

# Git, Emacs org-mode and mlflow: making applied machine learning research fully reproducible

An extension of Stanislav et al, 2015

Alex Seltmann <sup>1,2</sup>

<sup>1</sup>Institute of Applied Optic and Biophysics, Friedrich-Schiller University Jena, Max-Wien-Platz 1, 07743 Jena, Germany

<sup>2</sup>Leibniz-Institute of Photonic Technologies, Albert-Einstein Strasse 9, 07745 Jena, Germany

June 21, 2022

# Scope

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

## Scope

Computational research life  
cycle

Definitions: reproducibility

## Tools

Org-mode

Mlflow

Git

Reproducible Workflow

Discussion

# The computational research life cycle<sup>[1]</sup>

Git, org-mode and mflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

## Individual exploration

- = single investigator tests idea, algorithm, question with small-scale dataset / simulation
- tools: Excel, Matlab, Mathematica, Sage, R, SPSS, ...

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

Mflow

Git

Reproducible Workflow

Discussion

---

<sup>[1]</sup>Millman Pérez (2018)

# The computational research life cycle<sup>[1]</sup>

Git, org-mode and mflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

## Individual exploration

- = single investigator tests idea, algorithm, question with small-scale dataset / simulation
- tools: Excel, Matlab, Mathematica, Sage, R, SPSS, ...

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

Mflow

Git

Reproducible Workflow

Discussion

## Collaboration

- = bring together complementary expertise from colleagues
- tools: email, VCS, Dropbox, Github, (paper-final-v2-REALLY-FINAL-john.doc)

---

<sup>[1]</sup>Millman Pérez (2018)

# The computational research life cycle<sup>[1]</sup>

## Individual exploration

- = single investigator tests idea, algorithm, question with small-scale dataset / simulation
- tools: Excel, Matlab, Mathematica, Sage, R, SPSS, ...

## Collaboration

- = bring together complementary expertise from colleagues
- tools: email, VCS, Dropbox, Github, (`paper-final-v2-REALLY-FINAL-john.doc`)

## Production-scale execution

- = large data sets, complex simulations, supercomputers...
- tools: compiled code (C, C++, ...) and parallel computing libraries (MPI, Hadoop)

Git, org-mode and mflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

Mflow

Git

Reproducible Workflow

Discussion

---

<sup>[1]</sup>Millman Pérez (2018)

# The computational research life cycle<sup>[1]</sup>

## Individual exploration

- = single investigator tests idea, algorithm, question with small-scale dataset / simulation
- tools: Excel, Matlab, Mathematica, Sage, R, SPSS, ...

## Collaboration

- = bring together complementary expertise from colleagues
- tools: email, VCS, Dropbox, Github, (`paper-final-v2-REALLY-FINAL-john.doc`)

## Production-scale execution

- = large data sets, complex simulations, supercomputers...
- tools: compiled code (C, C++, ...) and parallel computing libraries (MPI, Hadoop)

## Publication & Education

- = paper, internal report, visualization → share with students and colleagues, cycle starts again
- tools: LATEX, Google Docs, Word, PowerPoint

<sup>[1]</sup>Millman Pérez (2018)

# The computational research life cycle<sup>[1]</sup>

Individual exploration

.

Collaboration

.

Production-scale execution

.

Publication & Education

.

Goal: Use one setup for everything! Reduce manual data transfer!

<sup>[1]</sup>Millman Pérez (2018)

# What are we talking about?<sup>[2]</sup>

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

Mlflow

Git

Reproducible Workflow

Discussion

## Methods reproducibility

- = get *same* results, use *same* data and tools
- topics: provide study protocols, reusable (meta)data, code, results, ...

Today: tools for methods reproducibility

---

<sup>[2]</sup>Goodman et al (2016)



# Tools

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

## Scope

Computational research life  
cycle

Definitions: reproducibility

## Tools

Org-mode

Mlflow

Git

Reproducible Workflow

Discussion

# What is org-mode

- **plain text-based** tool for outlining, note-taking, spreadsheets, project planning, ...

```
File Edit Options Buffers Tools Org Tbl Text Help
[Icons] [Icons] [Icons] [Icons] [Icons] [Icons] [Icons] [Icons] [Icons]
** Prepare the trip
*** Plan it
| Day | Main item |
|-----|
| 1 | Flight to berlin |
| 2 | Explore city |
| 3 | Flight back |

*** Essential purchases [2/4]
1. [X] tickets [1/2]
   + [X] to Berlin [2011-06-28 Tue]
   + [ ] from Berlin [2011-06-30 Thu]
2. [-] Insurance
3. [ ] Book hotel
4. [X] Panama hat

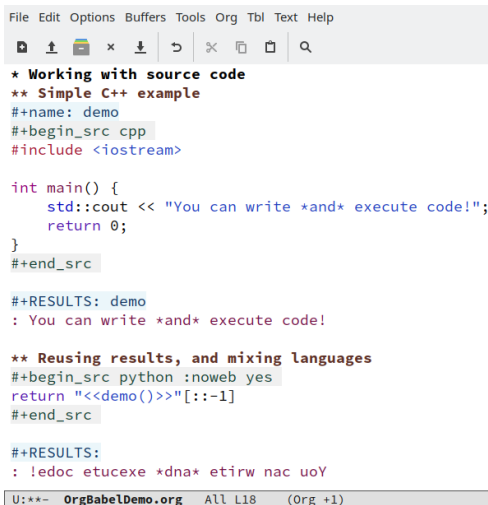
** The trip itself
*** Taxi to airport
We looked at Wikipedia Berlin, found some interesting
facts and *places to visit*.

U:**- TripPlanning.org All L20 (Org +1)
```

[3]

# What is org-mode

- **plain text-based** tool for outlining, note-taking, spreadsheets, project planning, ...
- **literate programming** via notebook-like environment for >70 programming languages



```
File Edit Options Buffers Tools Org Tbl Text Help
[Icons]
* Working with source code
** Simple C++ example
#+name: demo
#+begin_src cpp
#include <iostream>

int main() {
    std::cout << "You can write *and* execute code!";
    return 0;
}
#+end_src

#+RESULTS: demo
: You can write *and* execute code!

** Reusing results, and mixing languages
#+begin_src python :noweb yes
return "<<demo()>>"[::-1]
#+end_src

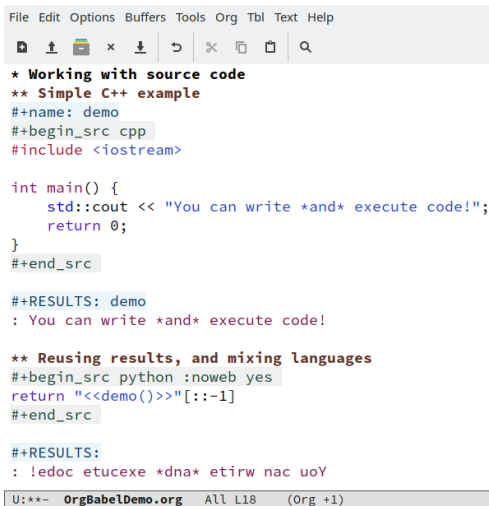
#+RESULTS:
: !edoc etucexe *dna* etirw nac uoY

U:*~ OrgBabelDemo.org All L18 (Org +1)
```

[3]

# What is org-mode

- **plain text-based** tool for outlining, note-taking, spreadsheets, project planning, ...
- **literate programming** via notebook-like environment for >70 programming languages
- **one-click publishing** as HTML (or full-fledged modern website), LaTeX, ODT, ...



```
File Edit Options Buffers Tools Org Tbl Text Help
[Icons]

* Working with source code
** Simple C++ example
#+name: demo
#+begin_src cpp
#include <iostream>

int main() {
    std::cout << "You can write *and* execute code!";
    return 0;
}
#+end_src

#+RESULTS: demo
: You can write *and* execute code!

** Reusing results, and mixing languages
#+begin_src python :noweb yes
return "<<demo()>>"[::-1]
#+end_src

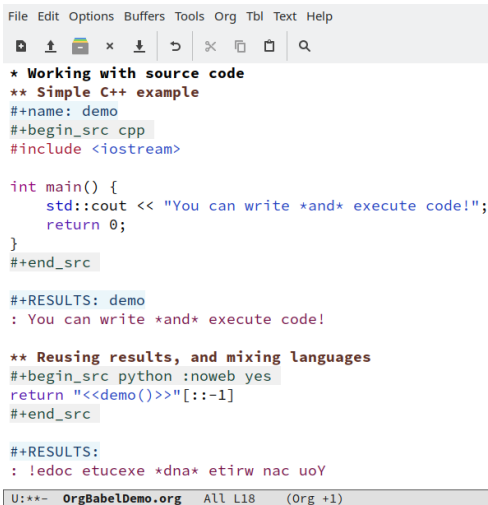
#+RESULTS:
: !edoc etucexe *dna* etirw nac uoY

U:*~ OrgBabelDemo.org All L18 (Org +1)
```

[3]

# What is org-mode

- **plain text-based** tool for outlining, note-taking, spreadsheets, project planning, ...
- **literate programming** via notebook-like environment for >70 programming languages
- **one-click publishing** as HTML (or full-fledged modern website), LaTeX, ODT, ...
- **seamless git compatibility**



```
File Edit Options Buffers Tools Org Tbl Text Help
[Icons]
* Working with source code
** Simple C++ example
#+name: demo
#+begin_src cpp
#include <iostream>

int main() {
    std::cout << "You can write *and* execute code!";
    return 0;
}
#+end_src

#+RESULTS: demo
: You can write *and* execute code!

** Reusing results, and mixing languages
#+begin_src python :noweb yes
return "<<demo()>>"[::-1]
#+end_src

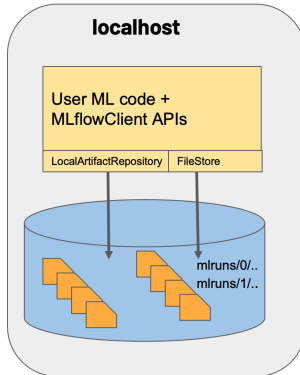
#+RESULTS:
: !edoc etucexe *dna* etirw nac uoY

U:*~ OrgBabelDemo.org All L18 (Org +1)
```

[3]

# What is Mlflow

- **Tracking:** API to log parameters, code, and results



**Scenario 1:** MLflow on the localhost

[4]

[4] From MLflow documentation, CC BY 4.0, <https://mlflow.org/docs/latest/projects.html>

# What is Mlflow

- **Tracking:** API to log parameters, code, and results
- **Projects:** code packaging format for reproducible runs using Conda and Docker

```
name: My Project

conda_env: my_env.yaml
# Can have a docker_env instead of a conda_env, e.g.
# docker_env:
#   image: mlflow-docker-example

entry_points:
  main:
    parameters:
      data_file: path
      regularization: {type: float, default: 0.1}
    command: "python train.py -r {regularization} {data_file}"
  validate:
    parameters:
      data_file: path
    command: "python validate.py {data_file}"
```

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

**Mlflow**

Git

Reproducible Workflow

[4]

Discussion

[4]From MLflow documentation, CC BY 4.0, <https://mlflow.org/docs/latest/projects.html>

# What is Mlflow

- **Tracking:** API to log parameters, code, and results
- **Projects:** code packaging format for reproducible runs using Conda and Docker
- **Models:** model packaging format and tools for deployment (from any ML library)

```
# Directory written by mlflow.sklearn.save_model(model, "my_model")
my_model/
├─ MLmodel
├─ model.pkl
├─ conda.yaml
├─ python_env.yaml
└─ requirements.txt
```

And its `MLmodel` file describes two flavors:

```
time_created: 2018-05-25T17:28:53.35

flavors:
  sklearn:
    sklearn_version: 0.19.1
    pickled_model: model.pkl
  python_function:
    loader_module: mlflow.sklearn
```

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

**Mlflow**

Git

Reproducible Workflow

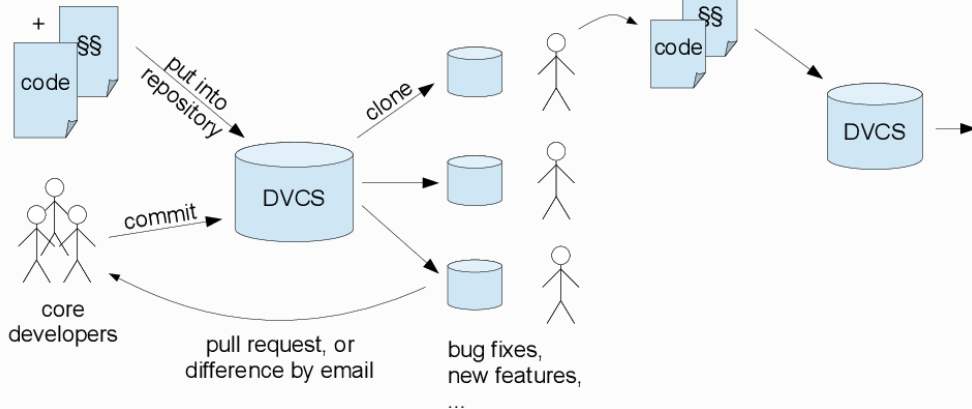
[4] Discussion

[4] From MLflow documentation, CC BY 4.0, <https://mlflow.org/docs/latest/projects.html>



# Distributed Version Control with Git

source code with license  
that allows reuse



Git, org-mode and mflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle  
Definitions: reproducibility

Tools

Org-mode  
Mflow  
Git

Reproducible Workflow

Discussion

[5]

[5]Vanschoren, J., Braun, M.L., Ong, C. (2018). Open science in machine learning. ArXiv, abs/1402.6013.

## Scope

Computational research life  
cycle

Definitions: reproducibility

## Tools

Org-mode

Mlflow

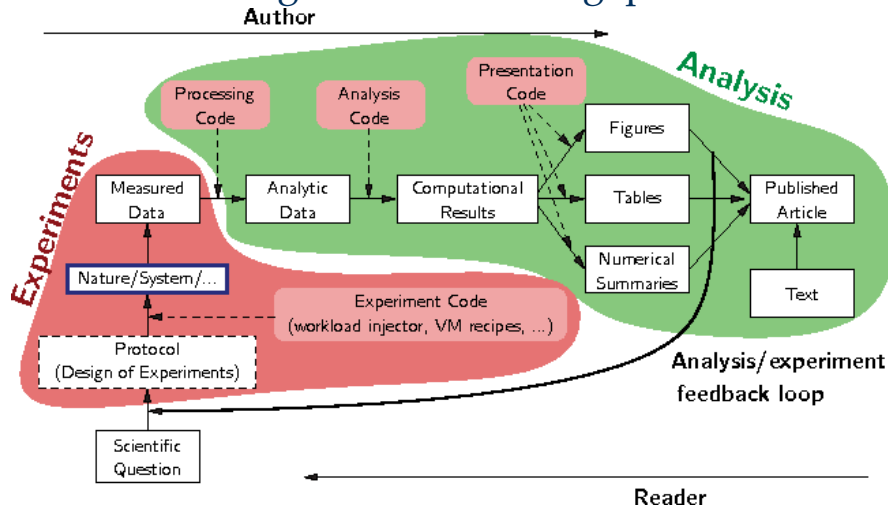
Git

## Reproducible Workflow

Discussion

# Reproducible Workflow

# Motivation: bridge author-reader gap<sup>[6]</sup>



Git, org-mode and mflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle  
Definitions: reproducibility

Tools

Org-mode  
Mflow  
Git

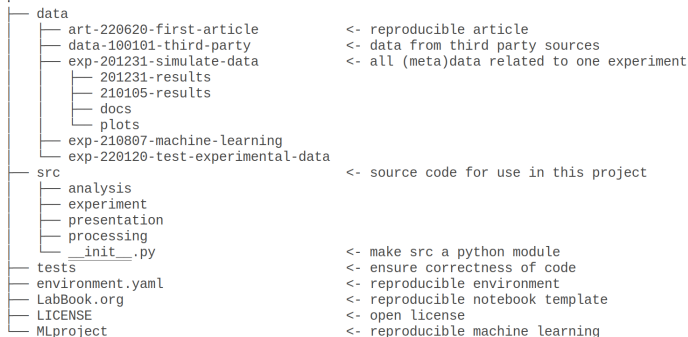
Reproducible Workflow

Discussion

<sup>[6]</sup>Stanisic, L., Legrand, A., Danjean, V. (2015). An Effective Git And Org-Mode Based Workflow For Reproducible Research. ACM SIGOPS Operating Systems Review, 49(1), 61–70. <https://doi.org/10/gfbx5x>

# Org-mode and Git for Reproducible Research<sup>[7]</sup>

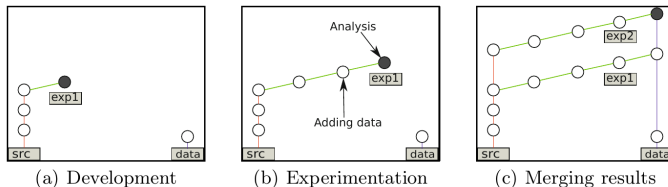
- **data file organization:**  
clear, coherent,  
hierarchical



<sup>[7]</sup>Stanisic, L., Legrand, A., Danjean, V. (2015). An Effective Git And Org-Mode Based Workflow For Reproducible Research. ACM SIGOPS Operating Systems Review, 49(1), 61–70. <https://doi.org/10/gfbx5x>

# Org-mode and Git for Reproducible Research<sup>[7]</sup>

- **data file organization:** clear, coherent, hierarchical
- **git branching structure;** version code, data, and results

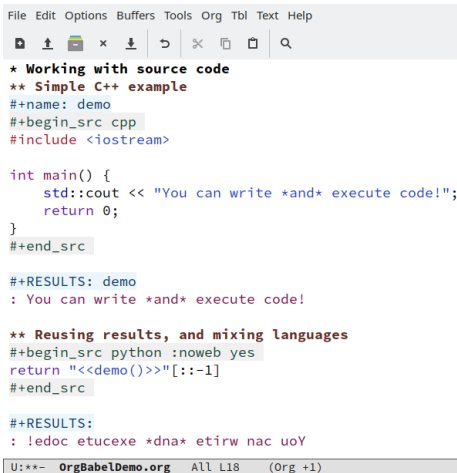


**Fig. 1.** Different phases in git workflow

<sup>[7]</sup>Stanisic, L., Legrand, A., Danjean, V. (2015). An Effective Git And Org-Mode Based Workflow For Reproducible Research. ACM SIGOPS Operating Systems Review, 49(1), 61–70. <https://doi.org/10/gfbx5x>

# Org-mode and Git for Reproducible Research<sup>[7]</sup>

- **data file organization:**  
clear, coherent,  
hierarchical
- **git branching  
structure;** version  
code, data, and results
- **org-mode LabBook:**  
key analysis details



```
File Edit Options Buffers Tools Org Tbl Text Help
[Icons] [Search]

* Working with source code
** Simple C++ example
#+name: demo
#+begin_src cpp
#include <iostream>

int main() {
    std::cout << "You can write *and* execute code!";
    return 0;
}
#+end_src

#+RESULTS: demo
: You can write *and* execute code!

** Reusing results, and mixing languages
#+begin_src python :noweb yes
return "<<demo()>>"[::-1]
#+end_src

#+RESULTS:
: !edoc etucexe *dna* etirw nac uoY

U:*-- OrgBabelDemo.org All L18 (Org +1)
```

<sup>[7]</sup>Stanisic, L., Legrand, A., Danjean, V. (2015). An Effective Git And Org-Mode Based Workflow For Reproducible Research. ACM SIGOPS Operating Systems Review, 49(1), 61–70. <https://doi.org/10/gfbx5x>

# Workflow Extension: mlflow + more

- **mlflow**: open source tool, that covers entire ML lifecycle and bridges the gap from ML research to application (e.g. model serving)

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

MLflow

Git

Reproducible Workflow

Discussion

# Workflow Extension: mlflow + more

- **mlflow**: open source tool, that covers entire ML lifecycle and bridges the gap from ML research to application (e.g. model serving)
- **more practices adapted from open source development**: single-click dependency setup (e.g. Docker), automated unit testing (e.g. tox), automated code documentation (e.g. sphinx)

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

Mlflow

Git

Reproducible Workflow

Discussion



# Workflow Extension: mlflow + more

- **mlflow**: open source tool, that covers entire ML lifecycle and bridges the gap from ML research to application (e.g. model serving)
- **more practices adapted from open source development**: single-click dependency setup (e.g. Docker), automated unit testing (e.g. tox), automated code documentation (e.g. sphinx)
- **leveraging org-mode single-click publishing**: host website documenting every step of your experiments (*Open Notebook Science*), easiest way to share experimental results including provenance

# Discussion

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

Mlflow

Git

Reproducible Workflow

Discussion

# Pros and Cons, Alternative Tools

- **Pros**

- combination of well-known, lightweight, open-source technologies
- facilitates reproducibility without taking away too much flexibility

- **Cons**

- some conventions not commonly used (git branching model)
- steep learning curve (org-mode preferably with Emacs)
- large files (possible solutions: git lfs, git-annex)

Git, org-mode and mflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

Mflow

Git

Reproducible Workflow

Discussion

# Pros and Cons, Alternative Tools

- **Pros**

- combination of well-known, lightweight, open-source technologies
- facilitates reproducibility without taking away too much flexibility

- **Cons**

- some conventions not commonly used (git branching model)
- steep learning curve (org-mode preferably with Emacs)
- large files (possible solutions: git lfs, git-annex)

- **Alternatives**

- jupyter notebooks instead org-mode - more commonly used, more intuitive, but not plain text (git integration), not the same flexibility
- R + knitr for literal programming
-

#### Scope

Computational research life  
cycle

Definitions: reproducibility

#### Tools

Org-mode

Mflow

Git

Reproducible Workflow

#### Discussion



# Questions?

# Acknowledgements



iaob  
Institute of Applied Optics  
and Biophysics

 Biophysical Imaging  
Eggeling Group

 FRIEDRICH-SCHILLER-  
UNIVERSITÄT  
JENA

Leibniz | ipht   
LEIBNIZ INSTITUTE OF  
PHOTONIC TECHNOLOGY



This work is licensed under a Creative Commons  
Attribution 4.0 International License.

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aselmann  
aselmann.github.io

Scope

Computational research life  
cycle

Definitions: reproducibility

Tools

Org-mode

Mlflow

Git

Reproducible Workflow

Discussion

# Thank you!

## Extra: Resources

- from comment after talk: DataLad as alternative to git
- Open online course to learn git and org-mode

Tools for research Digital and technology

### Reproducible research: methodological principles for transparent science

Ref. 41016

⌚ Effort: 24 hours ⚙️ Pace: Self paced

This Mooc proposes methodological principles for open and transparent science. It deals in a practical way with note-taking, computational documentation, replicability of analyses.



france-universite-numerique

**FUN-MOOC : Recherche reproducible : principes méthodologiques pour une science transparente**

[8][9]

[8]<http://handbook.datalad.org/en/latest/index.html#>

[9]<https://www.fun-mooc.fr/en/courses/reproducible-research-methodological-principles-transparent-scie/>

# Extra: Mlflow UI

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io

	Date	User	Source	Version	Parameters		Metrics		
					alpha	l1_ratio	mae	r2	rmse
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:10</a>	mlflow	train.py	05e956	1	1	0.649	0.04	0.862
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:10</a>	mlflow	train.py	05e956	1	0.5	0.648	0.046	0.859
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:10</a>	mlflow	train.py	05e956	1	0.2	0.628	0.125	0.823
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:09</a>	mlflow	train.py	05e956	1	0	0.619	0.176	0.799
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:09</a>	mlflow	train.py	05e956	0.5	1	0.648	0.046	0.859
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:09</a>	mlflow	train.py	05e956	0.5	0.5	0.628	0.127	0.822
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:09</a>	mlflow	train.py	05e956	0.5	0.2	0.621	0.171	0.801
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:09</a>	mlflow	train.py	05e956	0.5	0	0.615	0.199	0.787
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:09</a>	mlflow	train.py	05e956	0	1	0.578	0.288	0.742
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:09</a>	mlflow	train.py	05e956	0	0.5	0.578	0.288	0.742
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:09</a>	mlflow	train.py	05e956	0	0.2	0.578	0.288	0.742
<input type="checkbox"/>	<a href="#">2018-06-04 23:00:08</a>	mlflow	train.py	05e956	0	0	0.578	0.288	0.742

[10]

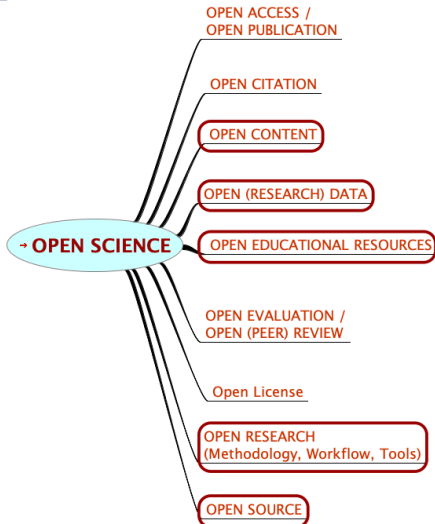
<sup>[10]</sup>From MLflow documentation, CC BY 4.0, <https://mlflow.org/docs/latest/projects.html>



# Extra: fields of Open Science covered

Git, org-mode and mlflow  
for reproducibility

Alex Seltmann  
Twitter @aseltmann  
aseltmann.github.io



This work by Peter Baumgartner is licensed under a  
→ Creative Commons Attribution-ShareAlike 4.0 International License.

# Extra: other definitions of reproducibility<sup>[11]</sup>

## Methods reproducibility

- = get *same* results, use *same* data and tools
- topics: provide study protocols, reusable (meta)data, code, results, ...

## Results reproducibility = replication

- = get *similar* results, use *similar* procedures and tools (maybe different data)
- topics: statistical significance, cumulative evidential weight, heterogeneity tests, effect sizes...

## Inferential reproducibility

- = get same scientific conclusions from independent study or re-analysis of the data, use *different* tools and methods
- topics: bayesian perspectives, avoiding multiplicity, HARKing, p-hacking

---

<sup>[11]</sup>Goodman et al (2016)