# Bonnes pratiques pour organiser vos projets en bioinfo

DUBii 2021

Hélène Chiapello & Pierre Poulain

# Organisation des données

# Deux références

**OPEN ACCESS Freely available online**

PLoS COMPUTATIONAL BIOLOGY

**Education**

## A Quick Guide to Organizing Computational Biology Projects

William Stafford Noble[1,2]*

Noble, PLoS Comput Biol, 2009
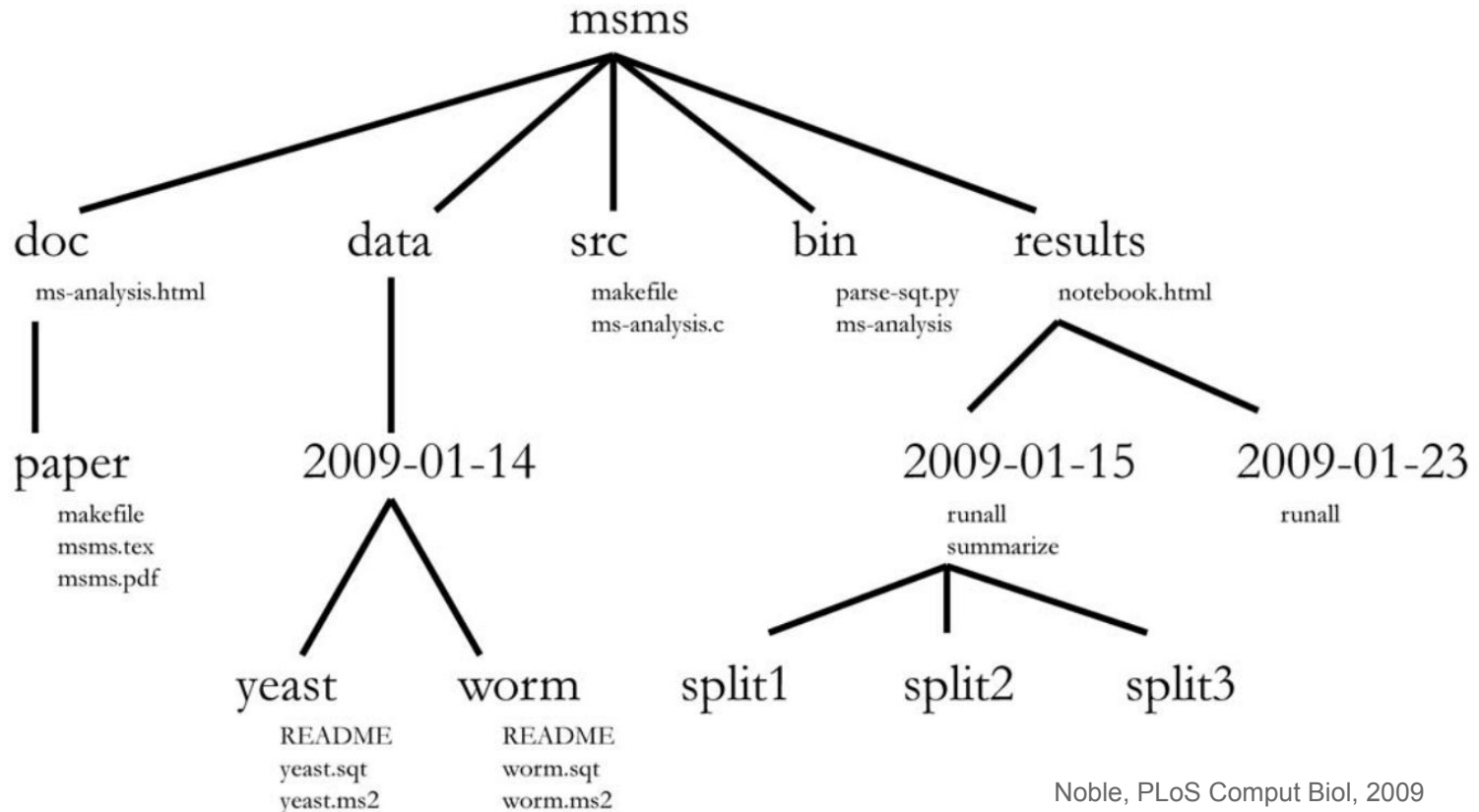DOI 10.1371/journal.pcbi.1000424

PLOS | COMPUTATIONAL BIOLOGY

PERSPECTIVE

## Good enough practices in scientific computing

Greg Wilson[1]*, Jennifer Bryan[2], Karen Cranston[3], Justin Kitzes[4], Lex Nederbragt[5], Tracy K. Teal[6]

Wilson, PLoS Comput Biol, 2017
DOI 10.1371/journal.pcbi.1005510

# Un exemple d'organisation



msms

doc — ms-analysis.html
data
src — makefile, ms-analysis.c
bin — parse-sqt.py, ms-analysis
results — notebook.html

doc → paper — makefile, msms.tex, msms.pdf

data → 2009-01-14 → yeast (README, yeast.sqt, yeast.ms2), worm (README, worm.sqt, worm.ms2)

results → 2009-01-15 (runall, summarize) → split1, split2, split3
results → 2009-01-23 (runall)

Noble, PLoS Comput Biol, 2009
DOI 10.1371/journal.pcbi.1000424

# Noms de fichiers et répertoires

**Pas d'espace**

_ ou - pour séparer les « mots »

Ex : `new_test`, `dubii-python`


Pas de caractères spéciaux

# Format de date

## ISO 8601 ?

# Format de date



So what's your idea of a perfect date?

YYYY-MM-DD

I find other formats a bit confusing

Mahdi Yusuf / @myusuf3
https://twitter.com/myusuf3/status/865722106071453696



PUBLIC SERVICE ANNOUNCEMENT:

OUR DIFFERENT WAYS OF WRITING DATES AS NUMBERS CAN LEAD TO ONLINE CONFUSION. THAT'S WHY IN 1988 ISO SET A GLOBAL STANDARD NUMERIC DATE FORMAT.

THIS IS *THE* CORRECT WAY TO WRITE NUMERIC DATES:

2013-02-27

THE FOLLOWING FORMATS ARE THEREFORE DISCOURAGED:

02/27/2013   02/27/13   27/02/2013   27/02/13
20130227   2013.02.27   27.02.13   27-02-13
27.2.13   2013. II. 27.   $27/2$-13   2013.158904109
MMXIII-II-XXVII   MMXIII $\frac{LVII}{CCCLXV}$   1330300800
$((3+3)\times(111+1)-1)\times 3/3-1/3^3$   2013
10/11011/1101   02/27/20/13

XKCD, ISO 8601
https://xkcd.com/1179/

7

# Un autre exemple d'organisation

```
Box 3. Project layout

.
|-- CITATION
|-- README
|-- LICENSE
|-- requirements.txt
|-- data
|    |-- birds_count_table.csv
|-- doc
|    |-- notebook.md
|    |-- manuscript.md
|    |-- changelog.txt
|-- results
|    |-- summarized_results.csv
|-- src
|    |-- sightings_analysis.py
|    |-- runall.py
```

```
.
|-- project_name
|    |-- current
|    |    |-- ...project content as described earlier...
|    |-- 2016-03-01
|    |    |-- ...content of 'current' on Mar 1, 2016
|    |-- 2016-02-19
|    |    |-- ...content of 'current' on Feb 19, 2016
```

Wilson, PLoS Comput Biol, 2017
DOI 10.1371/journal.pcbi.1005510

# Gestion des données

# Gestion des données

Source : PhD Comics

# Gestion des données : git / GitHub

- Garder une mémoire des modifications de fichiers
- Travailler collaborativement
- Partager des fichiers

- Git est un logiciel
- GitHub est un site internet (une plateforme d'échange)

# Gestion des données : git / GitHub



Sandve, PLOS Comput Biol, 2013
DOI 10.1371/journal.pcbi.1003285

J. Perkel, Nature Index, 2018

# Gestion des données : git / GitHub

Débuter avec Git et Github en 30 min



https://www.youtube.com/watch?v=hPfgekYUKgk
La capsule, 2017

D'autres ressources :

- https://cupnet.net/git-github/

- https://swcarpentry.github.io/git-novice/

intro : 10/03 après-midi 😆

# Quelques conseils

# Quelques conseils

This leads to the second principle, which is actually more like a version of Murphy's Law: Everything you do, you will probably have to do over again. Inevitably, you will discover some flaw in your initial preparation of the data being analyzed, or you will get access to new data, or you will decide that your parameterization of a particular model was not broad enough. This means that the experiment you did last week, or even the set of experiments you've been working on over the past month, will probably need to be redone. If you have organized

# Quelques conseils

**Record all the steps used to process data** (1e). Data manipulation is as integral to your analysis as statistical modeling and inference. If you do not document this step thoroughly, it is impossible for you or anyone else to repeat the analysis.

The best way to do this is to write scripts for *every* stage of data processing. This might feel frustratingly slow, but you will get faster with practice. The immediate payoff will be the ease with which you can redo data preparation when new data arrive. You can also reuse data

# Des conseils, encore !

Adopter des pratiques **robustes** et **reproductibles**

- Code
  - Lisible
  - Documenté
  - Utiliser des librairies existantes dès que c'est possible
  - Versionné et partagé
- Données
  - Versioning
  - Plans de Gestion de Données (PGD)
- Code + données + résultats
  - Gestionnaires de workflows
  - Notebooks