# Workflow Tools
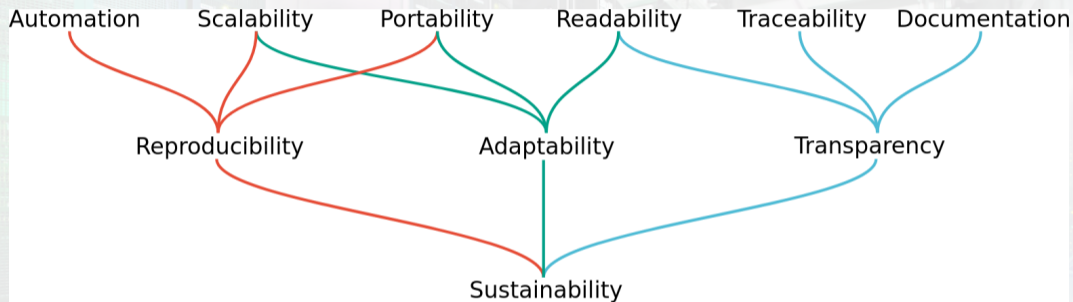
Snakemake, Galaxy, Aiida

Martin Paleico
martin-leandro.paleico@gwdg.de

June 11, 2024

hpc@gwdg.de
GWDG – Gesellschaft für wissenschaftliche Datenverarbeitung mbH Göttingen

(From: "Sustainable data analysis with Snakemake", 2021, official Snakemake Paper; and Snakemake Teaching Alliance repo)

- Snakemake paper has a good overview and categorization of different workflow approaches: f1000research.com/articles/10-33/v2
- Five "niches", three today:
  - Snakemake → Domain Specific Language
  - Galaxy → GUI based
  - Aiida → Gneric Programming Language based (Python)
- github.com/pditommaso/awesome-pipeline

**<u>Note:</u>** Snakemake and Galaxy are popular for bioinformatics, AiiDA for materials science and simulations, but any tool can in principle be used for any topic

# Snakemake

# Concept

- Snakefile scripts with defined tasks, inputs and outputs
- Snakemake identifies order of tasks and which tasks need to be rerun (*a lá* make)
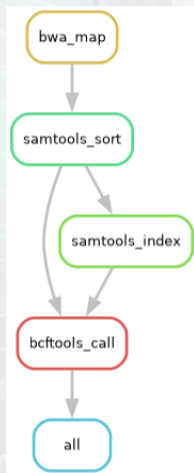
# Example

```
1  SAMPLES = ["A", "B"]
2
3  rule all:
4      input:
5          "calls/all.vcf"
6
7  rule bwa_map:
8      input:
9          "data/genome.fa",
10         "data/samples/{sample}.fastq"
11     output:
12         "mapped_reads/{sample}.bam"
13     shell:
14         "bwa mem {input} | samtools view -
     Sb - > {output}"
15 # more rules follow
```

- inputs and outputs are "promises/placeholders"
- Snakemake uses the chain of IOs to reconstruct what jobs to run and in which order
- {wildcards} are placeholders for text, variables, file names, etc.
- notice the Python syntax!

# Directed Acyclic Graph

- Good for understanding the flow and getting an overview of your job
- Useful for debugging
- Nice for publications

- SLURM (and generic) plug-ins available
- Currently need to run from a log-in node ($+$ nohup and &)
- Plug-in will submit jobs with the requested resources, submit independent jobs simultaneously, check for job completion, cleanup for failed jobs, etc.

```
1 executor: slurm
2 jobs: 2
3 default-resources:
4     mem_mb: 200
5     runtime: 100
6     slurm_partition: 'medium'
```

```
1 > snakemake —executor slurm —
      default—resources runtime
      =60 —j 5 —F —snakefile 05
      _Snakefile
2
3 Building DAG of jobs...
4 SLURM run ID: 09c94694—7362—4fe3
      —a8c5—b5c5c1eeb2f3
5 Using shell: /usr/bin/bash
6 Provided remote nodes: 5
7 Job stats:
8 job                      count
9 ——————————      ———————
10 all                      1
11 bcftools_call            1
12 bwa_map                  2
13 samtools_index           2
14 samtools_sort            2
15 total                    8
16
17 Select jobs to execute...
18 Execute 2 jobs...
19
20 [Mon Jun 10 12:08:11 2024]
21 rule bwa_map:
22     input: data/genome.fa, data/
          samples/B.fastq
23     output: mapped_reads/B.bam
24     jobid: 5
25     reason: Forced execution
26     wildcards: sample=B
27     resources: mem_mb=1000,
          mem_mib=954, disk_mb=1000,
          disk_mib=954, tmpdir=<TBD>,
          runtime=60
28
29 Job 5 has been submitted with
      SLURM jobid 247894 (log: /
      path/snakemake—tutorial—
      data—solution/.snakemake/
      slurm_logs/rule_bwa_map/B
      /247894.log).
30
31 [Mon Jun 10 12:08:12 2024]
32 rule bwa_map:
33     input: data/genome.fa, data/
          samples/A.fastq
34     output: mapped_reads/A.bam
35     jobid: 3
36     reason: Forced execution
37     wildcards: sample=A
38     resources: mem_mb=1000,
          mem_mib=954, disk_mb=1000,
          disk_mib=954, tmpdir=<TBD>,
          runtime=60
39
40 Job 3 has been submitted with
      SLURM jobid 247895 (log: /
      path/snakemake—tutorial—
      data—solution/.snakemake/
      slurm_logs/rule_bwa_map/A
      /247895.log).
41
42 [Mon Jun 10 12:11:32 2024]
43 Finished job 5.
44 1 of 8 steps (12%) done
45 [Mon Jun 10 12:11:32 2024]
46 Finished job 3.
47 2 of 8 steps (25%) done
48 Select jobs to execute...
49 Execute 2 jobs...
50
51 #etc.
```

# Pros and Cons

✓ Simple to understand

✓ Repository of workflows

✓ Widely used

✓ Plain text → easy to read at a glance

✓ Shell commands → easy to reconstruct

✓ Easy to install

✓ Easy to start with

✓ Many other features: mark temp and permanent files, mark local non-node tasks, various levels of config files, per task resource assignment, generate reports, etc.

✗ Hard to master

✗ Confusing formats (Python vs. YAML, : vs =, quote marks or not)

✗ Can be hard to debug

- Working directly on the cluster (or your own PC!)
- Some scripting, HPC, command line knowledge required
- Fast onboarding
- Similar approach: Nextflow

  More information: Snakemake Teaching Alliance (`s.gwdg.de/hYz6BZ`),
  GWDG Academy Course

# Galaxy

- Browser, GUI based
- Build a workflow by adding individual tool-based steps on a browser
- Internal DB for files (no on-system files directly)
- Any command line based tool can be turned into a Galaxy tool with a rich XML tool definition language
- Support for user quotas, workflow sharing, executing containers, HPC, etc.

# Example: Building a history

# Example: Exporting to a workflow

The following list contains each tool that was run to create the datasets in your current history. Please select those that you wish to include in the workflow.

Tools which cannot be run interactively and thus cannot be incorporated into a workflow will be shown in gray.

**Workflow name**

Workflow constructed from history 'Unnamed history'

[Create Workflow] [Check all] [Uncheck all]

| Tool | History items created |
|---|---|
| **Data Fetch** — *This tool cannot be used in workflows* | **1 genome.fa** — ☑ Treat as input dataset — genome.fa |
| **Data Fetch** — *This tool cannot be used in workflows* | **2 A.fastq** — ☑ Treat as input dataset — A.fastq |
| **Map with BWA-MEM** — ☑ Include "Map with BWA-MEM" in workflow | **3 Map with BWA-MEM on data 2 and data 1 (mapped reads in BAM format)** |
| **Samtools sort** — ☑ Include "Samtools sort" in workflow | **4 Samtools sort on data 3** |

History menu:
- You have 2 histories.
- Show Histories Side-by-Side
- Resume Paused Jobs
- Copy this History
- Delete this History
- Export Tool Citations
- Export History to File
- Archive History
- Extract Workflow
- Show Invocations
- Share or Publish
- Set Permissions
- Make Private

# Example: Editing a Workflow

# Pros and Cons

✓ GUI

✓ Tool repositories

✓ Easy to create GUI for your own tool

✓ Long development (15+ years)

✓ Easy to start with

✓ Constellation of Galaxy servers, intercompatible

✓ Links for sharing workflows, results, etc.

✓ Simple but useful Admin GUI panel

× GUI

× "Galaxy" is a very bad name for searches...

× Getting files to and from server is not so simple

× Might be hard to reconstruct DB and server if needed...

× HPC setup possible but not trivial

× Lots of configuration options

× Users can consume a lot of storage if not careful (100's of TBs)

- Can install on own PC and use instead of terminal
- Ideal for a group server setup
- Similar approach: KNIME

# AiiDA

## Concept

- Python based
- Works from your local PC or on-server, with its on quasi-scheduler daemon and command line tool (called "verdi")
- From there, jobs can run locally or on a remote HPC cluster (or another PC)
- Focused on "data provenance": Can query and filter jobs and files, create a single file archive of a group of jobs
- Can pack jobs into a workflow ("work chain")
- Local DB keeps track of jobs and files
- Use the verdi shell or script directly on Python
- Very complex

# Example: verdi shell

```
 1  # GROMACS plug-in
 2  > gmx_pdb2gmx -f 1AKI_clean.pdb -ff oplsaa -water spce -o 1AKI_forcefield.gro -p 1AKI_topology.top -i 1
        AKI_restraints.itp
 3
 4  > verdi process list -a
 5    PK  Created       Process label          Process State      Process status
 6  ----  ----------    -------------------    ----------------   ----------------
 7     4  11s ago       Pdb2gmxCalculation     Finished [0]
 8
 9  Total results: 1
10
11  # after more plug-in calls
12  > verdi process list -a
13    PK  Created       Process label          Process State      Process status
14  ----  ----------    -------------------    ----------------   -------------------------------------------
15     4  4m ago        Pdb2gmxCalculation     Finished [0]
16    12  2m ago        EditconfCalculation    Finished [0]
17    18  2m ago        SolvateCalculation     Finished [0]
18    26  1m ago        GromppCalculation      Finished [0]
19    33  53s ago       GenionCalculation      Finished [0]
20    41  9s ago        GromppCalculation      Finished [0]
21    47  4s ago        MdrunCalculation       Waiting            Monitoring scheduler: job state RUNNING
22  Total results: 7
23  Report: last time an entry changed state: 4s ago (at 12:40:52 on 2024-06-11)
24  Report: Checking daemon load... OK
25  Report: Using 0%% of the available daemon worker slots.
```

```
 1  > verdi node show 4
 2  Property       Value
 3  _____   _____
 4  type           Pdb2gmxCalculation
 5  state          Finished [0]
 6  pk             4
 7  uuid           2ab77e55-0936-47df-a6cf-9df4619ae831
 8  label
 9  description    record pdb2gmx data provenance via
        the aiida_gromacs plugin
10  ctime          2024-06-11 12:36:23.178336+02:00
11  mtime          2024-06-11 12:36:24.564897+02:00
12  computer       [1] localhost
13
```

```
14  Inputs         PK   Type
15  _____    ____ _____
16  code           1    InstalledCode
17  parameters     3    Pdb2gmxParameters
18  pdbfile        2    SinglefileData
19
20  Outputs        PK   Type
21  _____    ____ _____
22  grofile        8    SinglefileData
23  itpfile        10   SinglefileData
24  remote_folder  5    RemoteData
25  retrieved      6    FolderData
26  stdout         7    SinglefileData
27  topfile        9    SinglefileData
```

Everything has its own PK (primary key), UUID, and label

# Example: Provenance Graph

```
 1  > verdi computer setup
 2  Computer label: scc
 3  Hostname: gwdu101.gwdg.de
 4  Transport plugin: core.ssh
 5  Scheduler plugin: core.slurm
 6  Work directory on the computer [/scratch/{username
       }/aiida/]: /scratch/users/{username}/aiida/
 7  Success: Computer<3> scc created
 8
 9  > verdi -p quicksetup computer configure core.ssh
       scc
10  User name []: myusername
11  Port number [22]:
12  Look for keys [Y/n]: Y
13  SSH key file []: ~/.ssh/mysshkey
14  ...
15  Success: scc successfully configured for test@test
       .com
```

```
16
17
18  > verdi computer test scc
19  ...
20  Success: all 6 tests succeeded
21
22  > verdi code create core.code.installed
23  Computer: scc
24  Absolute filepath executable: /opt/sw/rev/23.12/
       linux-scientific7-cascadelake/gcc-11.4.0/
       gromacs-2023.3-4ehgu3/bin/gmx_mpi
25  Label: gmxmpi
26  Description: Remote GROMACS
27  Default 'CalcJob' plugin: gromacs
28  Success: Created InstalledCode<78>
29
30  > gmx_pdb2gmx --code gmxmpi@scc ...
```

# Pros and Cons

✓ Extremely flexible

✓ Launch jobs remotely

✓ Provenance

✓ Queryable databases

× Complex

× Difficult onboarding

× Have to set up plug-ins for all your tools, probably

× Not very popular/active

- Advanced users
- Handful of tools
- Long time use
- Similar approach: Covalent `covalent.readthedocs.io`

# Discussion

- Takeaways:
  - Many different paradigms
  - Pick the one you like the most
  - Local only or also remote?
  - How easily can I migrate away or retrieve my scripts and data?
  - But keep an eye on practicalities such as onboarding difficulty, availability of support, active development, etc.
- Do you know other unique workflow paradigms?
- Do you know other workflow tools within these example frameworks?
- Would you like to know more about any of these tools (workshop, articles, etc.)?