# An Analysis of Long-tailed Network Latency Distribution and Background Traffic on Dragonfly+

Majid Salimi Beni
Department of Computer Science
University of Salerno, Salerno, Italy
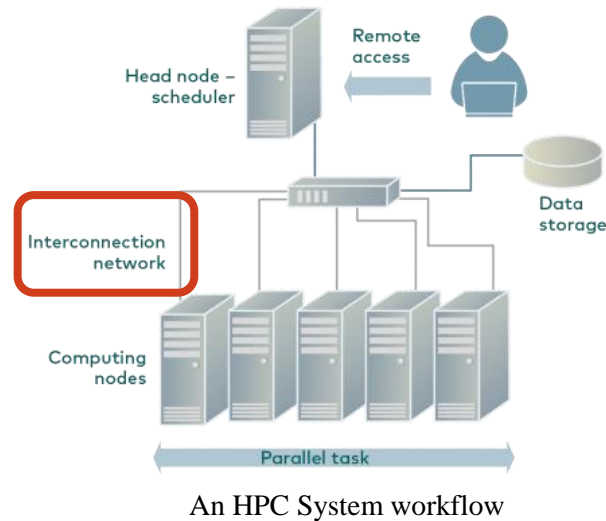msalimibeni@unisa.it

Biagio Cosenza
Department of Computer Science
University of Salerno, Salerno, Italy
bcosenza@unisa.it

November 2022

UNIVERSITÀ DEGLI STUDI DI SALERNO

DIPARTIMENTO DI INFORMATICA
DIPARTIMENTO DI ECCELLENZA

Bench 2022

# HPC Components and Some Challenges

❑ Does an HPC program always finish at the same time if we repeat the experiment?
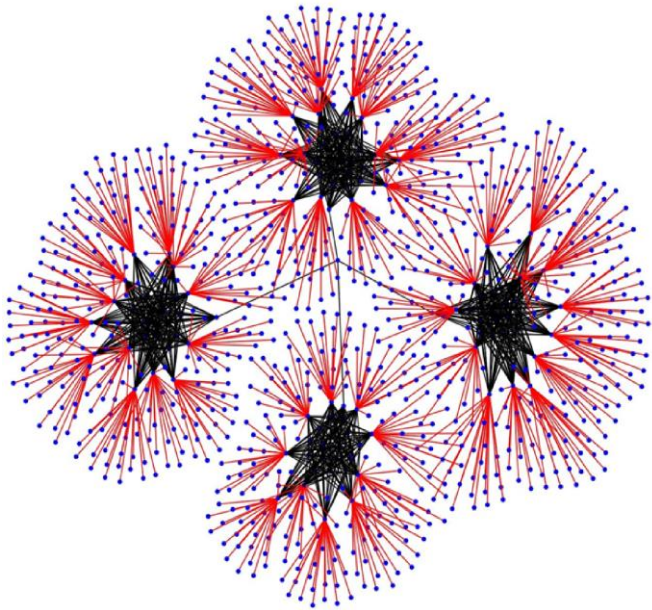


An HPC System workflow

> **Performance Variability:**
> The difference in the performance of an individual program in consecutive executions

❑ There are different sources:
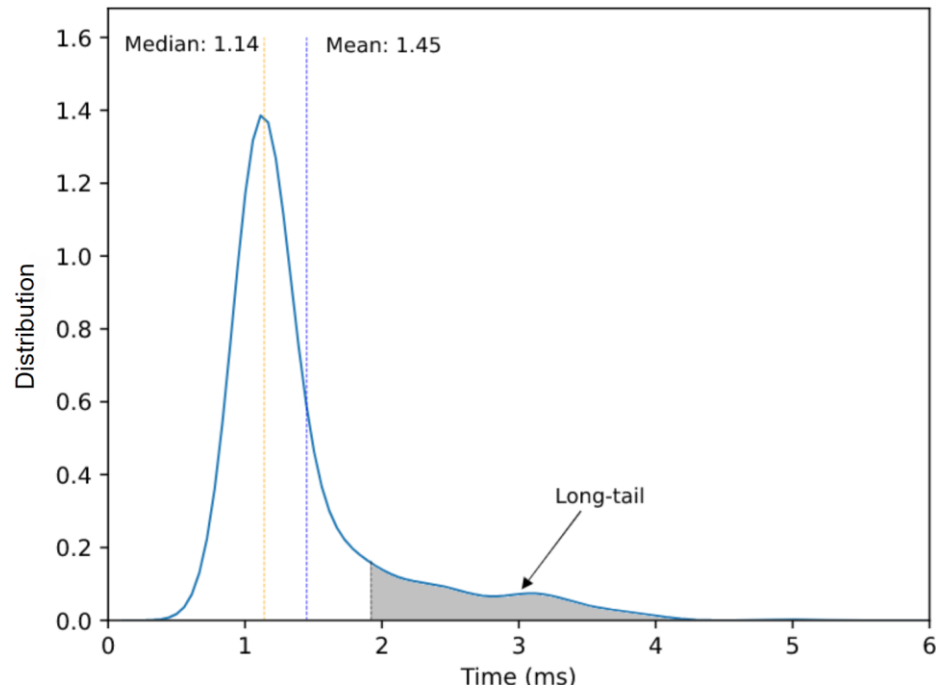  ❑ OS, I/O and file system, MPI, routing, **network,** etc.

# Performance Variability and Long Tail

❑ Distribution of latencies on Marconi100 when we repeat for 1000 times
   ❑ Some runs are taking longer than the majority
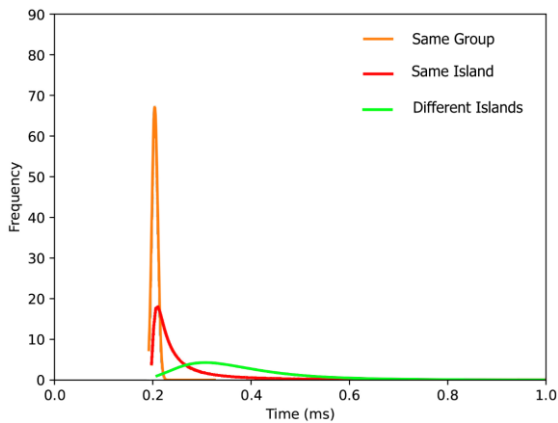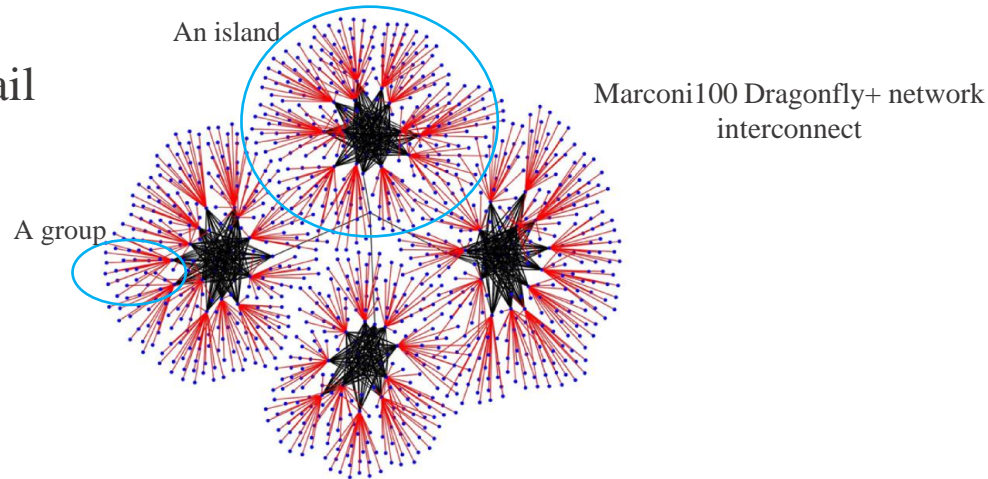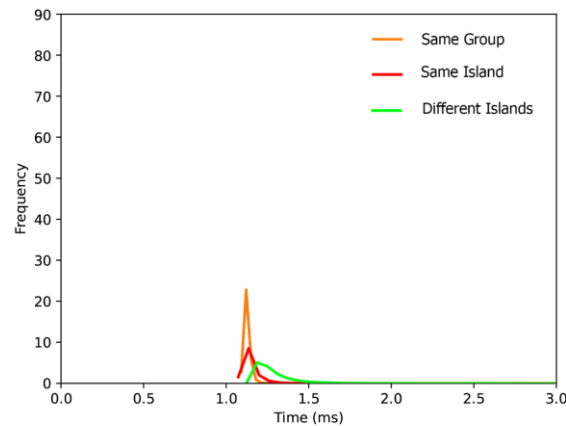


Dragonfly+ of Marconi100 @ Cineca
Supercomputing Center



Performance variability (Long-tail of the latency) distribution on Dragonfly+
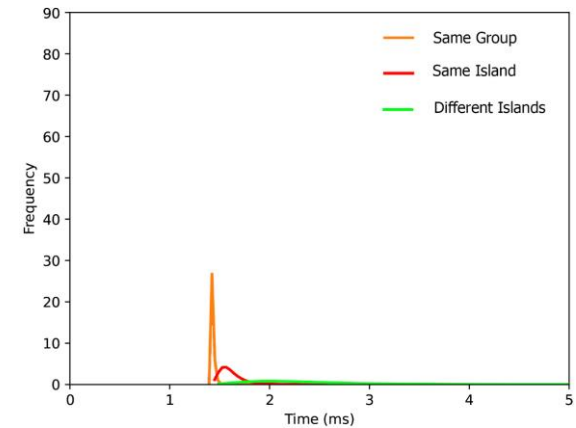
# Performance Variability and Job Placement Locality

❑ Locality and long-tail

An island

Marconi100 Dragonfly+ network interconnect

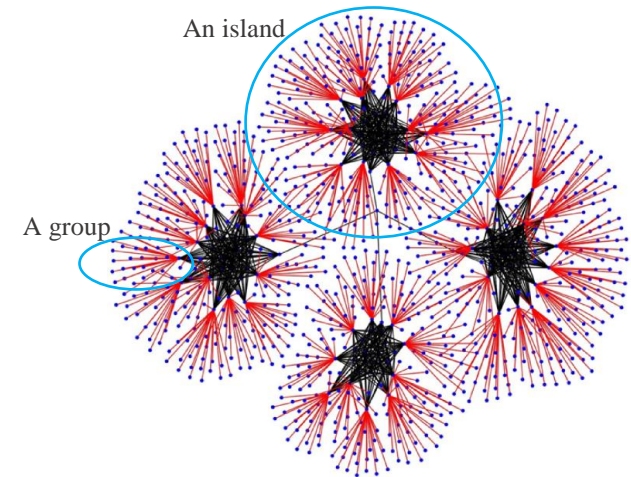A group



(a) Broadcast

(b) Reduce

(c) AlltoAll

Communication time frequency distribution of collective communications for 1000 iterations, with different allocation locality scenarios

# Performance Variability and Job Placement Locality

❑ Why don't we allocate all the nodes to the **Same Group**?
  ❑ Limited nodes in each group
  ❑ Long waiting time in the job queue for free groups

❑ What makes the "Different Islands" allocation more variable?
  ❑ Network is a shared resource
  ❑ There might be other users running communication-intensive jobs



An island

A group

Marconi100 Dragonfly+ network interconnect

Collecting information of other users from the **Job Scheduler**

3 months of data collection from the job scheduler of Marconi100

# Network Congestion

(Background traffic) $$b = \frac{N_c}{N_t} * \frac{N_c'}{N_a} * 100$$

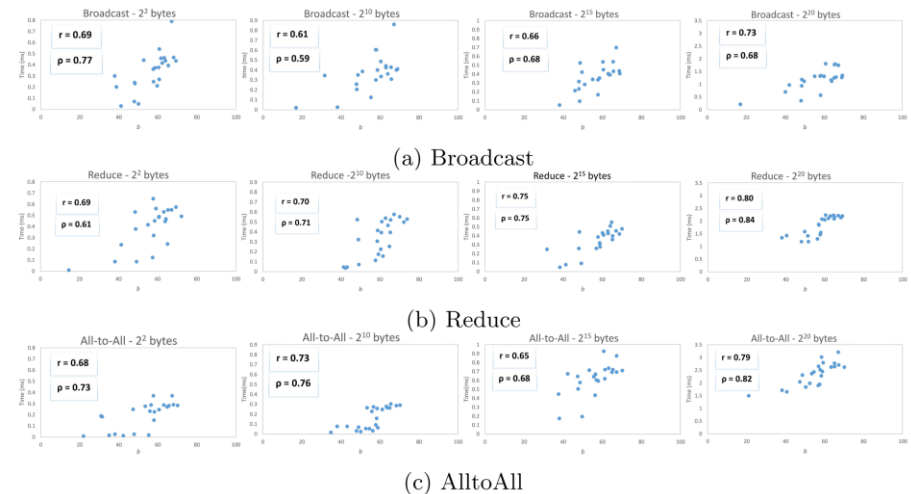$N_c$ : number of unique nodes contributing to communication
$N_t$ : total number of cluster physical nodes
$N_c'$ : the number of nodes contributing to communication (containing duplication)
$N_a$ : all allocated running nodes (containing duplication)

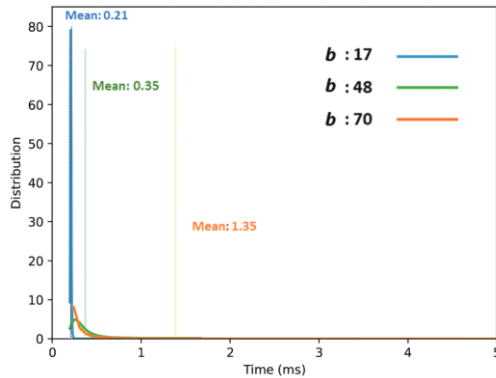> The ratio of nodes contributing to communication to all the physical cluster nodes.

- Using Pearson Correlation Coefficient (r) and Spearman Rank Correlation (ρ):

  - The heuristic is around 80 percent accurate
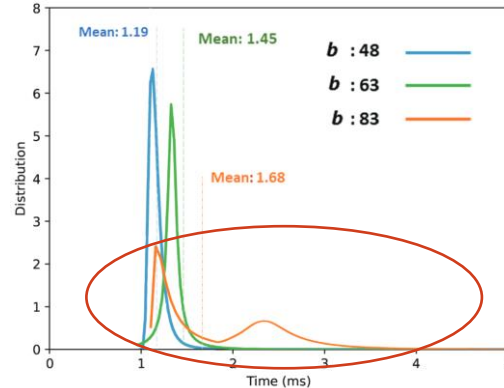
  - The correlations become stronger for larger data sizes



(a) Broadcast

(b) Reduce

(c) AlltoAll

The relation between background traffic ($b$) and the average communication time of different collectives with different message sizes
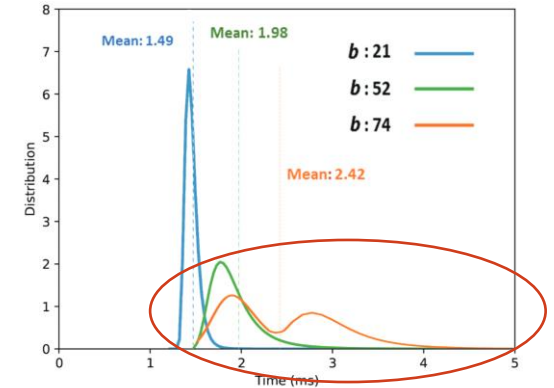
# The Impact of Background Traffic on Long-tail



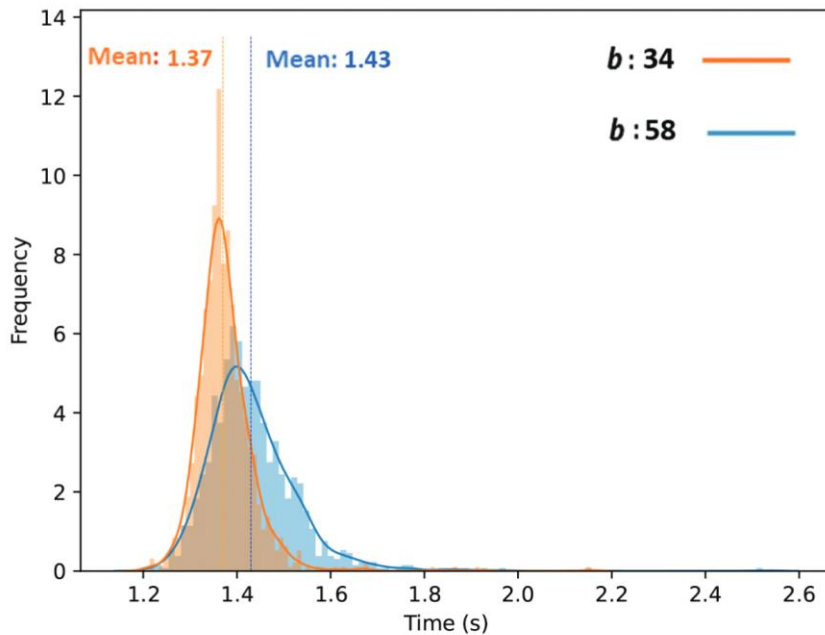(a) Broadcast      (b) Reduce      (c) All-to-All

Frequency distribution of communication times of 1000 iterations of Broadcast, Reduce, and All-to-All with different background traffics
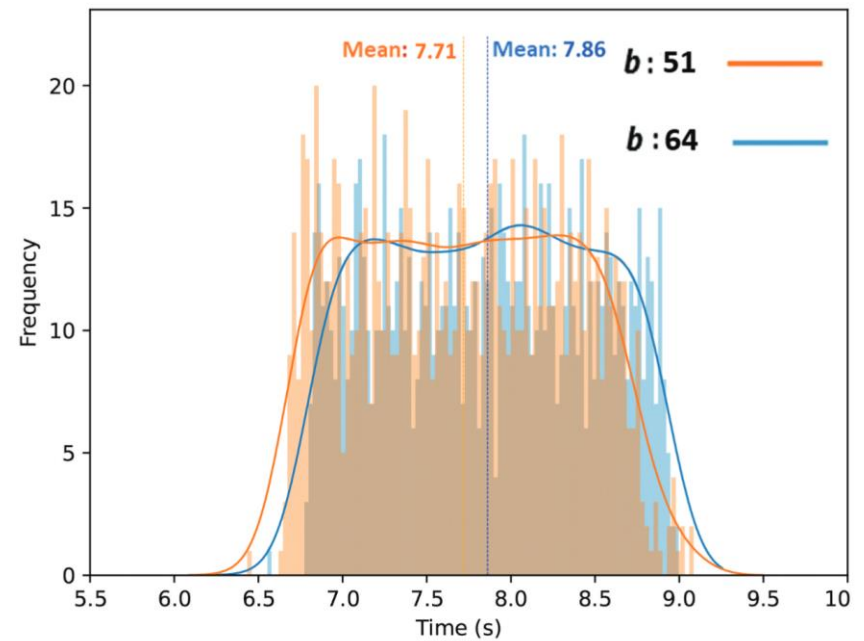
❑ The bigger the **b**, the longer the tail, the lower the peak!

❑ AlltoAll has the longest tail among all: It's more communication-intensive

❑ The Bimodal distribution is because of **Adaptive Routing, c**hoosing the non-minimal path while there is congestion on the shortest path

# Background Traffic and Applications



(a) HACC



(b) miniAMR

- ❏ Mini-Application analysis

    - ❏ Communication-intensive applications

    - ❏ Each consist of different communication patterns

# Conclusion and Future Work

❑ Performance variability study

    ❑ Communication patterns

    ❑ Message sizes

    ❑ Job placement locality

    ❑ Background traffic

❑ Future Work

    ❑ Gathering more network info such as: InfiniBand counters, job information, I/O, etc.

    ❑ Using ML models to make our heuristic more accurate

    ❑ Apply our findings to the job scheduler (SLURM)

THANK YOU

An Analysis of Long-tailed Network Latency Distribution and Background Traffic on Dragonfly+

Majid Salimi Beni, Biagio Cosenza

The 14th BenchCouncil International Symposium On Benchmarking, Measuring And Optimizing (Bench 2022) Nov 7-9, 2022

✉ Reach me at: msalimibeni@unisa.it