# Launching parallel EUTERPE job with srun vs mpiexec

## N. Shukla, A. Marani & D. Molinari

High-Performance Computing Department

CINECA, Casalecchio di Reno

Bologna Italy

**Email:** n.shukla@cineca.it

# Launching parallel EUTERPE job with srun vs mpiexec

**Methodology**

- Using srun to launch euterpe on sparse nodes

- Using mpiexec on sparse nodes

- Using mpiexec on sparse nodes but with the flag

  --hostfile ./nodelist

For each case we run **99** jobs each one running Euterpe on **32** nodes up to fill the entire cluster and check if the job starts correctly or is found in hang after few minutes of run.

# Results: Using srun to launch euterpe on sparse nodes

## Total jobs analized: 99

✓ Total jobs successful : 99

✳ Total nodes con job in hang: 0

# Results: Using mpiexec on sparse nodes

**Run1: Total jobs analized: 99**

✓ Total jobs successful : 16

✳ Total nodes con job in hang: 83

**Run 2: Total jobs analized: 99**

✓ Total jobs successful : 51

✳ Total nodes con job in hang: 48

**Run 3: Total jobs analized: 99**

✓ Total jobs successful : 92

✳ Total nodes con job in hang: 7

## Summary:

Total jobs analized: 297
Total jobs successful : 159
Total nodes con job in hang: 138

# Results: Using mpiexec on sparse nodes



Resolving The Problem: Wait for the parent process to complete

# Results: Using mpiexec but with a sleep of 30 seconds before launching mpiexec

**Run 4:** Total jobs analized: 99

✓ Total jobs successful : 99

✳ Total nodes con job in hang: 0

# Results: Using mpiexec on sparse nodes but with the flag --hostfile ./hostfile

**Run 4:** Total jobs analized: 99

✓ Total jobs successful : 99

✳ Total nodes con job in hang: 0

# Summary run

We repeated the mpiexec run to gain statistic, because it was the only one showing jobs in hang state

| Case | Total jobs | Successful | Hang |
|------|------------|------------|------|
| mpiexec | 297 | 159 | 138 |
| srun | 99 | 99 | 0 |
| mpiexec --hostfile ./nodelist | 99 | 99 | 0 |
| sleep(30); mpiexec | 99 | 99 | 0 |

# Summary

- Launching with srun things goes fine.

- launching with mpiexec, the job may go in hang due to ssh processes denied because of the concurrency between mpiexec and slurm notifying the nodes to be in the job.

- The effect are ssh processes in state "defunct" the master node as in the snapshot attached.

## Please use srun

### Regarding the mpiexec runs we also found some workarounds

- Launching mpiexec with the flag --hostfile ./nodelist

- adding a brief sleep before the mpiexec command