

Electronic Supporting Information:  
PerQueue: Managing Complex and Dynamic Workflows

Benjamin H. Sjølin, William S. Hansen, Armando A. Morin-Martinez, Martin H. Petersen,  
Laura H. Rieger, Tejs Vegge, Juan M. García-Lastra, and Ivano E. Castelli<sup>‡</sup>

*Department of Energy Storage and Conversion, Technical University of Denmark, Anker  
Engelunds Vej 301, DK-2800 Kongens Lyngby, Denmark.*

<sup>‡</sup>Corresponding author email: [ivca@dtu.dk](mailto:ivca@dtu.dk)

# 1 Submission Example Codes

Below are the code scripts for the workflow presented in Lst. 1 in the main text. The workflow uses a Monte Carlo method to estimate  $\pi$  and uses this value to calculate the proposed diameter of a perfect sphere made of silicon-28 weighing exactly 1 kilogram.

## 1.1 Estimation of $\pi$

This code uses a Monte Carlo method to estimate the value of  $\pi$ . This is done by randomly generating points in the 2D plane such that the x- and y-coordinate is in the range  $[0;1[$ . Dividing the fraction of points where  $x^2 + y^2 \leq 1$  by the total number of generated points and multiplying by 4 gives the estimated value of  $\pi$ .

This script estimates  $\pi$  using `N_samples` generated points and returns the estimate for each order of magnitude of points.

---

```
1 from typing import Tuple
2
3 import numpy as np
4
5
6 def main(N_samples: int = 1, **kwargs) -> Tuple[bool, dict]:
7     """Compute an estimate of pi using Monte Carlo sampling
8
9     Parameters
10    -----
11    N_samples : int
12        The number of samples to generate for the estimate
13    """
14    # Seed the random number generator
15    np.random.seed(1023)
16
17    # Generate N points in the x-y plane
18    X = np.random.random_sample((2, N_samples))
19
20    # Calculate the squared distance from origo of each point and determine
21    # which points are inside the unit circle
22    inside_mask = (X**2).sum(axis=0) <= 1
23
24    # Calculate the cumulative number of points inside the unit circle
25    cs = np.cumsum(inside_mask)
26
27    # Calculate the estimate of pi in logarithmic intervals
28    estimates = np.array([
29        4 * cs[10**i - 1] / 10**i
30        for i in range(np.log10(N_samples).astype(int) + 1)
31    ])
32
33    return True, {"pi_estimates": estimates, 'N_samples': N_samples}
```

---

## 1.2 Plotting of the Absolute Error to the Value of $\pi$

This script finds the absolute error between the estimates of  $\pi$  and the established value of  $\pi$ . It then generates a log-log plot of the error versus the number of points. Depending on the value of the parameter `savefig_filename`, the plot is either saved to file or shown to the user.

---

```
1 from typing import List, Tuple, Union
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 # Increase resolution of the saved figure
7 plt.rcParams["savefig.dpi"] = 300
8
9
10 def main(
```

```

11     pi_estimates: Union[List[float], np.ndarray] = None, N_samples: int = 1,
12     savefig_filename: str = None, **kwargs
13 ) -> Tuple[bool, dict]:
14     """Plot the absolute error of the given estimates of pi
15
16     Compute the absolute error between the given values of pi and the value
17     given by 'numpy.pi'.
18
19     Parameters
20     -----
21     pi_estimates : list of float or 1D 'numpy.ndarray'
22         A list of estimates for the value of the constant pi
23     N_samples : int
24         The number of samples used to generate the last estimate
25     savefig_filename : str, optional
26         If this parameter is not 'None', the generated graph is saved to the
27         filename and path given. Otherwise, the graph is shown as an
28         interactive window.
29     """
30     if pi_estimates is None or len(pi_estimates) == 0:
31         raise ValueError("Empty list of pi estimates given")
32
33     # Compute the x-tick locations
34     n = [10**i for i in range(np.log10(N_samples).astype(int) + 1)]
35
36     # Create figure and plot the absolute error as a log-log plot
37     fig, ax = plt.subplots()
38     ax.loglog(n, np.abs(np.array(pi_estimates) - np.pi))
39     ax.set(
40         xlabel="Number of generated points",
41         ylabel="Absolute error from value of \"numpy.pi\""
42     )
43     ax.grid()
44
45     # Either save or show the graph
46     if savefig_filename is not None:
47         fig.savefig(savefig_filename)
48     else:
49         plt.show()
50
51     return True, None

```

---

### 1.3 Calculation of the Diameter of the Silicon Sphere

The script below calculates the diameter of a perfect sphere made of silicon-28 using the estimated values of  $\pi$ . The sphere is assumed to weigh 1 kilogram. It returns both the diameters that use the estimates of  $\pi$  and one that uses the established value of  $\pi$ .

```

1 from typing import List, Tuple, Union
2
3 import numpy as np
4
5
6 def main(
7     pi_estimates: Union[List[float], np.ndarray] = None, **kwargs
8 ) -> Tuple[bool, dict]:
9     """Compute the diameter of the newest prototype for the kilogram
10
11     Using the estimates for pi that we've computed earlier, we want to compute
12     the diameter of the silicon sphere that works as a prototype for the
13     kilogram after the 2018 redefinition of the SI-system. Given that we know
14     the weight of the sphere, we can compute its diameter from the basic
15     definitions and constants of the revised SI-system.
16
17     We shall assume that the silicon in the sphere is only silicon-28, but the
18     amount of other impurities in the bulk and on the surface are still

```

```

19 included. These values and the method for computation is taken from
20 https://doi.org/10.1016/j.crhy.2018.12.005.
21
22 Parameters
23 -----
24 pi_estimates : list of float or 1D 'numpy.ndarray'
25     A list of estimates for the value of the constant pi
26     """
27 if pi_estimates is None or len(pi_estimates) == 0:
28     raise ValueError("Empty list of pi estimates given")
29
30 # Define physical constants
31 planck = 6.62607015e-34 # Planck's constant [J/Hz]
32 Rydberg = 10973731.568157 # Rydberg constant [1/m]
33 fine_structure = 7.2973525643e-3 # fine structure constant [-]
34 c = 299792458 # speed of light [m/s]
35 Ar_Si_28 = 27.9769265 # Relative molecular weight of Si-28 [-]
36 Ar_e = 5.489e-4 # Relative molecular weight of an electron [-]
37 lattice_parameter = 5.43102051e-10 # Si-28 lattice parameter [m]
38
39 # We know the target mass of the silicon sphere.
40 target_mass = 1.0 # Target mass of the sphere [kg]
41
42 # Define the mass deficits of the bulk and surface layer
43 m_SL = (7.1 + 12.0 + 58.2) * 1e-9 # Surface layer mass
44 m_deficit = (17.1 - 2.3 - 0.5 + 6.0) * 1e-9 # Bulk mass
45
46 # Compute the rest energy of an electron
47 m_e = 2*planck*Rydberg / (c*fine_structure**2)
48
49 # Compute the core volume of the silicon sphere
50 V = (target_mass - m_SL + m_deficit) * Ar_e * \
51     lattice_parameter**3 / (m_e * Ar_Si_28 * 8) # [m^3]
52
53 # Add the 'numpy' value for pi to the end of the array
54 pi_x = np.array([*pi_estimates, np.pi])
55
56 # Compute the core diameter of the sphere using the estimates of pi
57 # Final sphere diameters [mm]
58 sphere_diameters = np.cbrt(6 * V / np.array(pi_x)) * 1e3
59
60 return True, {
61     "sphere_diameters": sphere_diameters[:-1],
62     "actual_sphere_diameter": sphere_diameters[-1]
63 }

```

---

## 2 Use Case Workflow Codes

Here, we present the submission scripts used for the four use cases presented in Section 3. The contents of the task scripts are not presented herein, but they are available upon reasonable request. As a consequence of this, the arguments to the tasks herein have generally been replaced by `{...}` for brevity.

### 2.1 High-Throughput Screening

The following submission script shows how the high-throughput screening workflow, described in the main paper, is set up in PerQueue terms. It should be noted that due to the size of the screening (more than 6000 instances of the workflow) the submission script takes two inputs,  $j$  and  $k$ . These limit the index into the search space for decorations, allowing the user to submit a subset of the full workflow, such that other projects can use resources in parallel to this study.

Running the central part of the workflow in three parallel sub-workflows is achieved by wrapping tasks `t2-t6` in a `StaticWidthGroup` with a `width` of 3.

Since the convex hull and the band gap are both post-processing steps in this study they are fast to compute and are run on `local`.

---

```
1 from sys import argv
2
3 from perqueue import PersistentQueue, StaticWidthGroup, Task, Workflow
4
5 j = int(argv[1])
6 k = int(argv[2])
7 WIDTH = 3
8
9 # Defining the tasks for a specific range of entries.
10 for i in range(j, k):
11     # Define tasks
12     t1 = Task("generation.py", {'index': i}, "local:10m")
13
14     t2 = Task("relax.py", None, "40:xeon40:1:50h")
15     t3 = Task("convex_hull.py", None, "local:5m")
16     t4 = Task("band_gap.py", None, "local:5m")
17     t5 = Task("preneb.py", {...}, "112:xeon56:1:50h")
18     t6 = Task("neb.py", {...}, "168:xeon56:1:50h")
19
20     # Define the subworkflow as a list of tasks - each depends on the previous
21     swf = Workflow([t2, t3, t4, t5, t6])
22
23     # Wrap subworkflow in width group to get parallel workflows
24     swg = StaticWidthGroup(swf, width=WIDTH)
25
26     # Define final workflow layer
27     wf = Workflow([t1, swg])
28
29     # Submit entire workflow through PerQueue
30     with PersistentQueue() as pq:
31         pq.submit(wf)
```

---

### 2.2 Active Learning for MLIPs

The script presented below is that used for submitting the workflow for training machine-learned interatomic potentials using active learning. As shown, the complex workflow is distilled down into four distinct tasks, and the dynamic nature arises from wrapping these in sub-workflows inside `Static-/DynamicWidthGroups` and a `CyclicalGroup`.

The training, simulation and selection tasks all utilize GPU resources (through the `sm3090` resource), and the labeling runs Density Functional Theory calculations on a 24-core CPU resource.

To control when to break out of the `CyclicalGroup`, the `t_train` task returns the PerQueue constant `CYCLICALGROUP_KEY` with a value of `True` for stopping the loop or `False` for continuing to iterate.

---

```

1 from perqueue import CyclicalGroup, DynamicWidthGroup, PersistentQueue,
   StaticWidthGroup, Task, Workflow
2
3 # Define tasks
4 t_train = Task('work_train.py', None, '1:sm3090:30m')
5 t_sim = Task('work_simulate.py', None, '1:sm3090:30m')
6 t_select = Task('work_select.py', None, '1:sm3090:30m')
7 t_label = Task('work_label.py', None, '24:xeon24:10m')
8
9 # Define groups for workflow width
10 swg_train = StaticWidthGroup(t_train, width=2)
11 dwg_ssl = DynamicWidthGroup([t_sim, t_select, t_label])
12 dwg_train = DynamicWidthGroup(t_train)
13
14 # Wrap width groups in a CyclicalGroup for looping
15 cg = CyclicalGroup([dwg_ssl, dwg_train], max_tries=10)
16
17 # Define and submit workflow
18 wf = Workflow([swg_train, cg])
19
20 # Submit the workflow through PerQueue
21 with PersistentQueue() as pq:
22     pq.submit(wf)

```

---

## 2.3 Cluster Expansion

Below, we present the submission script used for the Cluster Expansion workflow explained in the main paper. Not much is different here from the previous use cases, which speaks to the simplicity for setting up workflows in PerQueue. Once the workflow structure is defined, converting it to PerQueue constructs yields rather simple code.

Here, we show the use of the optional `name` parameter to a `Task`, which gives it a different name for visualization and searching.

---

```

1 from perqueue import CyclicalGroup, DynamicWidthGroup, PersistentQueue,
   StaticWidthGroup, Task, Workflow
2
3 # Define tasks
4 CE_task_init = Task('CE_model.py', {...}, 'local:10m')
5 CE_task = Task('CE_model.py', {...}, '24:xeon24:10m', name='train')
6 relax_task = Task('relaxation.py', {...}, '8:sm3090:3h', name='optimize')
7 MC_task = Task('MC.py', {...}, '24:xeon24:10m')
8 KMC_task = Task('KMC.py', {...}, '24:xeon24:10m')
9
10 # Wrap up subworkflow for loop
11 dwg = DynamicWidthGroup(relax_task)
12 cg = CyclicalGroup([dwg, CE_task], max_tries=10)
13
14 # Wrap (kinetic) Monte Carlo in width group
15 swg = StaticWidthGroup([MC_task, KMC_task], width=5)
16
17 # Package up workflow
18 wf = Workflow([CE_task_init, cg, swg])
19
20 # Submit the workflow through PerQueue
21 with PersistentQueue() as pq:
22     pq.submit(wf)

```

---

## 2.4 Active Learning for Image Segmentation

While the image segmentation workflow should be the hardest to express due to the human-in-the-loop constraint, the submission script is the most simple with only three `Tasks` that are connected through a `CyclicalGroup` and wrapped in a workflow. Currently, human-in-the-loop is achieved by purposefully failing

the `t3` task after writing output to the human and restarting the `Entry` once data has been written to a pre-specified file for the task to read from. In future versions of `PerQueue` the intention is for this to be a more graceful maneuver.

The key to starting with the training step and have it be at the end of the `CyclicalGroup` is the ordering of the `Task` objects in line 9.

---

```
1 from perqueue import CyclicalGroup, PersistentQueue, Task, Workflow
2
3 # Define tasks
4 t1 = Task("train.py", {...}, "1:sm3090:2h")
5 t2 = Task("predict_select.py", None, "1:sm3090:30m")
6 t3 = Task("label.py", None, "1:sm3090:10m")
7
8 # Wrap tasks in CyclicalGroup for looping
9 cg = CyclicalGroup([t2, t3, t1], max_tries=10)
10
11 # Define full workflow
12 wf = Workflow([t1, cg])
13
14 # Submit the workflow through PerQueue
15 with PersistentQueue() as pq:
16     pq.submit(wf)
```

---