

Reproducible projects with python and snakemake – a subjective user perspective

Bi-weekly ESR meetings

25th Feb 2022 - David

Why interesting?

Research today:

- many projects in parallel, from first steps to publication, several years can pass
 - people leave their lab after few years
 - different labs work on the same research topic and need to share their data analysis workflows with the results to enable collaboration
- Reproducibility is needed
 - Things must work out of the box to be re-used

Why interesting?


- Reproducibility is needed
- Things must work out of the box to be re-used

One solution: Snakemake workflows

- description via a human readable, Python based language
- creation of reproducible and scalable data analyses
- scaling from laptop to cluster computing without modification of the workflow description
- optionally manages the software environment to execute individual steps of the workflow (via conda, docker)

Example project

- Neural data from a Utah array (96 electrodes) in monkey motor cortex



scientific **data** [View all journals](#) [Search](#) [Login](#)

[Explore content](#) [About the journal](#) [Publish with us](#)

[nature](#) > [scientific data](#) > [data descriptors](#) > [article](#)

[Open Access](#) | [Published: 10 April 2018](#)

Massively parallel recordings in macaque motor cortex during an instructed delayed reach-to-grasp task

[Thomas Brochier](#) , [Lyuba Zehl](#) , [Yaoyao Hao](#), [Margaux Duret](#), [Julia Sprenger](#), [Michael Denker](#), [Sonja Grün](#) & [Alexa Riehle](#)

[Scientific Data](#) **5**, Article number: 180055 (2018) | [Cite this article](#)

4822 Accesses | **17** Citations | **3** Altmetric | [Metrics](#)

Example project

- Neural data from a Utah array (96 electrodes) in monkey motor cortex

Table 3 Overview of files (names, size, and content) for each provided dataset of monkey L and N.

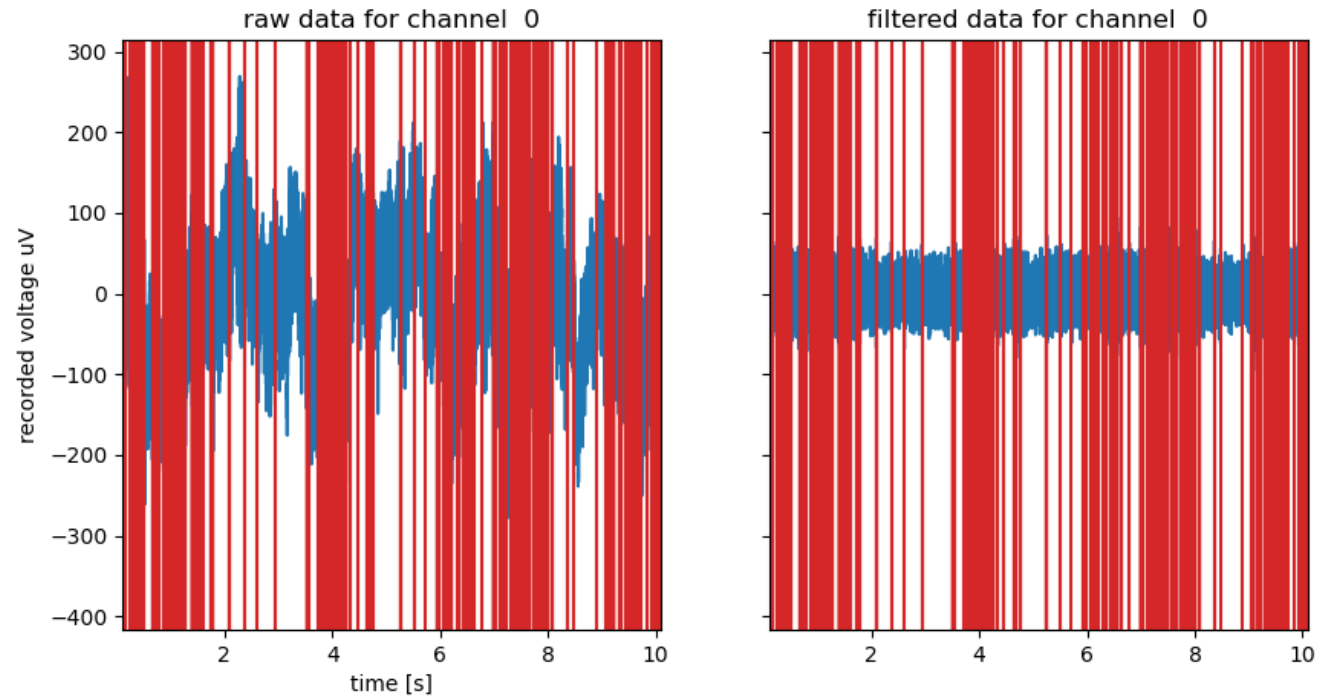
From: [Massively parallel recordings in macaque motor cortex during an instructed delayed reach-to-grasp task](#)

monkey L		
file names	file size	file content
I101210-001.ccf	108.2 kB	cerebus configuration file
I101210-001.nev	287.7 MB	digital events, unsorted spikes times & waveforms
I101210-001-02.nev	287.7 MB	digital events, sorted spikes times & waveforms
I101210-001.ns2	8.5 MB	analog signals of object sensors
I101210-001.ns5	4.1 GB	raw neuronal signal
I101210-001.odml	2.7 MB	metadata
monkey N		
file names	file size	file content
i140703-001.ccf	187.1 kB	cerebus configuration file
i140703-001.nev	168.3 MB	digital events, unsorted spikes times & waveforms
i140703-001-03.nev	168.3 MB	digital events, sorted spikes times & waveforms
i140703-001.ns2	204.7 MB	analog signals of object sensors and LFP signals
i140703-001.ns6	5.8 GB	raw neuronal signal
i140703-001.odml	2.3 MB	metadata

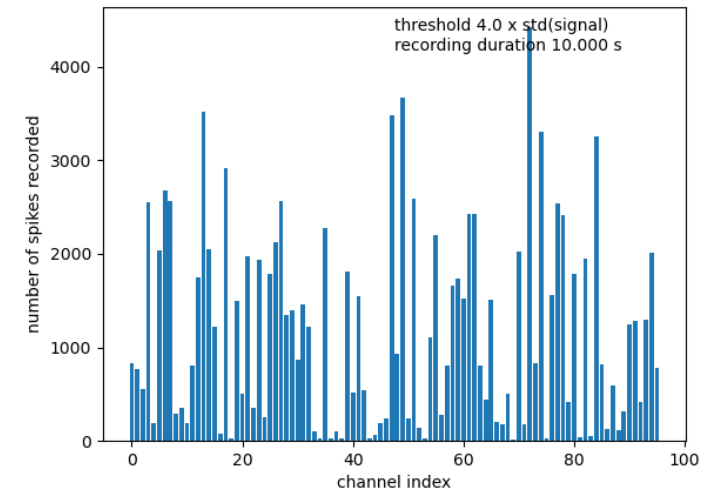
Outline

- Neural data from a Utah array (96 electrodes) in monkey motor cortex
- Workflow steps to perform on the dataset:
 - Download
 - Available on online repository
 - Preprocess
 - Bandpass filter the voltages recorded at all electrodes
 - Analyze
 - Extract spikes based on thresholding the filtered voltage signals
 - Count the number of spikes recorded per electrode
 - Plot results
 - Plot raw signal and filtered signal with spikes
 - Plot histogram of the number of spikes per channel
 -

Target



- Figure 1: raw signal, filtered signal, spikes
- Figure 2:
- Parameters:
- E.g. filename, filter frequencies, threshold for spike detection, file-format of the plots



Workflow Parameters

- Neural data from a Utah array (96 electrodes) in monkey motor cortex
 - Workflow steps to perform on the dataset:
 - Download
 - Available on online repository
 - **filename**
 - Preprocess
 - Bandpass filter the voltages recorded at all electrodes
 - **lower & upper filter frequency**
 - Analyze
 - Extract spikes based on thresholding the filtered voltage signals
 - Count the number of spikes recorded per electrode
 - **threshold for spike detection**
 - Plot results
 - Plot raw signal and filtered signal with spikes
 - Plot histogram of the number of spikes per channel
 - **plot parameters, e.g. file format**
- (many more parameters not listed here)

Workflow

Parameters

program to
perform the step

- Neural data from a Utah array (96 electrodes) in monkey motor cortex
 - Workflow steps to perform on the dataset:
 - Download **curl using bash shell command**
 - Available on online repository
 - **filename**
 - Preprocess **python script**
 - Bandpass filter the voltages recorded at all electrodes
 - **lower & upper filter frequency**
 - Analyze **python script**
 - Extract spikes based on thresholding the filtered voltage signals
 - Count the number of spikes recorded per electrode
 - **threshold for spike detection**
 - Plot results **python script**
 - Plot raw signal and filtered signal with spikes
 - Plot histogram of the number of spikes per channel
 - **plot parameters, e.g. file format**
- (many more parameters not listed here)

workflow recipe

in snakemake:

- Definition of target files
 - here we want to produce 2 plots:
 - Raw signal, filtered signal with spikes
 - Histogram of spike count per electrode
- Download **curl using bash shell command**
 - Available on online repository
 - **filename**
- Preprocess **python script**
 - Bandpass filter the voltages recorded at all electrodes
 - **lower & upper filter frequency**
- Analyze **python script**
 - Extract spikes based on thresholding the filtered voltage signals
 - Count the number of spikes recorded per electrode
 - **threshold for spike detection**
- Plot results **python script**
 - Plot raw signal and filtered signal with spikes
 - Plot histogram of the number of spikes per channel
 - **plot parameters, e.g. file format**

(many more parameters not listed here)

```
Snakefile_heuristic.txt
~/Documents/meeting_presentations/2022-02-23_Entrain_ESR_meeting/my_example
Save

1 configfile: "config.yml"
2
3 rule all:
4     input:
5         -> list of all target-filenames with their name including the employed parameters
6         in snakemake syntax:
7             "wildcard"1          "wildcard"2          "wildcard"3
8             results/hist__file_{file_name}__filter{freq-range}__spike-threshold_{thresh}.png
9             e.g. results/hist__file_i140703-001__filter [500,7500] __spike-threshold_ 3 .png
10
11 rule plot_figures:
12     input:
13         "{file_name}.ns6",
14         "data/filtered_signal_RELEVANT-PARAMETERS-HERE.npy",
15         "data/spiketrains_RELEVANT-PARAMETERS-HERE.npy",
16
17     output:
18         "results/signal__RELEVANT-PARAMETERS-HERE}.png",
19         "results/hist__file_{file_name}__filter{freq-range}__spike-threshold_{thresh}.png"
20     script:
21         "scripts/plot.py"
22
23 rule analyze:
24     input:
25         "data/filtered_signal__RELEVANT-PARAMETERS-HERE.npy"
26     output:
27         "data/spiketrains_RELEVANT-PARAMETERS-HERE.npy"
28     script:
29         "scripts/analyze.py"
30
31 rule preprocess:
32     input:
33         "{file_name}.ns6"
34     output:
35         "data/filtered_signal_RELEVANT-PARAMETERS-HERE.npy"
36     script:
37         "scripts/preprocess.py"
38
39 rule download_data:
40     output:
41         "{file_name}.ns6"
42     params:
43         link = 'download-data-here.com'
44     shell:
45         "curl -L {params.link} > {output}"
```

workflow configuration

Select the parameters to execute your workflow

- **Filename** to download
- **lower & upper filter frequency** for the bandpass filtering
- **threshold** for spike detection
- **plot parameters**, e.g. file format

(many more parameters not listed here)

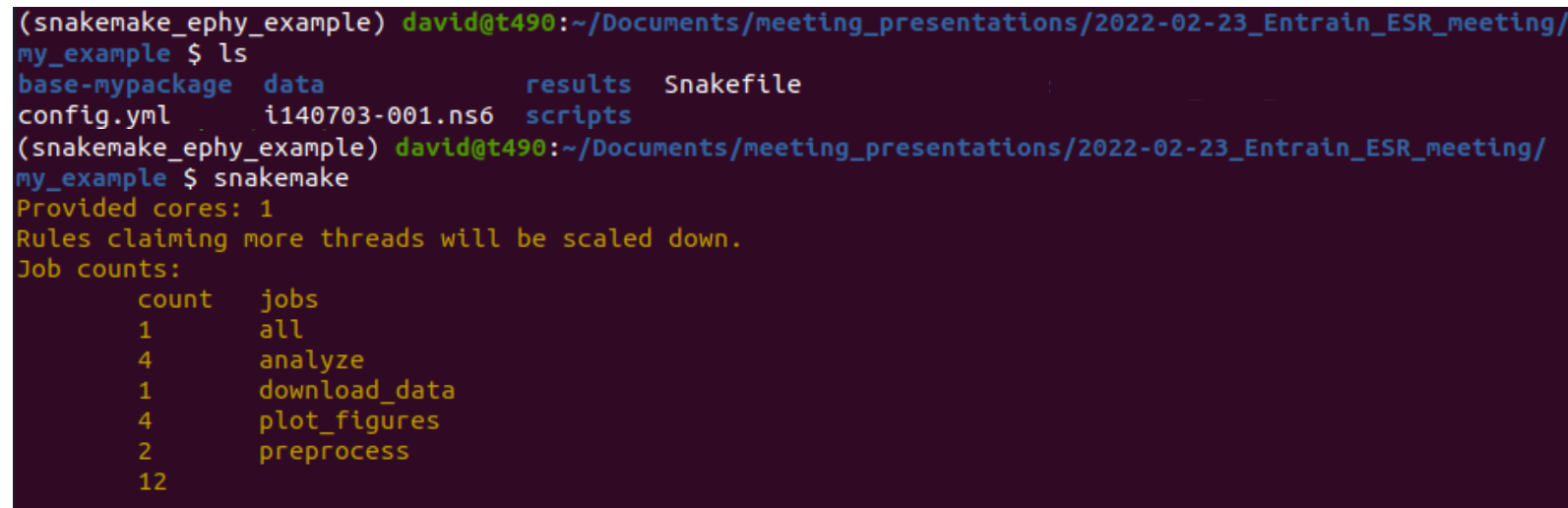


The screenshot shows a code editor window with two tabs: 'Snakefile_heuristic.txt' and 'config.yml'. The 'config.yml' tab is active, displaying the following configuration parameters:

```
1 file_name: ['i140703-001']
2
3 # slice loaded for analysis
4 tstart: [0.1]
5 tstop: [10.1]
6
7 # bandpass filter frequency bounds
8 lowcut: [500,500]
9 highcut: [5000,7500]
10
11 # threshold for spike detection (factor of std(filtered_signal))
12 spike_threshold_std: [3,4]
```

workflow execution

- Generation of all combinations of figures related to the parameter range
- When expanding the parameter range, snakemake automatically manages which parameter-sets have been already calculated and only performs the necessary steps

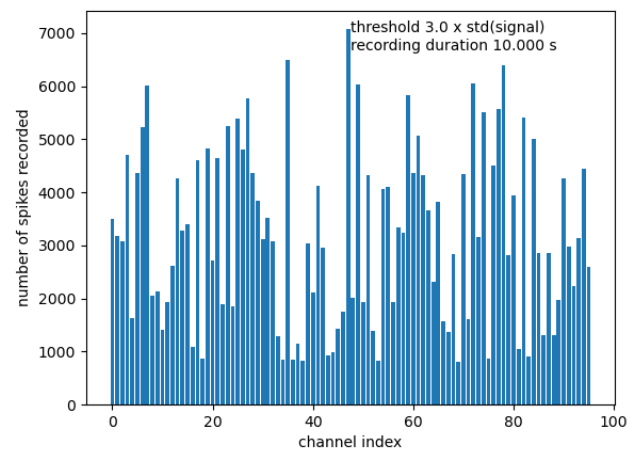


The screenshot shows a terminal window with the following output:

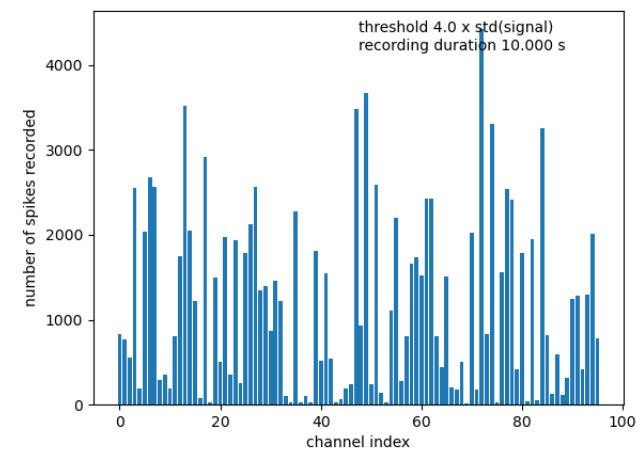
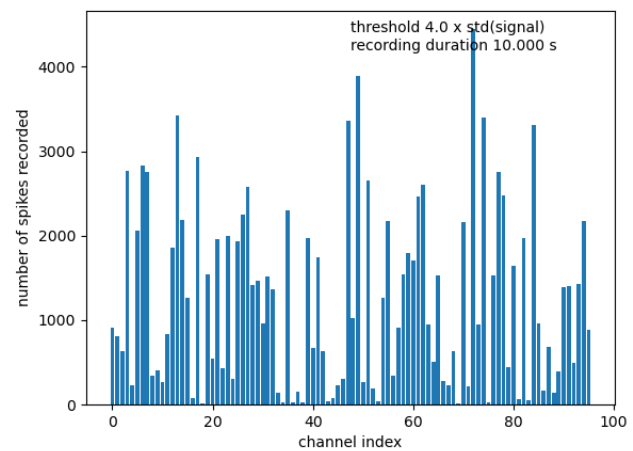
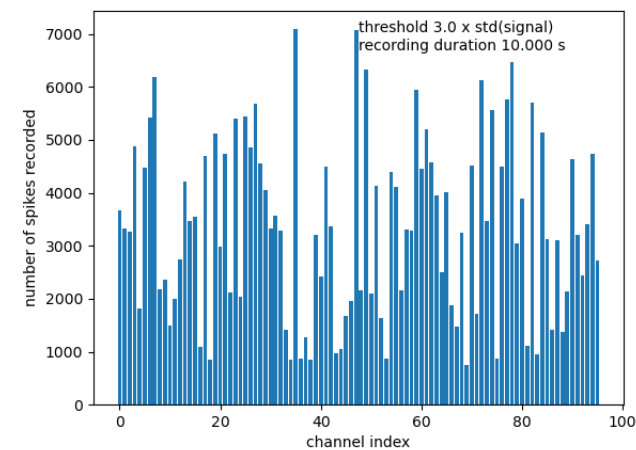
```
(snakemake_ephy_example) david@t490:~/Documents/meeting_presentations/2022-02-23_Entrain_ESR_meeting/my_example $ ls
base-mypackage  data          results  Snakefile
config.yml      i140703-001.ns6  scripts
(snakemake_ephy_example) david@t490:~/Documents/meeting_presentations/2022-02-23_Entrain_ESR_meeting/my_example $ snakemake
Provided cores: 1
Rules claiming more threads will be scaled down.
Job counts:
   count  jobs
     1    all
     4  analyze
     1 download_data
     4 plot_figures
     2 preprocess
    12
```

workflow results (only histograms)

Bandpass filter range 500Hz-5kHz



Bandpass filter range 500Hz-7.5kHz



Additional features

- software management can be integrated
 - execution of workflow steps with individually defined conda-environments or docker containers
- Scalable from your laptop to a cluster via execution flags
 - E.g. execute steps in parallel on X cores:
`snakemake -cores X`
 - Run your workflow on a computing cluster with job scheduling systems, e.g. slurm:
`snakemake --cluster "sbatch -J snake" --jobs 95`

Thanks for
your patience
and interest 😊

For those interested, we can go into the details now ;-)