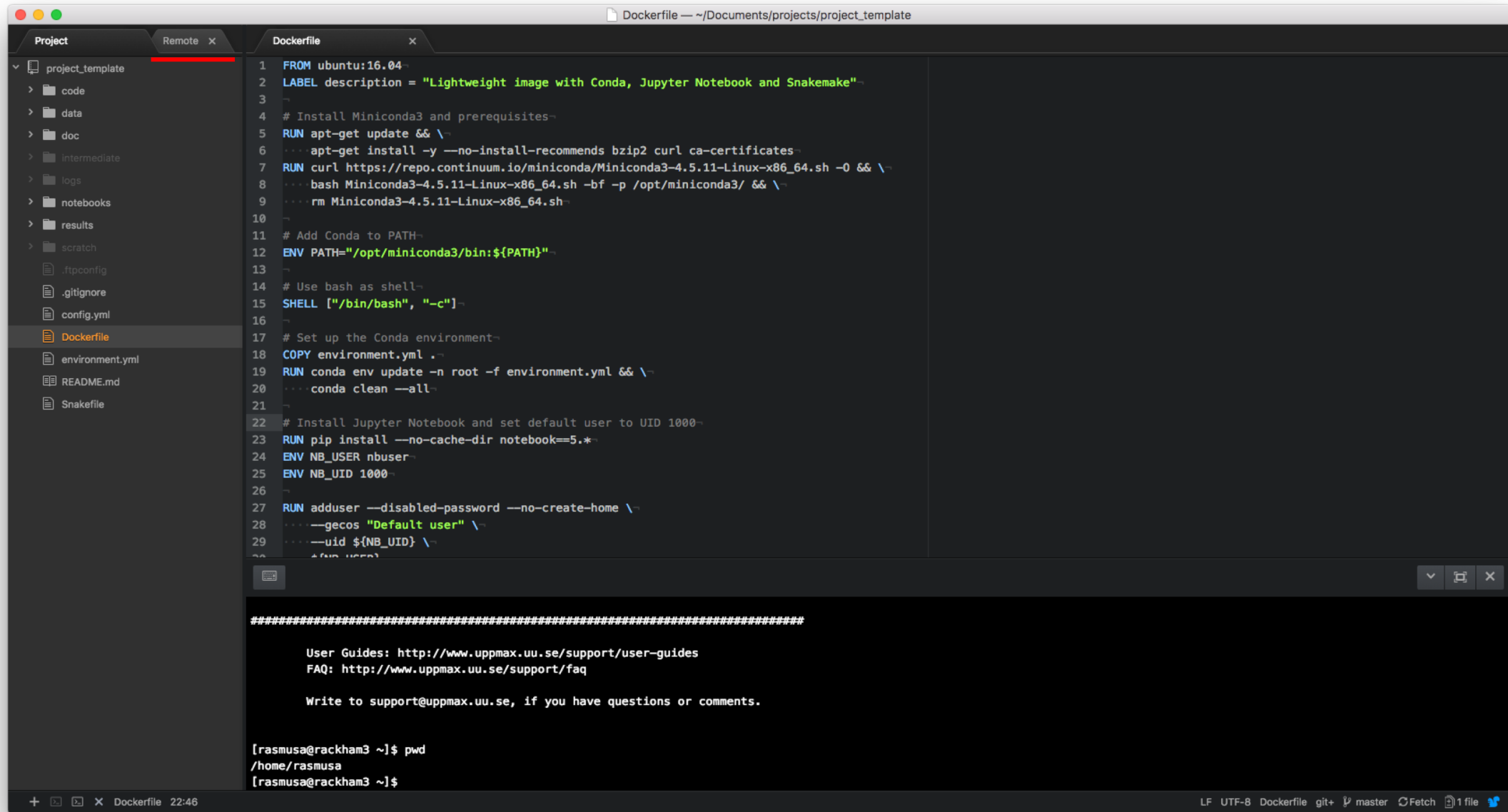
Syntax highlighting, indentation, and autocomplete

A project in Atom



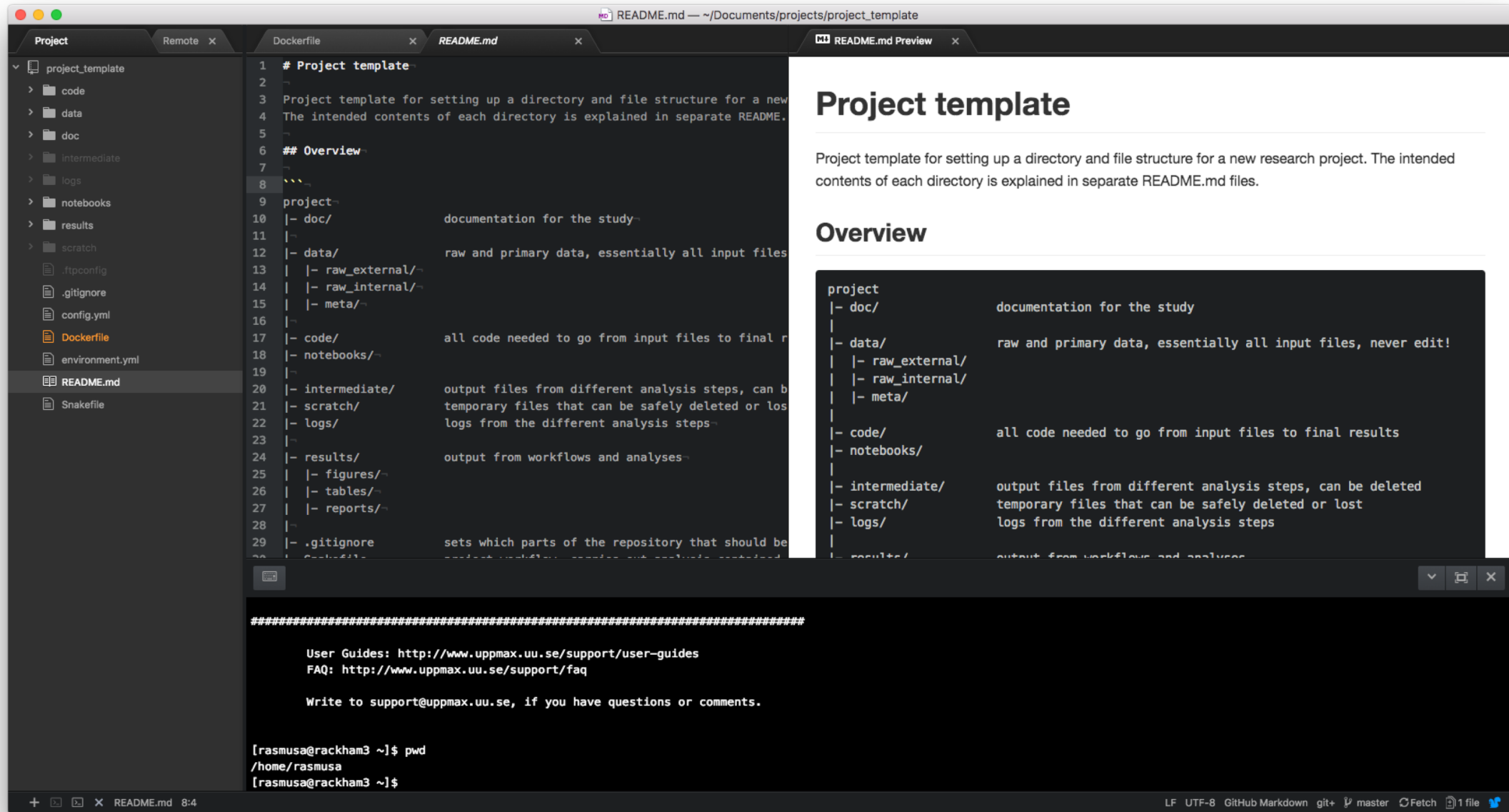
Integrated version control with Git

A project in Atom



Automatically sync files between local/remote

A project in Atom



Tons of plugins, e.g. for viewing different file formats

A project in RStudio

The screenshot displays the RStudio interface for a project named "project_template". The top toolbar shows the "Files" tab selected, with a red circle highlighting the "New Folder" button. The "Files" pane on the left shows the project structure, including directories like "code", "data", "doc", "logs", "notebooks", "results", "scratch", and "Snakefile", along with files like ".gitignore", "config.yml", "environment.yml", "README.md", and "README.html". The "Console" pane at the bottom left shows the output of a terminal session, including the command "ll" and the output of "git status". The "Source" pane on the right shows the content of "README.md", which describes the project template and its intended contents. The "Environment" pane at the bottom right shows the project's environment.

Files

Git

New Folder

Delete

Rename

More

Home > Work > Teaching > Reproducible_research > project_template

Name	Size	Modified
..		
.gitignore	256 B	Nov 23, 2018, 9:14 AM
code		
config.yml	20 B	Mar 19, 2018, 10:33 AM
data		
doc		
Dockerfile	818 B	Mar 20, 2018, 1:23 PM
environment.yml	61 B	Mar 19, 2018, 10:33 AM
intermediate		
logs		
notebooks		
README.md	1.2 KB	Mar 20, 2018, 1:23 PM
results		
scratch		
Snakefile	262 B	Mar 19, 2018, 10:33 AM
README.html	213.5 KB	Nov 23, 2018, 9:16 AM

Console

Terminal

R Markdown

Terminal 1

```
[base][master] >> ll
total 48
drwxr-xr-x  4 varemo  NET\Domain Users  128B Nov 23 09:11 .Rproj.user/
drwxr-xr-x 15 varemo  NET\Domain Users  480B Nov 23 09:14 .git/
-rw-r--r--  1 varemo  NET\Domain Users  256B Nov 23 09:14 .gitignore
-rw-r--r--  1 varemo  NET\Domain Users  818B Mar 20 2018 Dockerfile
-rw-r--r--  1 varemo  NET\Domain Users  1.2K Mar 20 2018 README.md
-rw-r--r--  1 varemo  NET\Domain Users  262B Mar 19 2018 Snakefile
drwxr-xr-x  3 varemo  NET\Domain Users   96B Mar 19 2018 code/
-rw-r--r--  1 varemo  NET\Domain Users  20B Mar 19 2018 config.yml
drwxr-xr-x  6 varemo  NET\Domain Users  192B Mar 20 2018 data/
drwxr-xr-x  4 varemo  NET\Domain Users  128B Mar 20 2018 doc/
-rw-r--r--  1 varemo  NET\Domain Users  61B Mar 19 2018 environment.yml
drwxr-xr-x  3 varemo  NET\Domain Users   96B Mar 20 2018 intermediate/
drwxr-xr-x  3 varemo  NET\Domain Users   96B Mar 19 2018 logs/
drwxr-xr-x  3 varemo  NET\Domain Users   96B Mar 19 2018 notebooks/
drwxr-xr-x  6 varemo  NET\Domain Users  192B Mar 20 2018 results/
drwxr-xr-x  3 varemo  NET\Domain Users   96B Mar 19 2018 scratch/

Fri Nov 23, 09:14:52 | MacBook ~/Work/Teaching/Reproducible_research/project_t
[base][master] >> git status
On branch master
Your branch is up to date with 'origin/master'.

nothing to commit, working tree clean
```

project_template

README.md

```
1 # Project template
2
3 Project template for setting up a directory and file structure for a new research project.
4 The intended contents of each directory is explained in separate README.md files.
5
6 ##Overview
7
8 ```
9 project
10 |- doc/           documentation for the study
11 |
12 |- data/         raw and primary data, essentially all input files, never edit!
13 |   |- raw_external/
14 |   |- raw_internal/
15 |   |- meta/
16 |
17 |- code/         all code needed to go from input files to final results
18 |- notebooks/
19 |
20 |- intermediate/ output files from different analysis steps, can be deleted
21 |- scratch/      temporary files that can be safely deleted or lost
22 |- logs/         logs from the different analysis steps
23 |
24 |- results/      output from workflows and analyses
25 |   |- figures/
26 |   |- tables/
27 |   |- reports/
28 |
29 |- gitignore     sets which parts of the repository that should be git tracked
```

Environment

History

Plots

Connections

Packages

Help

Viewer

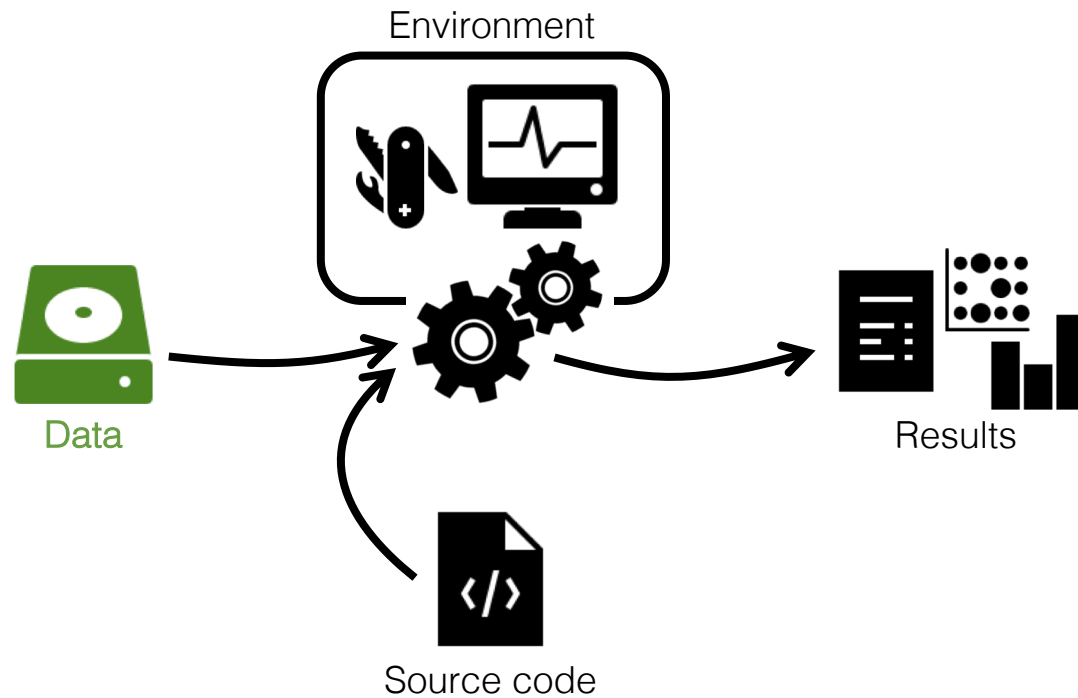
Publish

Project template

Project template for setting up a directory and file structure for a new research project. The intended contents of each directory is explained in separate README.md files.

Overview

```
project
|- doc/           documentation for the study
|
|- data/         raw and primary data, essentially all input files, never edit!
|   |- raw_external/
|   |- raw_internal/
|   |- meta/
|
|- code/         all code needed to go from input files to final results
|- notebooks/
|
|- intermediate/ output files from different analysis steps, can be deleted
|- scratch/      temporary files that can be safely deleted or lost
|- logs/         logs from the different analysis steps
```



```
project
|- doc/
|
|- data/
|   |- raw_external/
|   |- raw_internal/
|   |- meta/
|
|- code/
|- notebooks/
|
|- intermediate/
|- scratch/
|- logs/
|
|- results/
|   |- figures/
|   |- tables/
|   |- reports/
|
|- Snakefile
|- config.yml
|- environment.yml
|- Dockerfile
```

Tools for Reproducible Research



Managing
dependencies



Managing and executing
analysis workflow



Versioning and
collaborating on code
(and some other files)



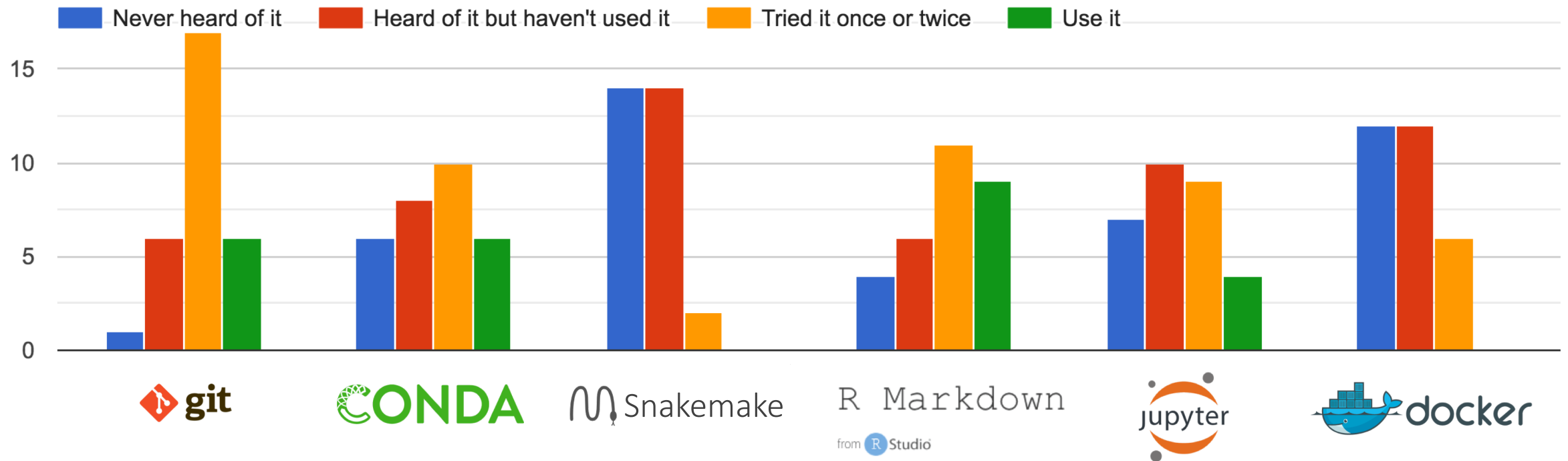
Connecting code
and reporting

and...



Isolating and exporting
environment

Student experience





Managing
dependencies



Managing and executing
analysis workflow



Versioning and
collaborating on code
(and some other files)

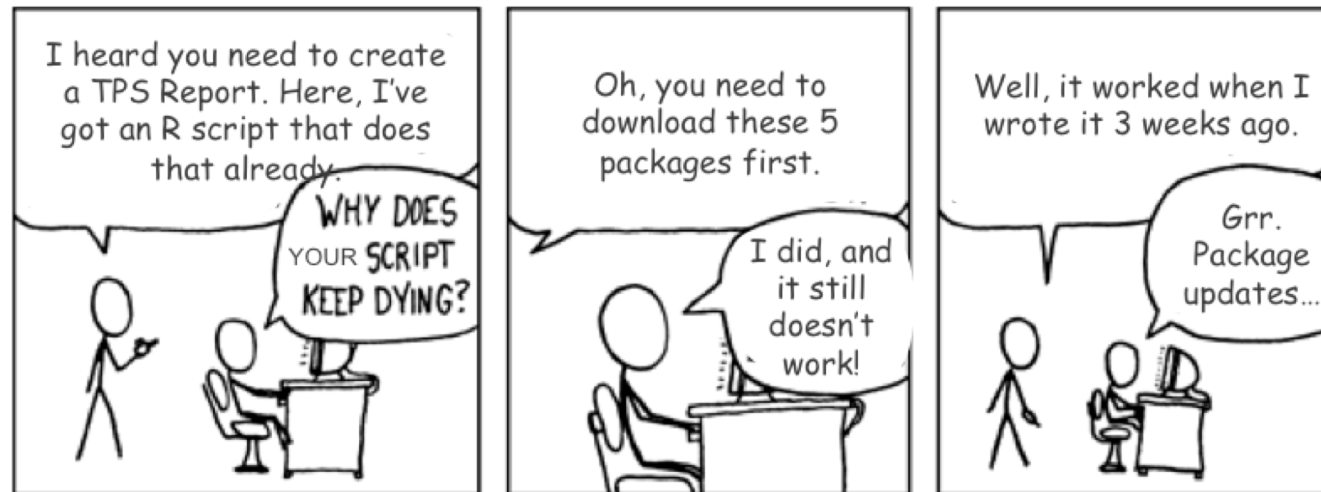


Connecting code
and reporting

and...



Isolating and exporting
environment

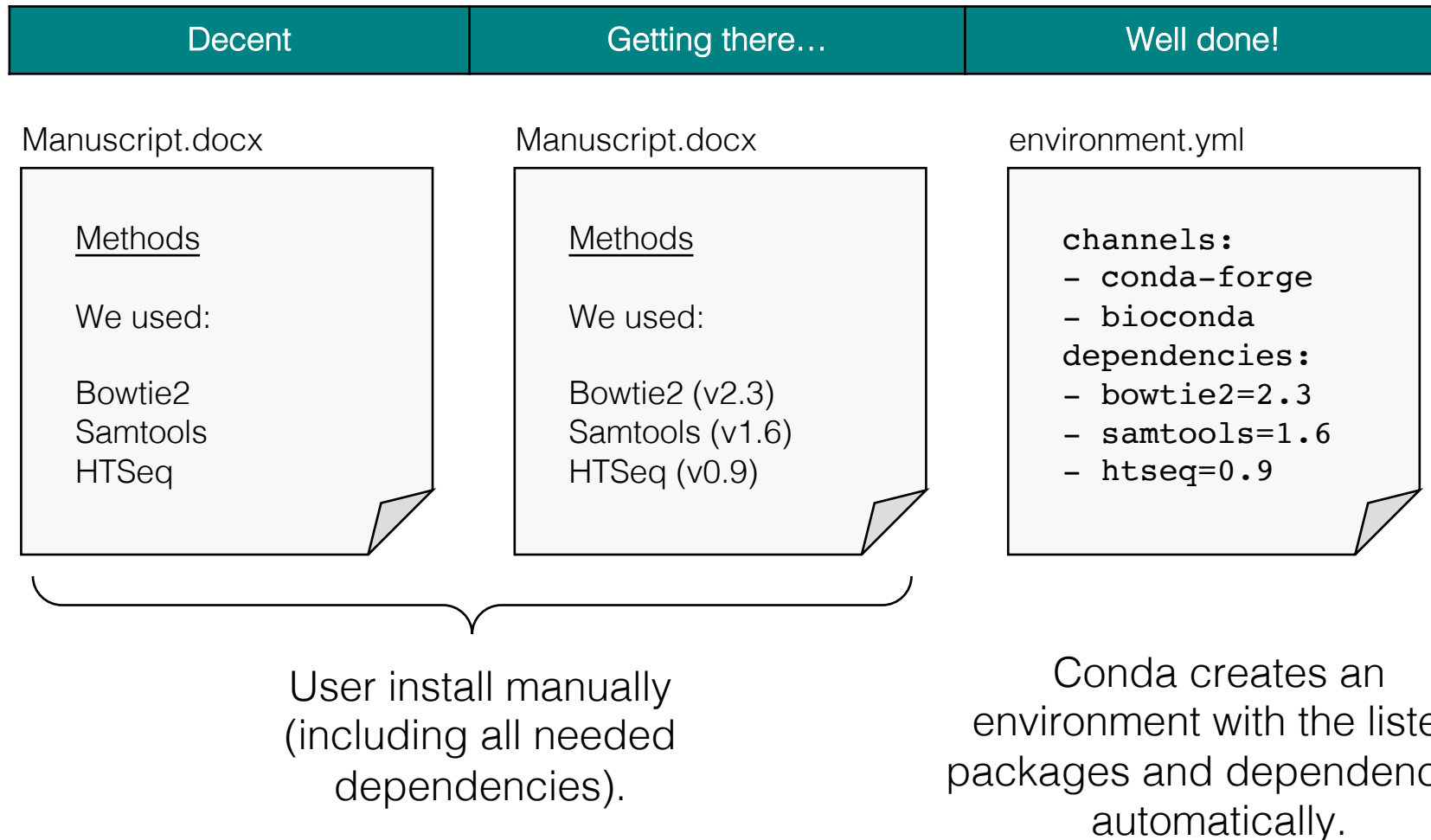


Full reproducibility requires the possibility to recreate the system that was originally used to generate the results.

What is Conda?



- Conda is a package, dependency, and environment manager.
- Works for software developed in any programming language.



Package manager



- Conda package: compressed tarball (system-level libraries, Python or other modules, executable programs, or other components).
- Conda keeps track of the dependencies between packages and platforms.
- Conda packages are downloaded from remote channels.

```
$ conda install -c conda-forge matplotlib
```

```
Fetching package metadata .....
```

```
Solving package specifications: .....
```

```
Package plan for installation in environment /Users/varemo/Applications/miniconda2/envs/test-r2:
```

```
The following packages will be downloaded:
```

package	build		
sqlite-3.13.0	1	1.4 MB	conda-forge
libpng-1.6.24	0	338 KB	conda-forge
python-2.7.12	1	11.8 MB	conda-forge
certifi-2016.8.31	py27_0	218 KB	conda-forge
freetype-2.6.3	1	782 KB	conda-forge
functools32-3.2.3.2	py27_1	16 KB	conda-forge
numpy-1.11.1	py27_0	3.1 MB	defaults
pyparsing-2.1.8	py27_0	89 KB	conda-forge
pytz-2016.6.1	py27_0	183 KB	conda-forge
six-1.10.0	py27_0	18 KB	conda-forge
cycler-0.10.0	py27_0	13 KB	conda-forge
python-dateutil-2.5.3	py27_0	236 KB	conda-forge
setuptools-26.1.1	py27_0	346 KB	conda-forge
matplotlib-1.5.3	np111py27_0	4.1 MB	conda-forge
wheel-0.29.0	py27_0	81 KB	conda-forge
pip-8.1.2	py27_0	1.5 MB	conda-forge

```
$ python
```

Total:	24.2 MB
--------	---------

```
The following NEW packages will be INSTALLED:
```

```
>>> import matplotlib
```

Environment manager



- Conda environment: directory that contains a specific collection of Conda packages that you have installed.
- Packages are symlinked between environments to avoid duplication.

```
$ conda create --name env1 -c bioconda fastqc
$ fastqc --version
-bash: fastqc: command not found
$ source activate env1
$(env1)fastqc --version
FastQC v0.11.5
$(env1)source deactivate
$ conda create --name env2 -c bioconda python=3 snakemake
$ python --version
Python 2.7.12 :: Continuum Analytics, Inc.
$ snakemake --version
-bash: snakemake: command not found
$ source activate env2
$(env2)python --version
Python 3.4.3 :: Continuum Analytics, Inc.
$(env2)snakemake --version
3.7.1
$(env2)
```

Defining and sharing environments



environment.yml

```
channels:  
- conda-forge  
- bioconda  
dependencies:  
- fastqc=0.11  
- sra-tools=2.8  
- snakemake=4.3.0  
- multiqc=1.3  
- bowtie2=2.3  
- samtools=1.6  
- htseq=0.9  
- graphviz=2.38.0
```

- Create an environment from specifications in a file.
- All additional dependencies will be included.
- The environment.yml file can be shared with others and used to recreate the environment on other systems.

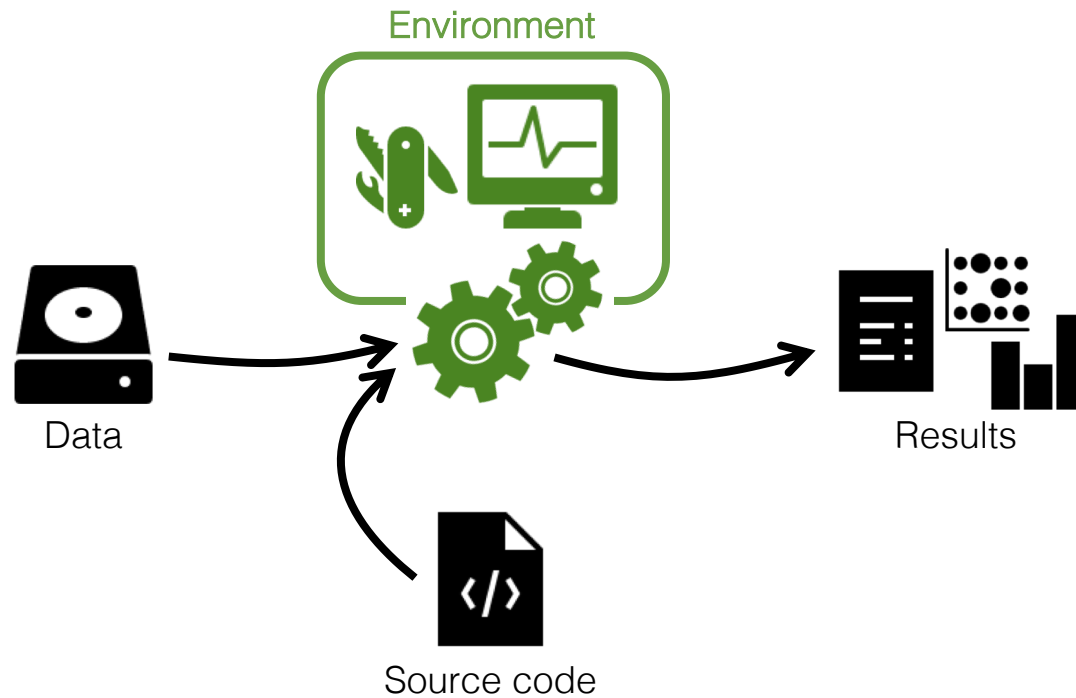
```
$ conda env create --name project_a -f environment.yml
```

- Update existing environment after adding new packages to environment.yml:

```
$ conda env update -f environment.yml
```

- Export existing environment as new yaml file (also includes dependencies):

```
$ conda env export > environment_full.yml
```



```
project
|- doc/
|
|- data/
|   |- raw_external/
|   |- raw_internal/
|   |- meta/
|
|- code/
|- notebooks/
|
|- intermediate/
|- scratch/
|- logs/
|
|- results/
|   |- figures/
|   |- tables/
|   |- reports/
|
|- Snakefile
|- config.yml
|- environment.yml
|- Dockerfile
```

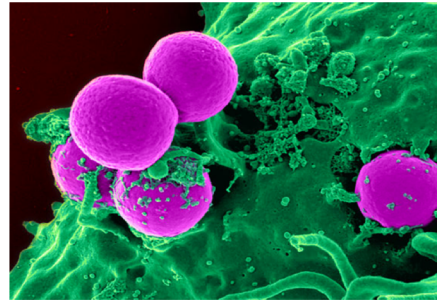
The tutorials

A few practical notes...

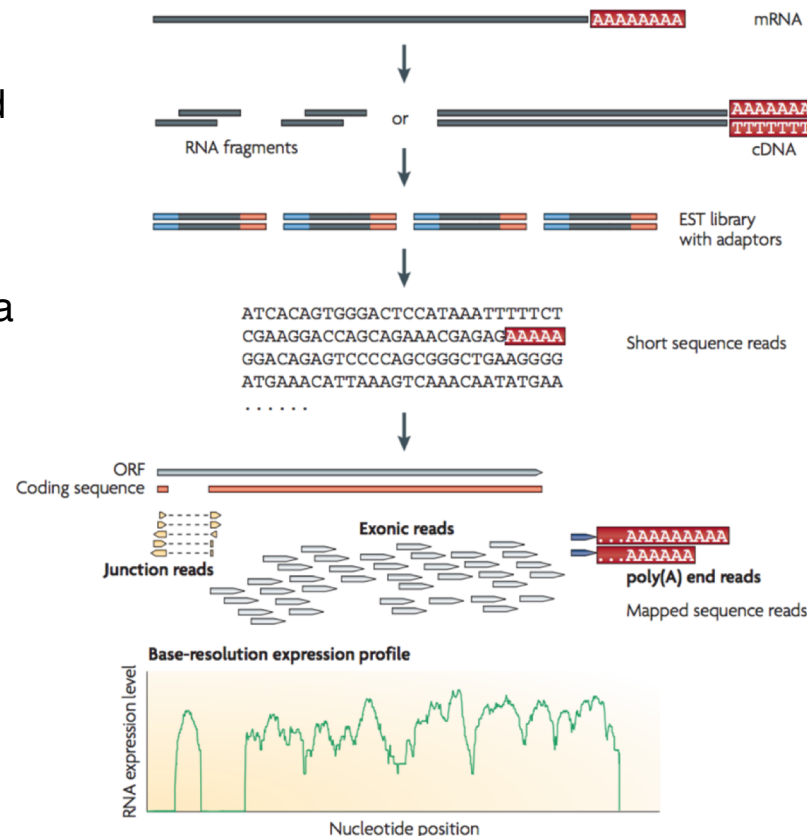
RNA-Seq Reveals Differential Gene Expression in *Staphylococcus aureus* with Single-Nucleotide Resolution

Joseph Osmundson^{1*}, Scott Dewell², Seth A. Darst¹

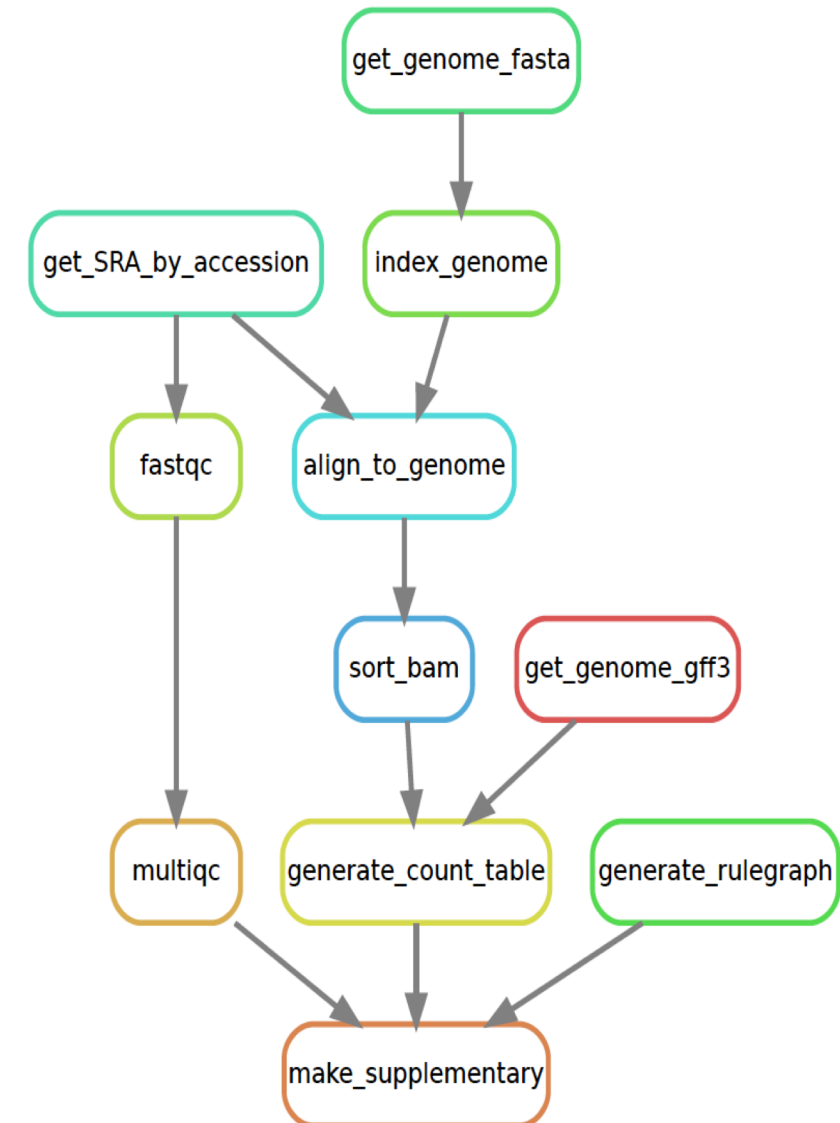
- Methicillin-resistant *Staphylococcus aureus* (MRSA):
 - is resistant to broad spectrum beta-lactam antibiotics
 - lead to difficult-to-treat infections in humans
- Lytic bacteriophages have been suggested as potential therapeutic agents, or as the source of novel antibiotic proteins or peptides.
- One such protein, gp67, was identified as a transcription-inhibiting transcription factor with an antimicrobial effect.
- To identify *S. aureus* genes repressed by gp67, the authors expressed gp67 in *S. aureus* cells.
- RNA-seq was performed on *S. aureus* strains:
 - RN4220 with pRMC2 with gp67
 - RN4220 with empty pRMC2
 - NCTC8325-4



Scanning electron micro-graph of a human neutrophil ingesting MRSA



The analysis workflow



The tutorials

Environment management

Set up and manage the project environment

CONDA

Start here!

Version control

Track and backup your project history



Workflow management

Move from separate scripts to a connected analysis

Snakemake

Reports

Connect code, output and text in fancy reports

R Markdown

from R Studio

Notebooks

Document your exploratory analysis

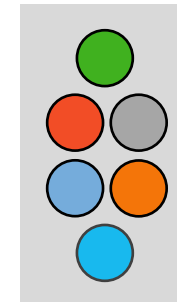
Jupyter

Containerization

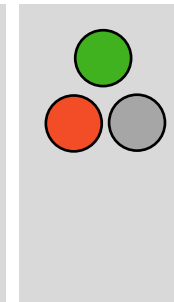
Make your project self-contained and distributable



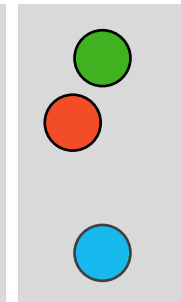
Do it all!



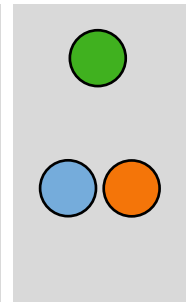
Workflow



Reproducible environment

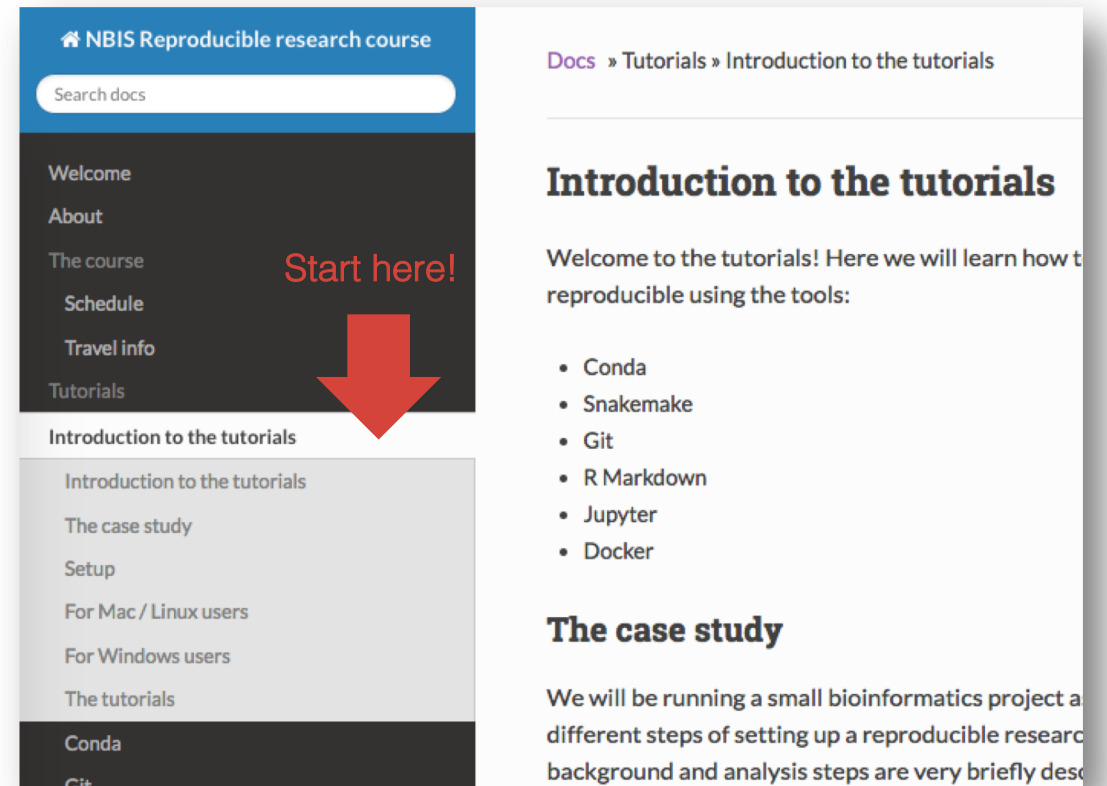


Interactive notebooks



Getting started

- Clone course git repository to get all files needed for tutorials!
- Each tutorial will run in a specific subdirectory within `reproducible_research_course`, make sure you are running from the right place!
- Exception: the git tutorial will be run in a user-created directory outside of `reproducible_research_course`.



<http://nbis-reproducible-research.readthedocs.io>



Managing
dependencies



Managing and executing
analysis workflow



Versioning and
collaborating on code
(and some other files)



Connecting code
and reporting

and...

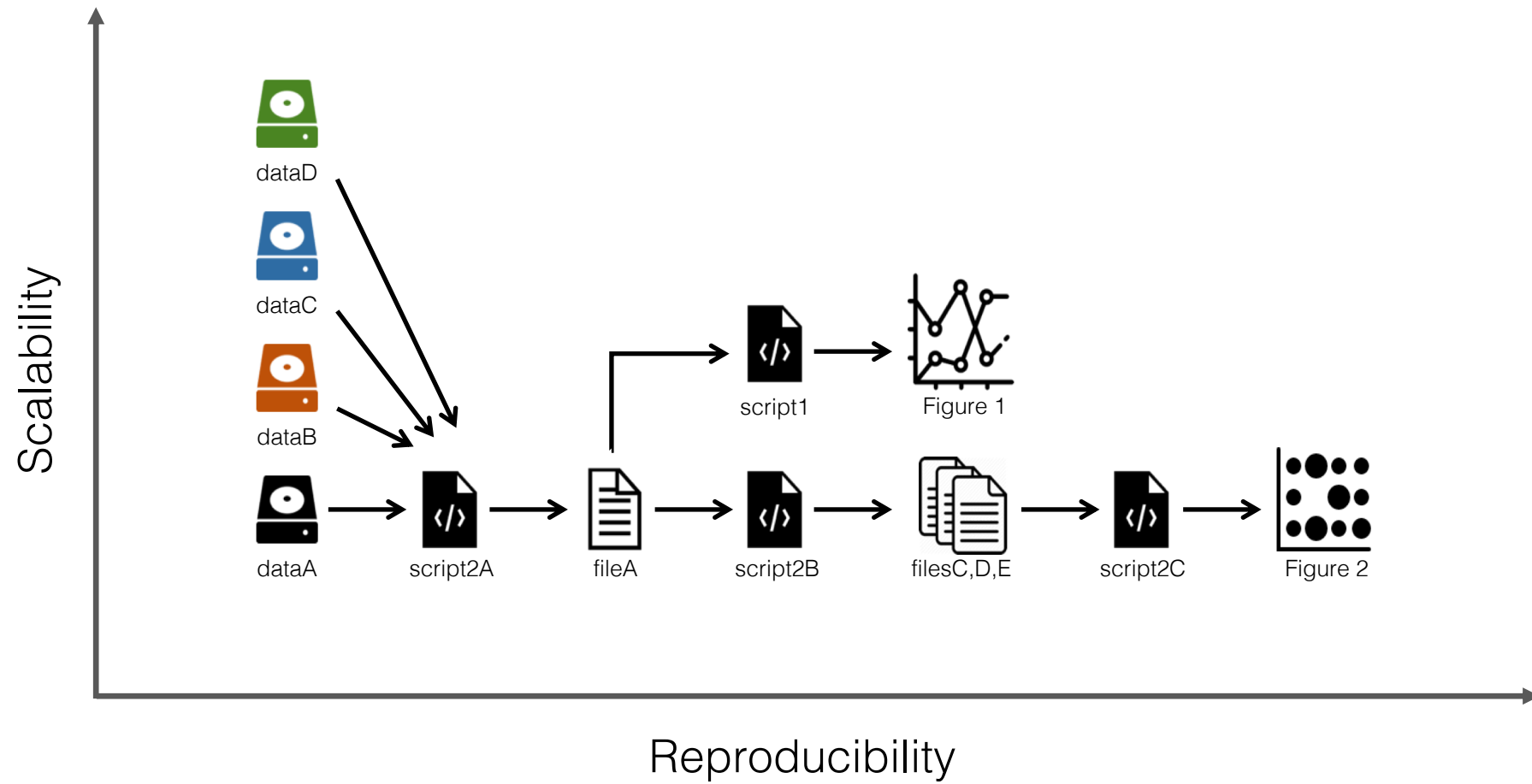


Isolating and exporting
environment

As projects grow or age, it becomes increasingly difficult to keep track of all the parts and how they fit together.



“Snakemake is a workflow management system that aims to reduce the complexity of creating workflows by providing a fast and comfortable execution environment, together with a clean and modern specification language in python style.”



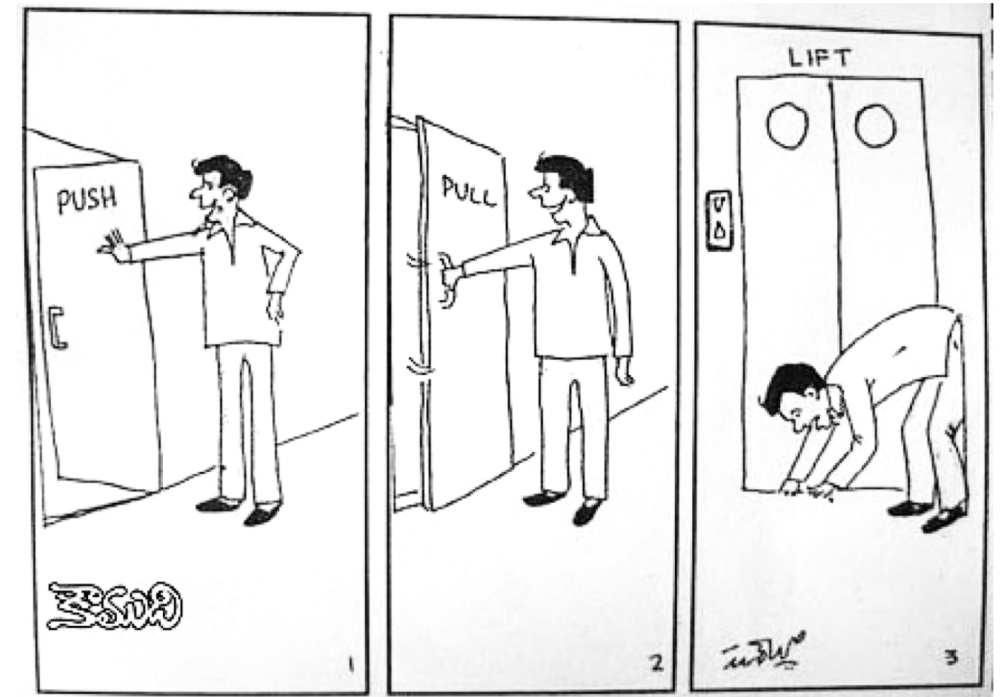
Workflow management systems come in different flavors.

Explicit syntax ("Push")

"Here are my inputs, please perform these operations in this order on them."

Implicit syntax ("Pull")

"I need this output, could you please figure out which operations to perform and in which order?"



Explicit approach using Bash

trim_and_zip.sh

```
for sample in *.fastq
do
    id=$(echo ${sample} | sed 's/\.fastq//')

    # Trim fastq file
    echo "Trimming ${id}"
    seqtk trimfq -b 5 -e 10 $sample > \
    ${id}.trimmed.fastq

    # Compress fastq file
    echo "Compressing ${id}"
    gzip -c ${id}.trimmed.fastq > \
    ${id}.trimmed.fastq.gz

    # Remove intermediate files
    rm ${id}.trimmed.fastq
done
```

```
$bash trim_and_zip.sh
Trimming sample: a
Compressing sample: a
Trimming sample: b
Compressing sample: b
```

Implicit approach using Snakemake

Snakefile

```
rule trim_fastq:
    input: "{prefix}.fastq"
    output: temp("{prefix}.trimmed.fastq")
    shell:
        "seqtk trimfq -b 5 -e 10 {input} > {output}"

rule gzip:
    input: "{prefix}"
    output: "{prefix}.gz"
    shell:
        "gzip -c {input} > {output}"
```

```
$snakemake {a,b}.trimmed.fastq.gz
Provided cores: 1
Rules claiming more threads will be
scaled down.
```

```
Job counts:
count  jobs
2      gzip
2      trim_fastq
4
```

```
rule trim_fastq:
    input: a.fastq
    output: a.trimmed.fastq
    wildcards: prefix=a
```

1 of 4 steps (25%) done

```
rule gzip:
    input: a.trimmed.fastq
    output: a.trimmed.fastq.gz
    wildcards: prefix=a.trimmed.fastq
```

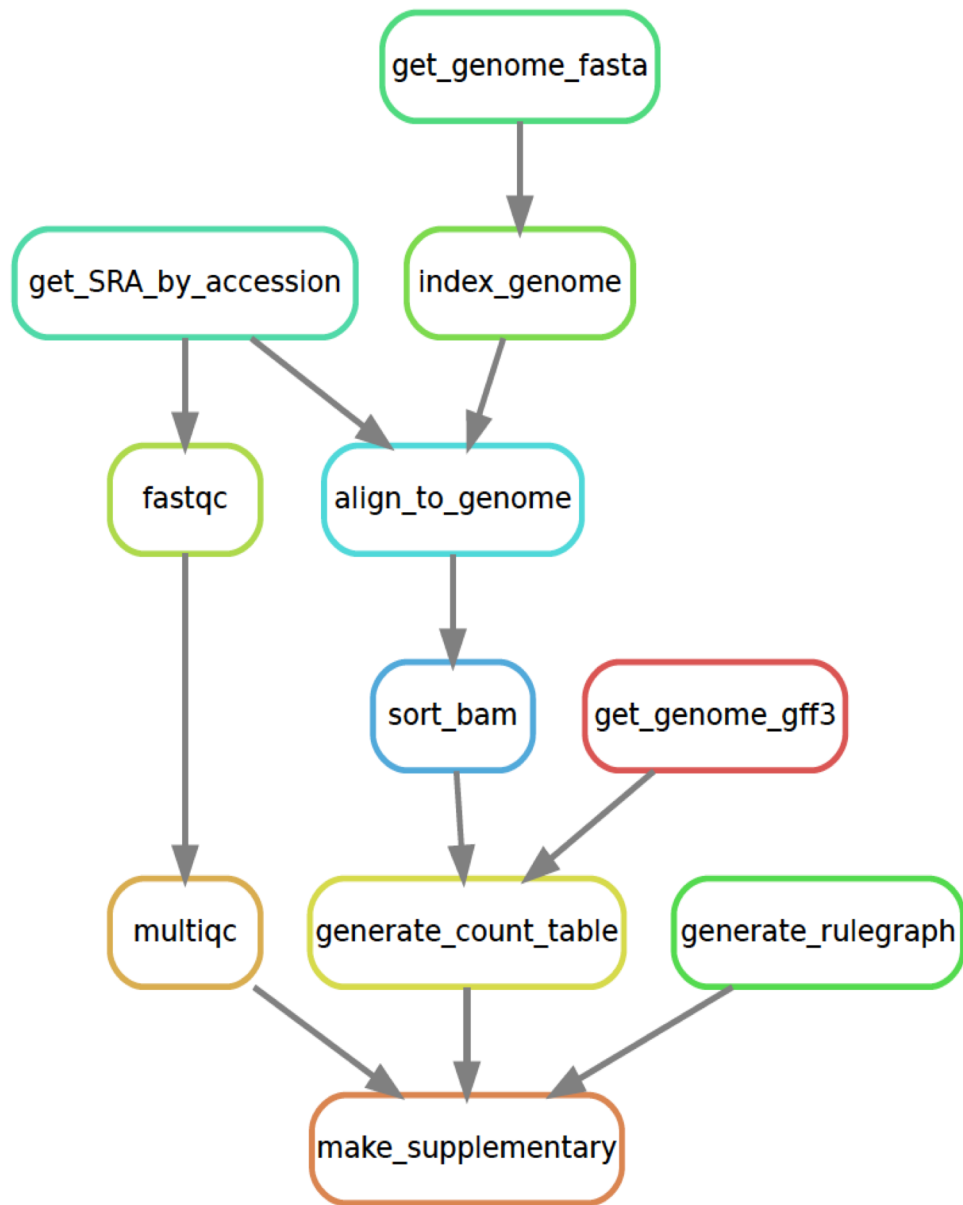
Removing temporary output file a.trimmed.fastq.
2 of 4 steps (50%) done

```
rule trim_fastq:
    input: b.fastq
    output: b.trimmed.fastq
    wildcards: prefix=b
```

3 of 4 steps (75%) done

```
rule gzip:
    input: b.trimmed.fastq
    output: b.trimmed.fastq.gz
    wildcards: prefix=b.trimmed.fastq
```

Removing temporary output file b.trimmed.fastq.
4 of 4 steps (100%) done



Snakemake figures out how rules can be pieced together to generate some requested output.

Here we ask for `supplementary.pdf`, which is an R Markdown report generated by the rule `make_supplementary`.

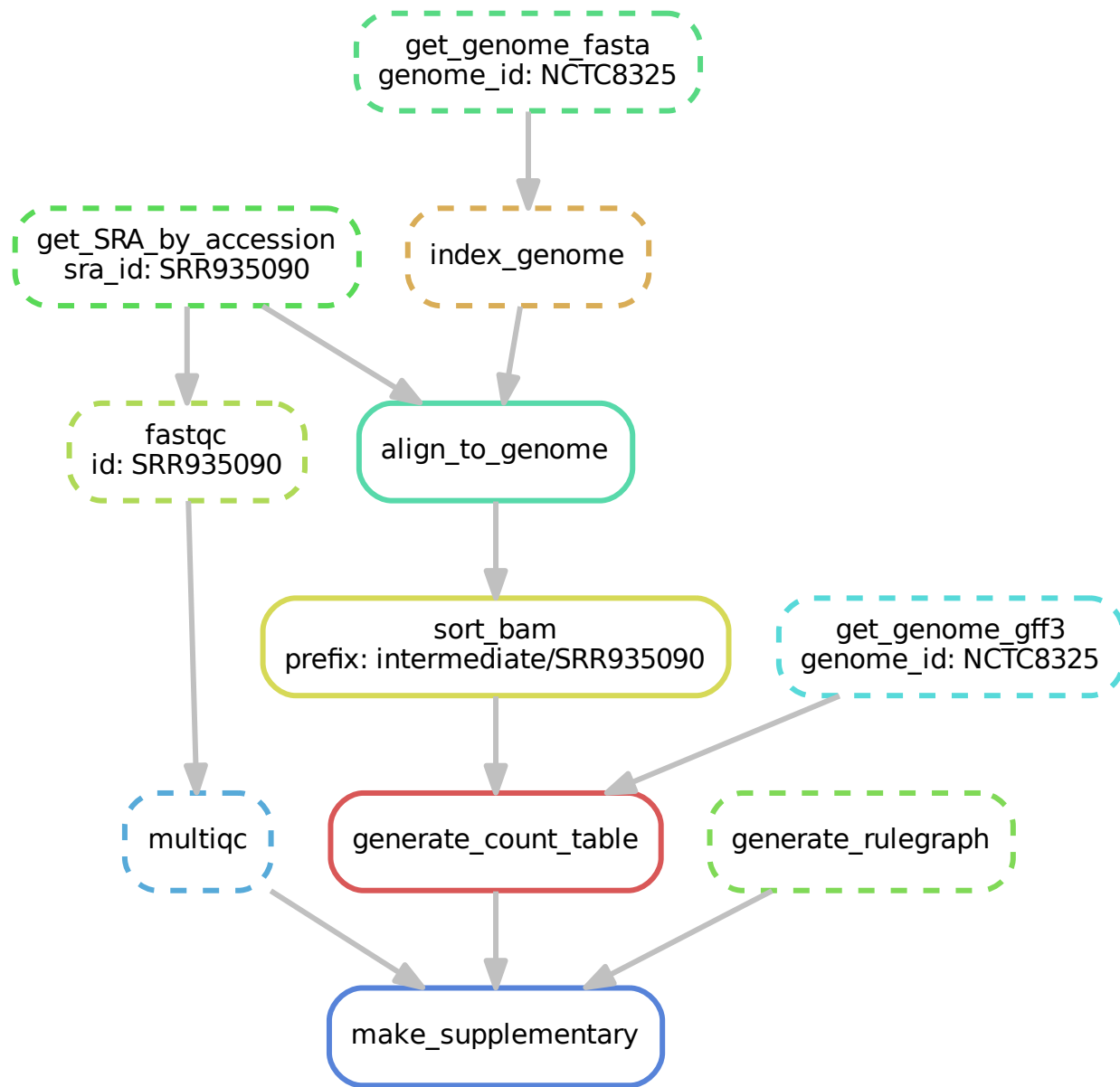
```
$snakemake supplementary.pdf --rulegraph | dot -Tpdf > rulegraph.pdf
```



Snakemake keeps track of when files were generated and by which rules.

Dotted rule boxes show that `supplementary.pdf` already exists and that it's newer than its dependencies (recursively).

```
$snakemake supplementary.pdf --dag | dot -Tpdf > dag.pdf
```

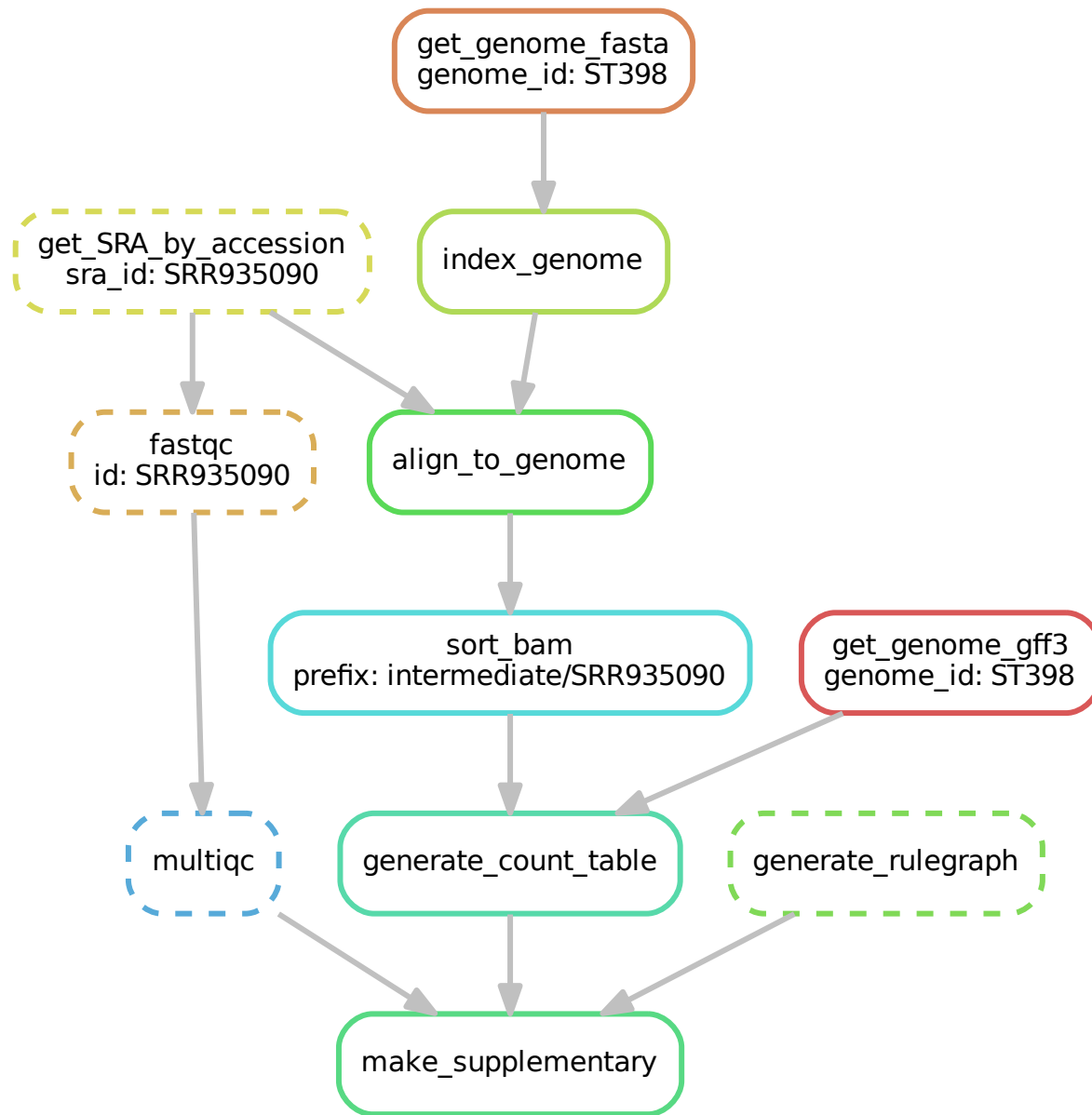



Here Snakemake detects that a file used in `align_to_genome` is newer than downstream files, so it reruns the necessary rules.

```

$touch intermediate/NCTC8325.1.bt2
$snakemake supplementary.pdf --dag | dot -Tpdf > dag.pdf

```



Forcing a rule (`get_genome_fasta` here) to be rerun also leads to rerunning all rules that depend on it.

Note that we also change the parameter "genome_id" to use another genome to align to. This causes `get_genome_gff3` to be rerun as well.

```
$snakemake supplementary.pdf --config genome_id=ST398
-f get_genome_fasta --dag | dot -Tpdf > dag.pdf
```

Anatomy of a Snakemake rule

```
import os

rule trim_fastq:
    input: "{prefix}.fastq"
    output: temp("{prefix}.trimmed.fastq")
    params:
        leftTrim=5,
        rightTrim=10
    log: "logs/trim_fastq.log"
    version: "0.1"
    message: "Trimming {input[0]}."
    shadow: True
    threads: 8
    priority: 90
    resources: mem=64
    conda: "envs/seqtk.yaml"
    singularity: "docker://quay.io/biocontainers/seqtk"
    run:
        if (os.stat(input[0]).st_size > 0):
            shell("seqtk trimfq -t {threads} -b {params.leftTrim}
                  -e {params.rightTrim} {input} > {output} 2> {log}")
        else:
            raise IOError(input[0]+" is empty.")
```

Command line interface

```
# execute the workflow with target a.trimmed.fastq.gz
snakemake a.trimmed.fastq.gz
```

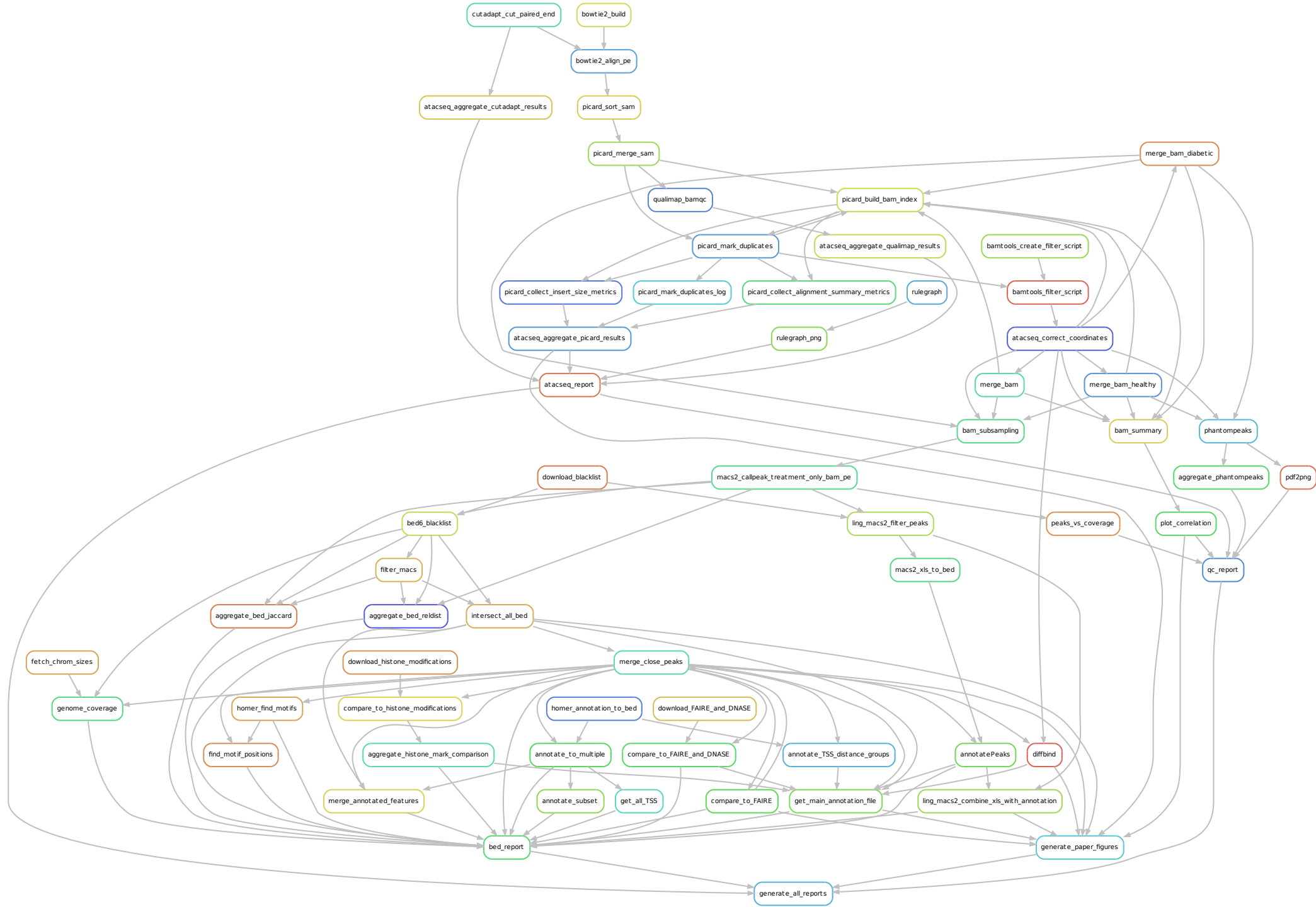
```
# execute the workflow with the first rule as target
snakemake
```

```
# dry-run, print shell commands and reason for execution
snakemake -n -p -r
```

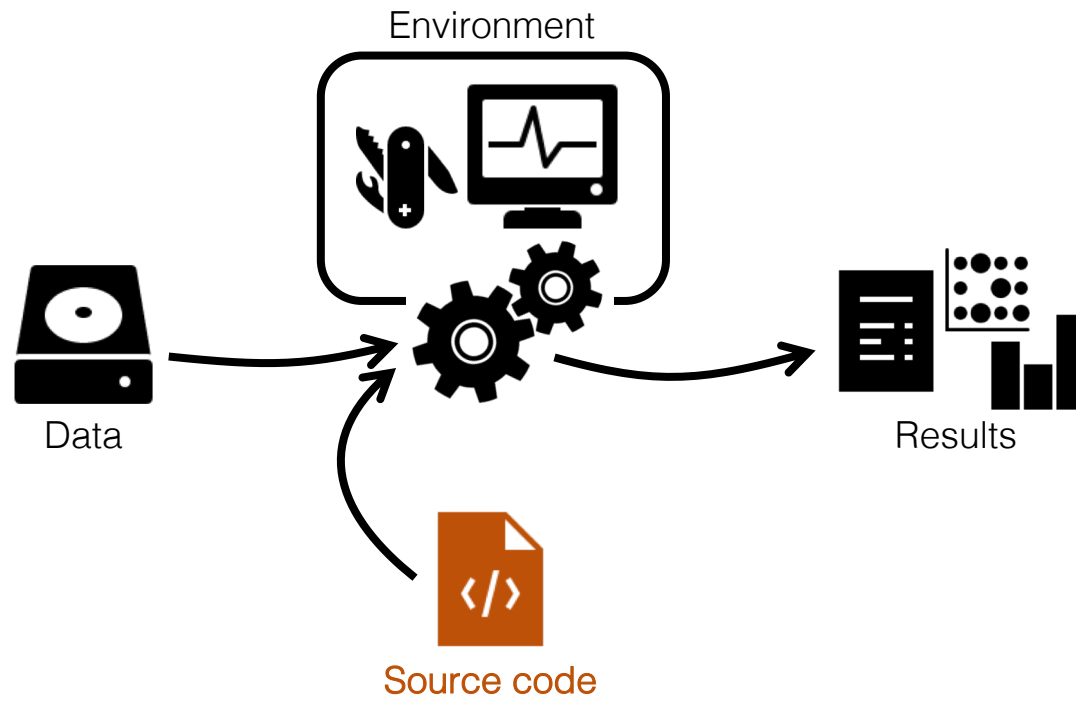
```
# visualize the DAG of jobs using the Graphviz dot command
snakemake --dag | dot -Tsvg > dag.svg
```

```
# execute the workflow with 8 cores
snakemake --cores 8
```

```
# run the workflow on a SLURM cluster
snakemake --cluster-config cluster.yml --cluster \
    "sbatch -A {cluster.account} -t {cluster.time}"
```



Things can get rather complex...



```
project
|- doc/
|
|- data/
|   |- raw_external/
|   |- raw_internal/
|   |- meta/
|
|- code/
|- notebooks/
|
|- intermediate/
|- scratch/
|- logs/
|
|- results/
|   |- figures/
|   |- tables/
|   |- reports/
|
|- Snakefile
|- config.yml
|- environment.yml
|- Dockerfile
```



Managing
dependencies



Managing and executing
analysis workflow



Versioning and
collaborating on code
(and some other files)



Connecting code
and reporting

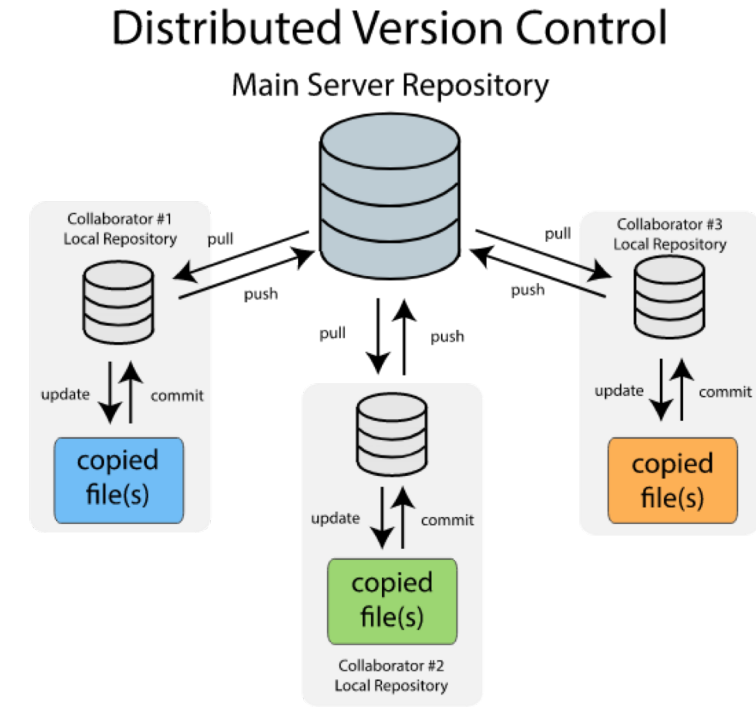
and...



Isolating and exporting
environment

What is Git?

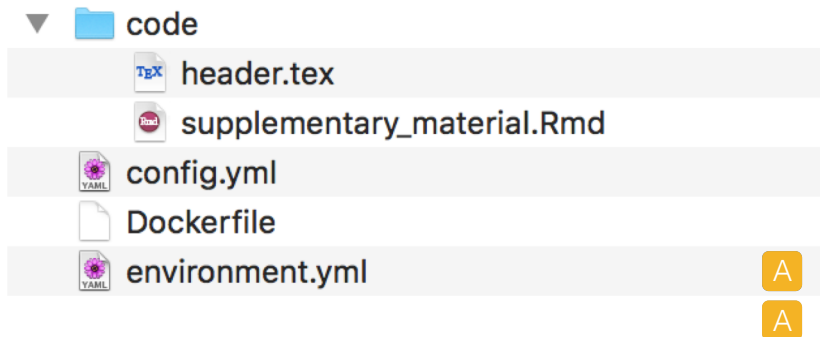
- Widely used system for file version control.
- Keeps copies of files and code from every stages in their lifecycle.
- Kind of like Dropbox, but you decide when each version is saved (and a lot more advanced features).
- Each file has a distinct history with specific incremental changes (each with a reference code).
 - Revert files to previous state.
 - Compare changes over time.
 - See who modified what.
- Makes you fearless.
- Runs on command line, but there also exists GUI and integration in e.g. text editors.
- Mainly for text files, not for binary files or large files.
- Versioning, backup, and sharing!



How does it work in practice?

Nomenclature

Repository	<i>Directory with all files, will include a .git folder</i>
Commit	<i>A specific version of the repository</i>
Push	<i>Upload local changes to remote repository</i>
Pull	<i>Download changes from remote repository</i>



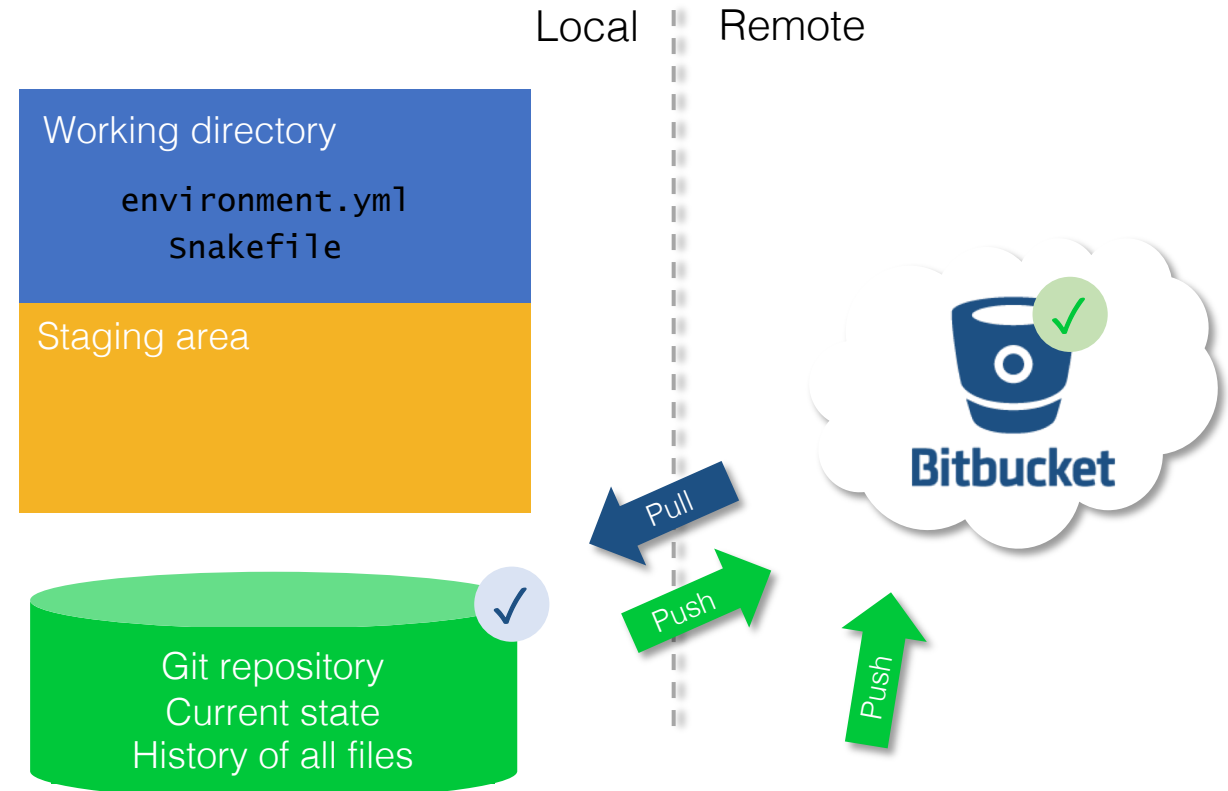
- Edit environment.yml
- Add a new file Snakefile
- `git add environment.yml`
- `git commit -m "Add snakemake 4.4.0"`
- `git add Snakefile`
- `git commit -m "Add Snakefile"`
- `git push`
- `git pull`


git log











```
a2c36bs Add heatmap figure
6152ff6 Format figure label
0abd0cb Update multiqc version
8dhfls8 Add snakemake 4.4.0
kfhs7s6 Add Snakefile
2kd7f0f Fix alignment command
```

git log

```
a2c36bs Add heatmap figure
6152ff6 Format figure label
0abd0cb Update multiqc version
8dhfls8 Add snakemake 4.4.0
kfhs7s6 Add Snakefile
2kd7f0f Fix alignment command
```









 reproducible_resear...


-  Overview
-  Source
-  Commits
-  Branches
-  Pull requests
-  Pipelines
-  Issues
-  Downloads
-  Boards
-  Settings






















SciLifeLab Bioinformatics LTS / Bioinformatics LTS / reproducible_research_course / Source

docker/


 master ▾  reproducible_research_course /

	..			
	code			
	Dockerfile	1.1 KB	5 hours ago	split git_jupyter_docker
	Snakefile	6.5 KB	4 hours ago	moved counts table
	config.yml	1.5 KB	5 hours ago	split git_jupyter_docker
	environment.yml	202 B	5 hours ago	split git_jupyter_docker



	Rasmus Ågren	232ffc0	language and spelling in git tutorial	5 hours ago
	varemo	a81f6e8	remove Where to next sections	6 hours ago
	varemo	d8e30df <small>M</small>	Merge branch 'master' of https://bitbucket.org/scilifelab-lts/reproducible_research...	6 hours ago
	varemo	449b4b9	remove todo	6 hours ago
	Rasmus Ågren	bfa0edd <small>M</small>	Merge branch 'master' of https://bitbucket.org/scilifelab-lts/reproducible_research...	7 hours ago
	Rasmus Ågren	851e43e	language and spelling in git tutorial	7 hours ago
	varemo	cc2bbcb	fix figure fonts	7 hours ago
	varemo	19c67d4 <small>M</small>	Merge branch 'master' of https://bitbucket.org/scilifelab-lts/reproducible_research...	7 hours ago
	varemo	3a1806f	fix nicer figure	7 hours ago
	Rasmus Ågren	1487f3b	test anchoring	8 hours ago
	Rasmus Ågren	3397711	split snakemake env for speed	8 hours ago
	varemo	8f265b0	tutorial text	9 hours ago
	varemo	414de52	typo	9 hours ago
	varemo	208e1d9	restructure pages	10 hours ago
	varemo	9fca2f7	remove take down part	10 hours ago
	varemo	831a44e	revert back to project a from project new	10 hours ago
	varemo	3575b9c <small>M</small>	Merge branch 'master' of https://bitbucket.org/scilifelab-lts/reproducible_research...	11 hours ago
	Rasmus Ågren	0e4296a	typo	11 hours ago
	varemo	aeffaa3	restructure pages	11 hours ago
	Rasmus Ågren	14c60cb <small>M</small>	Merge branch 'master' of https://bitbucket.org/scilifelab-lts/reproducible_research...	11 hours ago
	Rasmus Ågren	9795b0f	language and spelling in the conda tutorial	11 hours ago

Diff from  8baf4a1 2017-11-22 ▼ to  d797810 2017-11-23 ▼

 snakemake/Snakefile_mrsa **MODIFIED**

Side-by-side diff



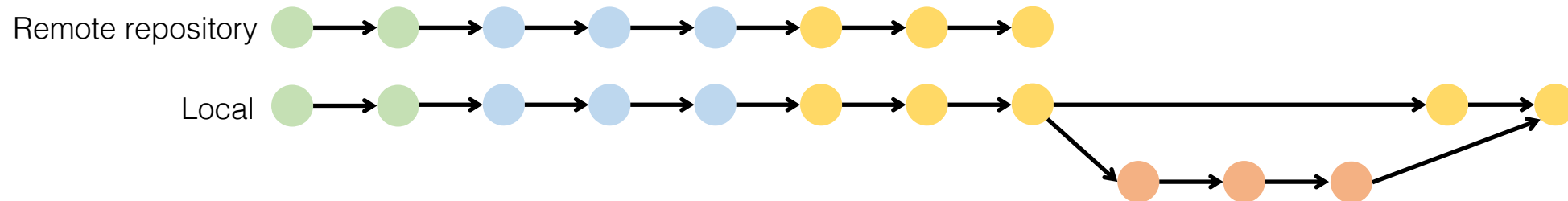
...

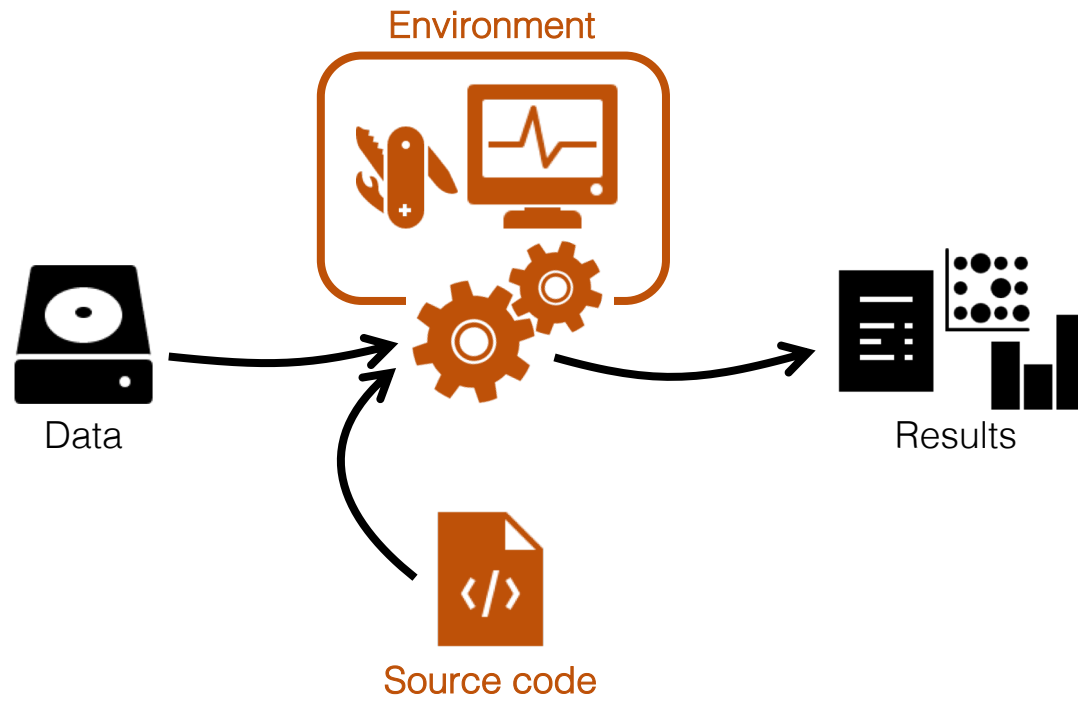
156	156	
157	157	# Save the count table as a temporary file and then prepend a header line
158	158	# with the sample names
159	-	htseq-count --format bam --type gene --idattr gene_id {input.bams} tempfile > tempfile2
	159	+ htseq-count --format bam --type gene --additional-attr description --idattr gene_id {input.bams} tempfile > tempfile2
160	160	echo '{input.bams}' tr ' ' '\t' cat - tempfile2 > {output}
161	161	
162	162	# Remove the temporary files

...

During the working day...

- Pull collaborator's latest work to get your local repository up to date.
- Carry on with your work and edit files.
- Commit often!
 - Each commit should be related to a distinct change/addition/task.
 - Write descriptive commit messages.
- Push your changes to the remote repository.
- If you know several people are actively working on the same repository, push and pull often!





```
project
|- doc/
|
|- data/
|   |- raw_external/
|   |- raw_internal/
|   |- meta/
|
|- code/
|- notebooks/
|
|- intermediate/
|- scratch/
|- logs/
|
|- results/
|   |- figures/
|   |- tables/
|   |- reports/
|
|- Snakefile
|- config.yml
|- environment.yml
|- Dockerfile
```



Managing
dependencies



Managing and executing
analysis workflow



Versioning and
collaborating on code
(and some other files)



Connecting code
and reporting

and...



Isolating and exporting
environment

Typical guidelines for keeping a notebook of wet-lab work

Record everything you do in the lab, even if you are following a published procedure.

Use a bound notebook so that tear-out would be visible.

When you finish a page, put a corner-to corner line through any blank parts that could still be used for data entry.

Write a title for each and every new set of entries.

The investigator and supervisor must sign each page.

If you make a mistake, put a line through the mistake and write the new information next to it.

If you're testing a specific hypothesis, write it down beforehand.

All pages must be pre-numbered.

Each page should be numbered and dated consistently.

Properly introduce and summarize each experiment.

Use a ball point pen so that marks will not smear nor will they be erasable.

Use a ball point pen so that marks will not smear nor will they be erasable.

It is critical that you enter all procedures and data directly into your notebook in a timely manner.

Typical guidelines for keeping a notebook of dry-lab work

Literate programming

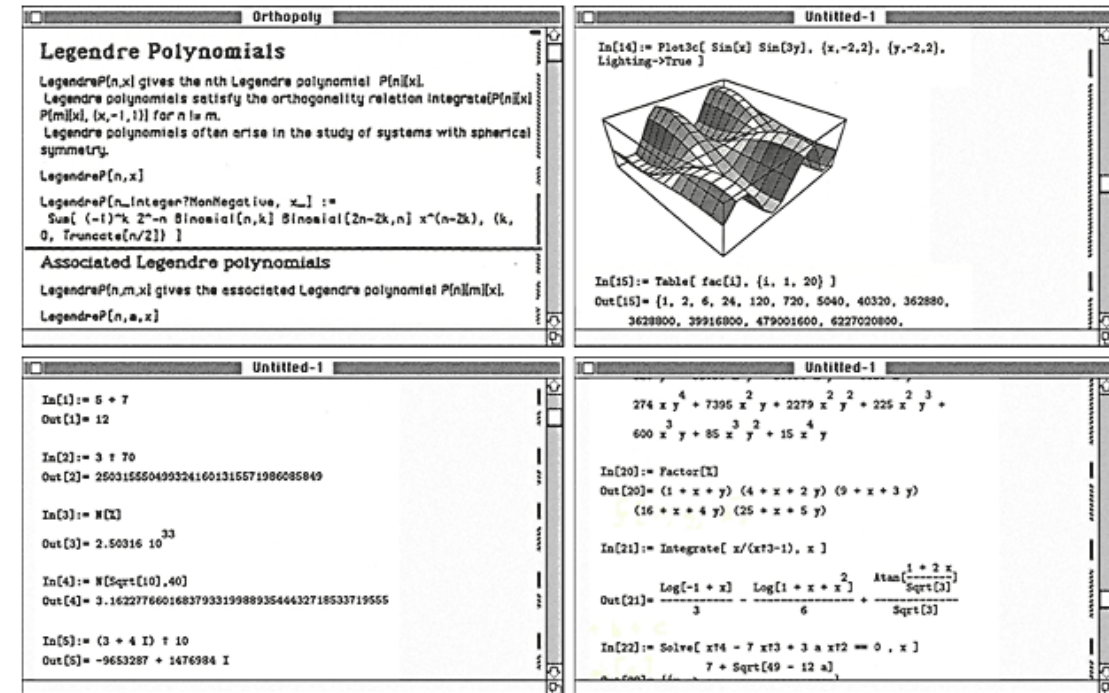
Instead of imagining that our main task is to instruct a computer what to do, let us concentrate rather on explaining to human beings what we want a computer to do.

Donald Knuth (1984)

Literate computing

A literate computing environment is one that allows users not only to execute commands interactively, but also to store in a literate document the results of these commands along with figures and free-form text.

Millman KJ and Perez F (2014)



Wolfgang Mathematica notebook (1988)

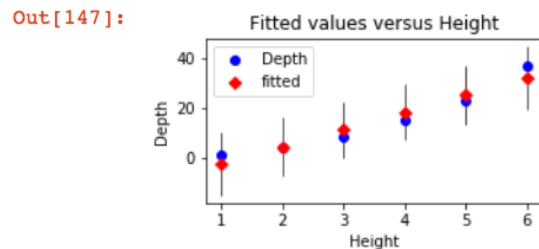
Study of velocity/energy relationship (1722-03-17, W. Gravesande)

If I drop brass balls from various heights and measure penetration depth in a block of clay, can I settle the dispute regarding conservation of energy?

```
In [146]: import statsmodels as sm; import matplotlib.pyplot as plt; import numpy
plt.rcParams["figure.figsize"] = (4,2)
data = {"Depth":[1.1, 4.3, 8.7, 15.5, 23.2, 37.0], "Height":[1, 2, 3, 4, 5, 6]}
```

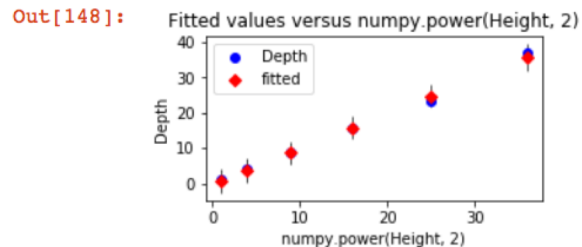
Test Bessler's hypothesis that kinetic energy, and thereby penetration depth, is linear with respect to height.

```
In [147]: res = sm.formula.api.ols(formula = 'Depth ~ Height', data = data).fit()
sm.graphics.api.plot_fit(res, 1)
```



Okayish, but maybe Bernoulli's quadratic model fits better?

```
In [148]: res = smf.ols(formula = 'Depth ~ numpy.power(Height,2)', data = data).fit()
sm.graphics.api.plot_fit(res, 1)
```



Neat!

- The Jupyter Notebook is a web application for *interactive* data science and scientific computing.
- In-browser editing for code, with automatic syntax highlighting, indentation, and tab completion/introspection.
- The ability to execute code from the browser, with the results of computations attached to the code which generated them.
- Mix and match languages to suit your needs (e.g. scikit-learn + ggplot2).

Runs as a local web server →

Load/save/manage notebooks →

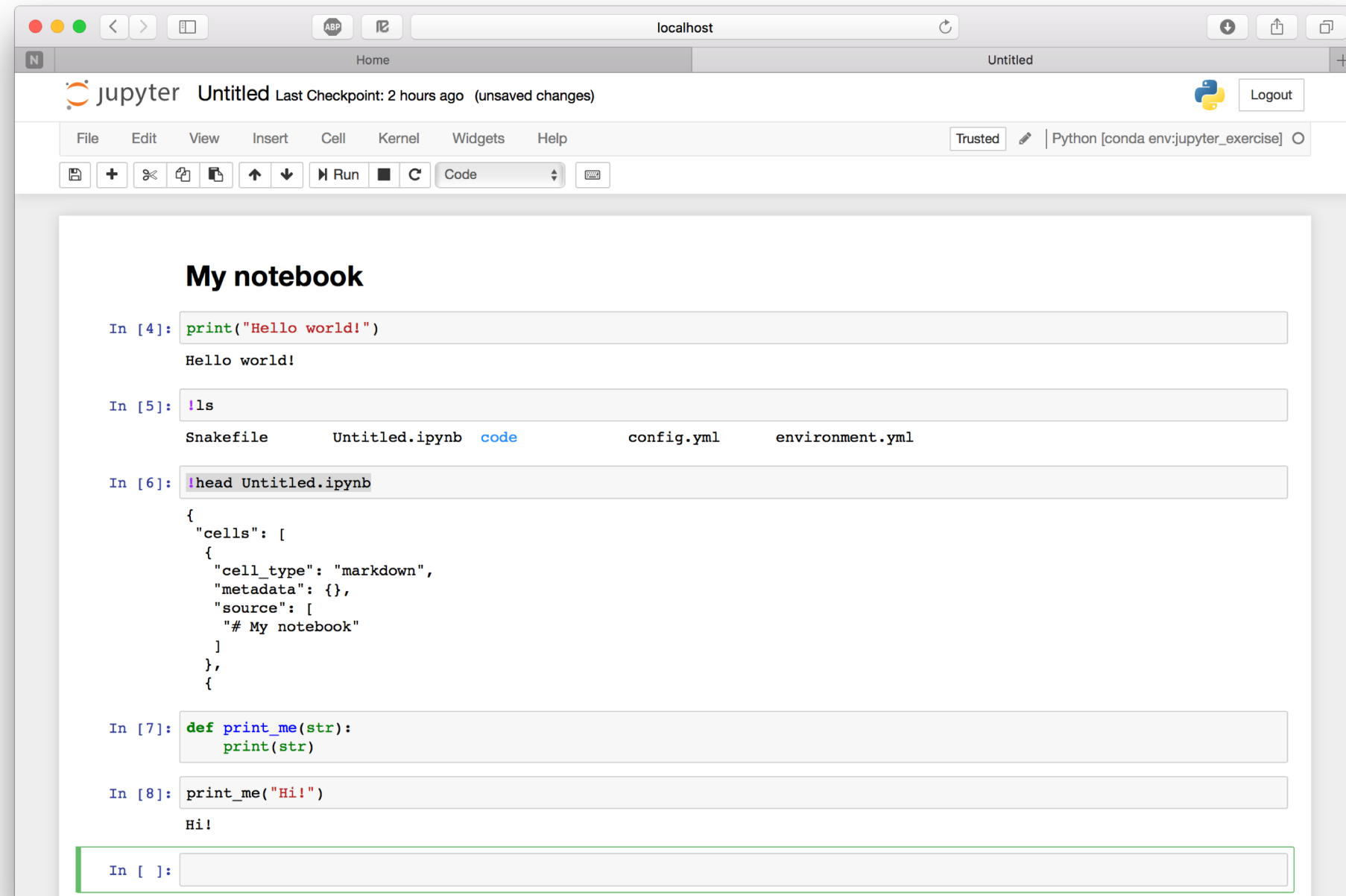
Markdown cell with a header →

Code cell with some Python code →

Run shell command to list files →

The notebook itself is a JSON file →

You can define and call functions →



localhost

Home Untitled

jupyter Untitled Last Checkpoint: 2 hours ago (unsaved changes) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python [conda env:jupyter_exercise]

Code

My notebook

```
In [4]: print("Hello world!")
Hello world!
```

```
In [5]: !ls
Snakefile      Untitled.ipynb  code            config.yml      environment.yml
```

```
In [6]: !head Untitled.ipynb
{
  "cells": [
    {
      "cell_type": "markdown",
      "metadata": {},
      "source": [
        "# My notebook"
      ]
    },
    {

```

```
In [7]: def print_me(str):
        print(str)
```

```
In [8]: print_me("Hi!")
Hi!
```

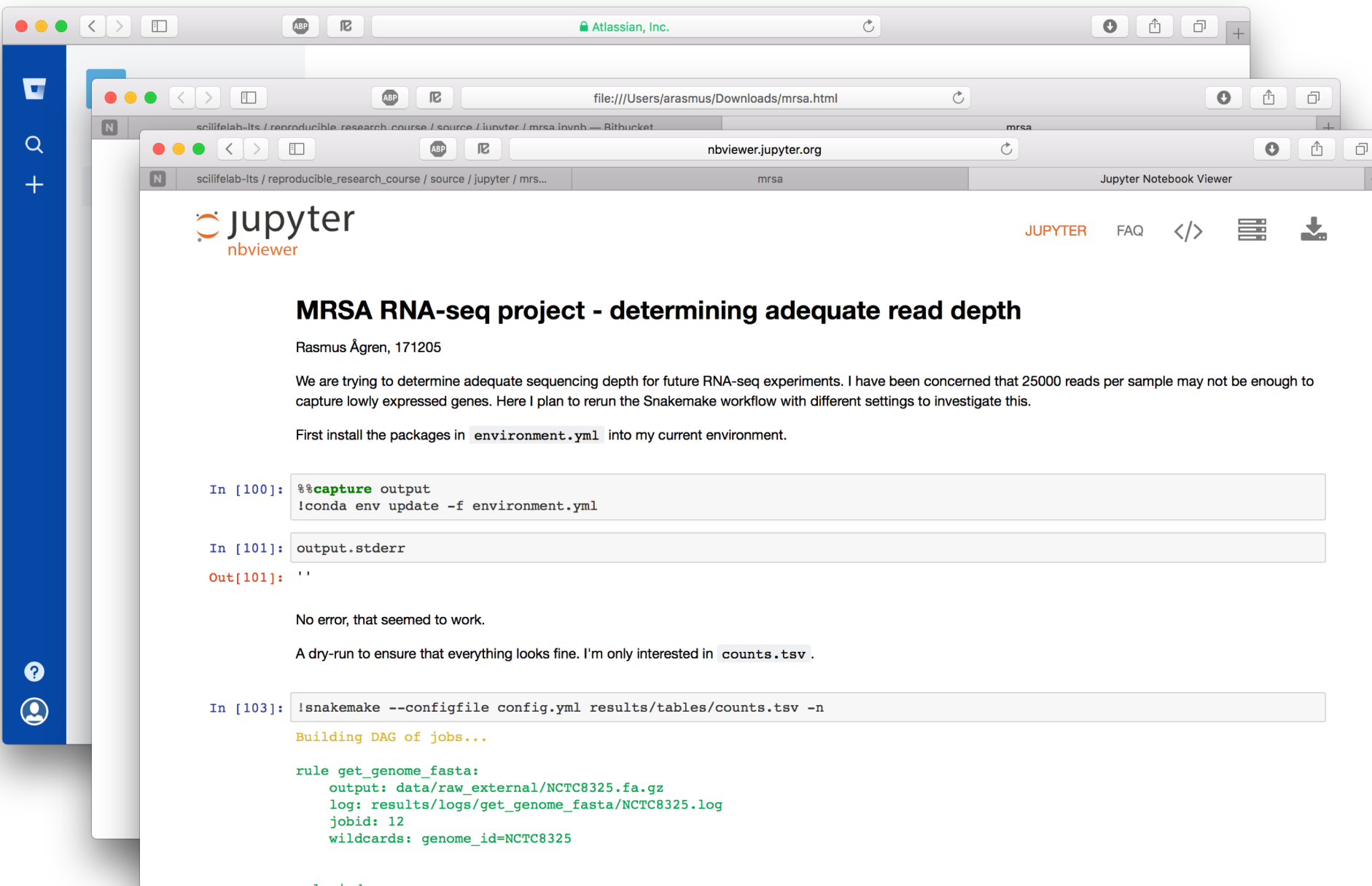
```
In [ ]:
```

Sharing is caring

Put the notebook on GitHub/Bitbucket and it will be rendered there..

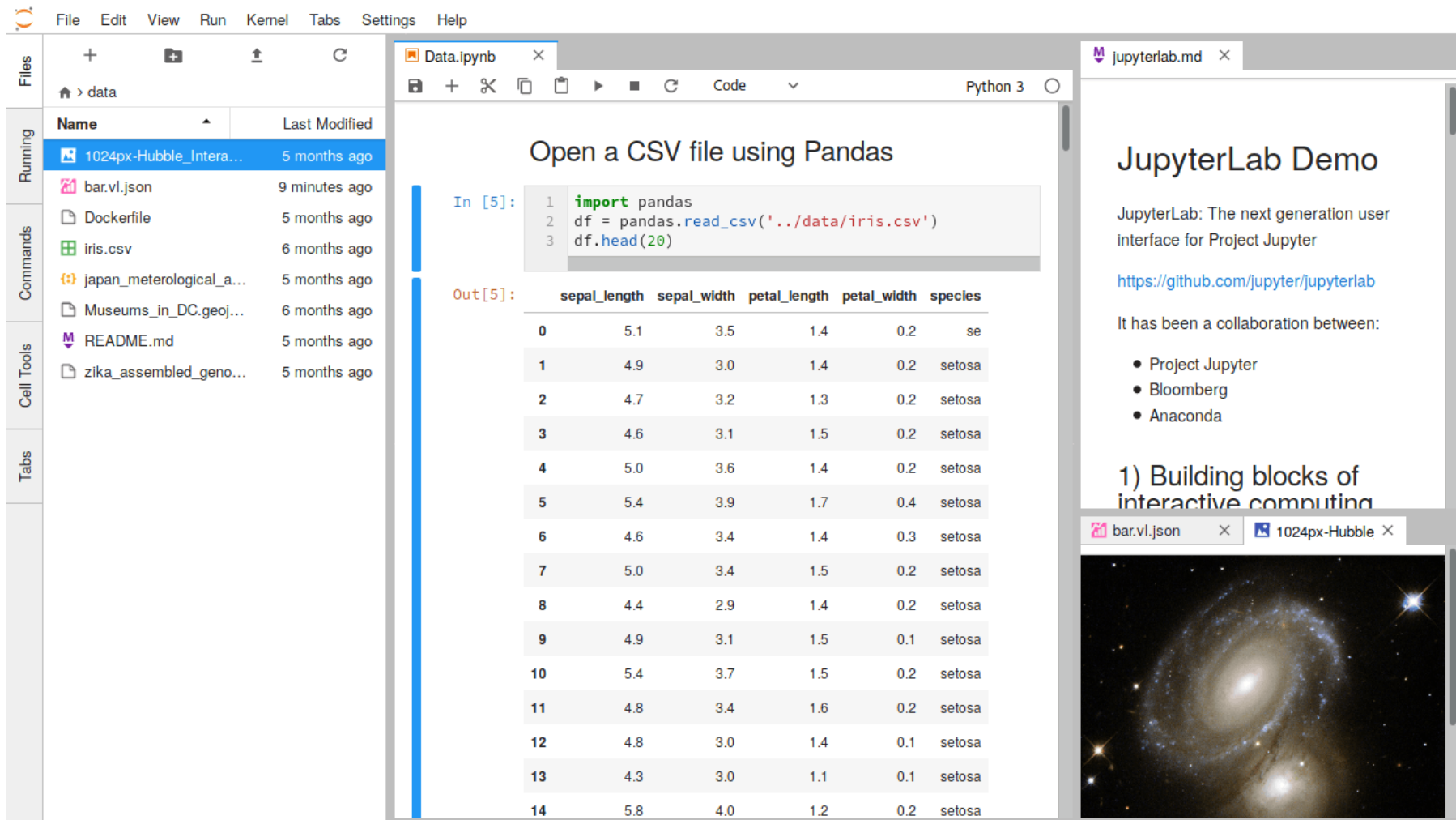
.. or export to one of many different formats, including HTML and PDF ..

.. or paste a link to any Jupyter notebook at nbviewer.jupyter.org and they will render it for you.



The screenshot shows a web browser displaying a Jupyter Notebook Viewer. The browser's address bar shows the URL `nbviewer.jupyter.org`. The notebook title is "MRSA RNA-seq project - determining adequate read depth" by Rasmus Ågren, 171205. The notebook content includes a text block explaining the project's goal to determine adequate sequencing depth for future RNA-seq experiments, followed by a code cell with the command `conda env update -f environment.yml`. The output of this cell is "No error, that seemed to work." and "A dry-run to ensure that everything looks fine. I'm only interested in counts.tsv". The notebook also includes a code cell with the command `snakemake --configfile config.yml results/tables/counts.tsv -n`.

JupyterLab



The screenshot displays the JupyterLab environment. On the left is a sidebar with a 'Files' view showing a directory structure under 'data'. The main area is divided into two panes. The top pane, titled 'Data.ipynb', contains a code cell with the following Python code:

```

In [5]: 1 import pandas
        2 df = pandas.read_csv('../data/iris.csv')
        3 df.head(20)

```

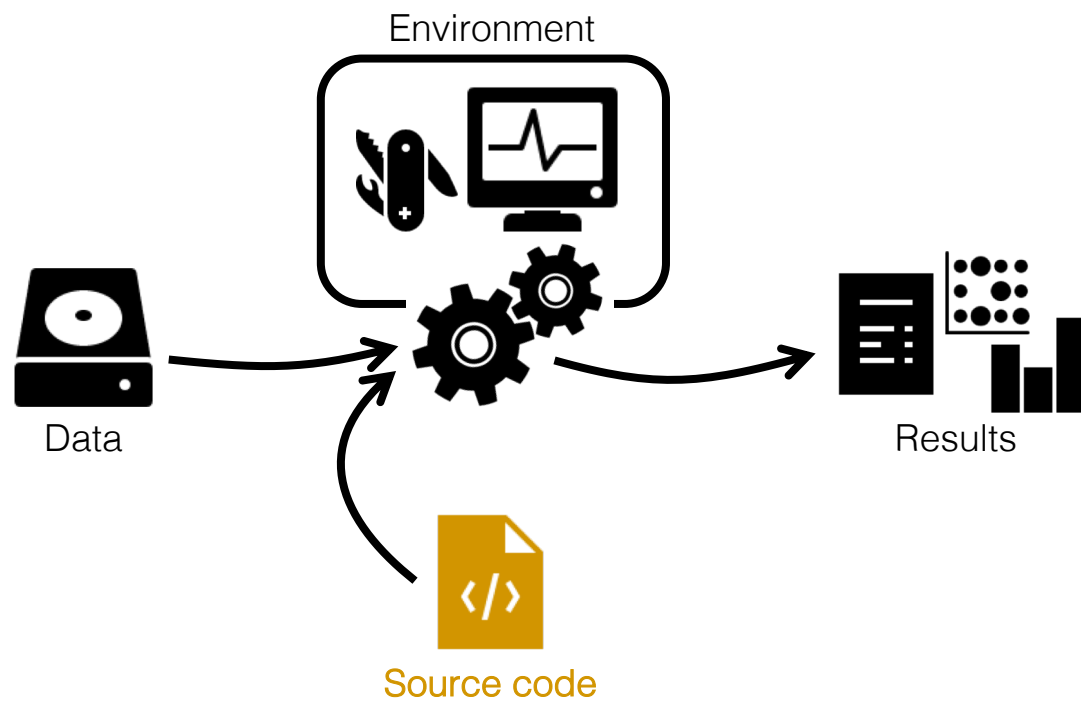
The bottom pane shows the output of the code as a table:

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	se
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
5	5.4	3.9	1.7	0.4	setosa
6	4.6	3.4	1.4	0.3	setosa
7	5.0	3.4	1.5	0.2	setosa
8	4.4	2.9	1.4	0.2	setosa
9	4.9	3.1	1.5	0.1	setosa
10	5.4	3.7	1.5	0.2	setosa
11	4.8	3.4	1.6	0.2	setosa
12	4.8	3.0	1.4	0.1	setosa
13	4.3	3.0	1.1	0.1	setosa
14	5.8	4.0	1.2	0.2	setosa

The right pane shows a 'JupyterLab Demo' page with the text: 'JupyterLab: The next generation user interface for Project Jupyter' and a link to <https://github.com/jupyter/jupyterlab>. Below this, it lists collaborators: Project Jupyter, Bloomberg, and Anaconda. At the bottom of the right pane, there is a preview of a file named '1024px-Hubble' showing a spiral galaxy.

JupyterLab is a full-fledged IDE, similar to e.g. Rstudio.

```
conda install -c conda-forge jupyterlab
```



```
project
|- doc/
|
|- data/
|   |- raw_external/
|   |- raw_internal/
|   |- meta/
|
|- code/
|- notebooks/
|
|- intermediate/
|- scratch/
|- logs/
|
|- results/
|   |- figures/
|   |- tables/
|   |- reports/
|
|- Snakefile
|- config.yml
|- environment.yml
|- Dockerfile
```



Managing
dependencies



Managing and executing
analysis workflow



Versioning and
collaborating on code
(and some other files)



Connecting code
and reporting

and...



Isolating and exporting
environment

Supplementary material

John Doe, Joan Dough, Jan Doh, Dyon Do

18 March, 2018

Read in the data

We have *count data* for three samples:

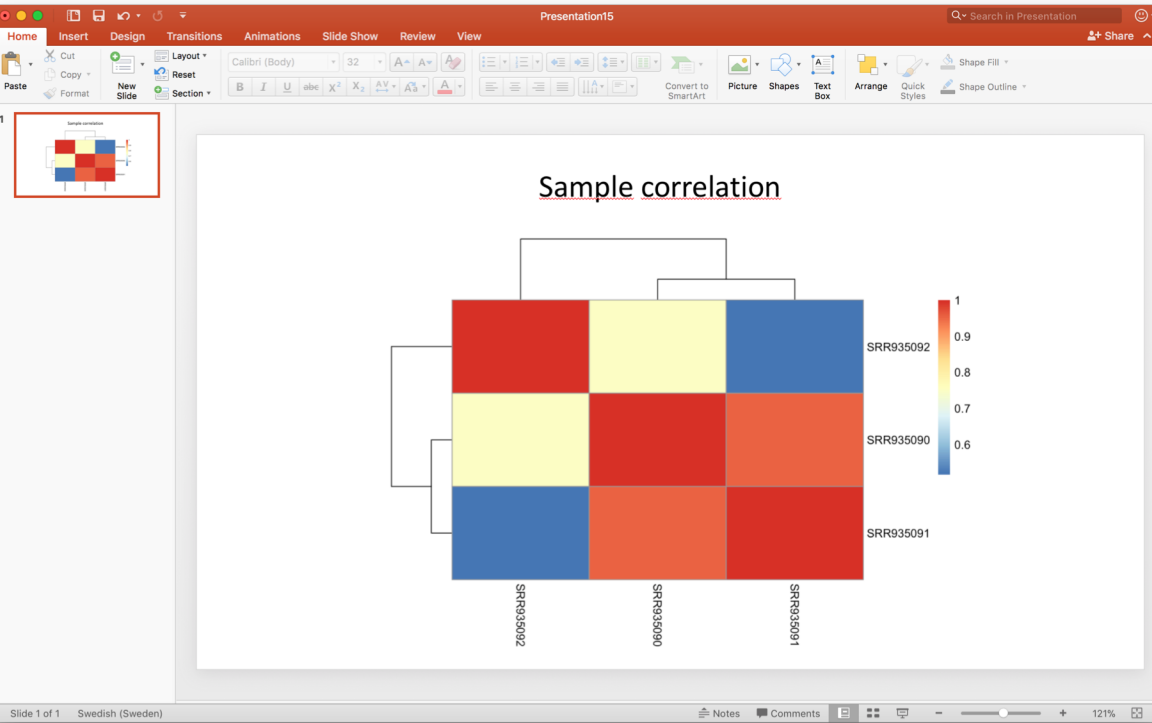
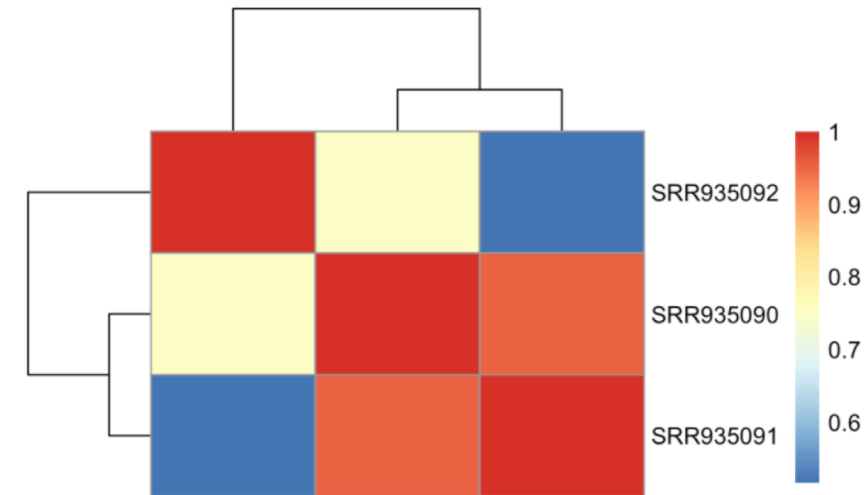
- SRR935090
- SRR935091
- SRR935092

```
# Read counts:
counts_file <- "results/tables/counts.tsv"
counts <- read.delim(counts_file, skip=1, header=F, row.names=1)
sample_names <- t(read.delim(counts_file, nrow=1, header=F))
colnames(counts) <- c("description", gsub(".*(SRR[0-9]*)\\.\\.\\.", "\\1", sample_names))
```

Plot sample correlation

Next, we will plot the sample correlation of the *count data*:

```
heatmap(cor((counts[, -1])), show_colnames=F)
```



- R Markdown makes your analysis more reproducible by connecting your code, figures and descriptive text.
- You can use it to make reproducible reports, rather than e.g. copy-pasting figures into a Word document.
- You can also use it as a notebook, in the same way as lab notebooks are used in a wet lab setting.

```
---
title: "Supplementary material"
author: John Doe, Joan Dough, Jan Doh, Dyon Do
date: "`r format(Sys.time(), '%d %B, %Y')`"
output: html_document
---
```

```
```{r, include=FALSE}
library("pheatmap")
```
```

```
# Read in the data
```

We have *count data* for three samples:

- SRR935090
- SRR935091
- SRR935092

```
```{r}
Read counts:
counts_file <- "results/tables/counts.tsv"
counts <- read.delim(counts_file, skip=1, header=F, row.names=1)
sample_names <- t(read.delim(counts_file, nrow=1, header=F))
colnames(counts) <- c("description",
 gsub(".*(SRR[0-9]*)\\.\\.\\.","\\1", sample_names))
```
```

```
# Plot sample correlation
```

Next, we will plot the sample correlation of the *count data*:

```
```{r, fig.height=3, fig.width=5}
pheatmap(cor((counts[, -1])), show_colnames=F)
```
```

Supplementary material

John Doe, Joan Dough, Jan Doh, Dyon Do

18 March, 2018

Read in the data

We have *count data* for three samples:

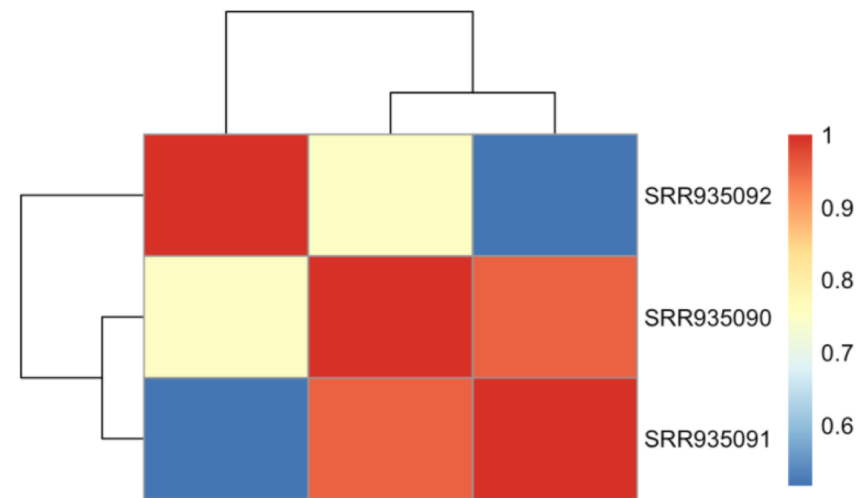
- SRR935090
- SRR935091
- SRR935092

```
# Read counts:
counts_file <- "results/tables/counts.tsv"
counts <- read.delim(counts_file, skip=1, header=F, row.names=1)
sample_names <- t(read.delim(counts_file, nrow=1, header=F))
colnames(counts) <- c("description", gsub(".*(SRR[0-9]*)\\.\\.\\.","\\1", sample_names))
```

Plot sample correlation

Next, we will plot the sample correlation of the *count data*:

```
pheatmap(cor((counts[, -1])), show_colnames=F)
```



```
```{r, fig.height=3, fig.width=5}
heatmap(cor((counts[,-1])), show_colnames=F)
```
```

- Document-wide options
- Output format
- Parameters

```
``{r, fig.height=3, fig.width=5}
heatmap(cor((counts[,,-1])), show_colnames=F)
````
```

- Evaluate R code and show output
- Also Bash, Python, Rcpp, SQL, Stan
- Chunk options

```
```{r, fig.height=3, fig.width=5}
pheatmap(cor((counts[,-1])), show_colnames=F)
```
```

- Freely add and format text using markdown

```

title: "Supplementary material"
author: John Doe, Joan Dough, Jan Doh, Dyon Do
date: "`r format(Sys.time(), '%d %B, %Y')`"
output: html_document

```

```
```{r, include=FALSE}  
library("pheatmap")  
```
```

```
Read in the data
```

We have *count data* for three samples:

- SRR935090
- SRR935091
- SRR935092

```
```{r}  
# Read counts:  
counts_file <- "results/tables/counts.tsv"  
counts <- read.delim(counts_file, skip=1, header=F, row.names=1)  
sample_names <- t(read.delim(counts_file, nrow=1, header=F))  
colnames(counts) <- c("description",  
                      gsub(".*(SRR[0-9]*)\\.\\.\\.","\\1", sample_names))  
```
```

```
Plot sample correlation
```

Next, we will plot the sample correlation of the *count data*:

```
```{r, fig.height=3, fig.width=5}  
pheatmap(cor((counts[, -1])), show_colnames=F)  
```
```

# Supplementary material

John Doe, Joan Dough, Jan Doh, Dyon Do

18 March, 2018

## Read in the data

We have *count data* for three samples:

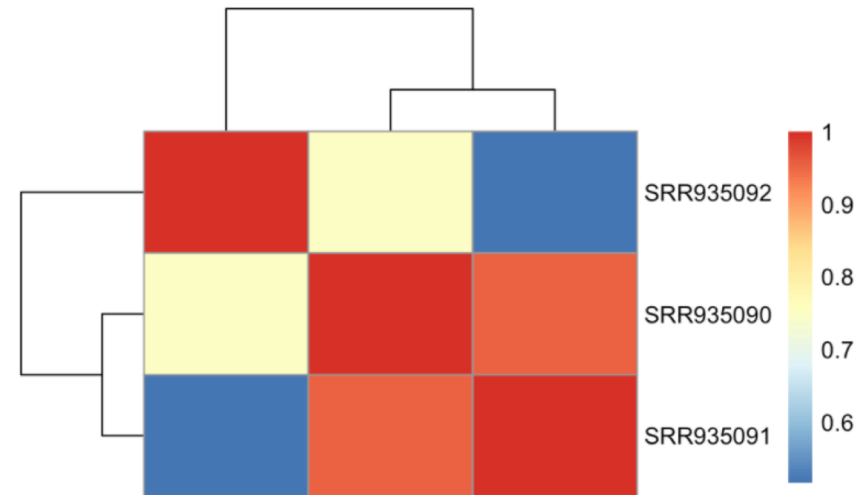
- SRR935090
- SRR935091
- SRR935092

```
Read counts:
counts_file <- "results/tables/counts.tsv"
counts <- read.delim(counts_file, skip=1, header=F, row.names=1)
sample_names <- t(read.delim(counts_file, nrow=1, header=F))
colnames(counts) <- c("description", gsub(".*(SRR[0-9]*)\\.\\.\\.","\\1", sample_names))
```

## Plot sample correlation

Next, we will plot the sample correlation of the *count data*:

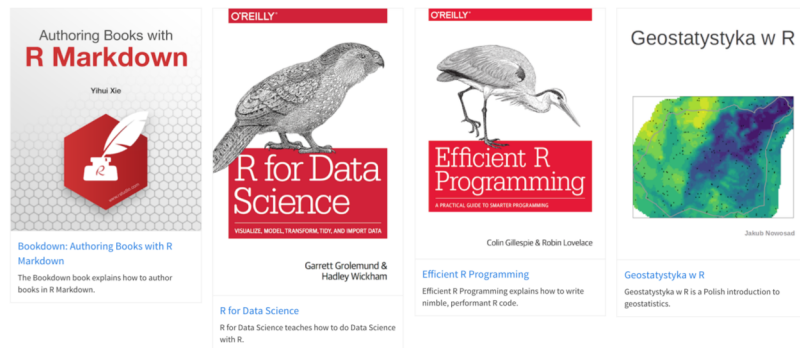
```
pheatmap(cor((counts[, -1])), show_colnames=F)
```



# Output formats

- Documents/reports (HTML, PDF, MS Word, Tufte handouts)
- Presentations (Powerpoint, Beamer, Slidy, ioslides, reveal.js)
- Interactive documents and dashboards (HTML widgets, Shiny)
- Books and websites
- Other templates...

Can require different markdown syntax depending on output!



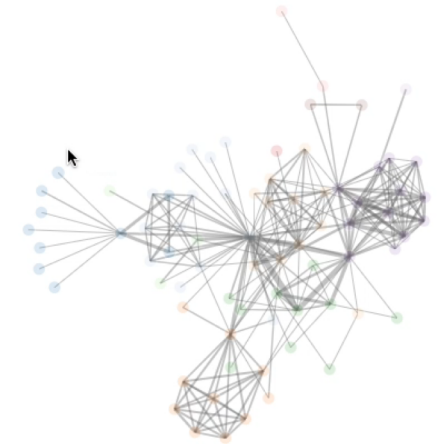
|                 |                                      |
|-----------------|--------------------------------------|
| Leaflet         | Geo-spatial mapping                  |
| dygraphs        | Time series charting                 |
| Plotly          | Interactive graphics with D3         |
| rbokeh          | R interface to Bokeh                 |
| Highcharter     | R interface to Highcharts            |
| visNetwork      | Graph data visualization with vis.js |
| networkD3       | Graph data visualization with D3     |
| d3heatmap       | Interactive heatmaps with D3         |
| DataTables      | Tabular data display                 |
| threejs         | 3D scatterplots and globes           |
| rglwidget       | Render scenes created with rgl       |
| DiagrammeR      | Diagrams and flowcharts              |
| MetricsGraphics | Scatterplots and line charts with D3 |

networkD3

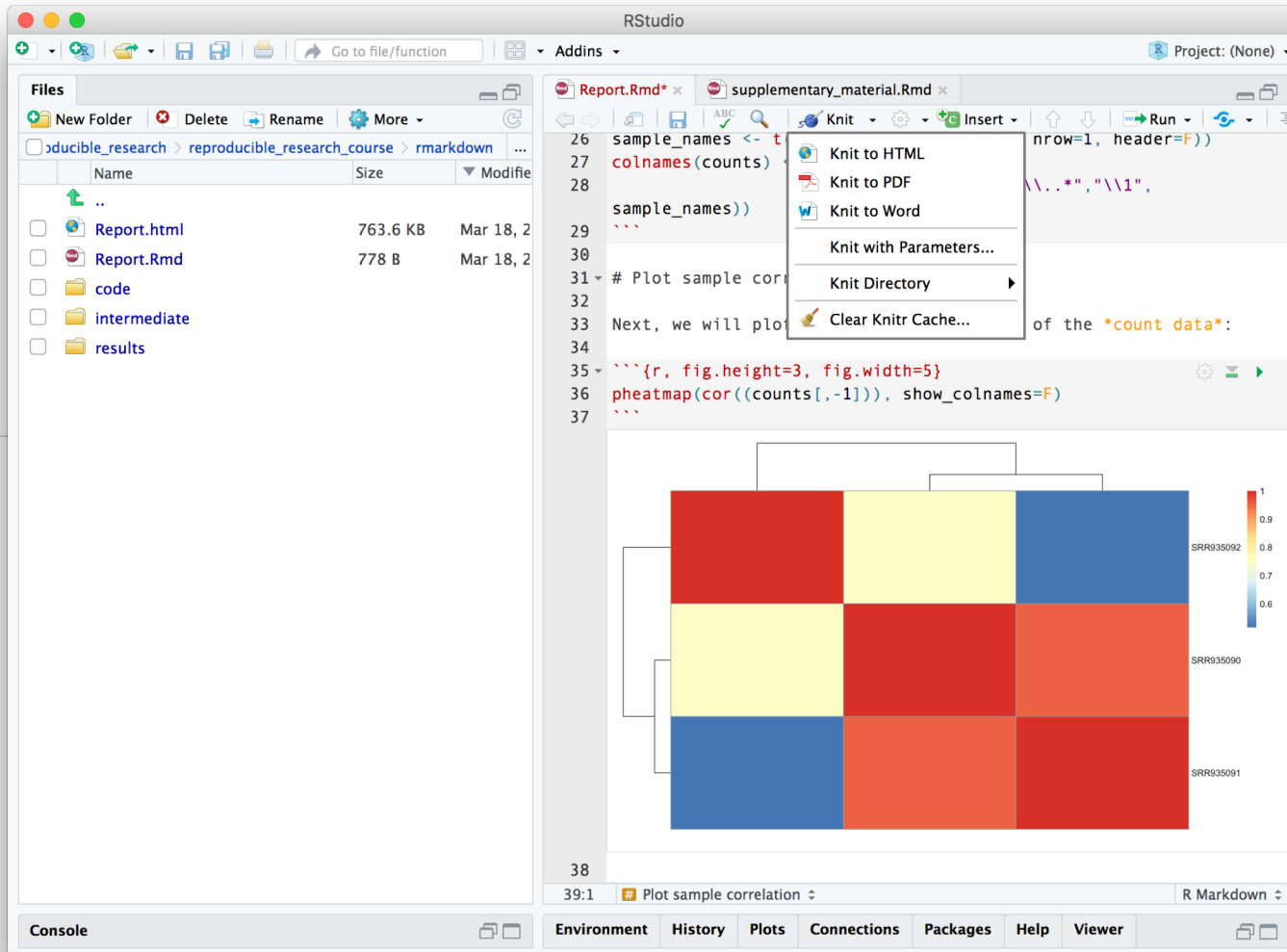
<http://christophergandrud.github.io/networkD3/>

networkD3 provides tools for creating D3 JavaScript network graphs from R.

```
library(networkD3)
data(MisLinks, MisNodes)
forceNetwork(Links = MisLinks, Nodes = MisNodes, Source = "source",
 Target = "target", Value = "value", NodeID = "name",
 Group = "group", opacity = 0.4)
```



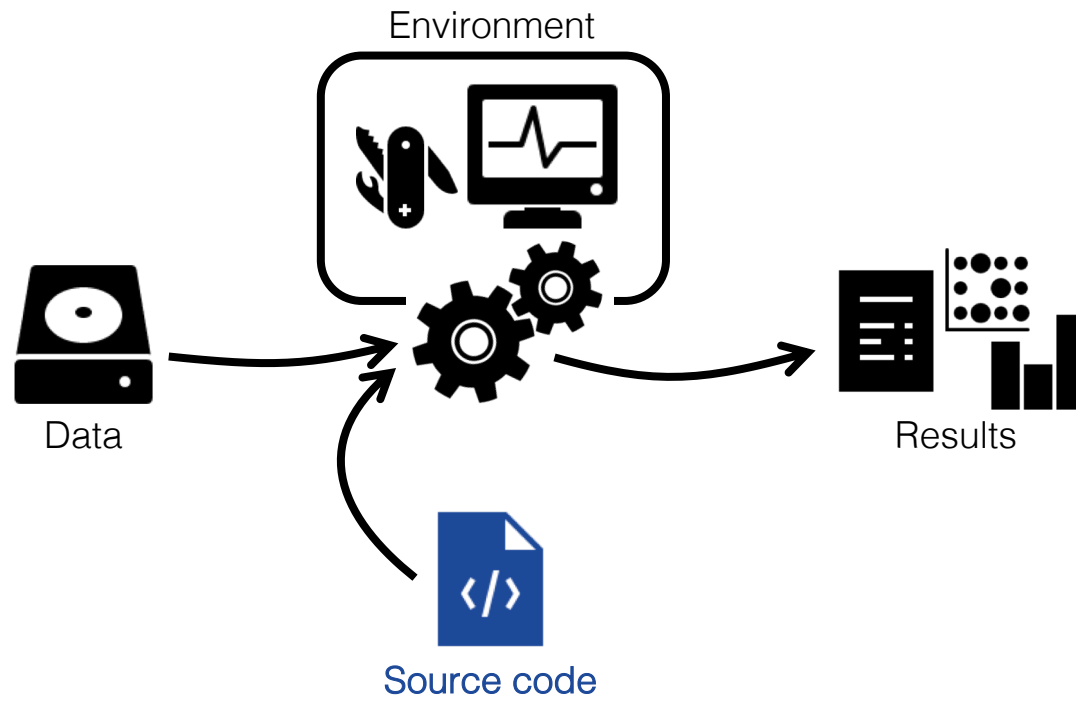
# R Markdown in RStudio



- Evaluate inline
- Render from menu
- Render from R console or terminal

```
$ R -e "rmarkdown::render('Report.Rmd')"
```





```
project
|- doc/
|
|- data/
| |- raw_external/
| |- raw_internal/
| |- meta/
|
|- code/
|- notebooks/
|
|- intermediate/
|- scratch/
|- logs/
|
|- results/
| |- figures/
| |- tables/
| |- reports/
|
|- Snakefile
|- config.yml
|- environment.yml
|- Dockerfile
```



Managing  
dependencies



Managing and executing  
analysis workflow



Versioning and  
collaborating on code  
(and some other files)



Connecting code  
and reporting

and...

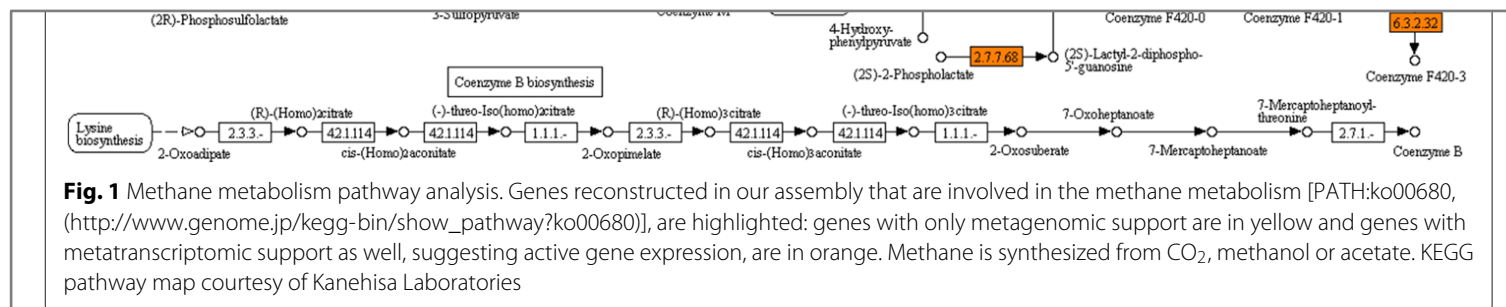


Isolating and exporting  
environment

Results should be possible to reproduce  
regardless of platform and with minimal effort.



“Docker provides a way to run applications securely  
isolated in a container, packaged with all its  
dependencies and libraries.”



## Discussion

We report extensive metagenomic and metatranscriptomic profiling of the microbial community from a production-scale biogas plant. Given the unprecedented sequencing depth and established bioinformatics, our data will be of great interest to the biogas research community in general and microbiologists working on biogas-producing microbial communities in particular. In a first applied study, our metagenome assembly was used to improve the characterization of a metaproteome generated from biogas plant fermentation samples and to investigate the metabolic activity of the microbial community [17].

By sharing our data, we want to actively encourage its reuse. This will hopefully result in novel biological and biotechnological insights, eventually enabling a more efficient biogas production.

## Availability of supporting data

### Data accession

Raw sequencing data are available in the European Nucleotide Archive (ENA) under study accession PRJEB8813 (<http://www.ebi.ac.uk/ena/data/view/PRJEB8813>). The datasets supporting the results of this article are available in *GigaScience's* GigaDB [2].

### Reproducibility

The complete workflow is organized in a single GNU Makefile and available on GitHub [18]. All data and results can be reproduced by a simple invocation of *make*. To further support reproducibility, we bundled all tools and dependencies into one Docker container available on DockerHub [19]. *docker run* executes the aforementioned Makefile inside the container. Reproduction

Bremges et al., "Deeply sequenced metagenome and metatranscriptome of a biogas-producing microbial community from an agricultural production-scale biogas plant", *GigaScience* (2015) 4:33, doi:10.1186/s13742-015-0073-6

```
$uname -a
```

```
Darwin dhcp-140-26.vpn.chalmers.se 15.6.0 Darwin Kernel Version 15.6.0:
Thu Sep 1 15:01:16 PDT 2016; root:xnu-3248.60.11~2/RELEASE_X86_64
x86_64
```

```
$docker pull ubuntu:16.04
```

```
16.04: Pulling from library/ubuntu
```

```
22dc81ace0ea: Pull complete
```

```
1a8b3c87dba3: Pull complete
```

```
91390a1c435a: Pull complete
```

```
07844b14977e: Pull complete
```

```
b78396653dae: Pull complete
```

```
Digest:
```

```
sha256:e348fbbbea0e0a0e73ab0370de151e7800684445c509d46195aef73e090a49bd6
```

```
Status: Downloaded newer image for ubuntu:16.04
```

```
$docker run -it ubuntu:16.04
```

```
root@407b0fd13fe5:/# uname -a
```

```
Linux 407b0fd13fe5 4.9.60-linuxkit-aufs #1 SMP Mon Nov 6 16:00:12 UTC
2017 x86_64 x86_64 x86_64 GNU/Linux
```

## Dockerfile

```
FROM ubuntu:16.04

Install prerequisites
RUN apt-get update && \
 apt-get install -y --no-install-recommends \
 bzip2 curl ca-certificates

Install Conda
RUN curl https://repo.continuum.io/miniconda.sh -O && \
 bash miniconda.sh -bf -p /opt/miniconda3/ && \
 rm miniconda.sh

Add conda to PATH
ENV PATH="/opt/miniconda3/bin:${PATH}"

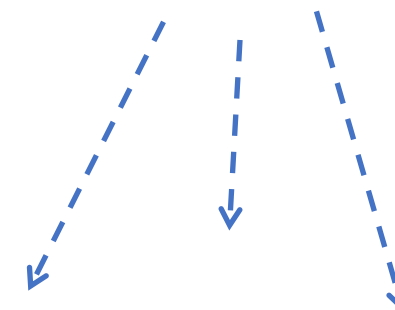
Install git and nano from conda-forge
RUN conda install -c conda-forge git nano

Use bash as shell
SHELL ["/bin/bash", "-c"]

Set workdir
WORKDIR /home
```

Build

Docker image



Docker container

Docker container

Docker container

# Mounting volumes

## Local project directory

```
project
|- doc/
|
|- data/
| |- raw_external/
| |- raw_internal/
| |- meta/
|
|- code/
|- notebooks/
|
|- intermediate/
|- scratch/
|- logs/
|
|- results/
| |- figures/
| |- tables/
| |- reports/
|
|- Snakefile
|- config.yml
|- environment.yml
|- Dockerfile
```

## Docker image file system

```
|- bin/
|- boot/
|- dev/
|- etc/
|- home/
| |- data/
| |- raw_external/
| |- raw_internal/
| |- meta/
|- lib/
|- lib64/
|- media/
|- opt/
|- proc/
|- root/
|- run/
|- sys/
|- tmp/
|- usr/
|- var/
```

```
$docker run -it -v $PWD/data:/home/data ubuntu:16.04
```

# What can I use Docker for?

As an advanced environment manager.



```
$docker run -it
-v $PWD:/home
my_image /bin/bash
```

To package your code with the environment it needs.



```
$docker run
-v $PWD/data:/home/data
-v $PWD/results:/home/results
my_image snakemake report.pdf
```

To package a whole workflow, e.g. to accompany a manuscript.



```
$docker run
-v $PWD/results:/home/results
my_image snakemake report.pdf
```

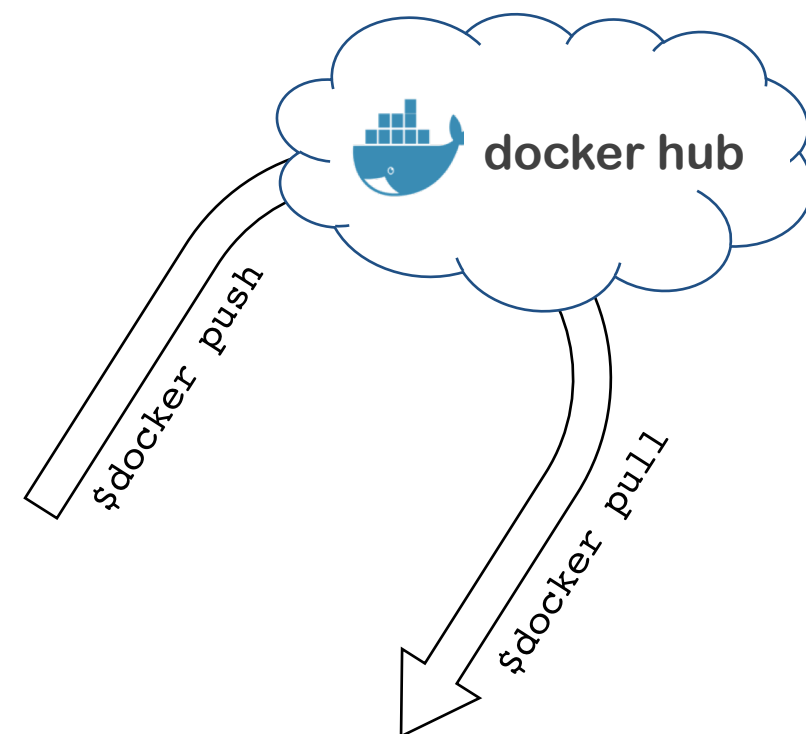
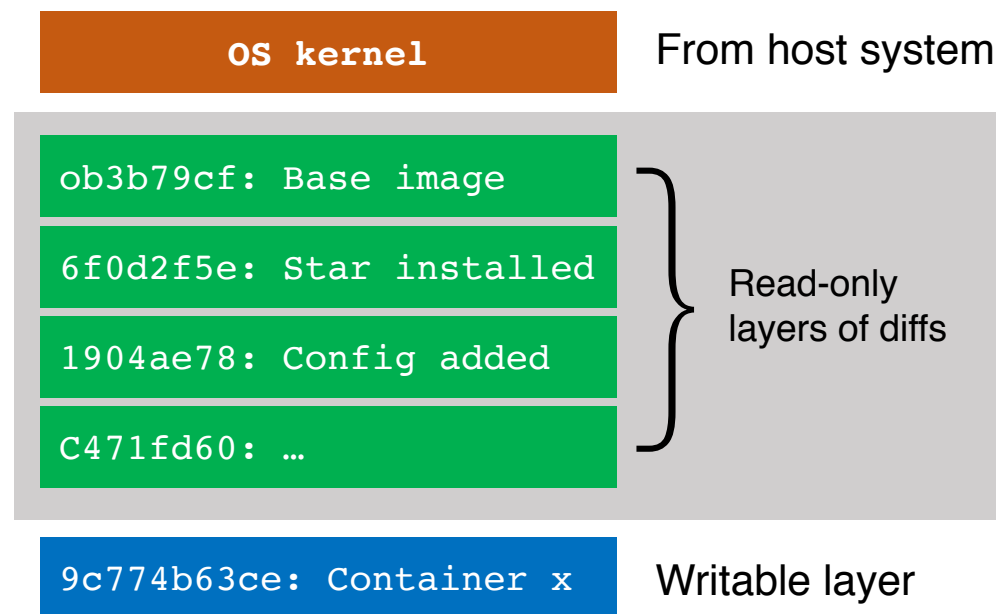


## A Docker image

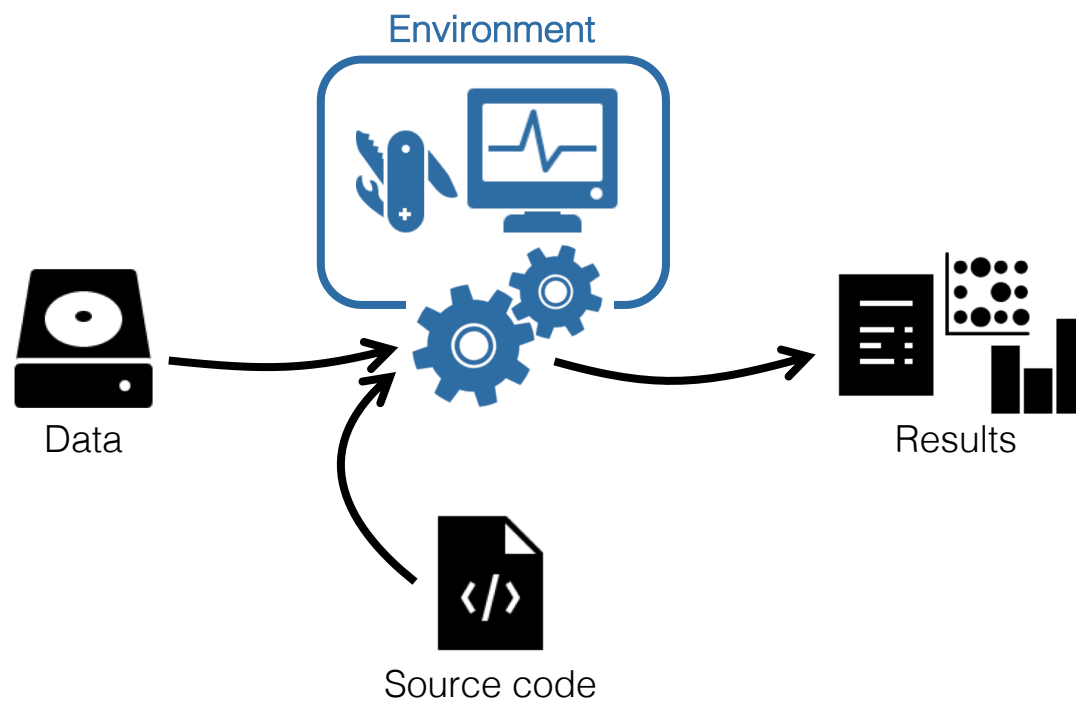
Dockerfile

```
FROM debian:latest
RUN conda install star
ADD config.yml
...
```

me\$ docker run image\_id

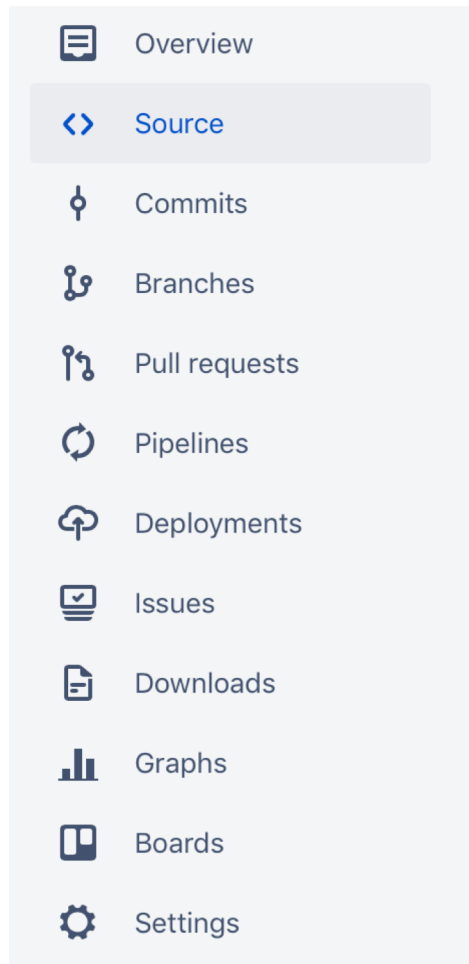


you\$ docker run image\_id

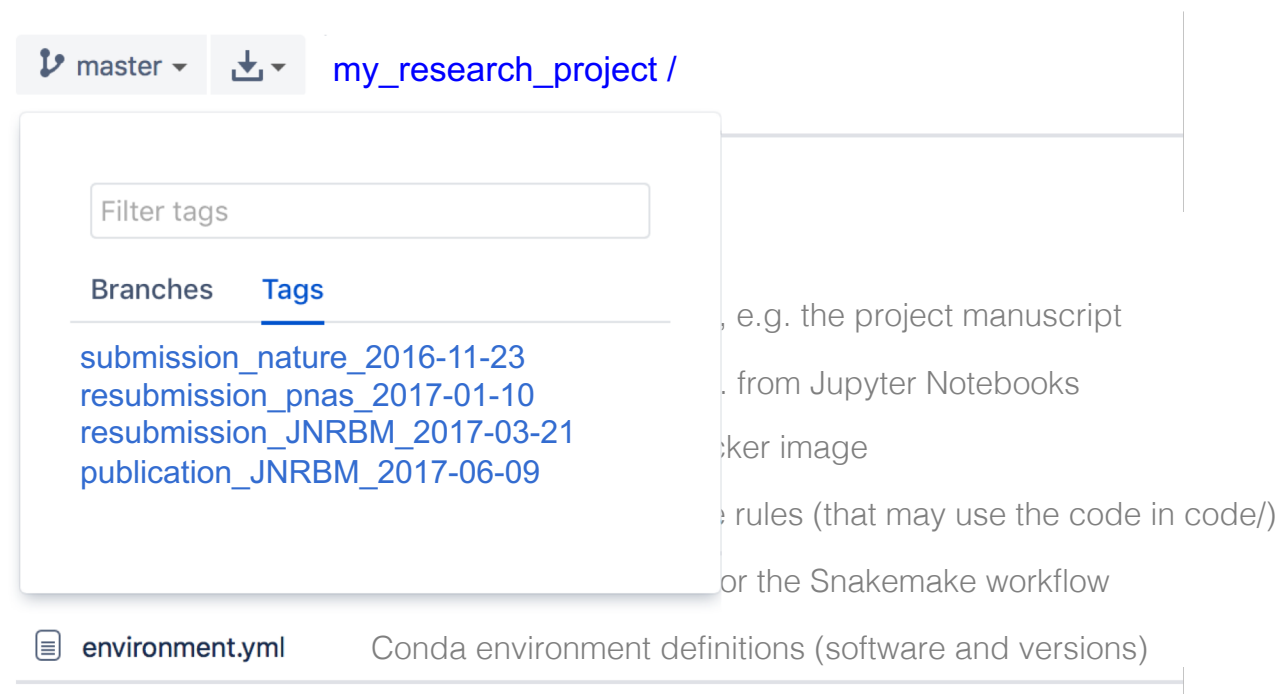


```
project
|- doc/
|
|- data/
| |- raw_external/
| |- raw_internal/
| |- meta/
|
|- code/
|- notebooks/
|
|- intermediate/
|- scratch/
|- logs/
|
|- results/
| |- figures/
| |- tables/
| |- reports/
|
|- Snakefile
|- config.yml
|- environment.yml
|- Dockerfile
```

Putting it all together



## my\_research\_project/



### Options for reproducing:

- Git clone and run workflow.
- Git clone, activate conda env, and run workflow.
- Git clone, docker build, and run workflow in container.
- Docker pull and run workflow in container.

# What is reasonable for your project?

Choose the right ambition level...

Minimal: code for reproducible results

 Snakemake

R Markdown

from  Studio

 jupyter

Good: versioned and structured repository

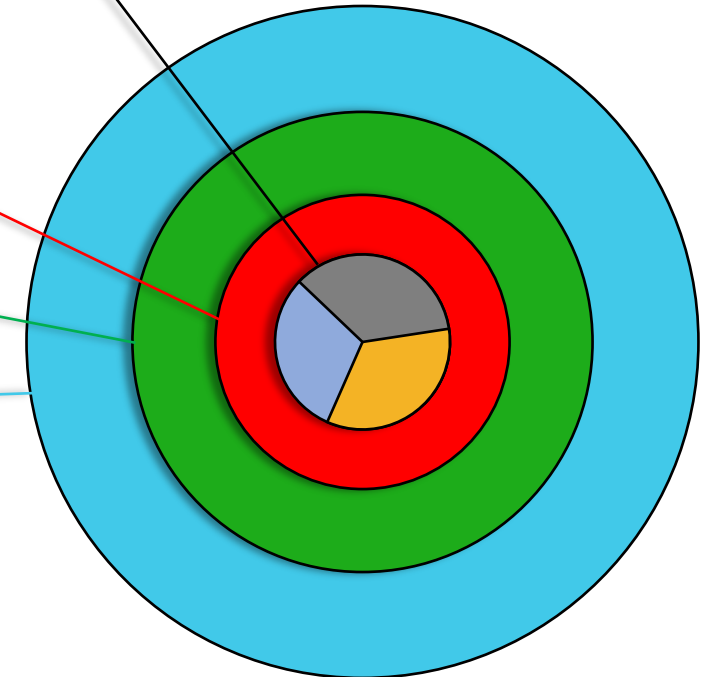
 git

Better: ambition to organize dependencies

 CONDA

Best: export everything!

 docker



# Reproducible research for bioinformatics projects

Leif Wigge (leif.wigge@scilifelab.se)

Rasmus Ågren (rasmus.agren@scilifelab.se)

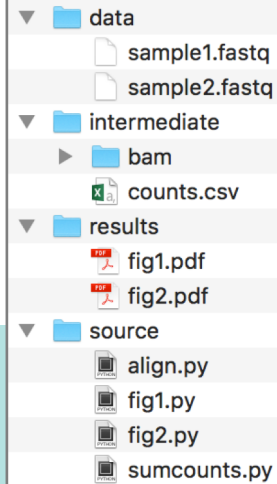
Bioinformatics long-term support (LTS)

## Everything can be a project

Divide your work into distinct projects and keep all files needed to go from raw data to final results in a dedicated directory with relevant subdirectories (see example).

Many software support the “project way of working”, e.g. Rstudio and the text editors Sublime Text and Atom.

**Tip!** Learn how to use git, a widely used system (both in academia and industry) for version controlling and collaborating on code.

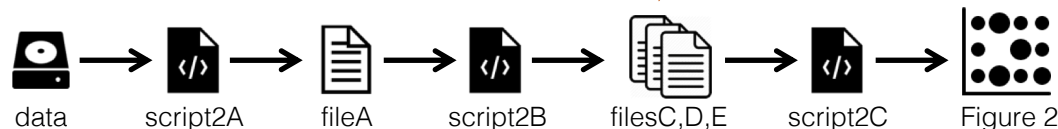


## Treasure your data

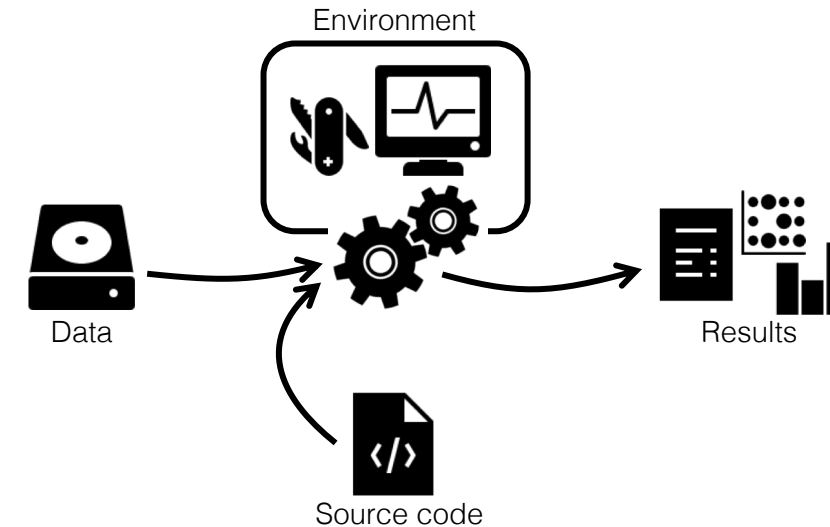
- Consider your input data static. Keep it read-only!
- Don't make *different* versions. If you need to preprocess it in any way, script it so you can recreate the steps (see box below).
- Backup! Keep redundant copies in different physical locations.
- Strive towards uploading it to its final destination already at the beginning of a project (e.g. specific repositories such as SRA, GEO, or GenBank, or general repositories such as Dryad or Figshare).

## Organize your coding

- Write scripts/functions/notebooks for specific tasks (connect raw data to final results)
- Keep parameters separate (e.g. top of file, or input arguments)



Avoid generating files interactively on the fly or doing things by hand (no way to track how they were made).



## Take control of your research by making it reproducible!

By moving towards a reproducible way of working you will quickly realize that you at the same time make your own life a lot easier! The added effort pays off by gain in control, organization and efficiency.

Below are all the components of a bioinformatics project that have to be reproducible.

## For the advanced

As projects grow, it becomes increasingly difficult to keep track of all the parts and how they fit together. Snakemake is a workflow management system that keeps track of how your files tie together, from raw data and scripts to final figures. If anything changes (script code, parameters, software version, etc) it will know what parts to rerun in order to have up to date and reproducible results.



## Connect your results with the code

Rmarkdown and Jupyter notebooks blur the boundaries between code and its output. They allow you to add non-code text (markdown) to your code. This generates a report containing custom formatted text, as well as figures and tables together with the code that generated them.

R Markdown

<http://rmarkdown.rstudio.com/> <http://jupyter.org/>



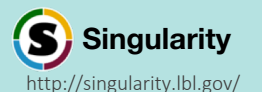
## Master your dependencies

- Full reproducibility requires the possibility to recreate the system that was originally used to generate the results.
- Conda is a package, dependency, and environment manager that makes it easy to install (most) software that you need for your project.
- Your environment can be exported in a simple text format and reinstalled by Conda on another system.



## For the advanced

- Conda cannot always *completely* recreate the system, which is required for proper reproducibility.
- A solution is to package your project in an isolated Docker container, together with all its dependencies and libraries.
- A vision is that every new bioinformatics publication is accompanied by a publically available Docker container!
- Singularity is an alternative to Docker which runs better on HPC clusters.



# alternatives

| Version control                                                                                                              | Environment / package managers                                                                                                         | Workflow managers                                                                                                               | Literate programming                                                                                                            | Containerization / virtualization                                                                                                           |
|------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Git</b> – Widely used and a lot of tools available + GitHub.                                                              | <b>Conda</b> – General purpose environment and package manager. Community-hosted collections of tools at bioconda or conda-forge.      | <b>Snakemake</b> – Based on Python, easily understandable format, relies on file names.                                         | <b>Jupyter</b> – Create and share notebooks in a variety of languages and formats by using a web browser.                       | <b>Docker</b> – Used for packaging and isolating applications in containers. Dockerhub allows for convenient sharing. Requires root access. |
| <b>Mercurial</b> – Distributed model just like Git, close to sourceforge.                                                    | <b>Pip</b> – Package manager for Python, has a large repository at pypi.                                                               | <b>Nextflow</b> – Based on Groovy, uses data pipes rather than file names to construct the workflow.                            | <b>Rmarkdown</b> – Developed by Rstudio, focuses on generating high-quality documents.                                          | <b>Singularity</b> – Simpler Docker alternative geared towards high performance computing. Does not require root.                           |
| <b>Subversion</b> – Centralized model unlike git/mercurial; no local repository on your computer and somewhat easier to use. | <b>Apt/yum/brew</b> – Native package managers for different OS. Integrated in OS and might deal with e.g. update notifications better. | <b>Make</b> – Used in software development and has been around since the 70s. Flexible but notoriously obscure syntax.          | <b>Zeppelin</b> – Developed by Apache. Closely integrated with Spark for distributed computing and Big Data applications.       | <b>Shifter</b> – Similar ambition as Singularity, but less focus on mobility and more on resource management.                               |
|                                                                                                                              | <b>Virtualenv</b> – Environment manager used to set up semi-isolated python environments.                                              | <b>Galaxy</b> - attempts to make computational biology accessible to researchers without programming experience by using a GUI. | <b>Beaker</b> – Newcomer based on Ipython, just as Jupyter. Has a focus on integrating multiple languages in the same notebook. | <b>VirtualBox/VMWare</b> – Virtualization rather than containerization. Less lightweight, but no reliance on host kernel.                   |