# Snakemake

Making data workflows easier and more reproducible

Johannes Hampp

16 March 2022

Center for International Development and Environmental Research (ZEU)
Justus-Liebig University Giessen

JUSTUS-LIEBIG-
UNIVERSITÄT
GIESSEN

Is it relevant for you and me?

Making research more FAIR?

Quick overview

Live presentation

# Is it relevant for you and me?

**Who am I?**

- B./M.Sc. in experimental physics (modelling + experiment data crunching)
- PhD student in Energy System Modelling
- Proponent of Open Source Software, Data, Access
- Co-maintainer of multiple OSS packages and models with 10+ core devs
- I've seen a lot of multi-generation models (physics, economics, energy systems)

<center>`snakemake` **turned my way of working upside down.**</center>

# Demand and supply: More than 1400 citations

Challenges (e.g. model or data pipeline for experimental data):

- Legacy work (previous PhD student, own work, student assistants, other researchers)

- Which data goes in, which data goes out? (documentation is nice, is it comprehensible and up-to-date?)

- How to execute the pipeline? (same data or new data)

- Something is not working, but what? (thousands of lines of monolithic code)

- How to extent on the work? ("I'll just put this in here…")

Source: devrant.com

5

| | |
|---|---|
| 📄 | Preprocess_Input.py |
| 📄 | Run10_Co2Constraint2019_SingleCty.py |
| 📄 | Run11_Co2Constraint2019_SingleCty.py |
| 📄 | Run12_MultiCtyMultiYear.py |
| 📄 | Run13_sensitivity_test.py |
| 📄 | Run14_sensitivity_test.py |
| 📄 | Run15_no_unmet_demand.py |
| 📄 | Run16_DAC.py |
| 📄 | Run17_PGP.py |
| 📄 | Run1_Co2Constraint2019_SingleCty.py |
| 📄 | Run2_Co2Constraint2019_SingleCty.py |
| 📄 | Run3_MultiCtyMultiYear.py |
| 📄 | Run4_From50To100.py |

"Documentation too lax; did not reproduce"?
A common workflow example

1. Order preserved by file names

2. When to run `PreProcess.py`?

3. No additional documentation (Except for paper: "We did so-and-so…")

4. Which parts do I need to run if I change external (input) data?

# Making research more FAIR?

What does that really entail?

- Repeat & Rerun
- Reproduce
- Replicate
- Reliable & Robust
- Rapport building

# Quick overview

Snakemake is a workflow management system.

It is a system to manage workflows.

# Core concept: Rules

```
rule do_research:
    input:
        # define input dependencies
        'raw_data.csv'
    output:
        # files created through this rule
        'research_results.csv'
    run:
        # your magic converting <input> to <output>
        'research.py'
```

Support for: Python, R, R Markdown, Julia, Rust, Jupyter notebooks and any shell command (!)

# Advantages (highly opiniated selection)

| What it does | How it helps |
| --- | --- |
| Human readable workflow definition | Easy and fast to learn |
| | Define (and implicitly document) dependencies |
| | Faster onboarding of new students & staff |
| Explicit dependencies | Reduces mishaps and mistakes |
| | from manual execution |
| "Rules" (Dependencies) defined and monitored | Automatic re-run if input or code is updated |
| Scales well | Independent rules run as such |
| | Rules can be kept small |
| | (good for collab., error tracking, re-running) |

# Live presentation

- Website, Docs, Tutorials, Videos, Best Practices: `https://snakemake.github.io`
- Rolling paper: `https://f1000research.com/articles/10-33/v1`
- Code from live demo: `https://github.com/euronion/snakemake-demo`
- Download and install (with Anaconda): `conda install -c bioconda snakemake`



Source: `https://snakemake.github.io`