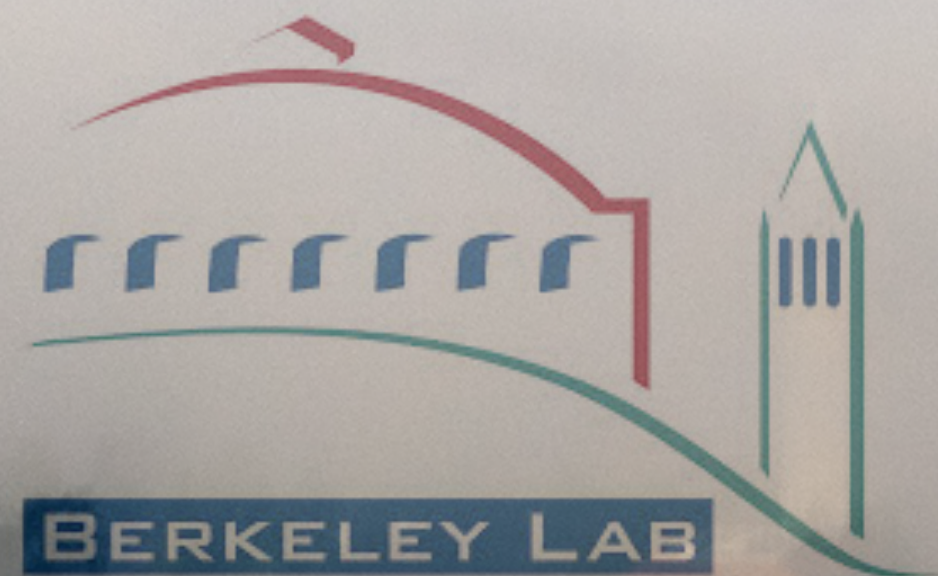


# MPI Job Manager (MPI\_JM)

Lattice 2022

André Walker-Loud





# MPI Job Manager

---

❑ [Ken McElvain](#) (architect), André Walker-Loud, [Evan Berkowitz]

❑ What is the issue that this code is addressing?

❑ How does it work? (high level)

❑ What are future features?

❑ How and when can you get it?



# What is the issue?

- ❑ The “measurement” stage of lattice QCD calculations requires the execution of millions of independent MPI tasks
- ❑ The tasks have mixed resource needs on heterogeneous systems
  - ❑ GPU-intense, CPU-only, heavy I/O
- ❑ The tasks often have chained/nested dependencies
  - ❑ A coherent sequential propagator for a 3-pt function can require 8 previous propagators to make a coherent sink
- ❑ A team may desire to run multiple projects with the same computing allocation with dynamic priorities (project A is higher priority generally, but you need results for project B for an upcoming seminar/conference...)
- ❑ Running this “swarm” of millions of MPI tasks as individual computations on “leadership computing facilities” is not tolerable by the supercomputing centers



---

# Jargon

---

## □ job

a resource allocation on the compute cluster

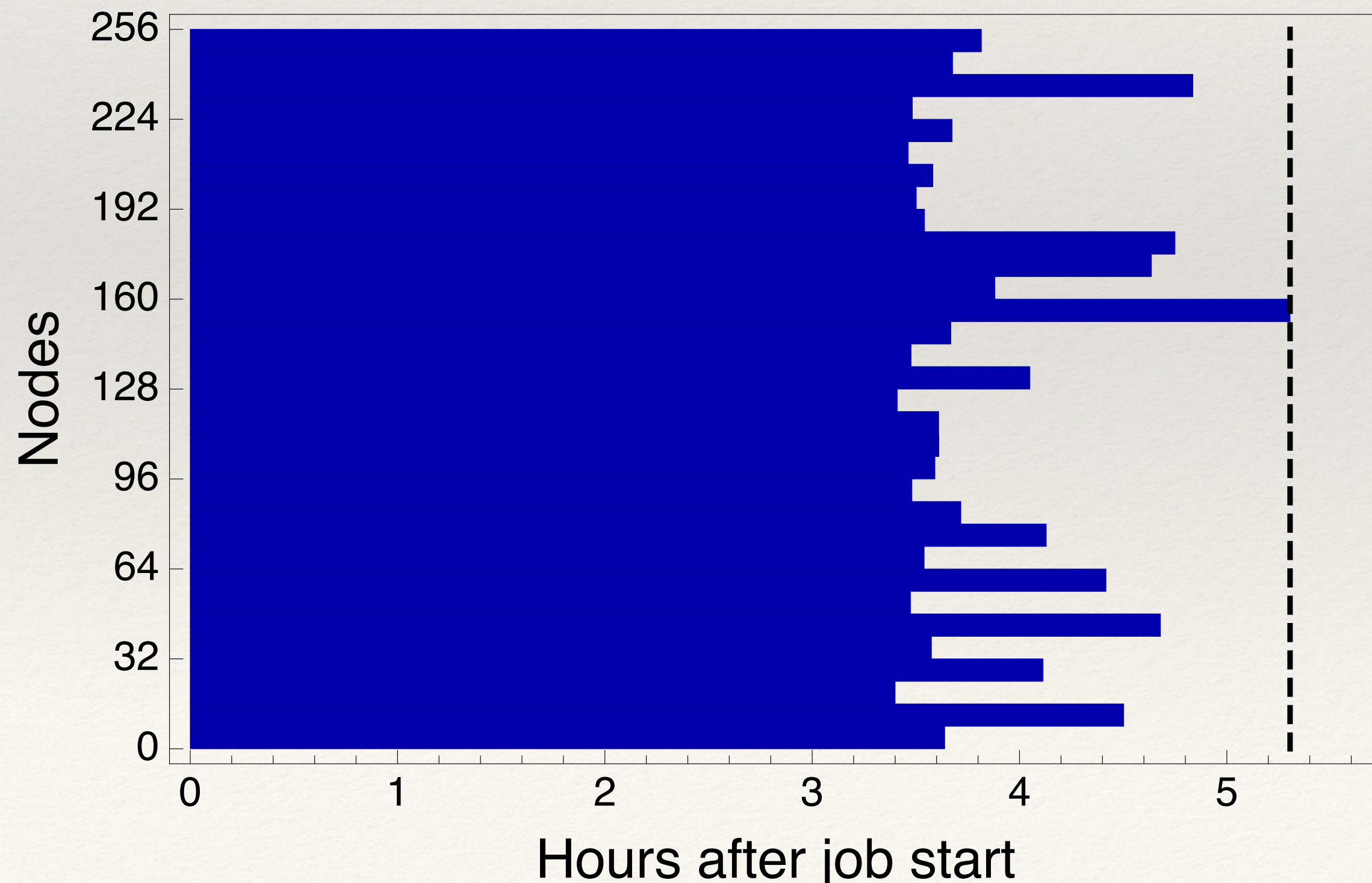
## □ task

an independent MPI task that could be run as a single job, or in a bundle on a larger job



# Managing millions of tasks on heterogeneous architectures

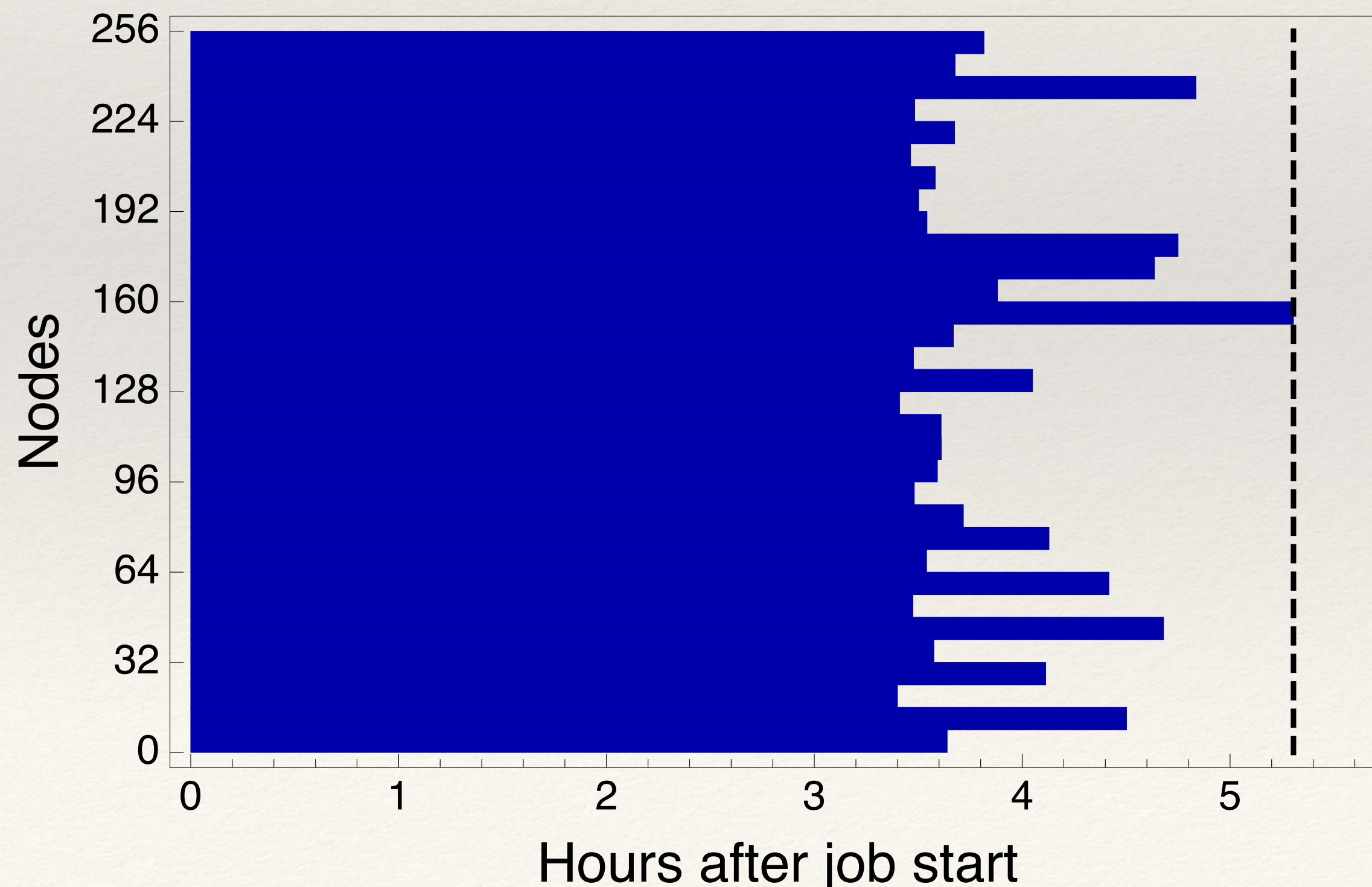
- ❑ Computing Centers often disfavor many small jobs
- ❑ We often bundle similar jobs into larger jobs
  - ❑ On heterogeneous architectures, different jobs may have very different resource needs
  - ❑ Even with similar jobs - significant waste can happen as fast/regular jobs wait for slow ones





# Managing millions of tasks on heterogeneous architectures

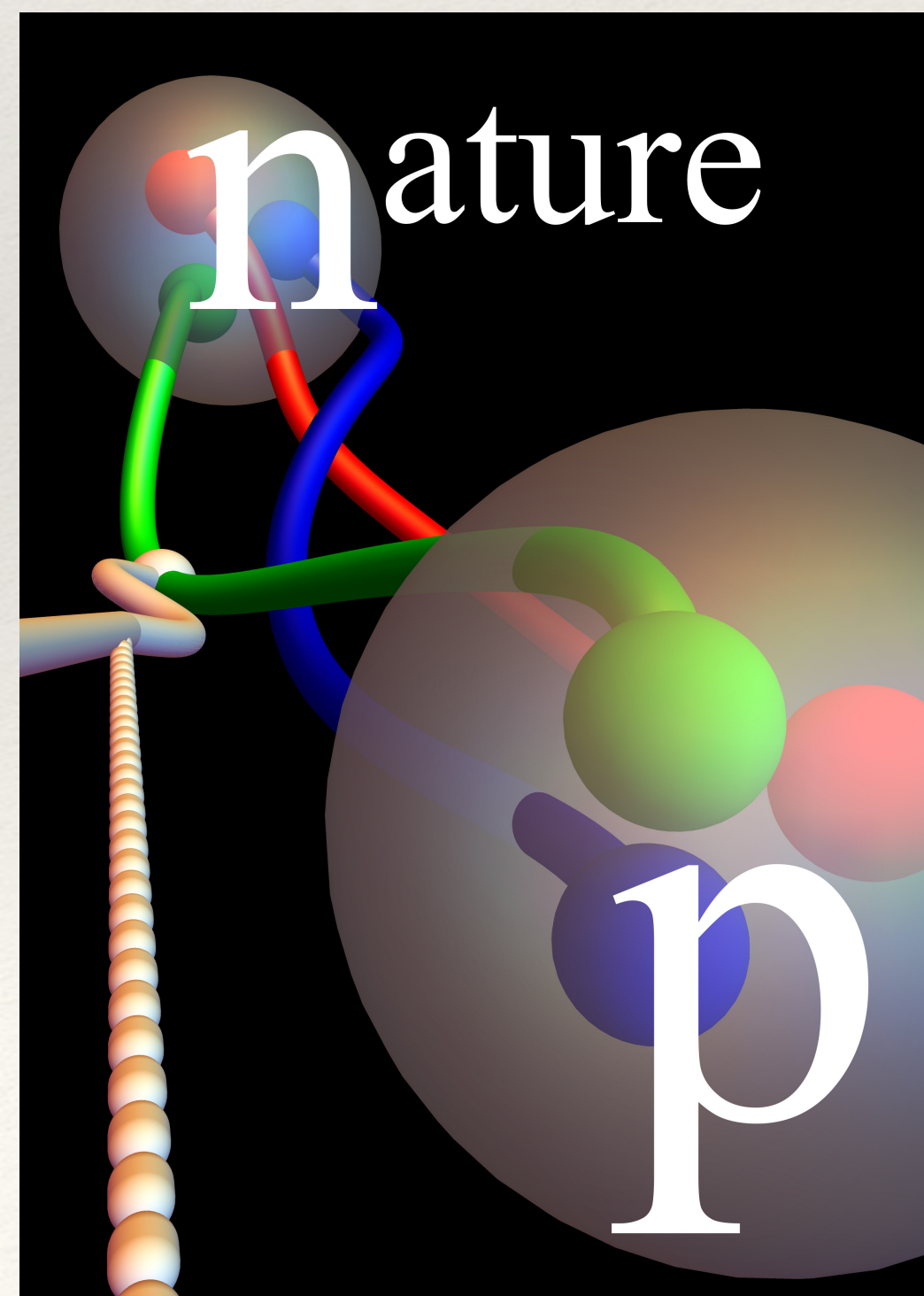
- ❑ METAQ - <https://github.com/evanberkowitz/metaq> - Evan Berkowitz
- ❑ Our first solution is a set of configurable **bash** scripts that can manage many different types of tasks and backfill idle nodes
- ❑ Different colors represent jobs with different resource needs  
(on average, 30% wasted idle nodes reduced to ~5% wasted idle nodes)





# Managing millions of tasks on heterogeneous architectures

- ❑ METAQ - <https://github.com/evanberkowitz/metaq> - Evan Berkowitz
- ❑ METAQ has served us very well - allowing us to very efficiently use Titan and its successor Summit (both at Oak Ridge Leadership Computing Facility (OLCF) @ ORNL)
- ❑ METAQ was used to simultaneously run computations for these two projects without METAQ - we would not have been able to run either, let alone both, at OLCF



## LETTER

<https://doi.org/10.1038/s41586-018-0161-8>

**A per-cent-level determination of the nucleon axial coupling from quantum chromodynamics**

**C. Chang et al (CalLat)**

**Nature 558 (2018) no.7708, 91-94 [arXiv:1805.12130]**

### ❑ additional requirements

- ❑ Ludicrously fast QUDA code (MDWF)
- ❑ Free configs from MILC
- ❑ Access to Titan/OLCF through INCITE



**Heavy Physics Contributions to Neutrinoless Double Beta Decay from QCD**

**A. Nicholson et al (CalLat)**

**Phys. Rev. Lett. 121, 172501 (2018) [arXiv:1805.02634]**



# Managing millions of tasks on heterogeneous architectures

- ❑ METAQ - <https://github.com/evanberkowitz/metaq> - Evan Berkowitz
- ❑ At the same time, METAQ has short comings that are not addressable with **bash**
- ❑ Managing tasks with **bash** means each MPI task requires an **mpirun** or equivalent call
  - ❑ This places significant strain on the service nodes that manage the compute nodes (we have crashed both Titan and Summit with too many tasks being simultaneously run)
- ❑ Desirable features:
  - ❑ Place tasks on neighboring nodes (don't communicate across racks)
  - ❑ Place ranks on specific cores/GPUs of node
  - ❑ Simultaneously run GPU-intense and CPU-only tasks on same nodes
  - ❑ Dynamically assign priority of tasks
  - ❑ Flexible run “configurations” for the same task (try this or that if this isn't possible)
  - ❑ Dynamically “shed” nodes in a job allocation
  - ❑ Collect node errors to identify problems (can't talk to disk, can't talk to GPU, slow fabric...)



# MPI\_JM

- ❑ MPI\_JM - [https://github.com/kenmcelvain/mpi\\_jm](https://github.com/kenmcelvain/mpi_jm) (temporarily private - license issue)
- ❑ To support these more desirable features, we have created MPI Job Manager (MPI\_JM)
  - ❑ C++/Python code
  - ❑ MPI\_JM uses a single `mpirun` call to start job-allocation and manage all tasks
  - ❑ MPI\_JM requires minimal interface to your code to work (normal running still works)
  - ❑ MPI\_JM uses a second Python Library (user written) that instructs MPI\_JM where to collect `yaml` files that describe resource requirements. This library enables the user to build generic types of jobs (propagator, baryon-block constructor, ...)
- ❑ MPI\_JM is configured to know about a machine
  - ❑ user only needs to specify resource requirements
  - ❑ MPI\_JM handles locking ranks to various parts of the node to optimize performance



# MPI\_JM

---

- ❑ A job is launched with MPI\_JM
  - ❑ MPI\_JM carves up the nodes in “lumps” of size specified by the user
  - ❑ The lumps are all initialized concurrently (can bring up entire Summit in < 5 minutes)
  - ❑ MPI\_JM begins to distribute tasks in each lump simultaneously
  - ❑ Tasks are located from user specified folders
  - ❑ User(s) can add more tasks during job (task list is refreshed regularly)
  - ❑ Multiple MPI\_JM instances can work on the same list of tasks
  - ❑ Tasks are given user-specified priorities to help MPI\_JM decide which jobs to run first
  - ❑ User can specify dependencies of tasks also
  - ❑ Tasks can have pre/post operations (check dependencies are ready - create dependent task)



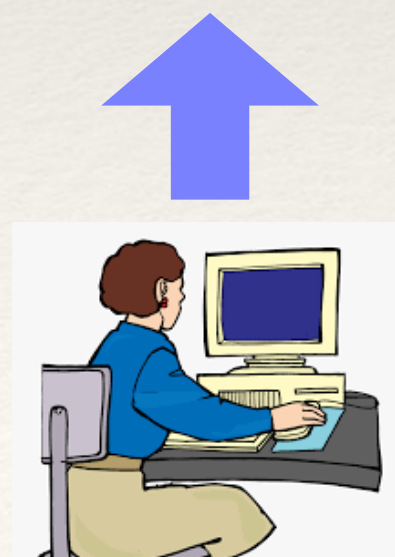
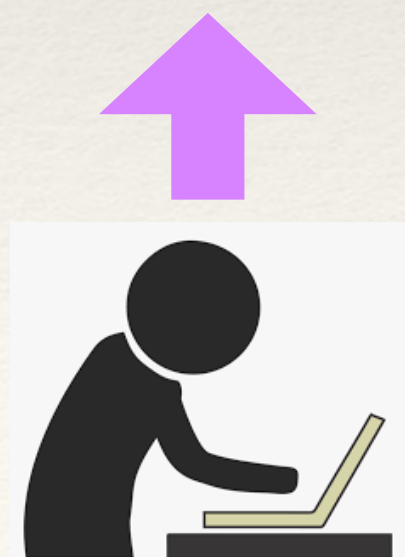
# MPI\_JM

Example: 4096 nodes, 1 mpirun, 32 “blocks” of 256 nodes,  
arbitrary number of “tasks” launched within blocks

256 nodes							

64 node CPU task ●  
32 node GPU task ●

128 node CPU+GPU task ●

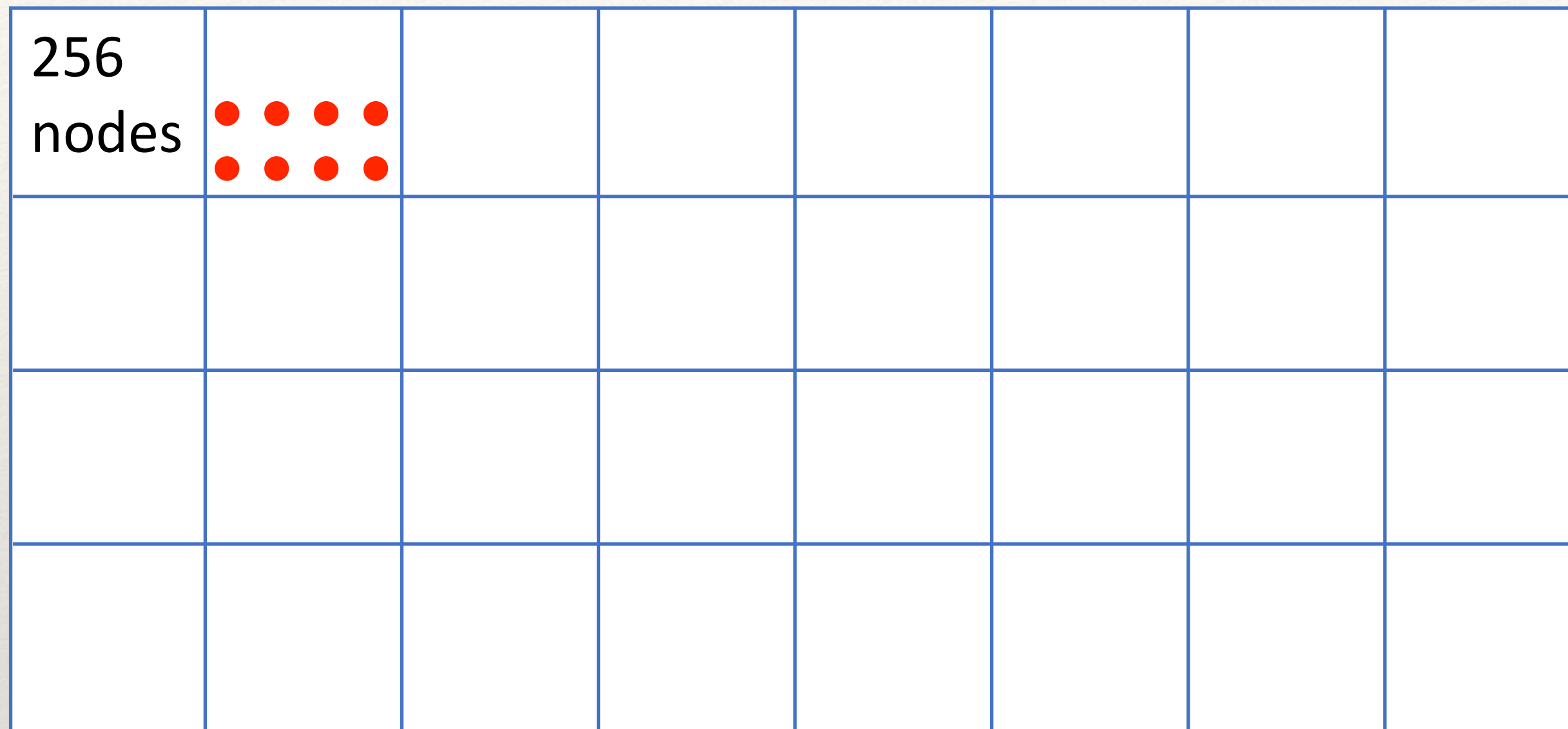


- ❑ Main communicator organizes nodes in sorted list to keep nodes physically local
- ❑ User specifies “block size” (256) to distribute master communicator over all nodes
- ❑ Multiple users can continuously feed tasks to MPI\_JM launch folder
- ❑ Tasks are launched via COMM\_SPAWN (failed task does not crash entire allocation)



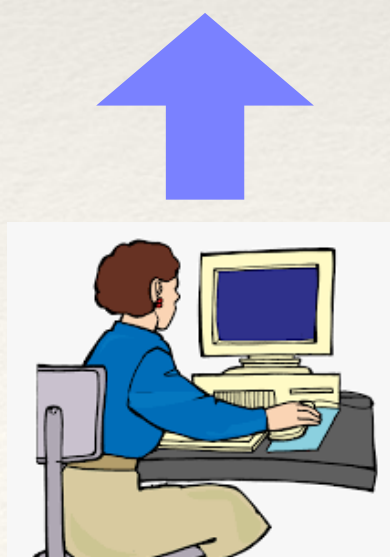
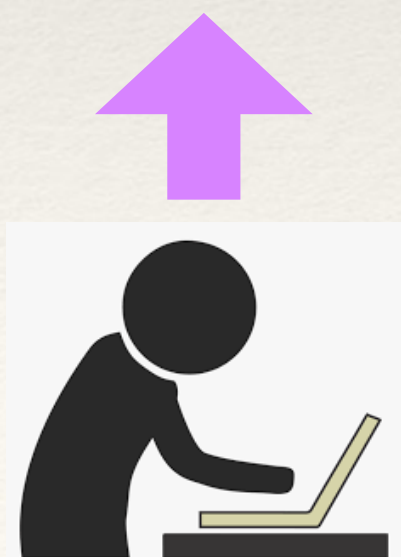
# MPI\_JM

Example: 4096 nodes, 1 mpirun, 32 “blocks” of 256 nodes,  
arbitrary number of “tasks” launched within blocks



64 node CPU task ●  
32 node GPU task ●

128 node CPU+GPU task ●

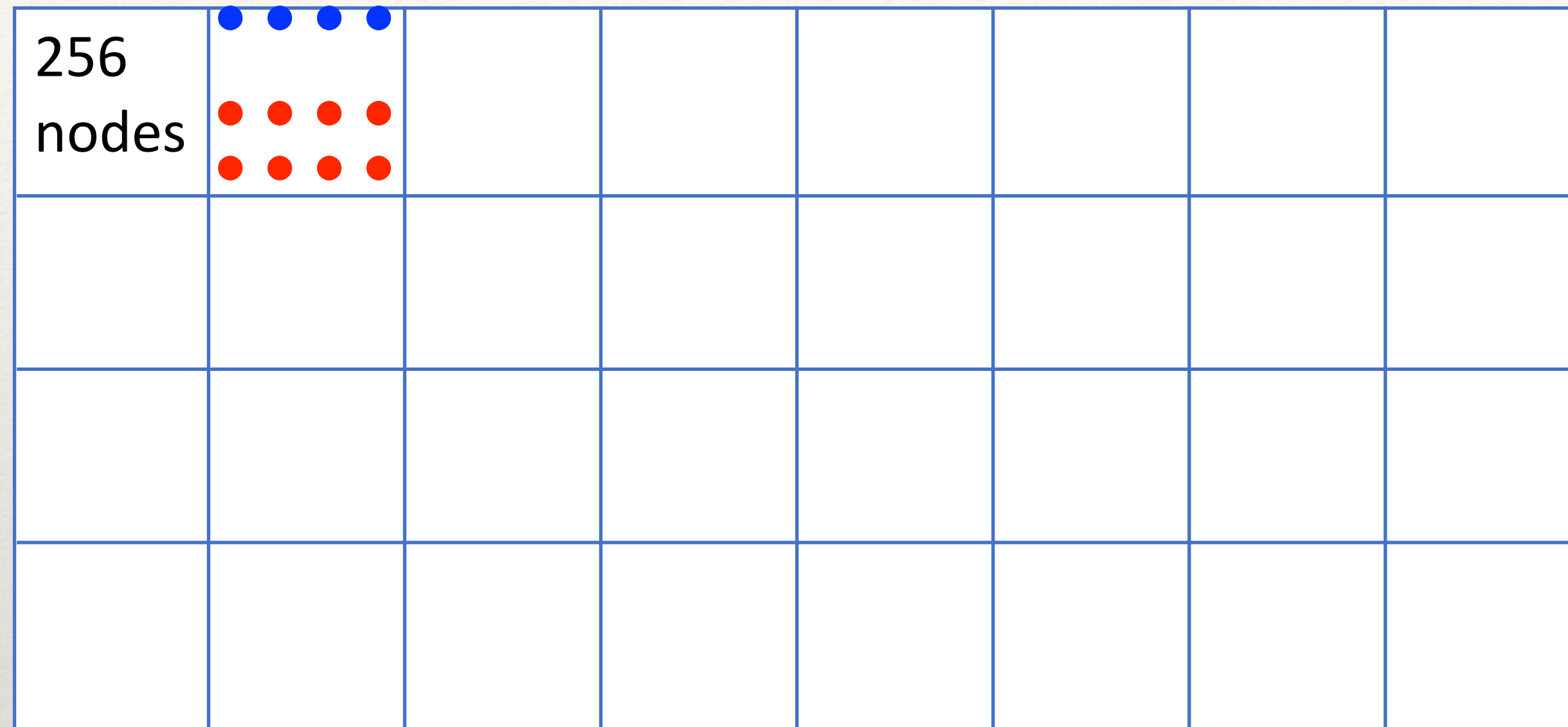


- ❑ Main communicator organizes nodes in sorted list to keep nodes physically local
- ❑ User specifies “block size” (256) to distribute master communicator over all nodes
- ❑ Multiple users can continuously feed tasks to MPI\_JM launch folder
- ❑ Tasks are launched via COMM\_SPAWN (failed task does not crash entire allocation)



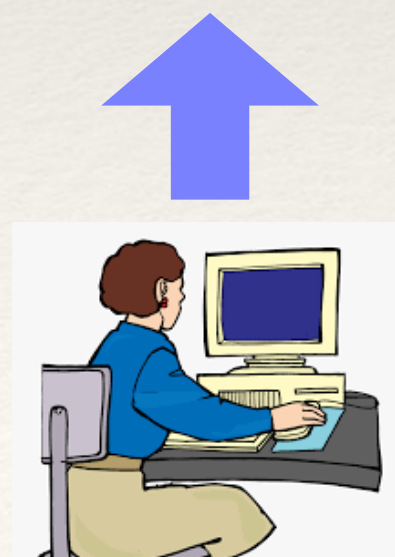
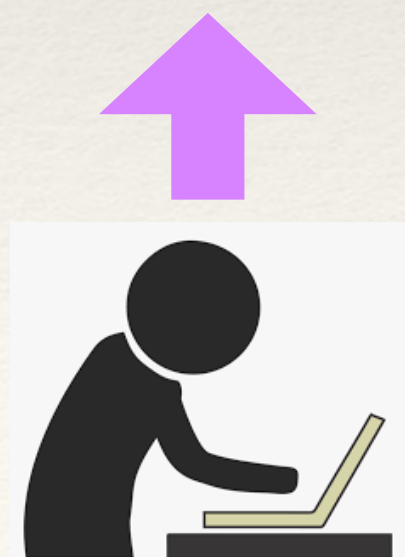
# MPI\_JM

Example: 4096 nodes, 1 mpirun, 32 “blocks” of 256 nodes,  
arbitrary number of “tasks” launched within blocks



64 node CPU task ●  
32 node GPU task ●

128 node CPU+GPU task ●

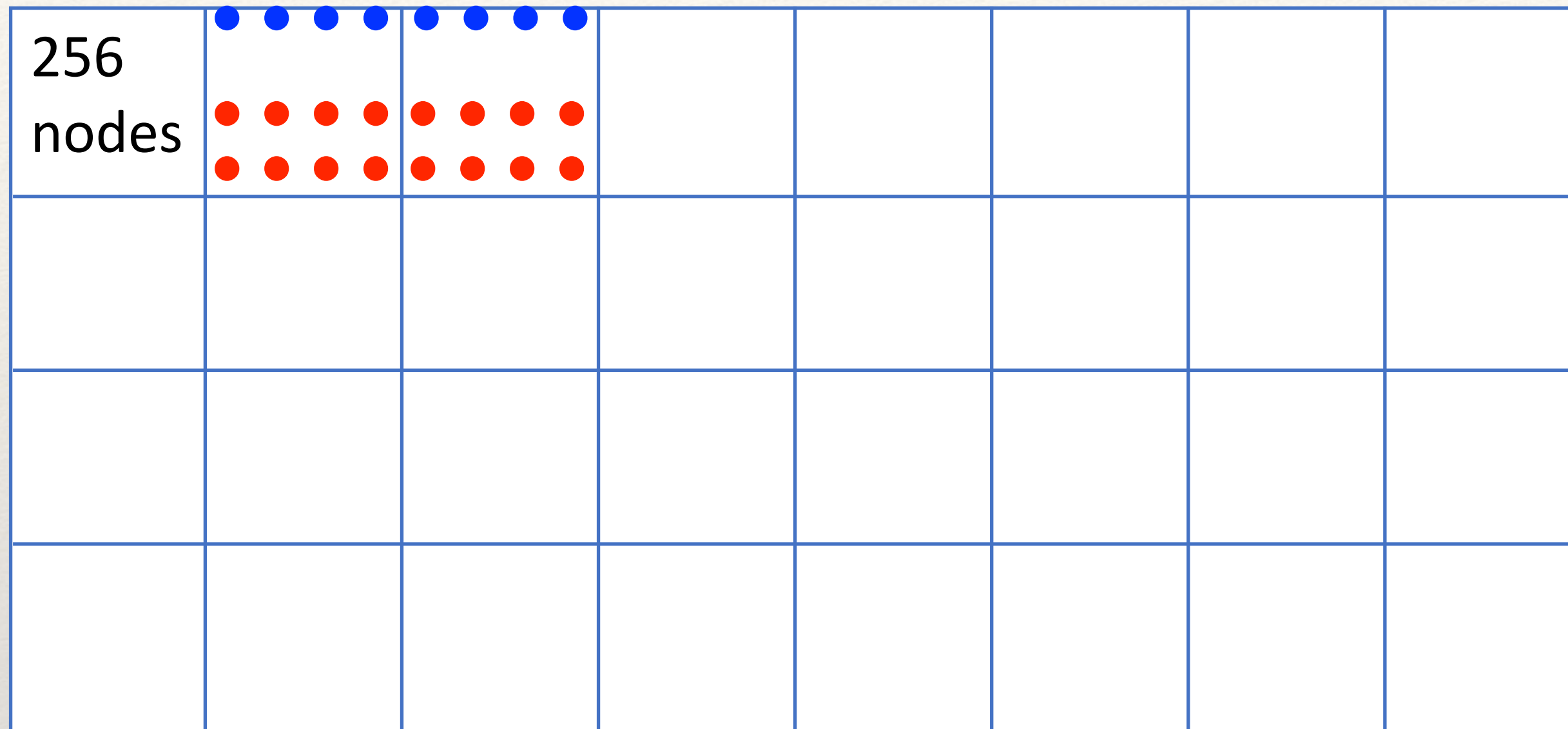


- ❑ Main communicator organizes nodes in sorted list to keep nodes physically local
- ❑ User specifies “block size” (256) to distribute master communicator over all nodes
- ❑ Multiple users can continuously feed tasks to MPI\_JM launch folder
- ❑ Tasks are launched via COMM\_SPAWN (failed task does not crash entire allocation)



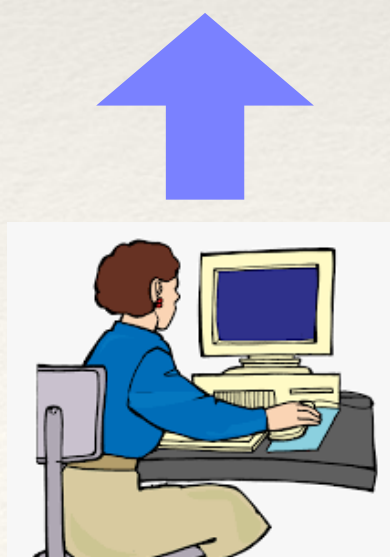
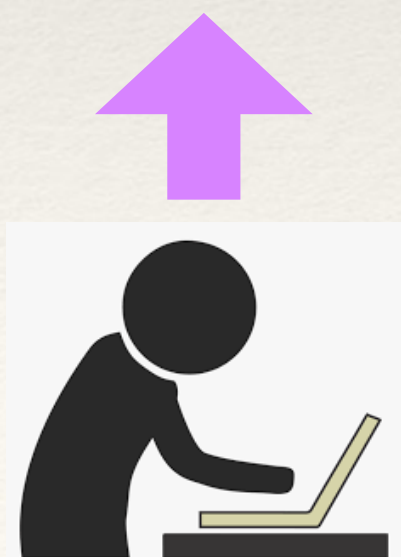
# MPI\_JM

Example: 4096 nodes, 1 mpirun, 32 “blocks” of 256 nodes, arbitrary number of “tasks” launched within blocks



64 node CPU task ●  
32 node GPU task ●

128 node CPU+GPU task ●

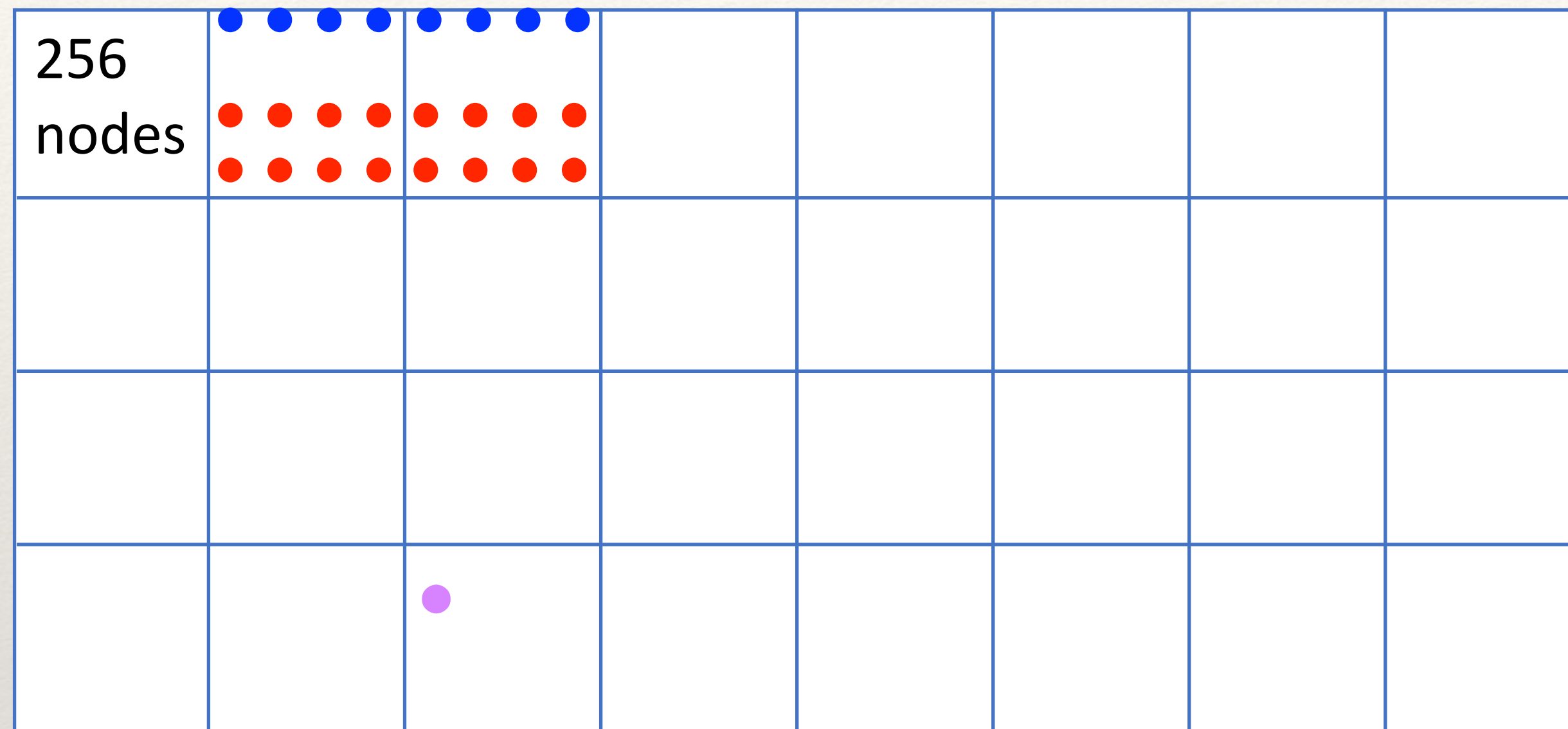


- ❑ Main communicator organizes nodes in sorted list to keep nodes physically local
- ❑ User specifies “block size” (256) to distribute master communicator over all nodes
- ❑ Multiple users can continuously feed tasks to MPI\_JM launch folder
- ❑ Tasks are launched via COMM\_SPAWN (failed task does not crash entire allocation)



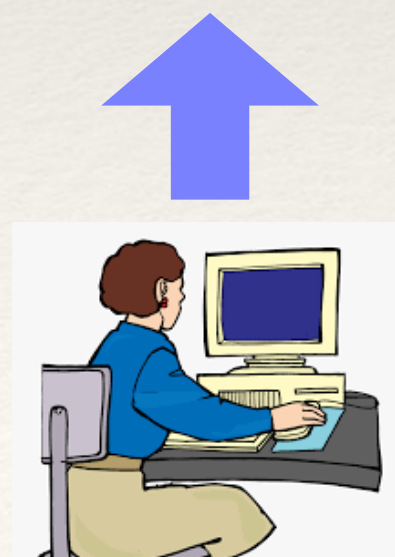
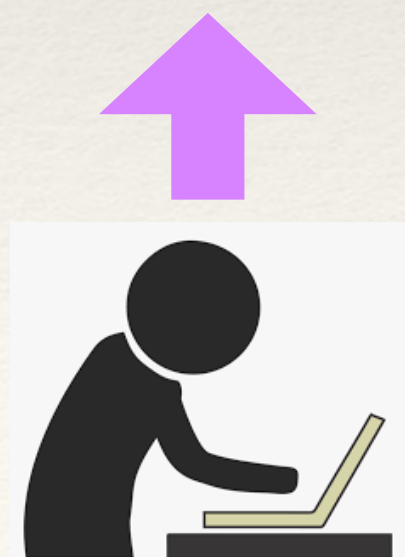
# MPI\_JM

Example: 4096 nodes, 1 mpirun, 32 “blocks” of 256 nodes,  
arbitrary number of “tasks” launched within blocks



64 node CPU task ●  
32 node GPU task ●

128 node CPU+GPU task ●

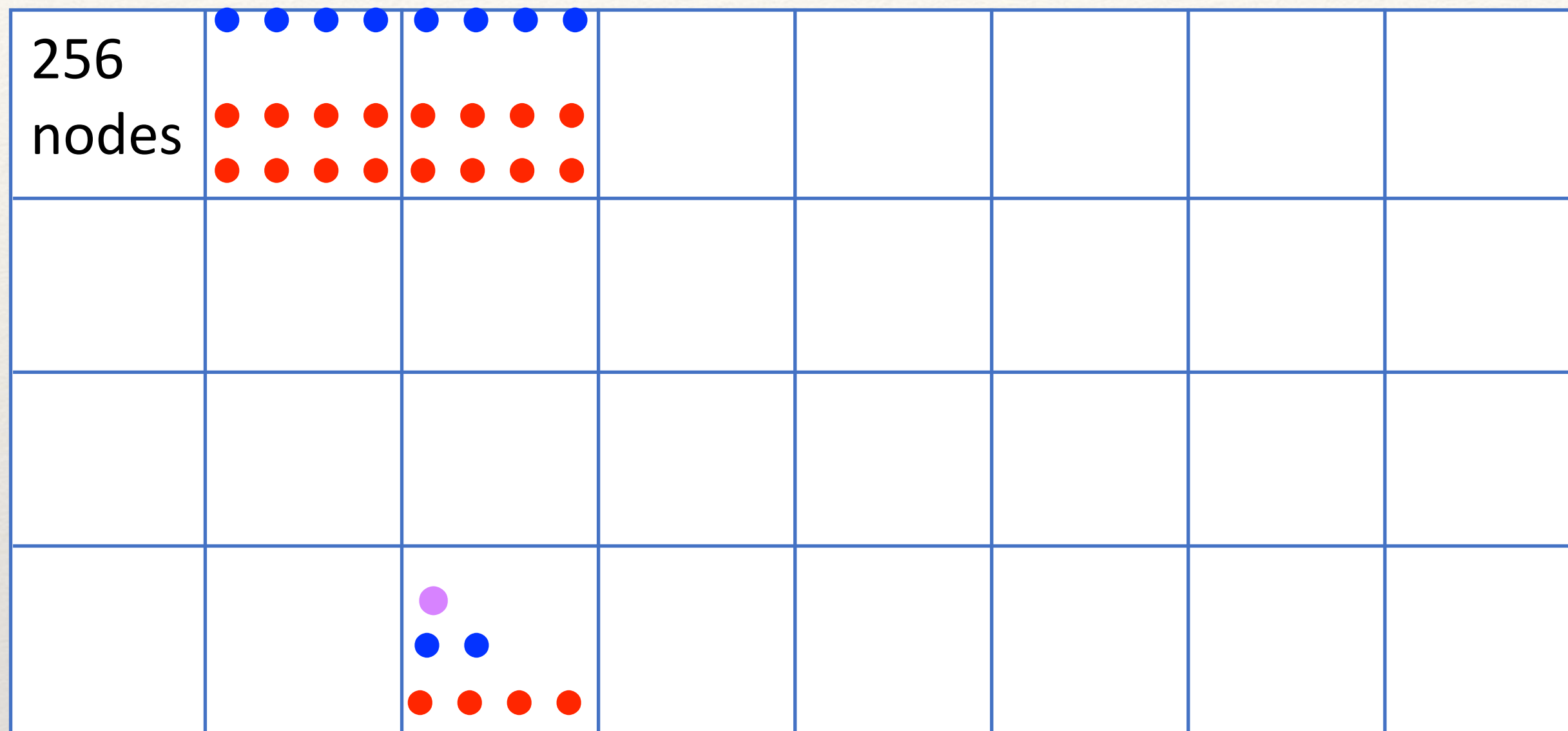


- ❑ Main communicator organizes nodes in sorted list to keep nodes physically local
- ❑ User specifies “block size” (256) to distribute master communicator over all nodes
- ❑ Multiple users can continuously feed tasks to MPI\_JM launch folder
- ❑ Tasks are launched via COMM\_SPAWN (failed task does not crash entire allocation)



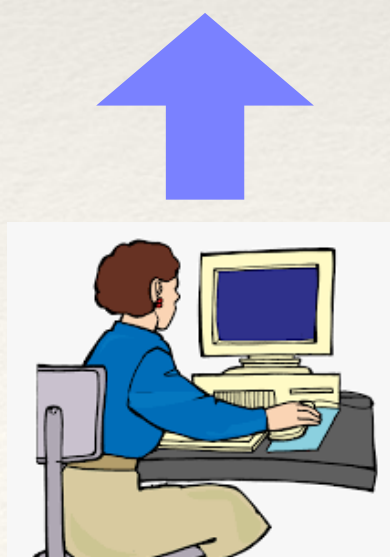
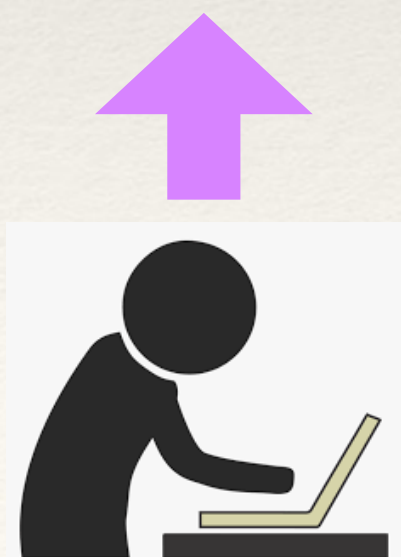
# MPI\_JM

Example: 4096 nodes, 1 mpirun, 32 “blocks” of 256 nodes,  
arbitrary number of “tasks” launched within blocks



64 node CPU task ●  
32 node GPU task ●

128 node CPU+GPU task ●

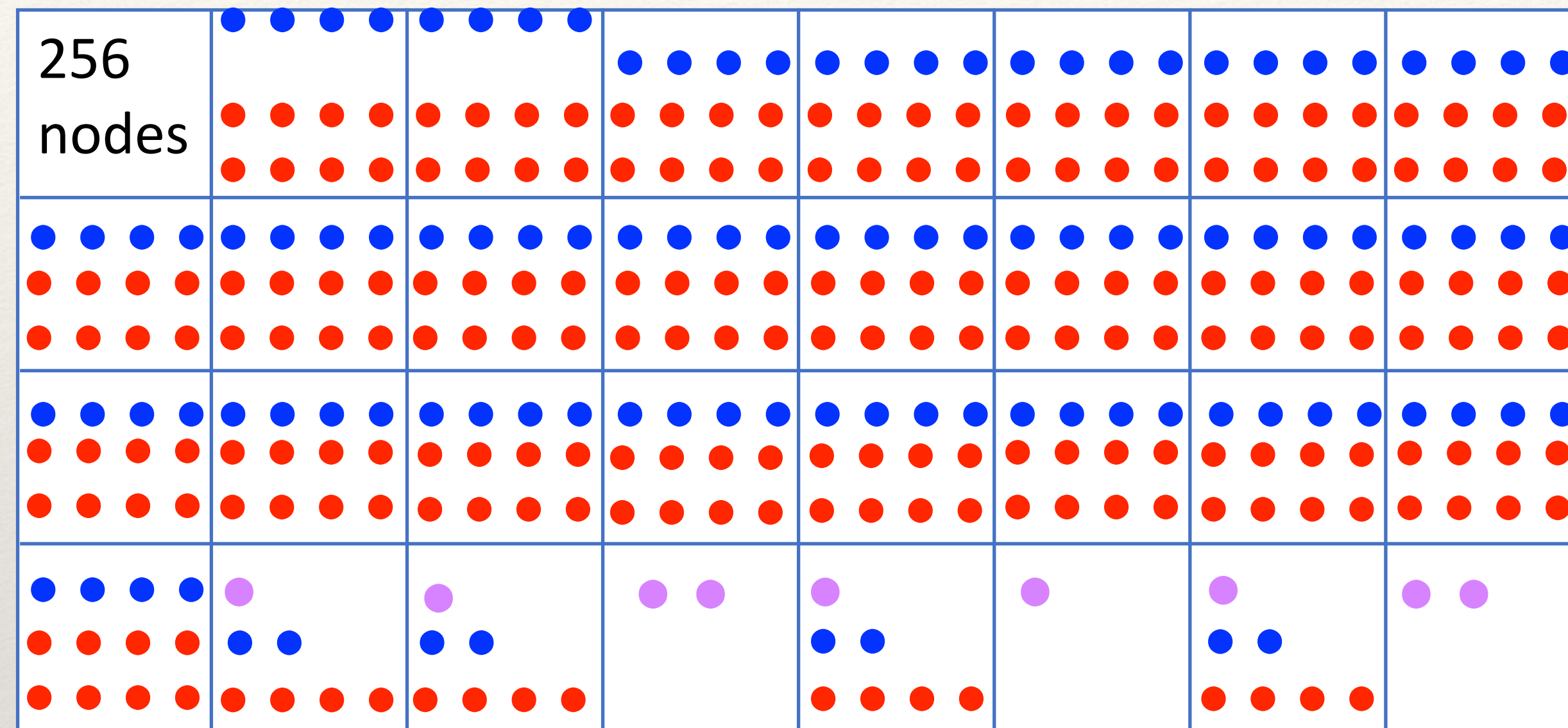


- ❑ Main communicator organizes nodes in sorted list to keep nodes physically local
- ❑ User specifies “block size” (256) to distribute master communicator over all nodes
- ❑ Multiple users can continuously feed tasks to MPI\_JM launch folder
- ❑ Tasks are launched via COMM\_SPAWN (failed task does not crash entire allocation)



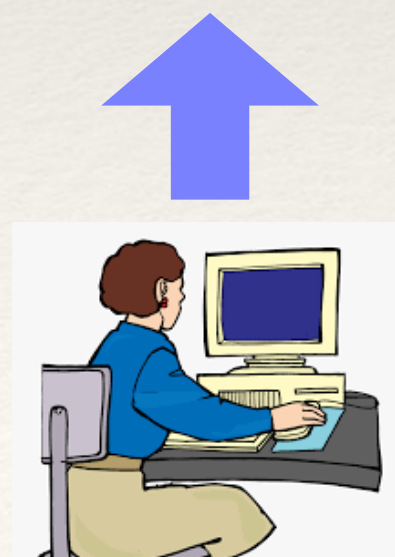
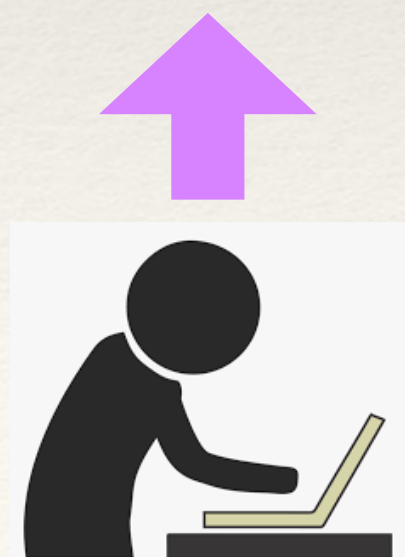
# MPI\_JM

Example: 4096 nodes, 1 mpirun, 32 “blocks” of 256 nodes,  
arbitrary number of “tasks” launched within blocks



64 node CPU task ●  
32 node GPU task ●

128 node CPU+GPU task ●



- ❑ Main communicator organizes nodes in sorted list to keep nodes physically local
- ❑ User specifies “block size” (256) to distribute master communicator over all nodes
- ❑ Multiple users can continuously feed tasks to MPI\_JM launch folder
- ❑ Tasks are launched via COMM\_SPAWN (failed task does not crash entire allocation)



# Minimal Interface

callat-qcd / qmp Public  
forked from usqcd-software/qmp

<> Code 🔗 Pull requests 🎮 Actions 📁 Projects 📖 Wiki 🛡 Security 📈 Insights ⚙ Settings

🔗 mpi\_jm.no\_call... qmp / lib / mpi / QMP\_init\_mpi.c Go to file ...

walkcloud merging... Latest commit c2505cb on Jan 21, 2019 🕒 History

🔍 6 contributors

135 lines (117 sloc) | 3.18 KB

Raw Blame

```
1  #include <stdio.h>
2  #include <string.h>
3  #include <stdlib.h>
4  #include <ctype.h>
5  #include <stdarg.h>
6
7  #include "qmp_config.h"
8  #include "QMP_P_COMMON.h"
9
10 #ifdef QMP_MPI_JM
11 #include "jm.h"
12 #endif
```

```
47
48 #ifdef QMP_MPI_JM
49     jm_parent_handshake(argc, argv);
50 #endif
51
```

```
102 void
103 QMP_finalize_msg_passing_mpi (void)
104 {
105     int flag;
106     MPI_Finalized(&flag);
107
108     if (!flag) {
109         MPI_Finalize();
110     }
111 #ifdef QMP_MPI_JM
112     jm_finish(0, "QMP MPI finalized.");
113 #endif
114 }
```

```
116 void
117 QMP_abort_mpi (int error_code)
118 {
119 #ifdef QMP_MPI_JM
120     /* try to do clean shutdown instead of ABORT */
121     /* BETTER: should really ask jm_master to kill job */
122     jm_finish(error_code, QMP_error_string(error_code));
123     MPI_Finalize();//this seems out of step with the change made to QMP with the &flag and finalize above
124     exit(error_code);
125 #else
126     MPI_Abort(MPI_COMM_WORLD, error_code);
127 #endif
128 }
```

- ❑ Code that uses QMP can utilize **MPI\_JM** (not yet in mainline QMP)



# Known Future Features

---

- ❑ Performance monitoring
- ❑ Improve overlapped use of distinct resources
- ❑ Inter-job communication tasks
- ❑ Dynamic task configuration
- ❑ Enhance fault-tolerance and diagnosis
- ❑ ...



# Known Future Features

---

## ❑ Performance monitoring

- ❑ User claims task will take 2 hours
- ❑ MPI\_JM will collect statistics on actual wall-clock time usage and use this to decide if a task can complete before end of allocation
- ❑ MPI\_JM can monitor performance of blocks, not just job type, in case some blocks of nodes are slower than others



# Known Future Features

---

- ❑ Improve overlapped use of distinct resources
  - ❑ After a task is finished, but needs to write I/O, the GPUs and CPUs for that task are otherwise idle.
  - ❑ MPI\_JM can be informed that the GPU-usage is complete, and a new GPU task can be placed on the nodes while I/O is being performed from the previous task



# Known Future Features

---

## ❑ Inter-job communication tasks

- ❑ Utilizing local NVME storage can be complicated for multi-node jobs with parallel I/O
- ❑ Users can create “follow up” tasks/work with MPI\_JM such that MPI\_JM will launch a helper task to collect the NVME from the individual nodes and build the complete file
- ❑ Task 37 may require the data files from Tasks 0-7, which were all written to NVME. MPI\_JM can collect these and place them on the NVME that Task 37 will run on (MPI\_JM can decide ahead of time which nodes Task 37 will run on)



# Known Future Features

## ❑ Dynamic task configuration

- ❑ The optimal number of nodes for task-A might be 4 (minimum required to fit job in memory)
- ❑ Towards the end of a job-allocation, there may be insufficient time to complete task-A on 4 nodes
- ❑ But, with a performance loss, the job could run on 8-nodes in the remaining wallclock time
- ❑ The user can pre-specify possible task-configurations that MPI\_JM can try to fit in the available resources, thus further reducing “wasted” time at the end of an allocation
- ❑ One often observes 1000-3000 nodes available for only 45 minutes - with the optional run-configurations, especially in bottom feeder mode - users could opt to utilize these nodes that would otherwise be idle as the scheduler makes room for large jobs



# Known Future Features

## ❑ Enhance fault-tolerance and diagnosis

- ❑ Even mature systems exhibit random trouble. eg. the GPUs on a node may take too long to initialize, so the task aborts/times out, but the very next task will run successfully
- ❑ MPI\_JM can have a “re-try” option to catch such failures
- ❑ MPI\_JM can perform pretests of node health before launching the first task
  - ❑ Can the nodes see the file system?
  - ❑ Can the nodes see their NVME?
  - ❑ Can the nodes communicate via MPI? Is the MPI slower than it should be?
  - ❑ MPI\_JM can detect if task is actually running (sometimes, failed I/O doesn't cause task to abort or time-out, MPI\_JM could trigger a “kill” to this job)
- ❑ All these errors/issues can be logged, ideally in some SQL database, that can be used to help users and the sys-admins identify problematic nodes



# Conclusions

---

- ❑ MPI\_JM has worked great for running our LQCD calculations on Summit
- ❑ While we haven't succeeded in getting \$ to bring it to full beta-stage - we are determined to finally do it this Fall - and to share with anyone interested in trying it out
- ❑ MPI\_JM is not specific to LQCD - it works with any compute problem that requires MANY individual tasks, and has dependencies etc.  
Ken has been running nuclear equation of state calculations with it  
Computation at one nuclear density can be used as initial guess for similar density
- ❑ There are many known features that will improve running as well as I am sure many not-yet-realized features
- ❑ It would be great to get users to get more feedback
- ❑ When it is ready - we'll email [latticenews](#) to get volunteers



*Thank You*