



**TECHNISCHE
UNIVERSITÄT
DRESDEN**



The Hebrew University
of Jerusalem

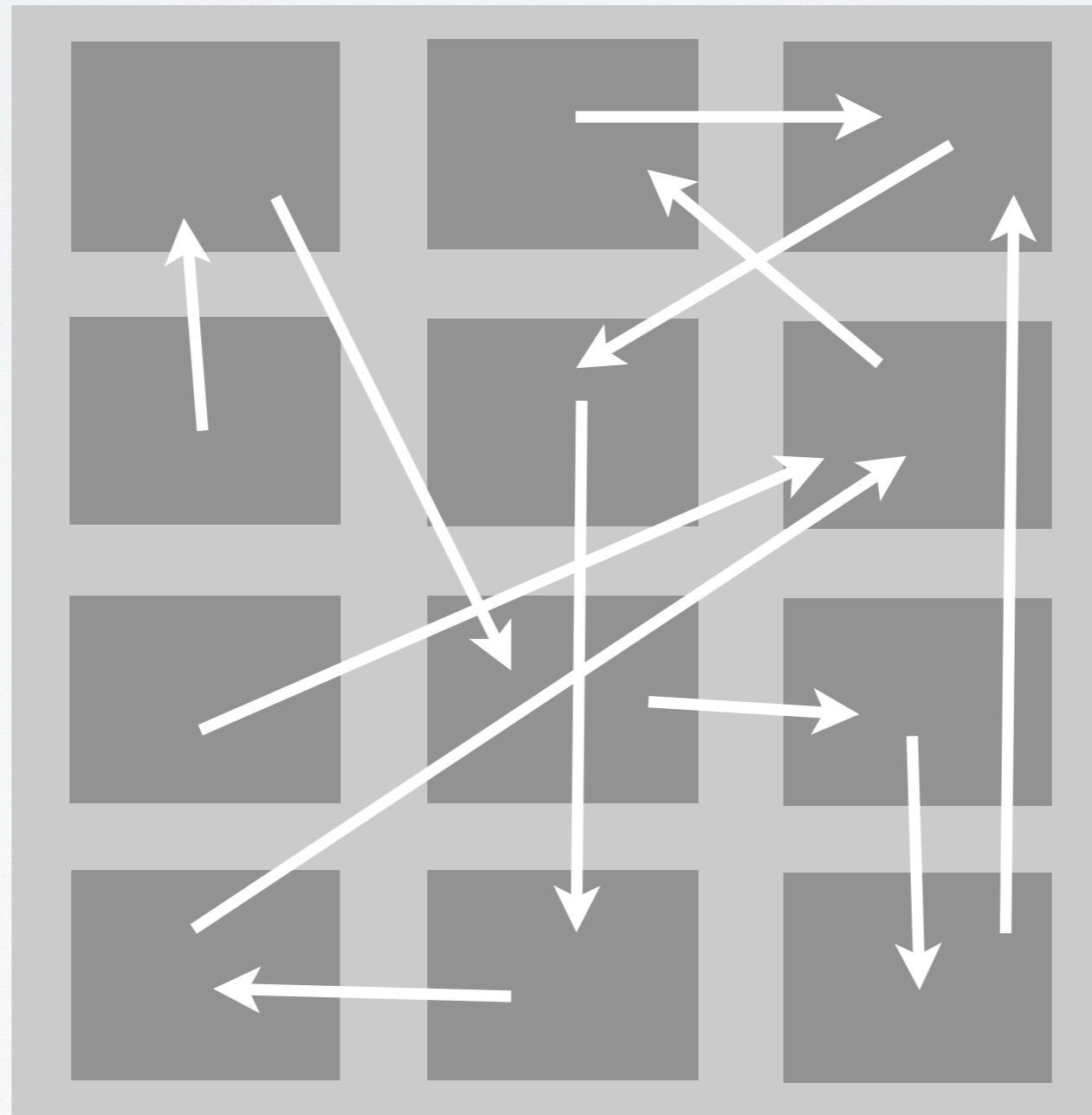
Faculty of Computer Science Institute of Systems Architecture, Operating Systems Group

OVERHEAD OF A DECENTRALIZED GOSSIP ALGORITHM ON THE PERFORMANCE OF HPC APPLICATIONS

**ELY LEVY, AMNON BARAK, AMNON SHILOH,
MATTHIAS LIEBER, CARSTEN WEINHOLD, HERMANN HÄRTIG**

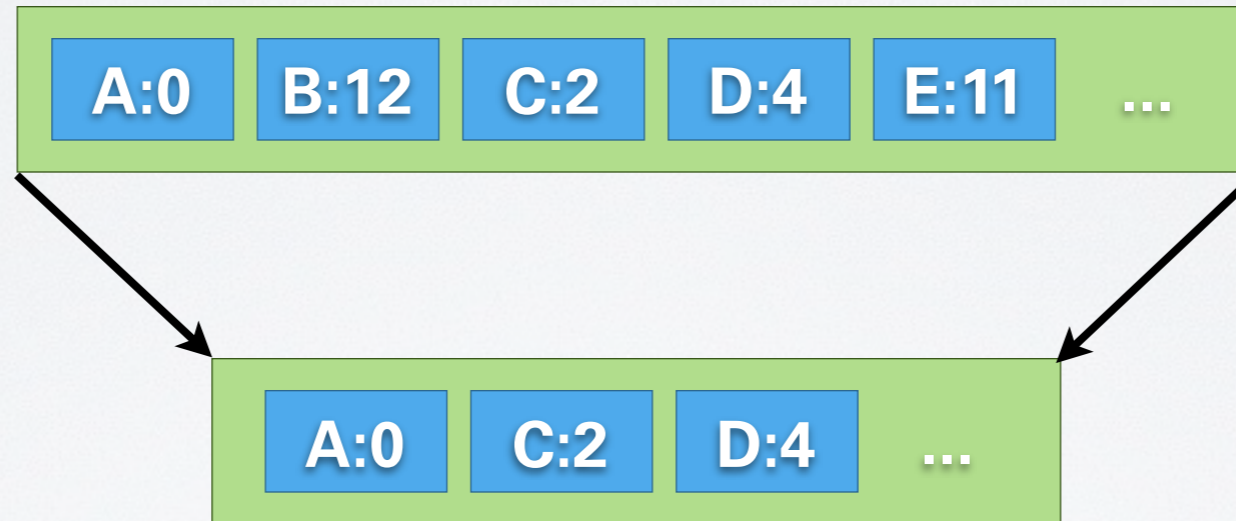
- Management tasks in supercomputers:
 - Process placement
 - Load management
 - System monitoring
- Up-to-date information required to make informed decisions

- Low overhead on application performance
- Scalability:
 - Decentralized information dissemination
 - Decentralized decision making
- Fault tolerance

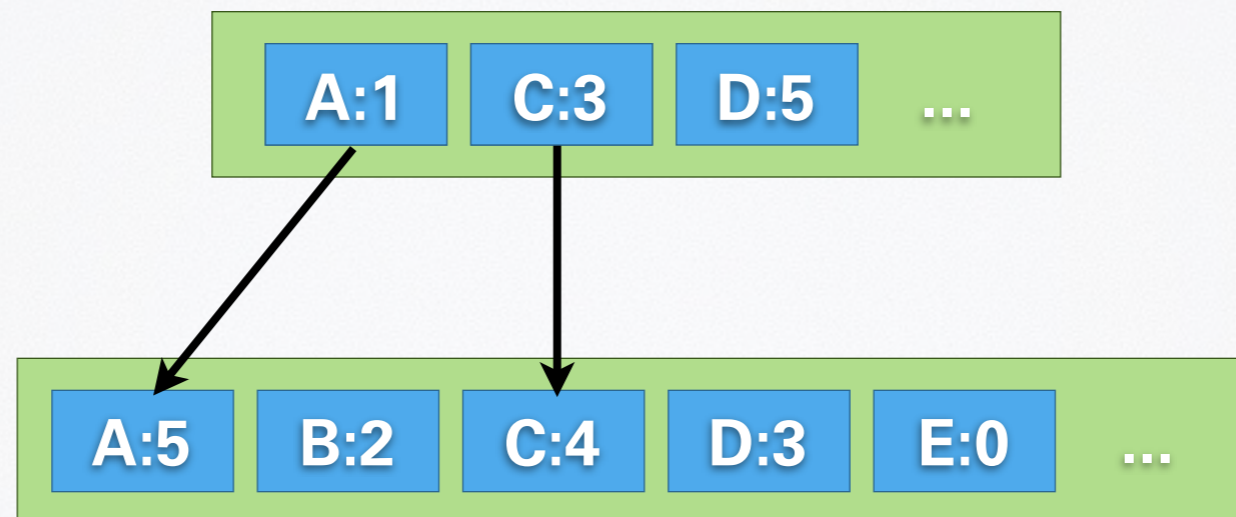


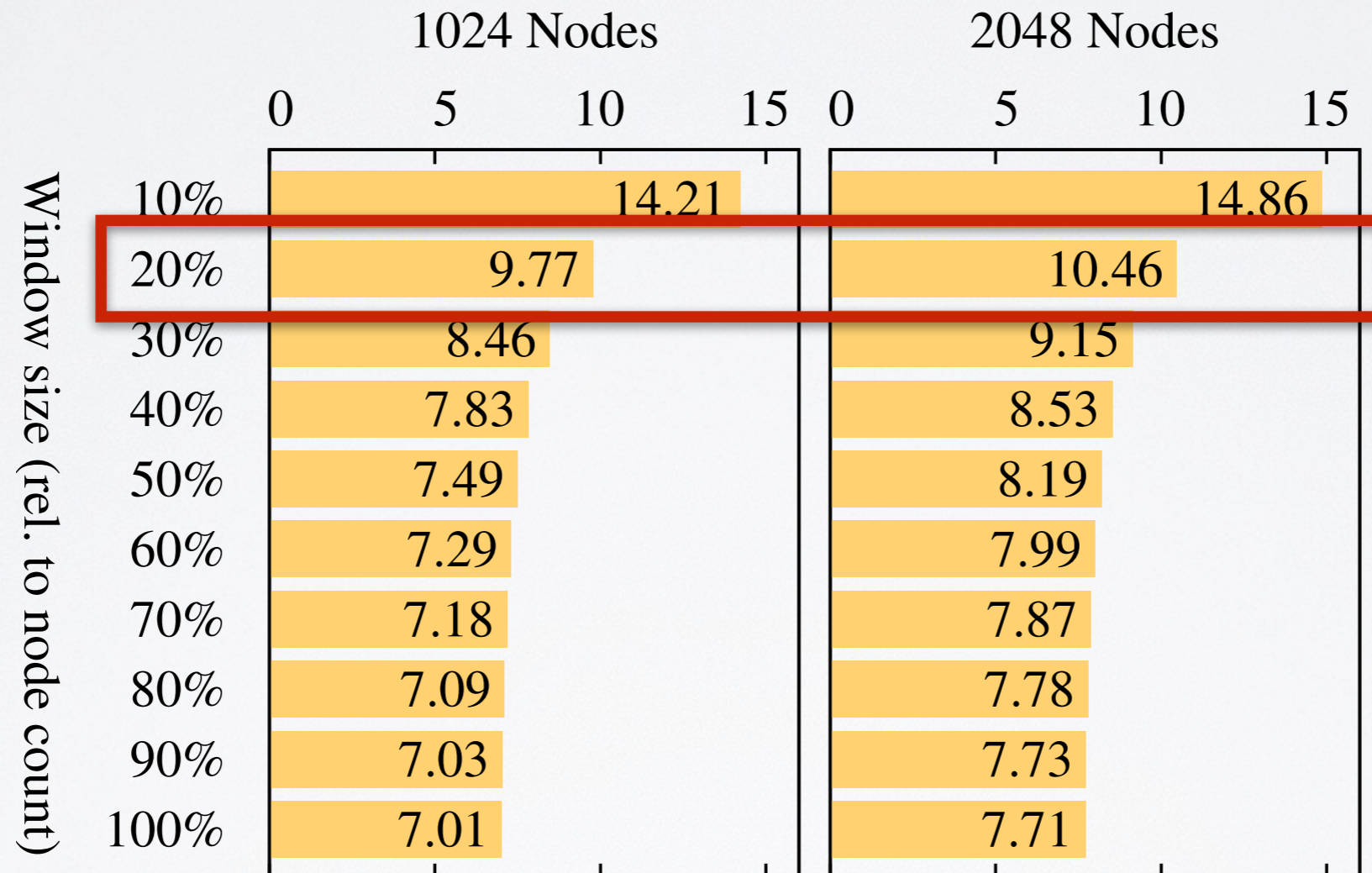
MERGING WINDOWS

Node A



Node E





How much data to send?

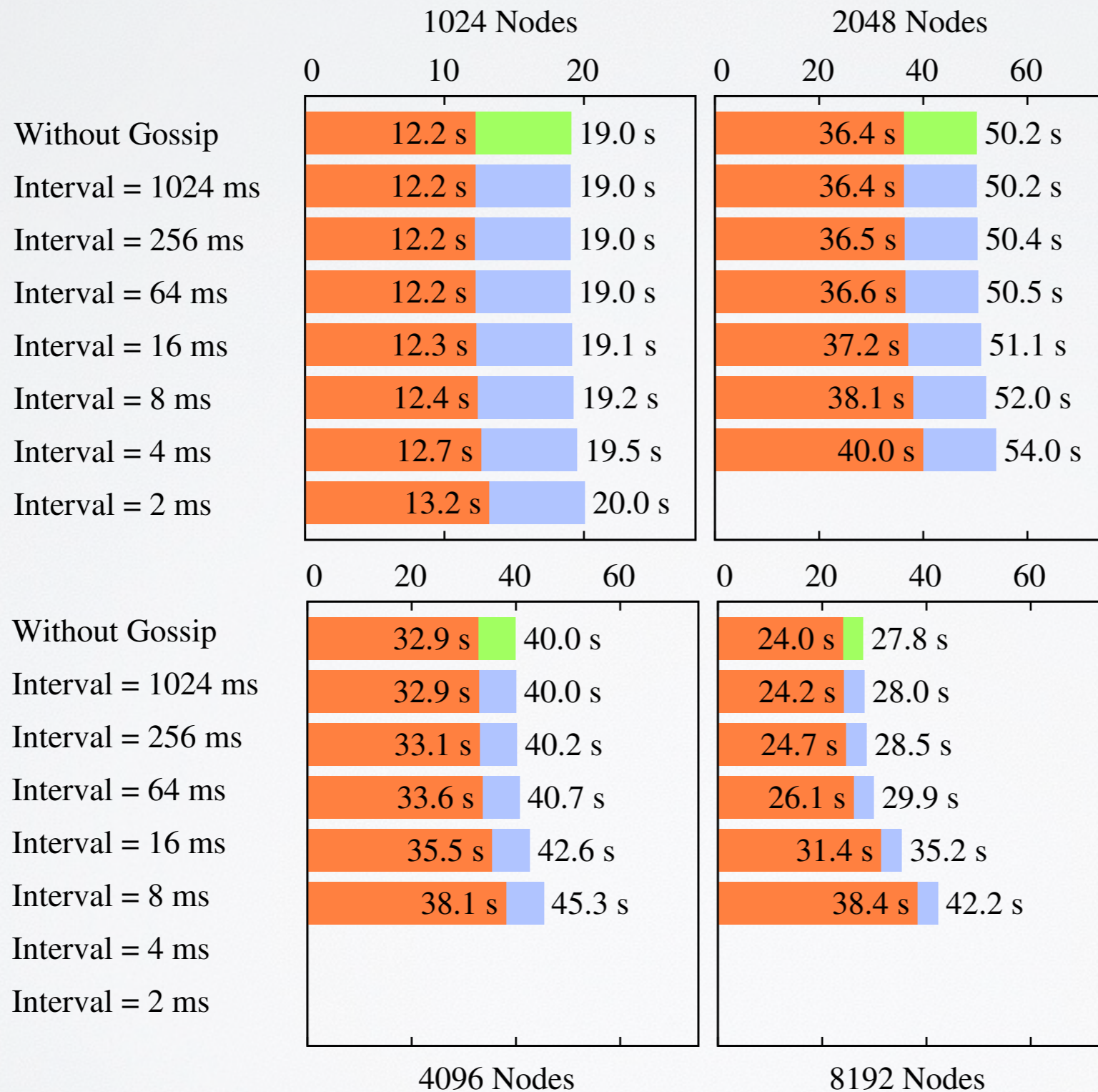
- Small window sizes already yield good average age
- Diminishing return for larger window sizes
- Example: 20% of 1024 nodes w/ 1KiB per node → 200 KiB

- BlueGene/Q at Jülich (JUQUEEN)
- 28.672 nodes total (used 1024 - 8192)
- 16 cores per node (PowerPC A2 @ 1.6 GHz)
- 5D torus network (10 links per node)
- 2 GB/s per link send + receive
- Total bandwidth per node: 40 GB/s
- 2.6 μ s worst-case latency

- MPI-based implementation (MPI_Bsend)
- Gossip algorithm runs on 1 core
- Application uses remaining 15 cores
- How to run two programs on BG/Q?
 - Gossip algorithm + application linked together
 - MPI communicators configured to hide every 16th core from application
 - Wrapped all uses of MPI_COMM_WORLD

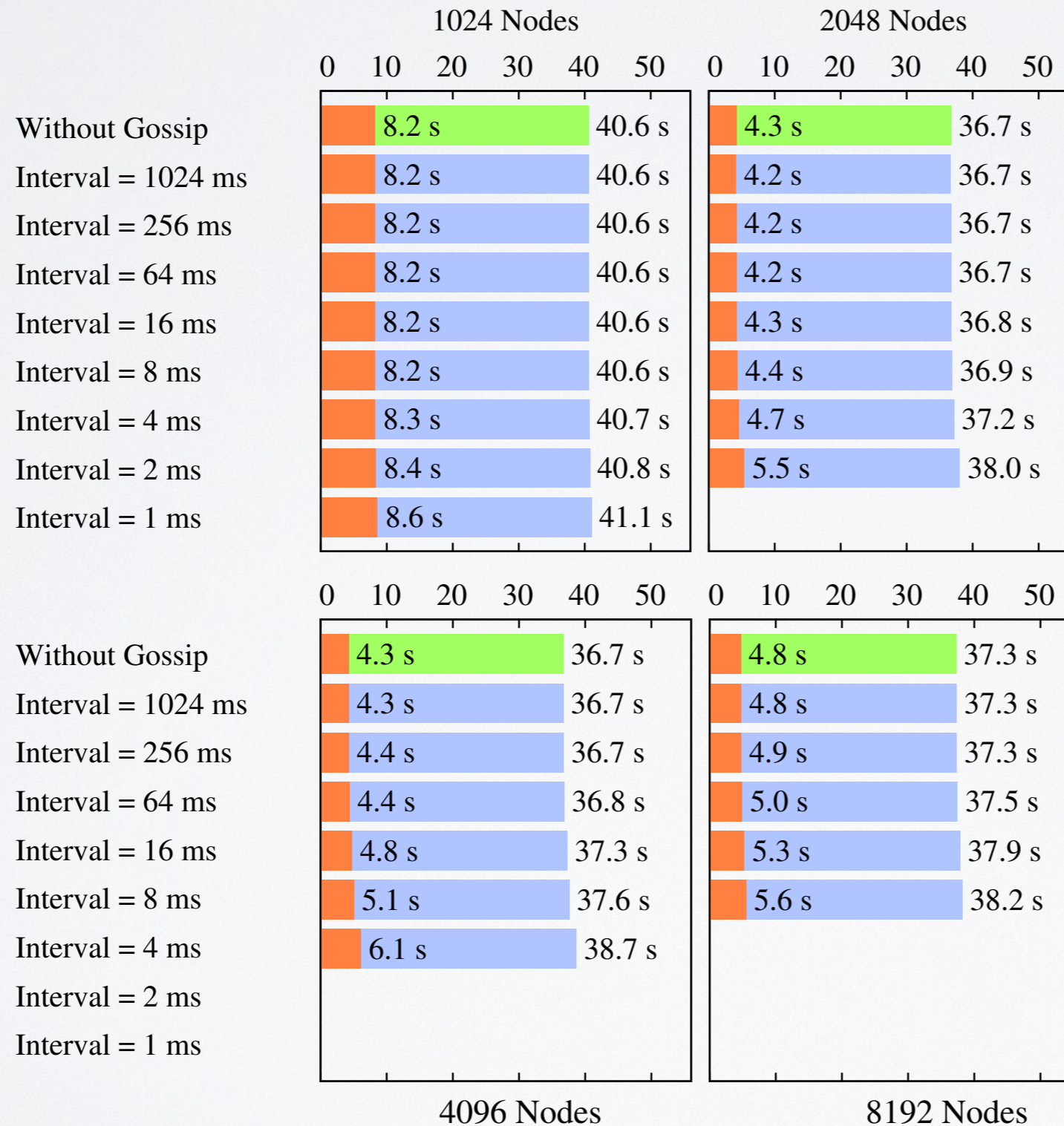
- HPCC suite: **MPI-FFT**  Heavy network usage
- HPCC suite: **PTRANS**
- Application: **COSMO-SPECS+FD4**  Moderate network usage

- HPCC suite: **MPI-FFT**
- HPCC suite: PTRANS
- Application: **COSMO-SPECS+FD4**



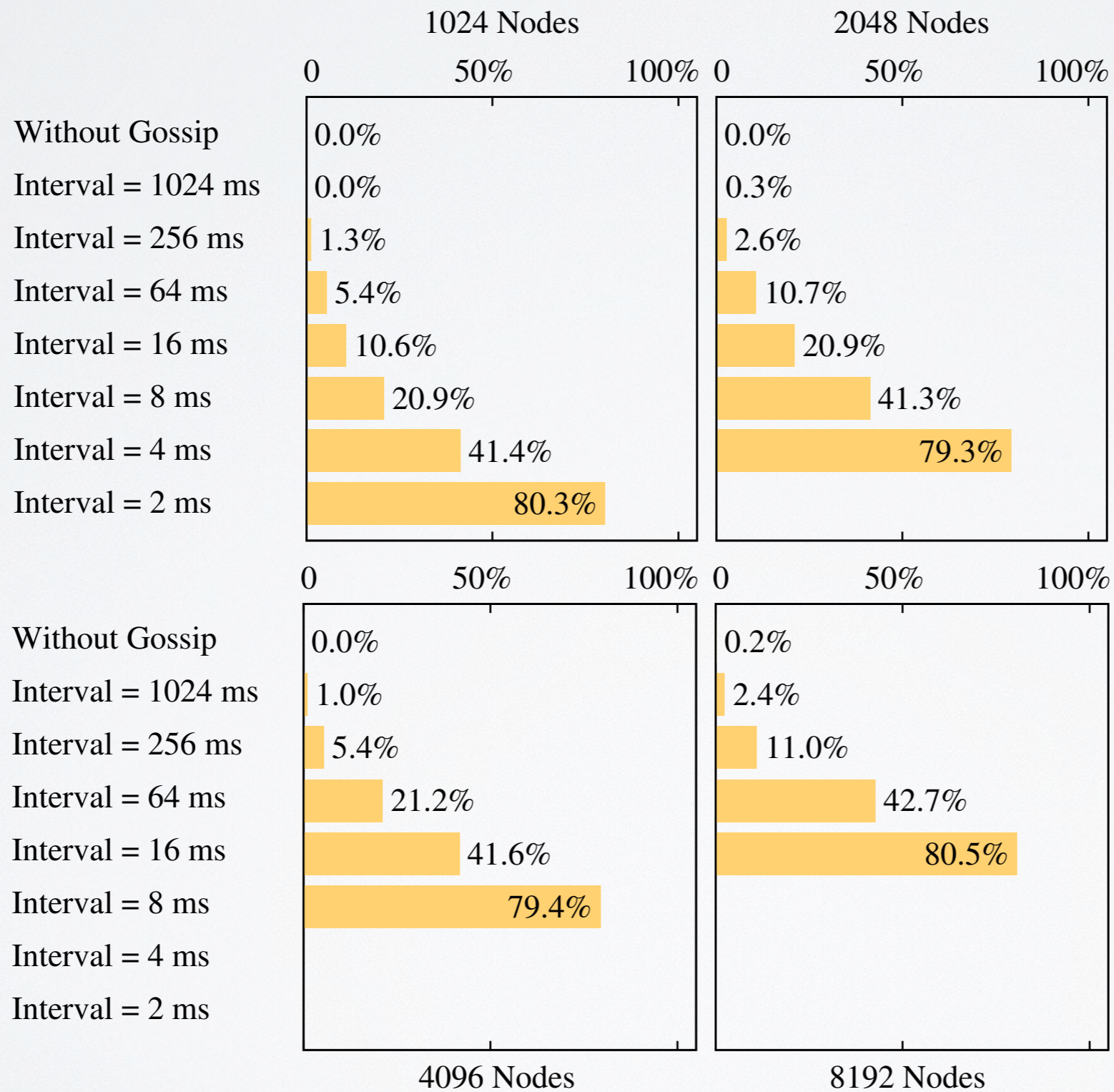
Benchmark: Fast Fourier Transform

- All-to-all communication pattern
- Stresses bisection bandwidth
- 1024 nodes: 136 million vector elements (2025 GiB)
- 2048+ nodes: 544 million vector elements (8100 GiB)



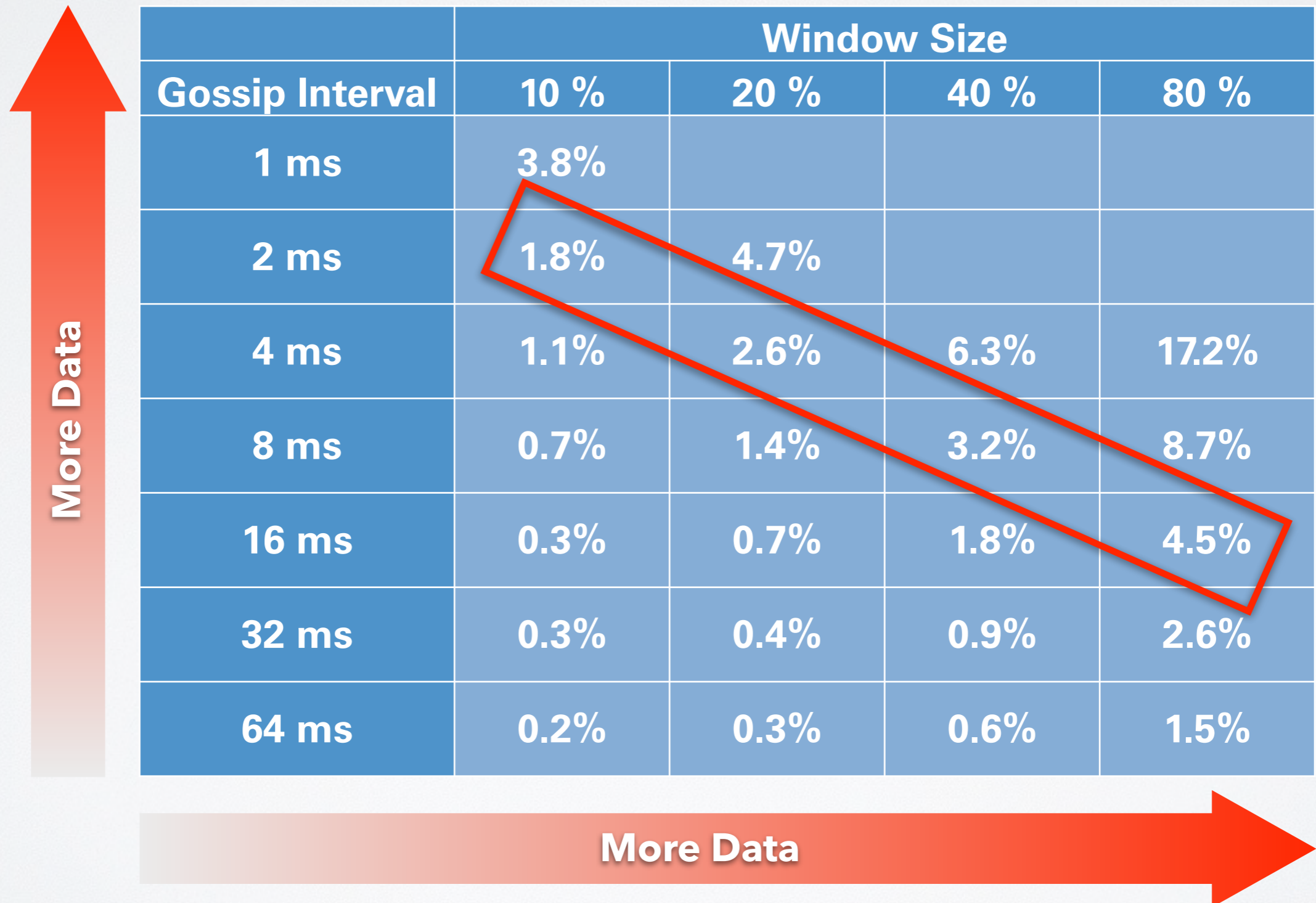
Benchmark: Atmospheric Simulation

- COSMO: static regular communication
- SPECS: dynamic, irregular communication
- Model coupling: dynamic, irregular, small volume
- Partitioning: collectives
- Migration: highly local, mostly between neighbors



Computational complexity:
 $O(n \cdot \log(n))$

RATE VS WINDOW SIZE



- Gossip intervals + average vector age:
 - 256 ms → **2-3 s**
 - 1024 ms → **10 s**
- Applicability for system services:
 - Global load information (allocation, ...)
 - Local load balancing (MOSIX-like, ...)
 - System monitoring (node health, ...)

- Other types of network (Infiniband, Cray, ...)
- Fault tolerance, loss of messages
- Must adapt for exascale systems:
 - Incomplete knowledge at each node
 - Groups of gossip nodes
 - Smaller vectors
 - Hierarchical gossip for global view

- Gossip algorithm scales to thousands of nodes
- Increasing window size causes more overhead than decreasing gossip interval
- Collective MPI communication most sensitive
- Gossip intervals of **256-1024 ms** with **no noticeable overhead** (in most cases)
- **Average age** of information at nodes in the order of **2-3 s** with gossip interval of 256 ms