# Open Software in Open Science

**Dr. Britta Westner**

Open Science and Reproducibility Workshop
March 12, 2019

AARHUS
UNIVERSITY

✉ britta@cfin.au.dk        🐦 @britta_wstnr        🐙 britta-wstnr

# Outline

**Let's find answers to the following questions:**

- Why is open source **essential for open science**?
- What are **best practices** for open tools?
- How does all this facilitate **reproducibility**?
- Is there an **open source crisis**?

# What is open source?

Software is *open source* if the source code

- is **freely available**

- may be **modified**

- may be **redistributed**

*What I cannot create, I do not understand.*

*Richard Feynman, 1988*

*Black boxes do not belong in science.*

*Fernando Pérez, 2017*

# Open source science

For **reproducibility** of results, the following things need to be considered:
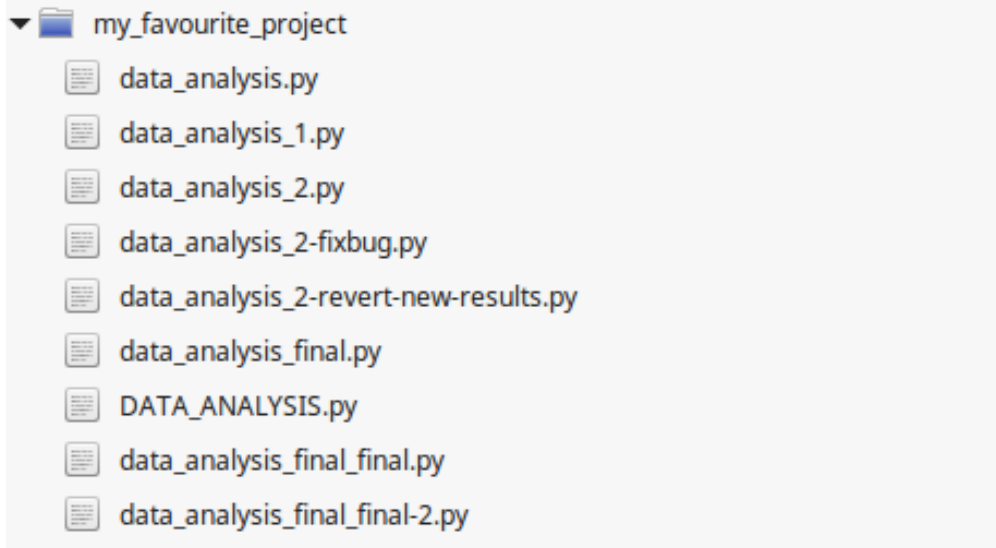
- **computational tools**: your scripts, toolboxes, programming language, operating system, . . .

- the **data**

- **sharing** of the work

- **communicating** the work

# Open source science

For **reproducibility** of results, the following things need to be considered:

- **computational tools**: **use open tools and share your code**

- the **data**: **share**

- **sharing** of the work: **in an easily accessible manner**

- **communicating** the work: **publish, tweet, … − and include links to code and data!**

# Making data analyses reproducible

**Reproducibility starts with you.**

*Looks familiar?*

```
▼ 📁 my_favourite_project
    📄 data_analysis.py
    📄 data_analysis_1.py
    📄 data_analysis_2.py
    📄 data_analysis_2-fixbug.py
    📄 data_analysis_2-revert-new-results.py
    📄 data_analysis_final.py
    📄 DATA_ANALYSIS.py
    📄 data_analysis_final_final.py
    📄 data_analysis_final_final-2.py
```

- can you reproduce your own results at a later stage?
- use **version control**
- **document** your code

# A word about version control

Using version control provides you
with your own **time machine**.

**Principle:**

- you are responsbile for time stamps
- file only exists in most recent version
- log of changes
- recommendation: **git**

photograph by Babbel1996 / CC-BY-2.5

# Making code public

**Where?**

GitHub     GitLab     GitBucket

**How? Etiquette for sharing code.**

- include a **license**
- share your code **formatted**: line width, coding stlyes (linters)
- **document** your code: comments, docstrings, project description
- note down **dependencies** and versions

# Got style?

## A demonstration how coding styles make things easier.



```
 1»import mne; import numpy as np
 2»from mne.beamformer import make_lcmv,apply_lcmv_epochs
 3»def run_lcmv_epochs(epochs,fwd,data_cov,reg, noise_cov=None,pick_ori='max-power', weigh
 4»    filters = make_lcmv(epochs.info, fwd, data_cov=data_cov,  noise_cov=noise_cov, pick
 5»    stcs=apply_lcmv_epochs(epochs,filters,return_generator = True,max_ori_out='signed',
 6»    stcs_mat=np.ones((epochs._data.shape[0], fwd['nsource'],len(epochs.times)))
 7     if verbose is False:
 8         mne.set_log_level('WARNING')
 9     for trial in range(epochs._data.shape[0]):
10»        if trial==0:
11»            stc = next(stcs);stcs_mat[trial, :, :]= stc.data
12         else:
13»            stcs_mat[trial,:,:]=next(stcs).data
14»    return stcs_mat,stc, filters
15
```

```
File        Line Col Level      ID      Message (Checker)
bad_e…       1   11  warning    E702    multiple statements on one line (semicolon) (python-flake8)
bad_e…       2   37  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…       5   27  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…       5   31  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…       5   40  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…       5   60  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…       5   80  warning    E501    line too long (116 > 79 characters) (python-flake8)
bad_e…       6   80  warning    E501    line too long (152 > 79 characters) (python-flake8)
bad_e…       7    9  warning    E225    missing whitespace around operator (python-flake8)
bad_e…       7   34  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…       7   42  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…       7   59  warning    E251    unexpected spaces around keyword / parameter equals (python-flake8)
bad_e…       7   61  warning    E251    unexpected spaces around keyword / parameter equals (python-flake8)
bad_e…       7   66  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…       7   80  warning    E501    line too long (104 > 79 characters) (python-flake8)
bad_e…       8   13  warning    E225    missing whitespace around operator (python-flake8)
bad_e…       8   60  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…      12   17  warning    E225    missing whitespace around operator (python-flake8)
bad_e…      13   29  warning    E702    multiple statements on one line (semicolon) (python-flake8)
bad_e…      13   29  warning    E231    missing whitespace after ';' (python-flake8)
bad_e…      13   51  warning    E225    missing whitespace around operator (python-flake8)
bad_e…      15   27  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…      15   29  warning    E231    missing whitespace after ',' (python-flake8)
bad_e…      15   32  warning    E225    missing whitespace around operator (python-flake8)
bad_e…      16   20  warning    E231    missing whitespace after ',' (python-flake8)
```

# Got style?

**A demonstration how coding styles make things easier.**

```python
1  import mne
2  import numpy as np
3  from mne.beamformer import make_lcmv, apply_lcmv_epochs
4
5
6  def run_lcmv_epochs(epochs, fwd, data_cov, reg, noise_cov=None,
7                      pick_ori='max-power', weight_norm='nai', verbose=False):
8      filters = make_lcmv(epochs.info, fwd, data_cov=data_cov,
9                          noise_cov=noise_cov, pick_ori=pick_ori,
10                          reg=reg, weight_norm=weight_norm, verbose=verbose)
11     stcs = apply_lcmv_epochs(epochs, filters, return_generator=True,
12                          max_ori_out='signed', verbose=verbose)
13     stcs_mat = np.ones((epochs._data.shape[0], fwd['nsource'],
14                          len(epochs.times)))
15     if verbose is False:
16         mne.set_log_level('WARNING')
17     for trial in range(epochs._data.shape[0]):
18         if trial == 0:
19             stc = next(stcs)
20             stcs_mat[trial, :, :] = stc.data
21         else:
22             stcs_mat[trial, :, :] = next(stcs).data
23     return stcs_mat, stc, filters
```

# Got style?

**A demonstration how coding styles make things easier.**

```python
1  import mne
2  import numpy as np
3  from mne.beamformer import make_lcmv, apply_lcmv_epochs
4
5
6  def run_lcmv_epochs(epochs, fwd, data_cov, reg, noise_cov=None,
7                      pick_ori='max-power', weight_norm='nai', verbose=False):
8
9      filters = make_lcmv(epochs.info, fwd, data_cov=data_cov,
10                         noise_cov=noise_cov, pick_ori=pick_ori,
11                         reg=reg, weight_norm=weight_norm, verbose=verbose)
12
13     stcs = apply_lcmv_epochs(epochs, filters, return_generator=True,
14                              max_ori_out='signed', verbose=verbose)
15
16     stcs_mat = np.ones((epochs._data.shape[0], fwd['nsource'],
17                        len(epochs.times)))
18
19     if verbose is False:
20         mne.set_log_level('WARNING')
21
22     for trial in range(epochs._data.shape[0]):
23         if trial == 0:
24             stc = next(stcs)
25             stcs_mat[trial, :, :] = stc.data
26         else:
27             stcs_mat[trial, :, :] = next(stcs).data
28
29     return stcs_mat, stc, filters
```

# Got style?

**A demonstration how coding styles make things easier.**

```python
1  import mne
2  import numpy as np
3  from mne.beamformer import make_lcmv, apply_lcmv_epochs
4
5
6  def run_lcmv_epochs(epochs, fwd, data_cov, reg, noise_cov=None,
7                      pick_ori='max-power', weight_norm='nai', verbose=False):
8      """Run LCMV on epochs."""
9
10     filters = make_lcmv(epochs.info, fwd, data_cov=data_cov,
11                         noise_cov=noise_cov, pick_ori=pick_ori,
12                         reg=reg, weight_norm=weight_norm, verbose=verbose)
13
14     # apply that filter to epochs
15     stcs = apply_lcmv_epochs(epochs, filters, return_generator=True,
16                              max_ori_out='signed', verbose=verbose)
17
18     # preallocate matrix
19     stcs_mat = np.ones((epochs._data.shape[0], fwd['nsource'],
20                         len(epochs.times)))
21
22     if verbose is False:
23         mne.set_log_level('WARNING')
24
25     # resolve generator
26     for trial in range(epochs._data.shape[0]):
27         # first time: also save stc
28         if trial == 0:
29             stc = next(stcs)
30             stcs_mat[trial, :, :] = stc.data
31         else:
32             stcs_mat[trial, :, :] = next(stcs).data
33
34     return stcs_mat, stc, filters
```

# Got style?

**A demonstration how coding styles make things easier.**

```python
1  import mne
2  import numpy as np
3  from mne.beamformer import make_lcmv, apply_lcmv_epochs
4
5
6  def run_lcmv_epochs(epochs, fwd, data_cov, reg, noise_cov=None,
7                      pick_ori='max-power', weight_norm='nai', verbose=False):
8      """Run LCMV on epochs.
9      Run weight-normalized LCMV beamformer on epoch data, will return matrix of
10     trials or stc object.
11
12     Parameters
13     ----------
14     epochs : MNE epochs
15         epochs to source reconstruct.
16     fwd : MNE forward model
17         forward model.
18     data_cov : MNE covariance estimateg
19         data covariance matrix
20     reg : float
21         regularization parameter
22     noise_cov : MNE covariance estimate
23         noise covariance matrix, optional
24     verbose : bool
25         overrides default verbose level, defaults to False, i.e., no logger
26         info.
27
28     Returns
29     -------
30     stcs_mat : numpy array
31         matrix with all source trials
32     stc : MNE stc
33         single trial stc object (last trial)
34     filters : dict
35         spatial filter used in computation
36     """
37     filters = make_lcmv(epochs.info, fwd, data_cov=data_cov,
38                         noise_cov=noise_cov, pick_ori=pick_ori, reg=reg,
39                         weight_norm=weight_norm, verbose=verbose)
40
41     # apply that filter to epochs
42     stcs = apply_lcmv_epochs(epochs, filters, return_generator=True,
43                              max_ori_out='signed', verbose=verbose)
44
45     # preallocate matrix
46     stcs_mat = np.ones((epochs._data.shape[0], fwd['nsource'],
47                        len(epochs.times)))
48
49     if verbose is False:
50         mne.set_log_level('WARNING')
51
52     # resolve generator
53     for trial in range(epochs._data.shape[0]):
54         # last time: also save stc
55         if trial == 0:
56             stc = next(stcs)
57             stcs_mat[trial, :, :] = stc.data
```

# How GitHub facilitates open science

**On GitHub\*/Lab/Bucket you can:**

- **share** code

- **follow** researchers and toolboxes to stay up-to-date

- **collaborate** on projects

- **fork** projects to make your own version of them

- **contribute** to projects, e.g., open source toolboxes

\* GitHub itself is **not** open source!

# Making data public

For full reproducibility, **data** is needed.
One possibility for sharing data: **The Open Science Framework**

# OSF: Keeping data and code together

# Technical vs. practical reproducibility

## How easy is it to re-run your analysis?



https://www.gw-openscience.org/tutorials/

# Practical reproducibility: Binder
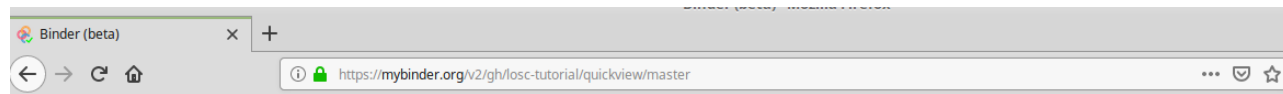
Notebooks are great, but:

- still need to download the **data**

- still need to create the right **environment**
  (software versions, operating system)



`https://www.gw-openscience.org/tutorials/`

**Wait, couldn't we write whole papers like this?**

# Practical reproducibility: eLife

**Why should I contribute to open source?**

- **solve a problem**
  1/2 of Github contributors contribute only once   Eghbal 2017

- for the **reputation**

- for the **community**
  *Came for the language, stayed for the community.*   Brett Cannon

# Contributing to open source: getting started

- Annoyed by that one **bug** in the toolbox? **Open an issue.**

- Know how to **fix** it? **Open a PR.**

- Most communities have a **how to contribute** wiki page.

- Most communities are very **welcoming**!

- Open source is **essential** for open science.

- Spans from **sharing code** to **using open source** toolboxes and software.

- **Practical reproducibility** is important.

- **Contributing** to open source toolboxes is fun!

# Is there an open source crisis?

## OpenSSL

The toolkit for internet connection security was used on 66% of all web servers worldwide (2014).

Prior to "Heartbleed", it was maintained by only a handful of volunteers.

Eghbal 2016; Klug & Miller 2018

## NumPy and scientific Python

Being one of the pillars of scientific Python, NumPy only secured stable funding in 2017.

The scientific Python world relied on an estimated 30 people in 2011.

NumFOCUS 2017; Pérez 2011

# Open source crisis — toolbox maintenance

2/3 of top projects on GitHub are maintained by only one or two people.

<div align="right">Avelino et al. 2017</div>

**The Truck Factor of toolboxes:**

*minimal number of developers that have to be hit by a truck before a project is lost.*

| Project | Truck Factor |
|---|---|
| git | 12 |
| scikit-learn | 7 |
| IPython | 4 |
| pandas | 2 |

<div align="right">Avelino et al. 2017</div>

- **funding**

- **needs of maintainers**: traditionally not considered in open source

  *Our goal should be to spread freedom and then defend it. That is more important than making our software popular, which would just be catering to our egos.*

  *Richard Stallman, 2005*

- **burning out on projects**: workload and toxic feedback

  *[T]he angry response has been overwhelming. Every single day I'm reading someone else rant about how awful of a job we're doing. It's been hard to stay motivated.*

  *James Kyle, 2016*

**Software work in science can be career suicide.**

Fernando Pérez, 2011

# Open source crisis — academia

**Chris Holdgraf**
@choldgraf

Something I didn't expect when working on @mybinderteam and @ProjectJupyter: not being attached to a specific scientific field is really scary! You realize how much the academic world doesn't understand how to value work that isn't attached to publications in "your domain."

🌐 Tweet übersetzen

13:35 - 11. Feb. 2019

27 Retweets  157 „Gefällt mir"-Angaben

💬 7      ⟲ 27      ♡ 157      ✉

# What can we do about it?

**The problems:**

- incentive structure of **modern academia** fits poorly with developers: *contributions instead of publications*

- **tradeoff**: expertise vs. time

**Possible solutions:**

- **critical mass**: sharing and contributing

- consider open source in **teaching** and supervising

- consider open source "sacrifices" in **hiring decisions** and with **grants**

# Conclusions

- Open source is **essential** for open science.

- Ways towards **higher reproducibility**.

- Ways towards **contributing to open source**.

- Awareness of the **open source dilemmas** and **ideas how to cope**.

# Acknowledgements

**CFIN @ Aarhus university**

Sarang Dalal

**MNE-Python**

Alexandre Gramfort

Denis A. Engemann

Eric Larson