# Reproducibility in Action

My experiences writing a reproducible paper

*Rasp, Selz and Craig, 2017. Variability and clustering of mid-latitude summertime convection: [...]. Submitted to JAS*

Stephan Rasp

- LMU Munich -

## Computation section in the manuscript

*e. Computational details and reproducibility*

This subsection closely follows the guidelines on publishing computational results proposed by Irving (2016). The analysis and plotting of model and observation data was done using Python. The Python libraries NumPy (Numerical Python; van der Walt et al. (2011)) and SciPy (Jones et al. 2001–) were used heavily. The raw data were read with the Python module cosmo_utils (code available upon request). The figures were plotted using the Python module Matplotlib (Hunter 2007). Plotting colors were chosen according to the Hue-Chroma-Luminance color space (Stauffer et al. 2015). Some plots were post-processed using the vector graphics program Inkscape.

To enable reproducibility of the results, this paper is accompanied by a version-controlled code repository (https://github.com/raspstephan/convective_variability_analysis) and a Figshare repository (Rasp 2017), which contains a snapshot of the code repository at the time of submission and supplementary log files for each figure. These log files contain information about the computational steps taken from the raw data to the generation of the plots. While the model code and initial data is not openly available, a detailed technical description of the model simulations can be found in the cosmo_runscripts directory of the code repository. The Jupyter notebooks (Kluyver et al. 2016) mentioned in the text are stored in the directory jupyter_notebooks of the repository. Links to non-interactive versions of the notebooks can be found on the front page of the Github repository; rendered PDF versions are also added to the supplement of this paper.
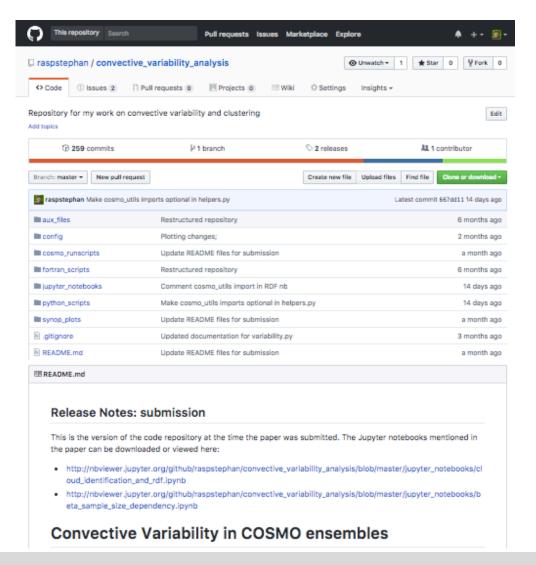
Citation of software tools
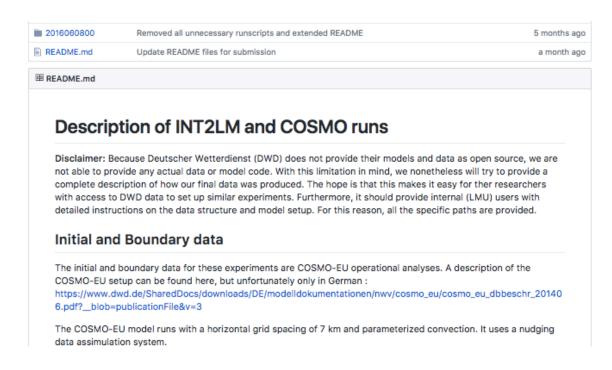
Github repository

Figshare repository with log files

Jupyter notebooks

The github repository contains all code used from the raw data to the final figures with adequate documentation.
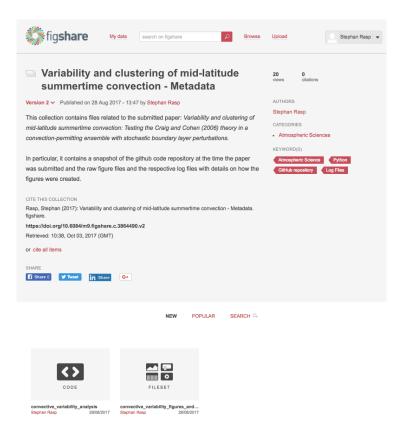




## Description of INT2LM and COSMO runs

Disclaimer: Because Deutscher Wetterdienst (DWD) does not provide their models and data as open source, we are not able to provide any actual data or model code. With this limitation in mind, we nonetheless will try to provide a complete description of how our final data was produced. The hope is that this makes it easy for ther researchers with access to DWD data to set up similar experiments. Furthermore, it should provide internal (LMU) users with detailed instructions on the data structure and model setup. For this reason, all the specific paths are provided.

### Initial and Boundary data

The initial and boundary data for these experiments are COSMO-EU operational analyses. A description of the COSMO-EU setup can be found here, but unfortunately only in German :
https://www.dwd.de/SharedDocs/downloads/DE/modelldokumentationen/nwv/cosmo_eu/cosmo_eu_dbbeschr_20140 6.pdf?__blob=publicationFile&v=3

The COSMO-EU model runs with a horizontal grid spacing of 7 km and parameterized convection. It uses a nudging data assimilation system.

Since driving data and model source code are not open source as much information as possible about the model version and how to obtain the data are given. Additionally, the location of data and model are listed for internal users.

A figshare repository, which is referenced in the paper, contains a snapshot of the Github repository and the all raw figures and log files.



Rasp, S., 2017: Variability and clustering of mid-latitude summertime convection - metadata. Figshare, accessed 28 August 2017, doi:10.6084/m9.figshare.c.3864490.
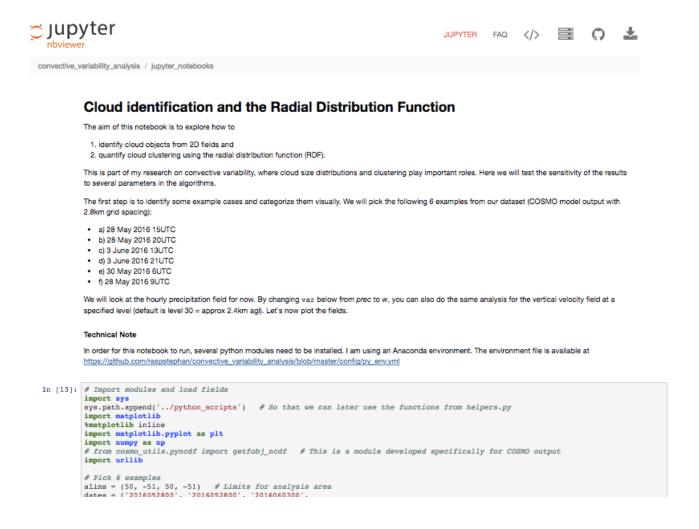
```
Plotting log

##################################################
Time: 2017-08-01T12:41:32

Executed command
------------------
python weather_time_series.py --date_start 2016052800 --date_end 2016060800 --time_start 1 --time_end 24 --nens 50
--pp_name all_days_50mem --plot_name all_days_50mem --radar_mask day --plot_type prec_cape_comp


in directory: /home/s/S.Rasp/repositories/convective_variability_analysis/python_scripts


Git hash: f20fb9f16b9a0a32569d1eceae8db271400e3c96


Full argparse parameters
----------------------------
--plot_type prec_cape_comp
--recompute False
--time_start 1
```

The log file lists the exact command executed and the computational environment. For an example of a simple log file creation in Python:
https://raspstephan.github.io/2017/08/24/reproducibility-hack.html

Jupyter notebooks allow for **literate programming**. Great for data exploration, code examples and simple analysis. In combination with a VCS literate programming can make decisions more reproducible. For R: R Markdown

Jupyter notebooks allow for **literate programming**. Great for data exploration, code examples and simple analysis. In combination with a VCS literate programming can make decisions more reproducible. For R: R Markdown
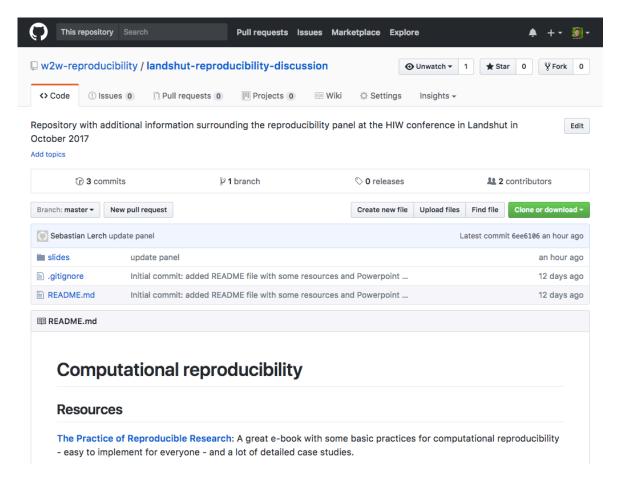
**Final thoughts from my side:**

- Few examples of reproducible research in climate and weather science.

- Initial (re)structuring of code was time-consuming.

- Striving to write reproducible code made me really think about my analysis.

- In the long run, reproducible code probably saves time, personally and for the community.

- Uncertainty how to deal with non-open-source data and model code.

Resources available at https://github.com/w2w-reproducibility/landshut-reproducibility-discussion

Link on HIW Webpage: https://hiw2017.wavestoweather.de

## Panel and open discussion:

Jenny Sun, NCAR - Julia Keller, WMO - Hannah Christensen, NCAR - Linus Magnusson, ECMWF

- Why is so little research currently computationally reproducible in the weather and climate sciences?

- What are the biggest problems related to reproducibility at the moment?

- How can these problems be tackled by the community?

- Should there be strict guidelines and requirements regarding reproducibility? Who enforces them (funding agencies, journals, universities, …)?