

Asynchronous execution of heterogeneous tasks in ML-driven HPC workflows

Paper by Pascuzzi et al.

Presentation by David Bunde

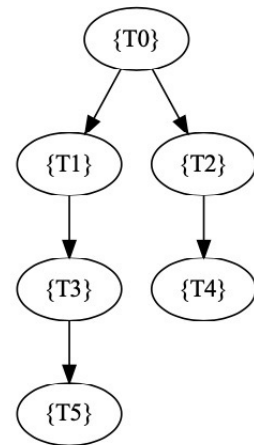
Setting

Types of tasks:

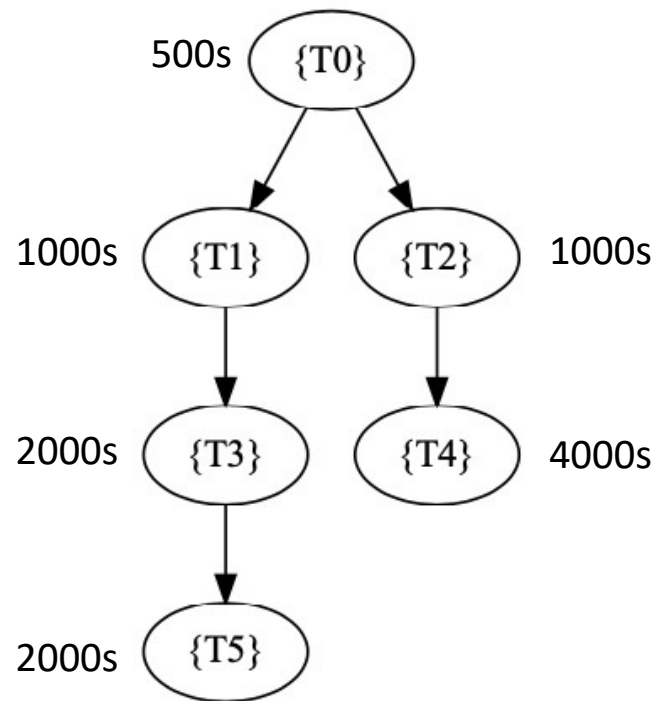
- Asynchronous: can execute independently of each other (top of pg 9)
- Synchronous: have dependencies or that interact (?)
- Concurrent: execute at the same time (top of pg 9)

Resources: bounded numbers of CPUs and GPUs

Application: ML pipelines with simulations, aggregation, training, and inference stages (and interactions)



Benefit of asynchronicity



Sequential time:

$$500 + 1000 + 2000 + 2000 + 1000 + 4000 = 7500s$$

Using asynchronicity:

$$500 + \max(1000 + 2000 + 2000, 1000 + 4000) = 5500$$

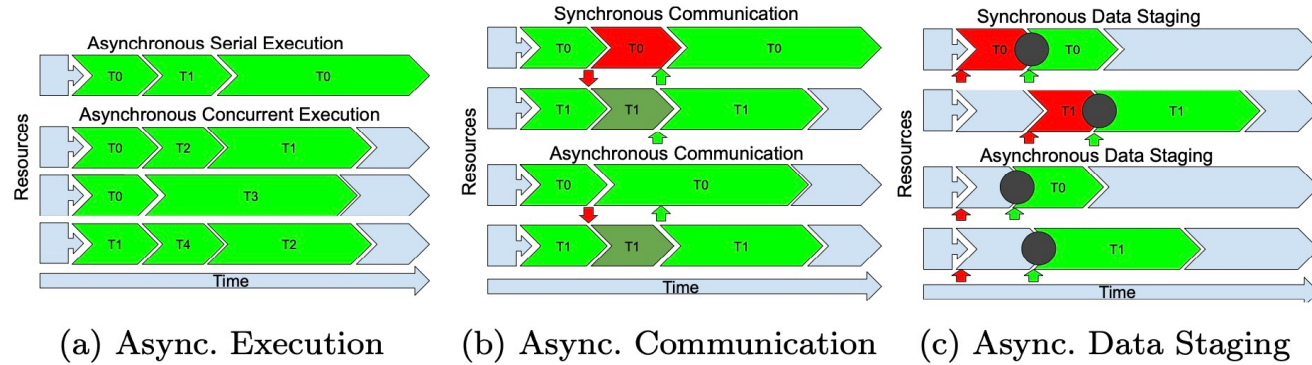


Fig. 1: Types of Asynchronicity. Fig 1a shows two asynchronous executions: Top, T0 and T1 execute asynchronously and serially; bottom, T0 and T1 execute asynchronously and concurrently and T0 and T2 execute synchronously. Note that we used DG in Figure 2c, assuming that T0 and T1 do not have any dependencies and T2 depends on T0. Fig 1b shows asynchronous and synchronous communication: **red** arrows are sent messages and **green** arrows are received messages; **red** bars indicate synchronous tasks that are waiting for their messages to be exchanged before resuming execution. and **dark green** bar is the time taken to execute incoming request before replying. Fig 1c shows asynchronous and synchronous data staging: **red** arrows show when data staging started and **green** arrows show when data is ready; **red** bars indicate synchronous tasks that cannot be executed before data become available.

Contributions

- An asynchronous implementation of DeepDriveMD, a framework to execute ML-driven molecular-dynamics workflows on HPC platforms at scale
- A performance evaluation of asynchronous DeepDriveMD
- A model of asynchronous behavior
- A general performance evaluation of that model for workflows with varying degrees of asynchronous execution