

Embarrassingly Parallel workflows and how to execute them

Embarrassingly parallel workflows

What?

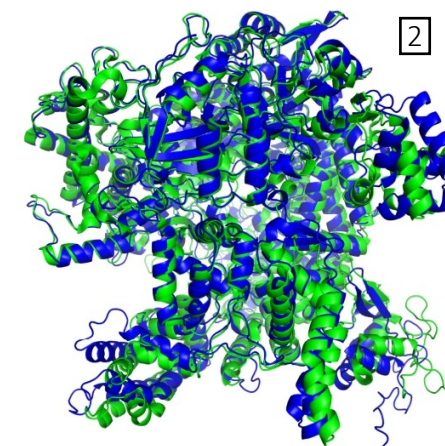
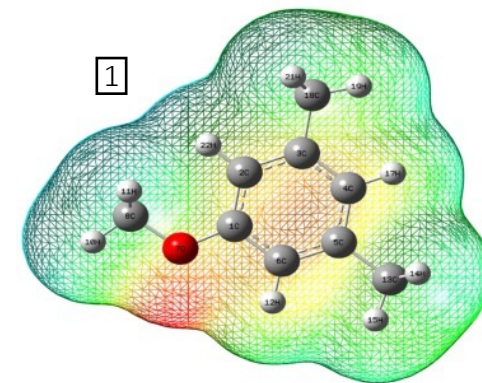
- Executing many independent small jobs simultaneously to use the hardware resource efficiently.

When?

- It is not possible to combine multiple runs/tasks into a single application.
- Large parameter sweep based on the same application or workflow with dependent pipeline steps.
- Workflow involves runs that are depending on several independent software which do not scale beyond a certain number of cores.

Why?

- Nodes are getting denser with respect to resources such as number of cores/threads, GPU cards, Memory.
- Applications do not necessarily scale even within a node: shared memory (OpenMP, MPI SHM, pthreads)



How?

Tools

- Linux background processes (within a node) along with **numactl**
 - Supported on Snellius and requires prior knowledge of bash scripting and process placement using **numactl**.
 - Gives most manual control if one knows what they are doing.
- SLURM arrays (to perform similar jobs on multiple allocations)
 - Supported on Snellius and requires prior knowledge of SLURM array options.
 - Is limited to cgroups supported by SLURM which means each allocation cannot be below the minimum allocation within SLURM. (Snellius -> 16 cores and 28 GiB of memory on “rome” partition)
 - One has to wait for each task in the array to get allocated and finish.
- STOPOS
 - Not supported in Snellius anymore. A tool to execute embarrassingly parallel jobs, requires prior knowledge of bash scripting.
- QCG-PilotJob (<https://qcg-pilotjob.readthedocs.io/en/develop/>)
 - Supported on Snellius and requires prior knowledge of Python scripting.



| There are many more tools



QCGPilot-Job

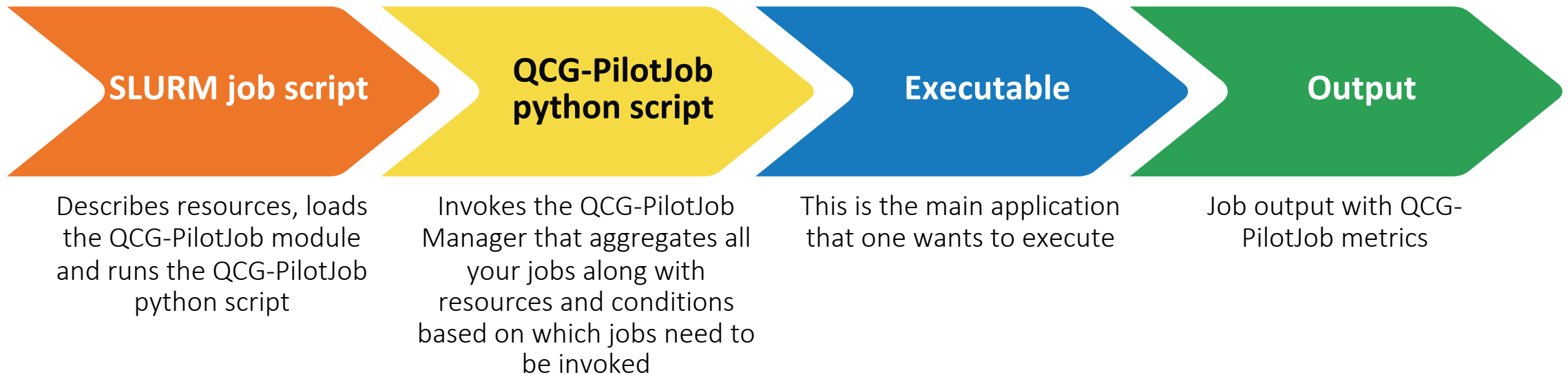
Job manager inside SLURM:

- Submits jobs within a single large allocation across multiple nodes.
- Does not have to wait for small individual allocations like in SLURM arrays.
- Tries to optimize resources for each job based on requirements such that the queue is always full.
- Possibility of restarting failed/timed out jobs.
- One can also specify complex result-based job dependencies.

Some caveats and drawbacks:

- One needs to estimate, the total time of all the jobs combined.
- Job/system state files get produced in the submission directory.
- Might bring some complexity into execution of a Singularity containers with embedded MPI-based codes.
- In the case of problems, the log files may not be very expressive.

Workflow

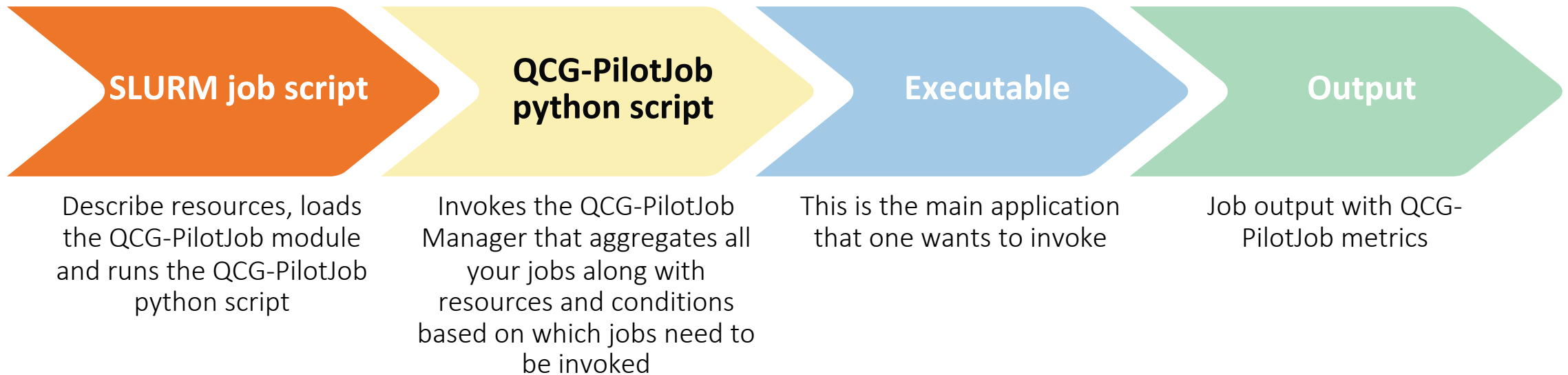


QCGPilot-Job: Set up the environment

- The **QCG-PilotJob** package is available as a module in the **2022** environment
- Note that the environment from the login node is ***not transferred*** to the compute nodes, unless defined in, e.g., the **.bashrc** file. Thus, you have to re-load the QCG-PilotJob package in your job script.

```
$ module purge  
$ module load 2023  
$ module load QCG-PilotJob/0.14.1-foss-2023a
```

Workflow



QCGPilot-Job: SLURM job script

run.sh

```
#!/bin/bash
```

```
#SBATCH -N 1
```

```
#SBATCH -n 16
```

```
#SBATCH -p rome
```

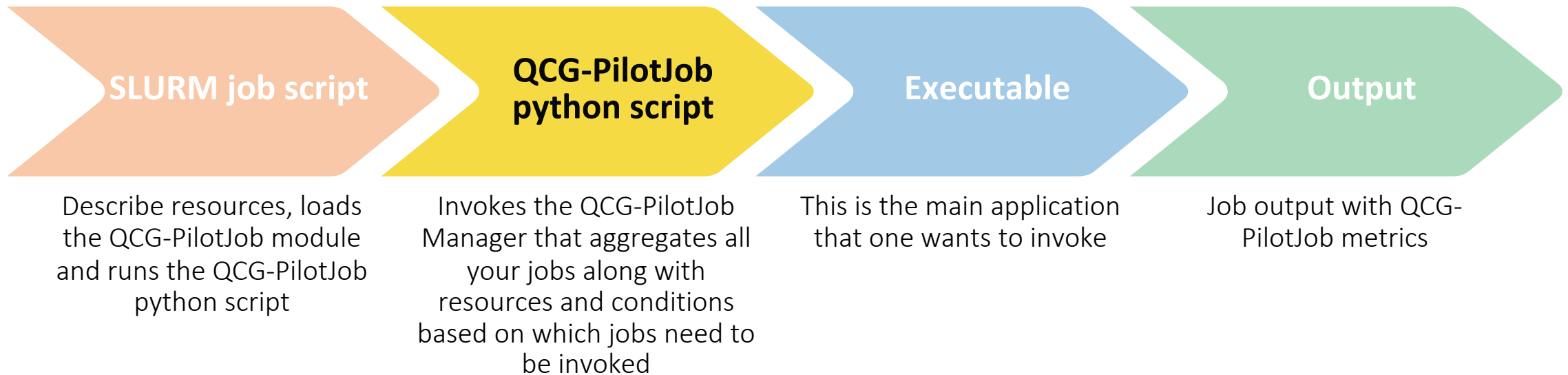
```
#SBATCH -t 10
```

```
module load 2023
```

```
module load QCG-PilotJob/0.14.1-foss-2023a
```

```
python test.py
```

Workflow



QCGPilot-Job: Python script

- Work with the **QCG-PilotJob** manager consists of a few mandatory steps:
 - Creation of an instance of **QCG-PilotJob Manager**
 - Creation of a class object
 - Waiting for all submitted jobs to finish
 - will wait until some or all submitted to the QCG-PilotJob Manager jobs finish
 - Call for the **finish()** method at the end of the script
 - this method terminates the background thread in which the instance of a QCG-PilotJob Manager has been run

QCGPilot-Job: Python script

test.py

```
from qcg.pilotjob.api.manager import LocalManager

manager = LocalManager()
print('available resources: ', manager.resources())
manager.wait4all()
manager.finish()
```

Example output

```
available resources: {'total_nodes': 1,
'total_cores': 96, 'used_cores': 0, 'free_cores': 96}
```

QCGPilot-Job: Python script

test.py

Create an object

Wait for all jobs to finish

Finalize the execution

```
from qcg.pilotjob.api.manager import LocalManager

manager = LocalManager()
print('available resources: ', manager.resources())
manager.wait4all()
manager.finish()
```

Example output

```
available resources: {'total_nodes': 1,
'total_cores': 96, 'used_cores': 0, 'free_cores': 96}
```

QCGPilot-Job: Python script

test.py

Create an object

Wait for all jobs to finish

Finalize the execution

```
from qcg.pilotjob.api.manager import LocalManager

manager = LocalManager()
print('available resources: ', manager.resources())
manager.wait4all()
manager.finish()
```

Example output

```
available resources: {'total_nodes': 1,
'total_cores': 96, 'used_cores': 0, 'free_cores': 96}
```

Accessible cores

Total number of nodes and cores within an allocation

QCGPilot-Job: Create jobs

- To execute multiple jobs within a single SLURM allocation we need to:
 - Create an object of the **Jobs** class
 - Populate it with some parameters (describe jobs)
 - Submit the created object to **QCG-PilotJob Manager**

```
...  
from qcg.pilotjob.api.job import Jobs  
  
jobs = Jobs()  
jobs.add(...)  
jobs_ids = manager.submit(jobs)  
print('submitted jobs: ', str(job_ids))  
...
```

QCGPilot-Job: Job description (main parameters)

- **name** - the job name
- **exec** - path to the executable program
- **args** - executable program arguments
- **script** - bash script content
- **stdin** - path to file which content should be passed to the standard input stream
- **stdout** - path to the file where standard output stream should be saved
- **stderr** - path to the file where standard error stream should be saved
- **modules** - list of modules that should be loaded before job start
- **numCores** - number of required cores specification
- **numNodes** - number of required nodes specification
- **iteration** - iterations definition
- **model** - model of execution

QCGPilot-Job: Job description (main parameters)

- **name** - the job name
- **exec** - path to the executable program
- **args** - executable program arguments
- **script** - bash script content
- **stdin** - path to file which content should be passed to the standard input stream
- **stdout** - path to the file where standard output stream should be saved
- **stderr** - path to the file where standard error stream should be saved
- **modules** - list of modules that should be loaded before job start
- **numCores** - number of required cores specification
- **numNodes** - number of required nodes specification
- **iteration** - iterations definition
- **model** - model of execution

QCGPilot-Job: execution models

- **default** - only a single process is launched within the allocation
- **threads** - is designed for running OpenMP tasks on a single node
- **openmpi** - the processes are started with the mpirun command
- **intelmpi** - same but for Intel MPI
- **srunmpi** - the processes are started with the srun command

Read more about execution models in the official docs:

https://qcg-pilotjob.readthedocs.io/en/develop/execution_models.html

QCGPilot-Job: Calling for an executable

test.py

```
from qcg.pilotjob.api.manager import LocalManager
from qcg.pilotjob.api.job import Jobs

manager = LocalManager()
jobs = Jobs()

jobs.add(script='echo "job ${it} executed at `date` @
`hostname`"',
         name='test_job',
         stdout='job.out.${it}',
         iteration=4)

job_ids = manager.submit(jobs)
print('submitted jobs: ', str(job_ids))
job_status = manager.status(job_ids)
print('job status: ', job_status)

manager.wait4all()
manager.finish()
```

| Hands-on

Hands-on

Clone the repository

- <https://github.com/sara-nl/course-module-qcg-pilotjob>

```
$ cd ~  
$ git clone https://github.com/sara-nl/course-module-qcg-pilotjob  
$ cd course-module-qcg-pilotjob
```

Hands-on #1

Task

- Execute **512 independent short jobs** on two nodes on the “rome” partition. Note that one node on the “rome” partition has 128 cores.

Steps

- Go to **hands-on/example_01**
- Take a look at **qcg_job.py** and **qcg_job.sh**
- Submit the SLURM job script, check the results

```
from qcg.pilotjob.api.job import Jobs
from qcg.pilotjob.api.manager import LocalManager

# create the QCG manager and the Jobs object
manager = LocalManager()
jobs = Jobs()

# create 512 independent jobs
for job in range(512):
    # define a job
    jobs.add(name="job_{}".format(job),
             exec='hostname',
             stdout='output/job.{}.out'.format(job),
             stderr='output/job.{}.err'.format(job),
             model='default',
             numCores={ "exact" : 1},
             iteration=1)

print("-- submit jobs --")
manager.submit(jobs)
print("-- wait for all jobs --")
manager.wait4all()
manager.finish()
print("-- finished")
```

Hands-on #2

Task

- Compute the average of one column in 1000 CSV files. Aggregate the results in one CSV file.

Steps

- Go to **hands-on/example_02**
- Take a look at **README.md**, **qcg_job_v1.py** and **qcg_job.sh**
- Extract the input data (**input.tar.gz**)
- Submit the SLURM job script, check the results

```
job_names = []

# scan input directory for CVS files
for filename in glob(os.path.join("input/*.csv"),
recursive = False):
    # determine the job name and append it to the list
    job_names.append(os.path.basename(filename))

# create the QCG manager and the Jobs object
manager = LocalManager()
jobs = Jobs()

# loop over all files and populate the list of jobs
for jname in job_names:
    print("submit {}".format(jname))
    jobs.add(name=jname,
            exec='python3',
            args=["average.py",
                os.path.join("input", jname)],
            stdout='average_{}'.format(jname),
            stderr='job.{}.err'.format(jname),
            modules=[
                "2023", "Python/3.11.3-GCCcore-12.3.0"],
            iteration=1
            )

# add a job for the aggregation of the results
jobs.add(name="aggregate",
        script='cat average_*.csv | sort',
        stdout='result.csv',
        stderr='aggregate.err',
        after=job_names)
```



SURF

Have fun! :)