

Workflow Tools

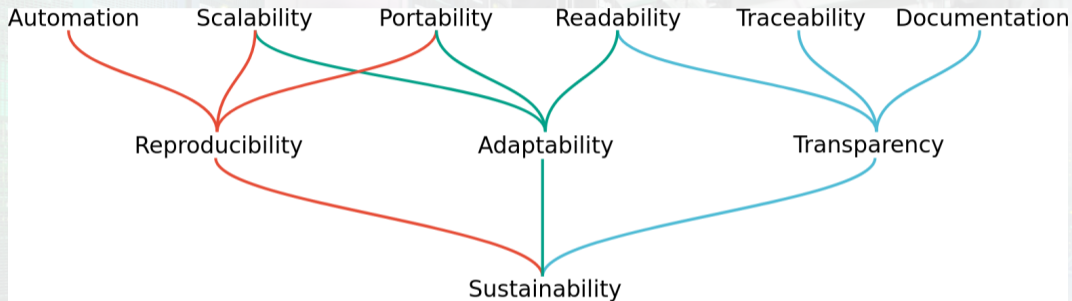
Covalent

Martin Paleico
martin-leandro.paleico@gwdg.de

July 10, 2024

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

Why Workflows?



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

More Workflows?

- Snakemake paper has a good overview and categorization of different workflow approaches: f1000research.com/articles/10-33/v2
- github.com/pditommaso/awesome-pipeline
- See previous presentation for more on Snakemake, Galaxy, Aiida

Covalent

Concept

- Similar in concept to Aiiida
- Python based - decorators
- Works from your local PC or on-server, with its own quasi-scheduler daemon and a browser based GUI
- From there, jobs can run locally or on a remote HPC cluster (or another PC)
- Python object based as the unit of communication between individual tasks
- Some rough edges

The Covalent logo is displayed in white text on a blue rectangular background. It features a stylized icon of three interconnected nodes to the left of the word "Covalent." data-bbox="379 734 625 805"/>

Covalent.

- Commercial side
- Python library is open source
- Server/cloud application for enterprises/computing centres which is not
- Never know if they'll drop support, migrate features, etc...

```
1 conda create -n covalent_env  
    python=3.8  
2 pip3.8 install covalent  
3 covalent start
```

Covalent

GitHub: <https://github.com/AgnostiqHQ/covalent>
Docs: <https://docs.covalent.xyz>

Starting Local Server

Covalent Server Configuration

Dispatcher Address	localhost
Port	48008
Develop	False
Disable Dask	False
Memory per Worker	None
Threads per Worker	None
Workers	None
Ignore Migrations	False
Disable Triggers	False
Triggers Only	False

Server Status: ✓ Running

Covalent UI can be accessed at <http://localhost:48008>

For a summary of the system status, use 'covalent status'

For additional help, use 'covalent --help'


```
1 import covalent as ct
2 # Construct manageable tasks by adding
  the @ct.electron decorator
3
4 @ct.electron
5 def add(x, y):
6     return x + y
7
8 @ct.electron
9 def multiply(x, y):
10    return x*y
11
12 @ct.electron
13 def divide(x, y):
14    return x/y
15
16
17 # Construct the workflow by stitching
18 # the electrons with the @ct.lattice
  decorator
19 @ct.lattice
20 def workflow(x, y):
21     r1 = add(x, y)
22     r2 = [multiply(r1, y) for _ in
  range(4)]
23     r3 = [divide(x, value) for value in
  r2]
24     return r3
25
26 # Dispatch the workflow
27 dispatch_id = ct.dispatch(workflow)(x
  =1, y=2)
28 result = ct.get_result(dispatch_id)
29 print(result)
```



```
1 Lattice Result
2 =====
3 status: RUNNING
4 result: None
5 input args:
6 input kwargs: 'x': 1, 'y': 2
7 error: None
8
9 start_time: 2024-07-05 11:18:06.310826
10 end_time: None
11
12 results_dir: /home/username/.cache/
13              covalent/results/bf3b6b6c-74eb-46
14              d8-9ce2-d70c77bffc09
15
16 Node Outputs
17 -----
18 add(0): None
19 :parameter:1(1): 1
20 :parameter:2(2): 2
21 multiply(3): None
22 :parameter:2(4): 2
23 multiply(5): None
24 :parameter:2(6): 2
25 multiply(7): None
26 :parameter:2(8): 2
27 (...)
28 :postprocess:reconstruct(19): None
29 :electron_list:(20): None
```

The screenshot displays a 'Dispatch list' interface with a dark theme. At the top, there are four summary cards: 'Total jobs running' (0), 'Total jobs done' (1), 'Latest running task status' (Completed), and 'Total dispatcher duration' (< 1min). Below these is a filter bar with 'All' (1), 'Running' (0), 'Completed' (1), and 'Failed' (0) buttons, along with a search input field. A table below shows a single job entry with columns for Dispatch ID, Lattice, Runtime, Started, Ended, and Status.

Dispatch ID	Lattice	Runtime	Started ↓	Ended	Status
bf3b6b6c-74eb-46d8-9ce2-d70c77bffc09	workflow	< 1min	Jul 05, 13:18:06	Jul 05, 13:18:07	21/21

bf3b6b6c...d70c77bffc09
Completed 21 / 21

Overview

Started - Ended
Jul 05, 13:18:06 - Jul 05, 13:18:07

Runtime
1s

Directory
/home/vlestrad2...9ce2-d70c77bffc09

Input
{'args': (), 'kwargs': {'x': 1, 'y': 2}}

Result
[0.16666666666666666, 0.16666666666666666, 1

Executor **dask**

```
@ct.lattice
def workflow(x, y):
    r1 = add(x, y)
    r2 = [multiply(r1, y) for _ in range(4)]
    r3 = [divide(x, value) for value in r2]
    return r3
```

Multiply

Status
Completed

Started - Ended
Jul 05, 13:18:06 - Jul 05, 13:18:07

Runtime
0s

Input
{'args': ('3', '2'), 'kwargs': {}}

Result
6

Executor **dask**

```
@ct.electron
def multiply(x, y):
    return x*y
```

```
1 # Redispatch and modify
2 @ct.electron
3 def weighted_sum_xy(x, y):
4     return 0.5 * (x + y)
5
6 redispatch_id = ct.redispatch(
7     dispatch_id ,
8     replace_electrons={'add':
9         weighted_sum_xy},
10    reuse_previous_results=True)()
11
12
13
14
```

```
15 # Cancel specific subtasks (and
16     dependent subtasks)
17 ct.cancel(dispatch_id , task_ids=[0, 2])
18 # Get and reuse results
19 result = ct.get_result(dispatch_id=
20     dispatch_id , wait=True)
21
22 # Combine decorators for Sublattices
23 @ct.electron
24 @ct.lattice
25 def cart_dist(x=0, y=0):
26     x2 = square(x)
27     y2 = square(y)
28     sum_xy = add(x2, y2)
29     return sqroot(sum_xy)
```

Other features...

- enforce software/environment dependencies
- dynamic responsive workflows
- workflow triggers

```
1 import covalent as ct
2 from pathlib import Path
3
4 #Get abs. paths
5 prot_name = "1aki"
6 init_pdb = Path('./input_files/{}.
    pdb'.format(prot_name)).resolve()
7 prot_name_upper = prot_name.upper()
8 clean_pdb = Path('./input_files/{}.
    _clean.pdb'.format(prot_name_upper
    )).resolve()
9
10 do_clean_pdb = ct.Lepton(
11     language="bash",
12     command="grep -v HOH {f_1} > {f_2}
13     ",
14     argtypes=[
15         (str, ct.Lepton.INPUT_OUTPUT),
16     ],
17     named_outputs=["f_1", "f_2"],
18 )
19
20 @ct.lattice
21 def run_gromacs(init_pdb, clean_pdb):
22     do_clean_pdb(
23         f_1=init_pdb,
24         f_2=clean_pdb,
25     )
26     return 1
27
28 dispatch_id = ct.dispatch(run_gromacs)
29 result = ct.get_result(dispatch_id,
30     wait=True)
31 print(result)
```

- Need to be careful with the location of files (absolute vs. relative paths)
- For remote execution, file transfer handlers are available
- All in all, not as suited for file handling as I/O to and from tasks
- Compare this with file-based Snakemake
- **Tasks should interact through Python objects**
- Which of course can replace files if you are smart about it

SLURM plug in

```
1 Local:
2 pip install covalent-slurm-plugin
3 # current asyncssh v. 2.15.0 breaks the plug in
4 pip install --force-reinstall -v "asyncssh==2.14.0"
5
6 Remote:
7 # need a conda env. with the same package versions:
8 conda create -n covalent_env python=3.8
9 conda activate covalent_env
10 # make sure this is the pip from the above env.
11 pip install covalent cloudpickle
```

SLURM Executor: Config

```
1 # Local: vi ~/.config/covalent/covalent.conf
2 [executors.slurm]
3 username = "remoteusername"
4 address = "glogin.hlrn.de"
5 ssh_key_file = ""
6 remote_workdir = "/home/remoteusername/covalent-workdir"
7 cache_dir = "/some/local/path"
8 conda_env = "covalent_env" #remote
9 create_unique_workdir = true
10 prerun_commands = []
11 postrun_commands = []
12 use_srun = false
13 srun_append = ""
14 bashrc_path = "$HOME/.bashrc" #remote, might error if this can't be loaded
    correctly
15 slurm_path = "/usr/bin" #remote
16 poll_freq = 30 #every this many seconds
17 cleanup = false # set to false for debugging and testing
18 ignore_versions = true # checks version compat. of remote and local packages
```

```
1 [executors.slurm.options]
2 nodes = 1
3 ntasks = 1
4 cpus-per-task = 1
5 partition = "standard96:test"
6 parsable = ""
7
8 [executors.slurm.variables]
9
10 [executors.slurm.srun_options]
```

(Thanks to Student Assistant Dhiraj Kumar)

- Can also be done at the Python level by defining an executor object
- Config file gets reprocessed by the covalent daemon
- Sometimes might need to stop and start the covalent daemon when modifying
- Config file can also be changed from the covalent browser GUI, and at the code level (`ct.get_config`, `ct.set_config`)
- Careful what option are local and which remote, particularly with paths

- Currently, can only use non-encrypted/passwordless SSH keys...
- **DO NOT** use SSH keys without passwords in our systems
- Two options: See appendix slides for more info:
 - Start an SSH tunnel with your usual encrypted SSH key, and connect to that one in the executor options (with an empty key)
 - Comment out some lines in the executor code, unlock your SSH key so it's available in ssh-agent, then use as normal

```
1 import covalent as ct
2
3 executor="slurm"
4
5 @ct.electron(executor=executor)
6 def task1(x):
7     return x + 1
8
9 @ct.electron(executor=executor)
10 def task2(x):
11     return x * 2
12
13 @ct.lattice
14 def workflow(x):
15     y = task1(x)
16     z = task2(y)
17     return z
18
19
20 dispatch_id = ct.dispatch(workflow)(2)
21 output = ct.get_result(dispatch_id=
22     dispatch_id, wait=True)
23 print(output)
```



```
1 /home/remoteusername/covalent-workdir/bd7a129d-027c-48d7-a527-ecd49a47c067/  
   node_0  
2 slurm-bd7a129d-027c-48d7-a527-ecd49a47c067 -0.sh  
3 func-bd7a129d-027c-48d7-a527-ecd49a47c067 -0.pkl  
4 result-bd7a129d-027c-48d7-a527-ecd49a47c067 -0.pkl  
5 script-bd7a129d-027c-48d7-a527-ecd49a47c067 -0.py  
6 stdout-bd7a129d-027c-48d7-a527-ecd49a47c067 -0.log  
7 stderr-bd7a129d-027c-48d7-a527-ecd49a47c067 -0.log
```

- slurm.sh: Actually submitted SLURM script (check that your SBATCH directives are correct, if the conda env. is being correctly loaded, etc.)
- script.py: Called from slurm.sh
- func.pkl, result.pkl: serialized Python objects and functions, your electron and its results are contained here
- stdout.log, stderr.log: Usual SLURM log files

Discussion

- Snakemake
 - Script based
 - Local only
 - Robust, developed
- Galaxy
 - GUI/browser based
 - Robust, developed
 - Best as a hosted server
- Covalent
 - Python-based
 - Remote execution possible
 - Best used for pure Python approaches

- Basics
 - <https://covalent.readthedocs.io/en/latest/concepts/basics.html>
 - <https://docs.covalent.xyz/docs/get-started/quick-start>
- Interesting applied tutorials
 - Electronic Structure Calculation
https://docs.covalent.xyz/docs/cloud/tutorials-cloud/simple_energy/
 - MNIST
<https://docs.covalent.xyz/docs/user-documentation/tutorials/mnist>
 - Many AI/ML, and Quantum Computing examples also available

Appendix

- Create 2 new hosts in your local `.ssh/config` as in the following slide
- SSH to the first host, the tunnel, don't close the connection; should ask you for your usual authentication method
- (There is also a way of getting the tunnel to go into the background instead of maintaining an open connection)
- Configure covalent SLURM plug in to connect to the second host, point to a fake unused key file (or make modifications to the plugin code similar to method 2)
- To test, the setup, open the tunnel, then attempt to connect to the second host with SSH, should open a session without needing to input any password:
`ssh arbitraryuname@localhost -p 1234`


```
1 Host covalenttunnel
2 Hostname glogin.hlrn.de
3 User remoteuname
4 IdentityFile ~/.ssh/your-usual-ssh-key
5 IdentitiesOnly yes
6 LocalForward 1234 localhost:22 #L option on the command line
7 ServerAliveInterval 120
8
9 Host covalentssh
10 Hostname localhost
11 User arbitraryuname
12 Port 1234
```


- Modify the installed plug in code (`/some/path/python3.8/site-packages/covalent_slurm_plugin/slurm.py`) at around line 215 as in next slide
- Write the plug in configuration file as you would correctly (correct log in node address, username, and your usual, encrypted SSH key)
- Unlock your encrypted SSH key (many ways, but a normal SSH into the remote server will unlock the key for the duration of your current session and pass it to ssh-agent)
- Test the covalent SLURM plug in normally
- This works because `asynccssh` by default can retrieve keys from ssh-agent, no need to force an SSH key like the plug in does, so just let `asynccssh` do its thing
- Tested on Ubuntu LTS 22.04, might not work in other OS's

```
1  #if not self.ssh_key_file:
2  #     raise ValueError("ssh_key_file is a required parameter in the Slurm plugin.")
3
4  if self.cert_file and not Path(self.cert_file).exists():
5      raise FileNotFoundError(f"Certificate file not found: {self.cert_file}")
6
7  #if not Path(self.ssh_key_file).exists():
8  #     raise FileNotFoundError(f"SSH key file not found: {self.ssh_key_file}")
9
10 if self.cert_file:
11     client_keys = [
12         (
13             asyncssh.read_private_key(self.ssh_key_file),
14             asyncssh.read_certificate(self.cert_file),
15         )
16     ]
17 else:
18     #client_keys = [asyncssh.read_private_key(self.ssh_key_file)]
19     client_keys=[] #blank key so it goes to ssh-agent
```