

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:

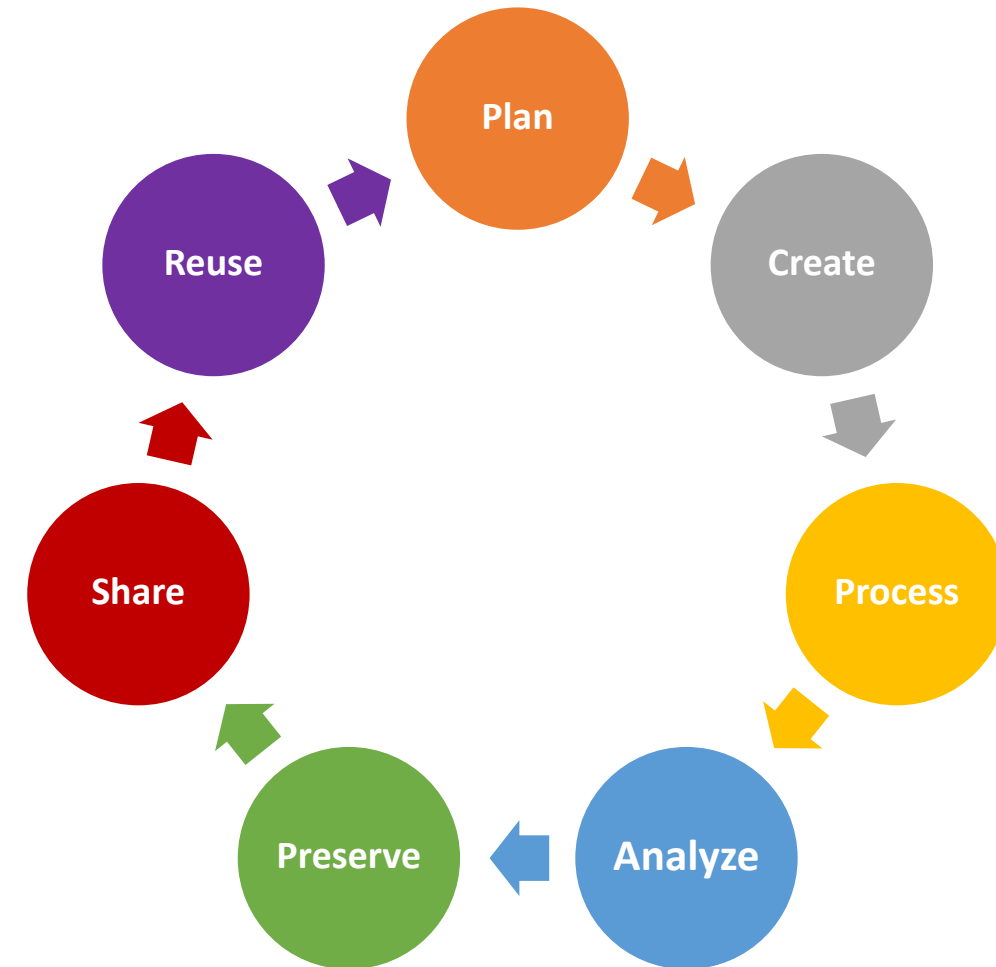
- R1. meta(data) are richly described with a plurality of accurate and relevant attributes
 - R1.1. (meta)data are released with a clear and accessible data usage license
 - R1.2. (meta)data are associated with detailed provenance
 - R1.3. (meta)data meet domain-relevant community standards



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
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.

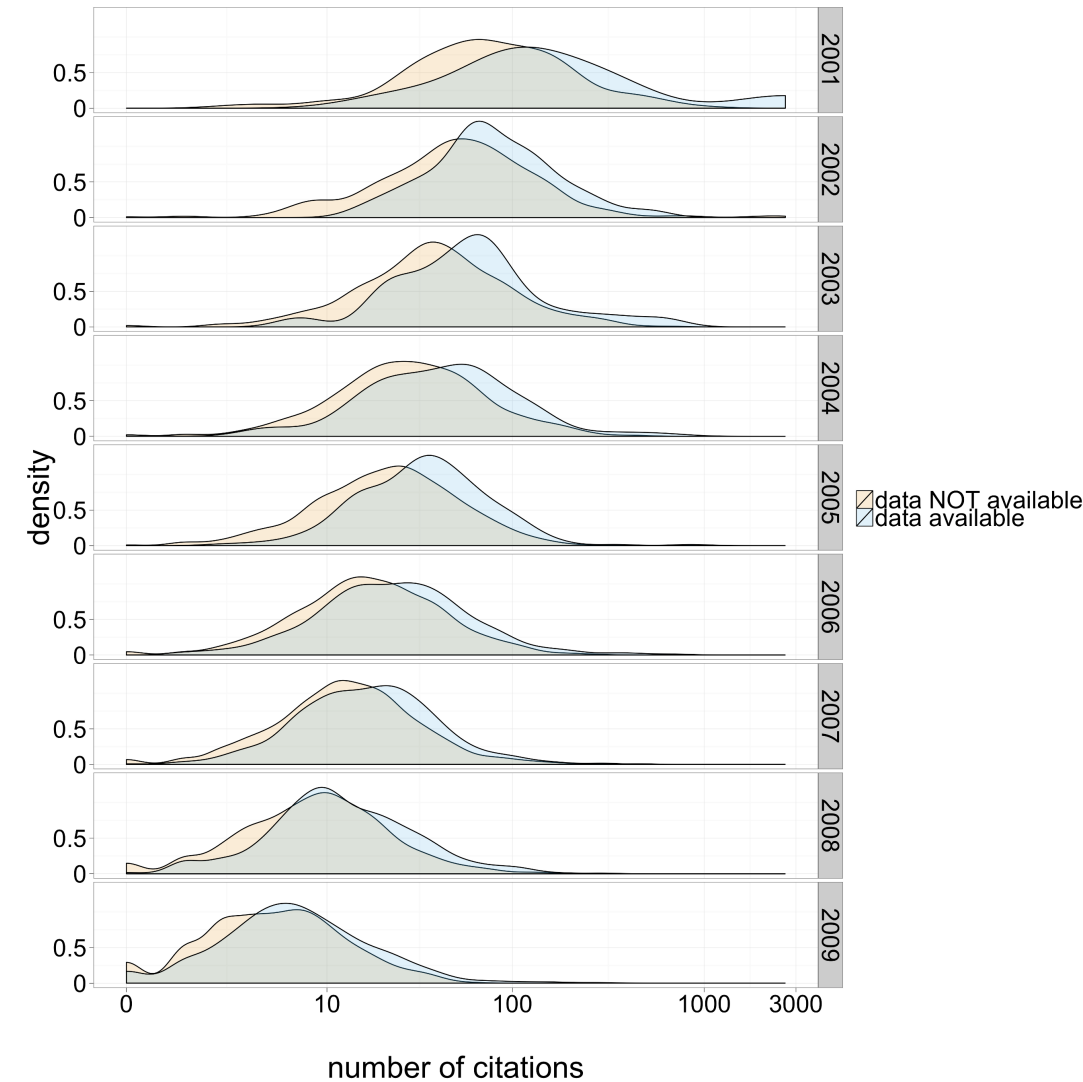
-	+
Binary	Text-based
Proprietary	Open
New kid on the block	Old as the hills
Compressed/encrypted	Uncompressed/unencrypted
Platform dependent	Interoperable
Complex	Simple

			
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.



Piowar 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

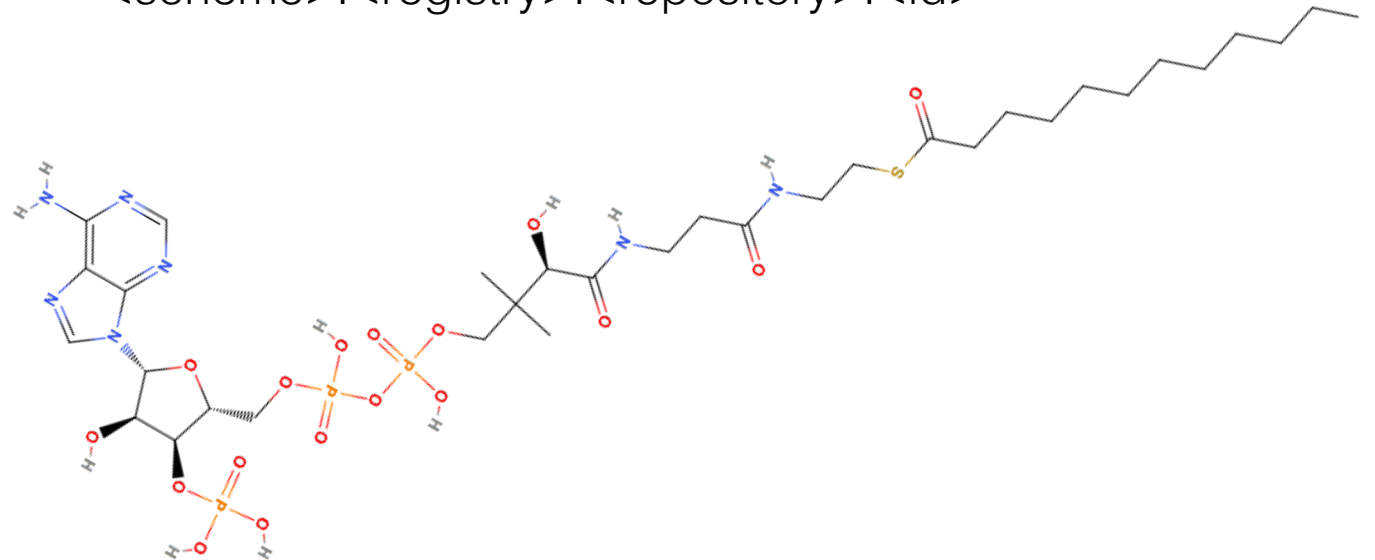
1-undecanecarboxylic acid CoA

vulvic acid CoA

Who am I?

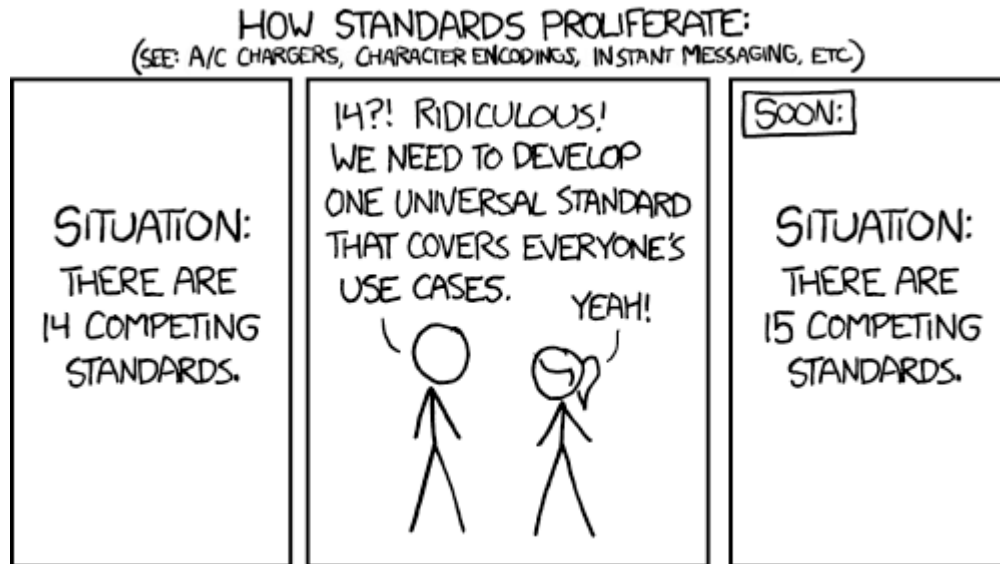
urn.miriam.chebi:15521 of course!

<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 Search

Advanced Search Search Wizard

Standards Databases Policies Collections/Recommendations

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

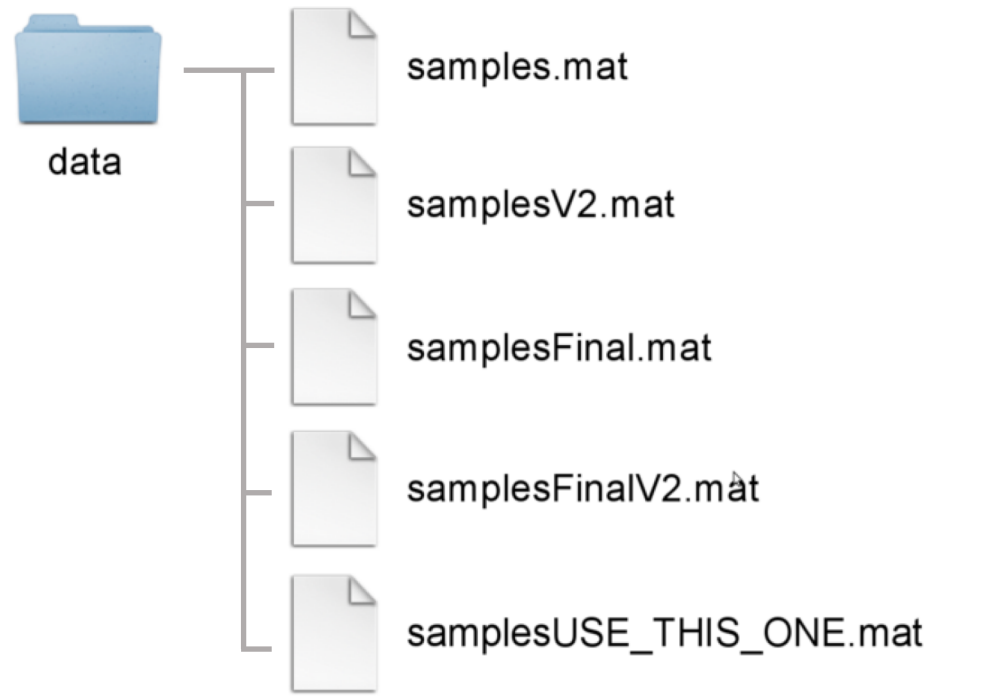
97 Policies

Funder	22
Journal	68
Society	3

View all

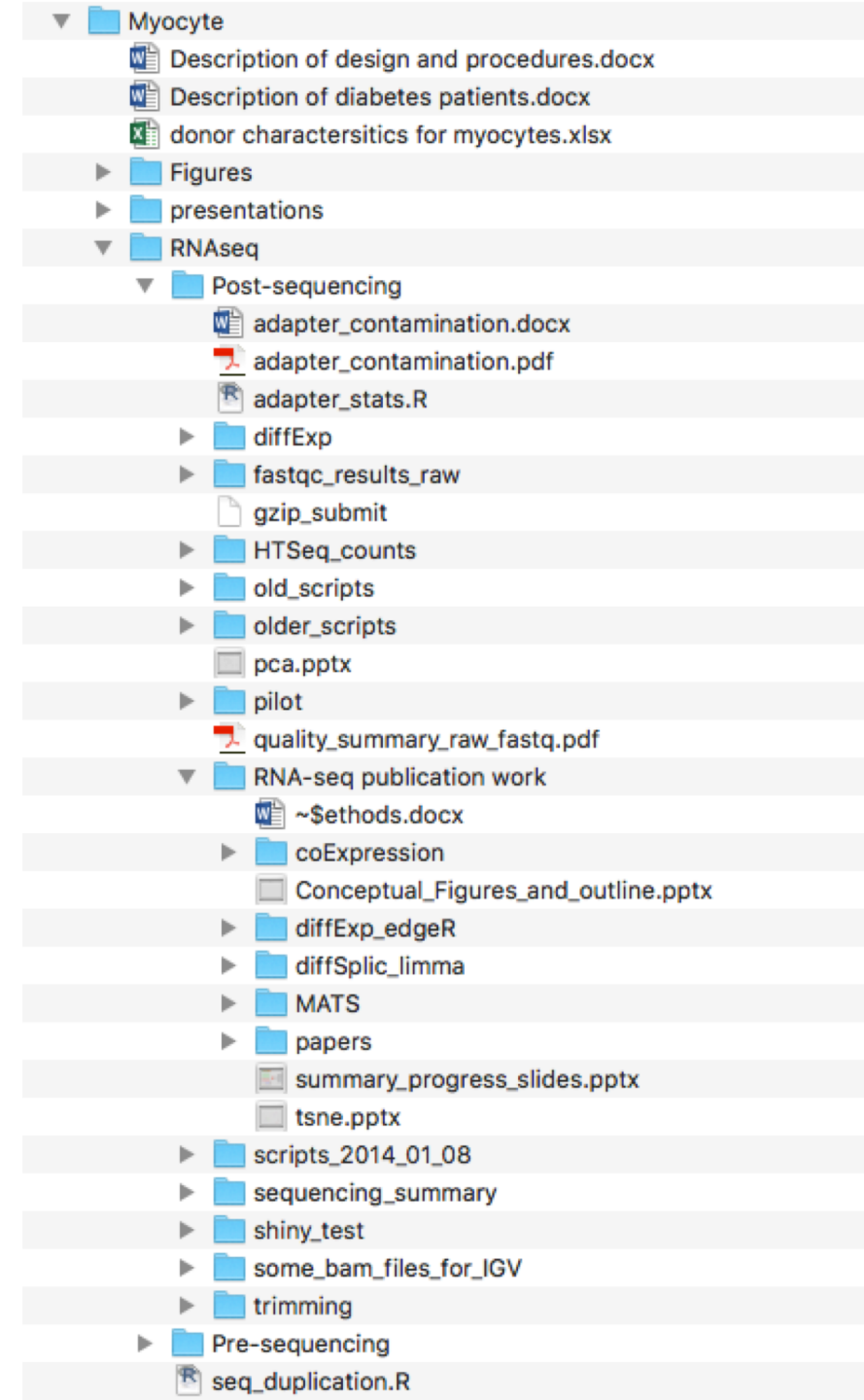
Project organization

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

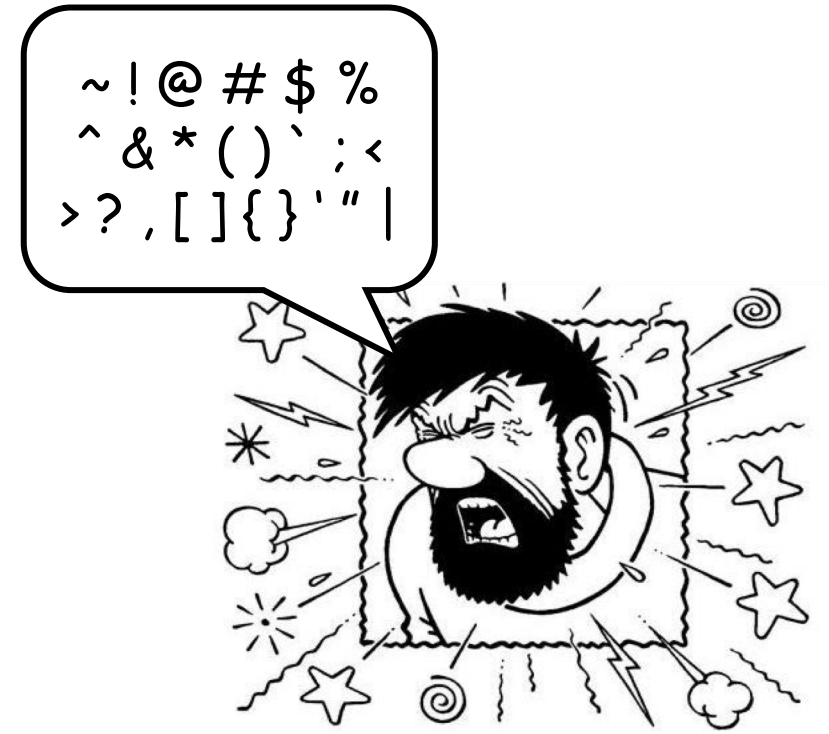
- Avoid special characters, e.g.: ~!@#\$%^&*() `; <>?, [] {} ' " |
- Avoid spaces, alternatives:
 - file_name.txt
 - file-name.txt
 - filename.txt
 - FileName.txt

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



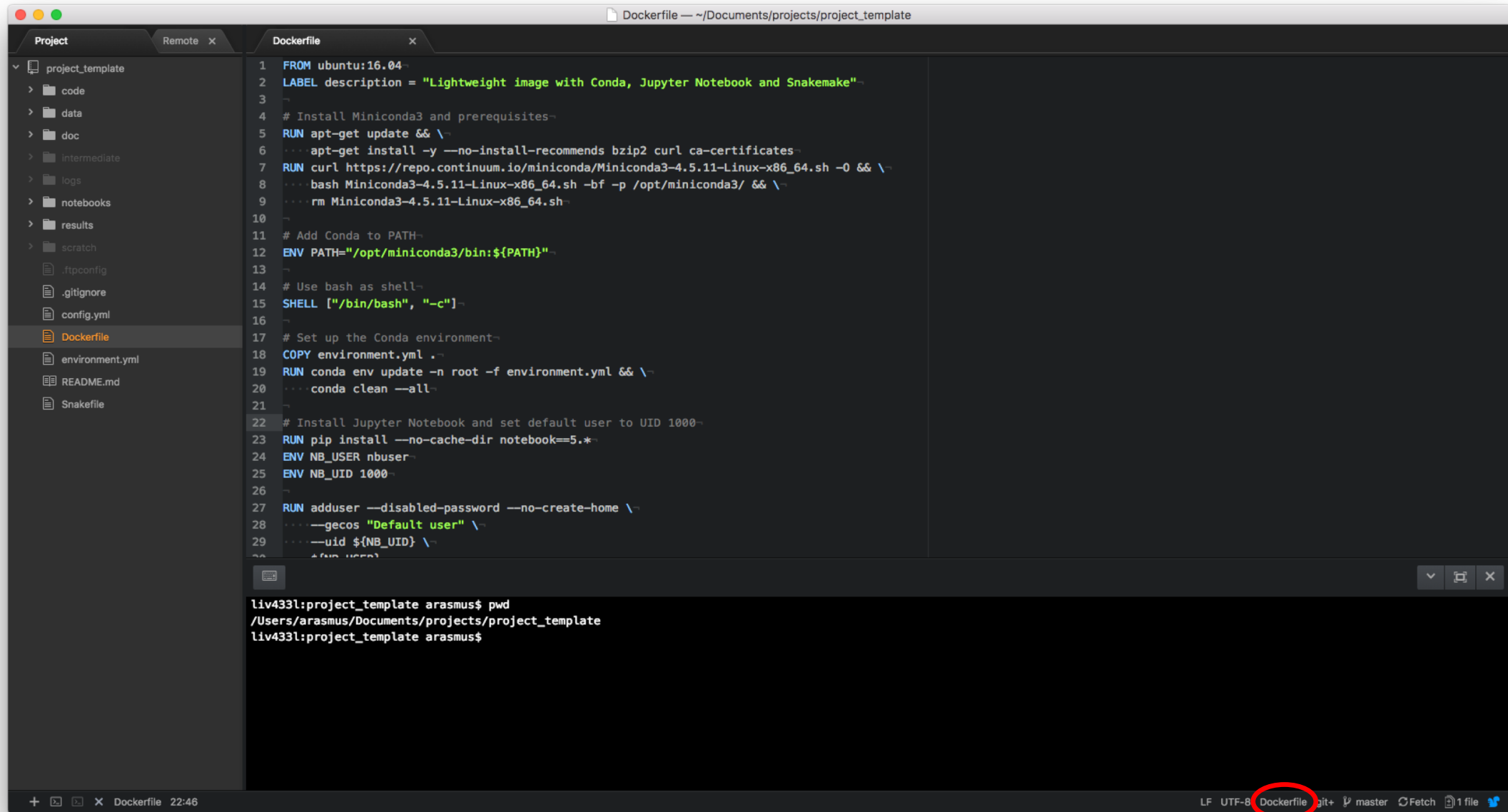
Bad examples:

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

Good examples:

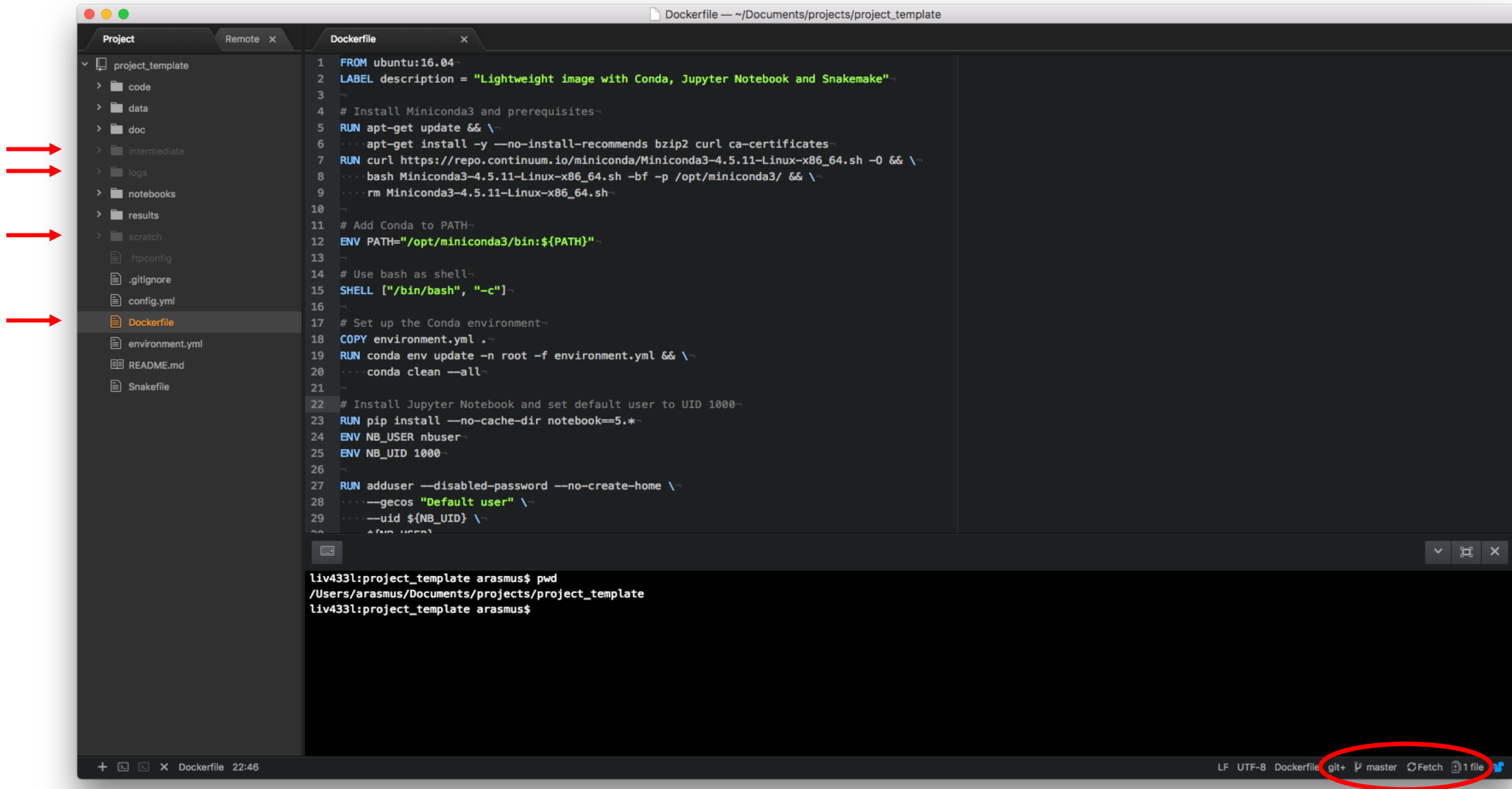
 2018-11-28_Gothenburg_Reproducible_research.pptx
 suppl_fig_09_barplot_alignment_stats.png
 suppl_fig_10_samples_PCA_count_data.png

A project in Atom



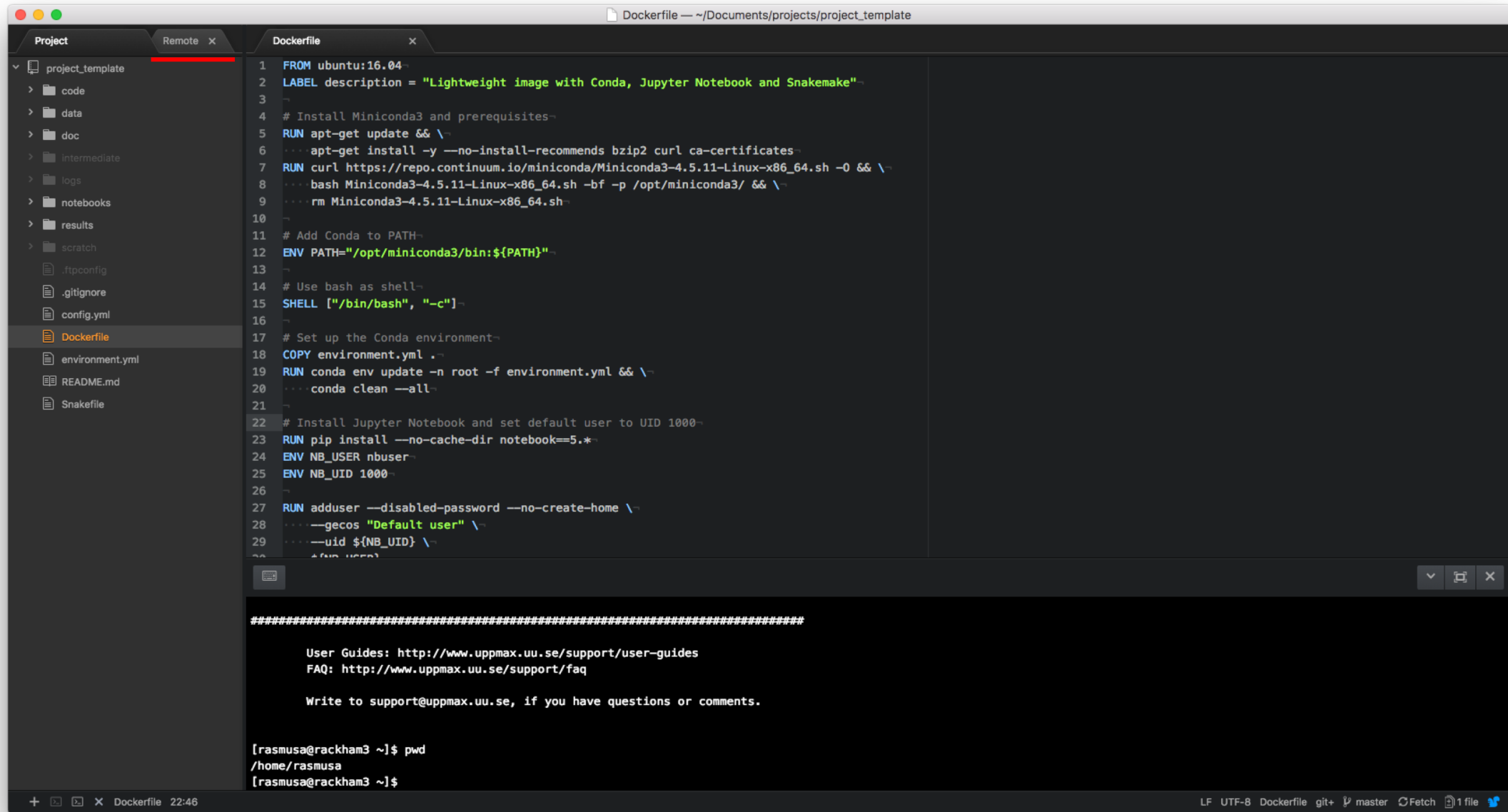
Syntax highlighting, indentation, and autocomplete

A project in Atom



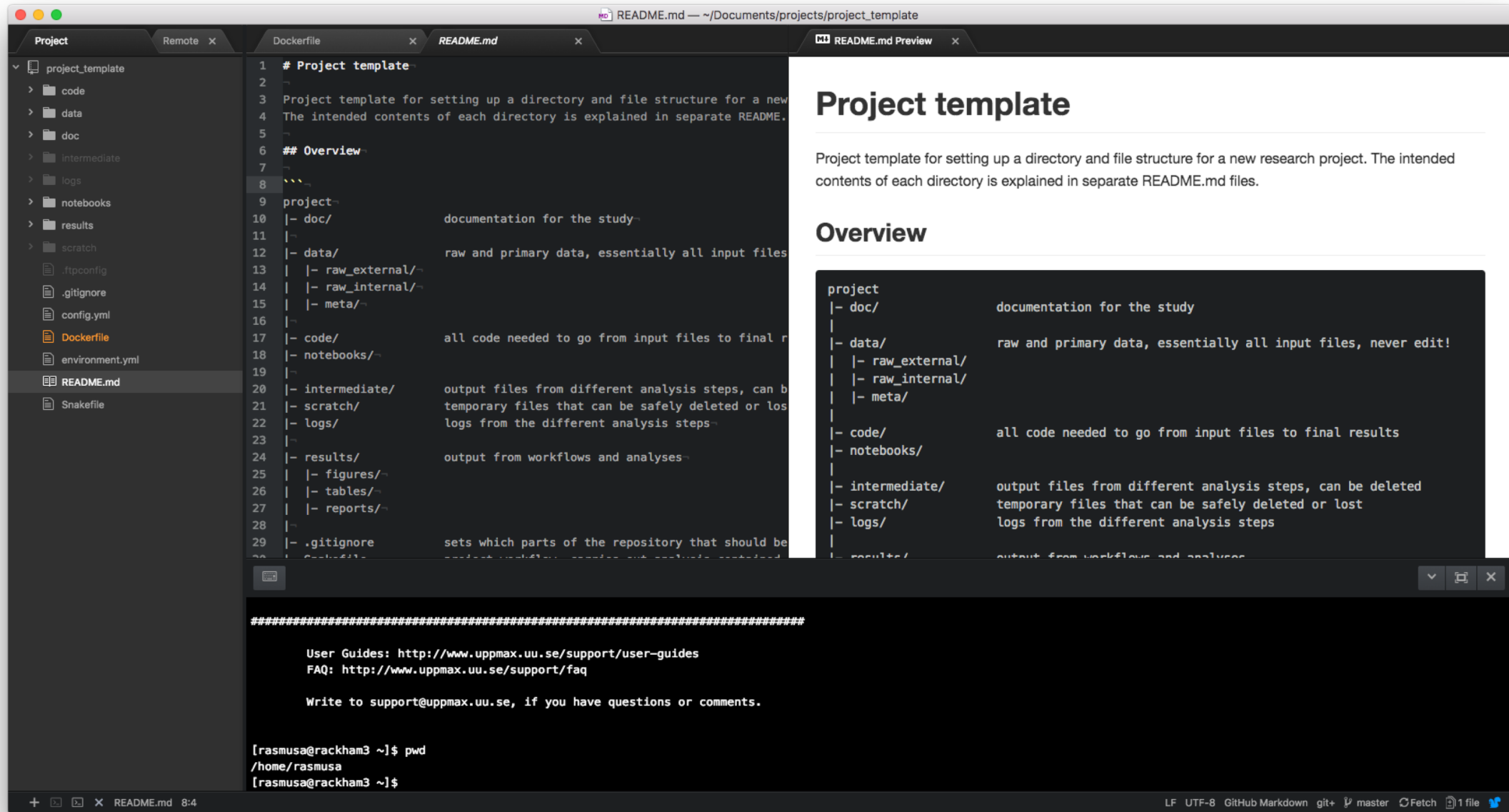
Integrated version control with Git

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, and the "project_template" dropdown menu is visible in the top right corner. The left pane shows the file explorer with a list of files and folders. The main editor pane displays the "README.md" file, which contains a project template structure. The bottom pane shows the terminal output of the "ls" command, listing the files and their sizes. The right pane shows a preview of the "README.md" file.

Files

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

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
30 ```
```

Terminal

```
[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'.

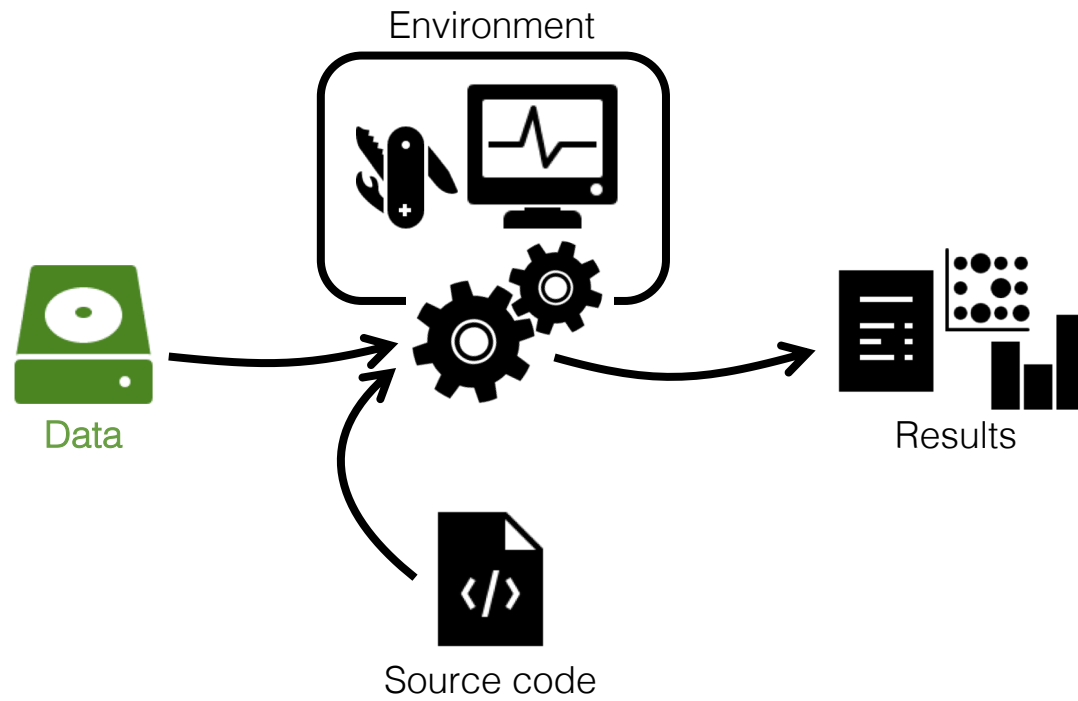
nothing to commit, working tree clean
```

Project template

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
```