

A Congestion Avoidance Technique Using Packet Pacing toward Exascale Interconnect

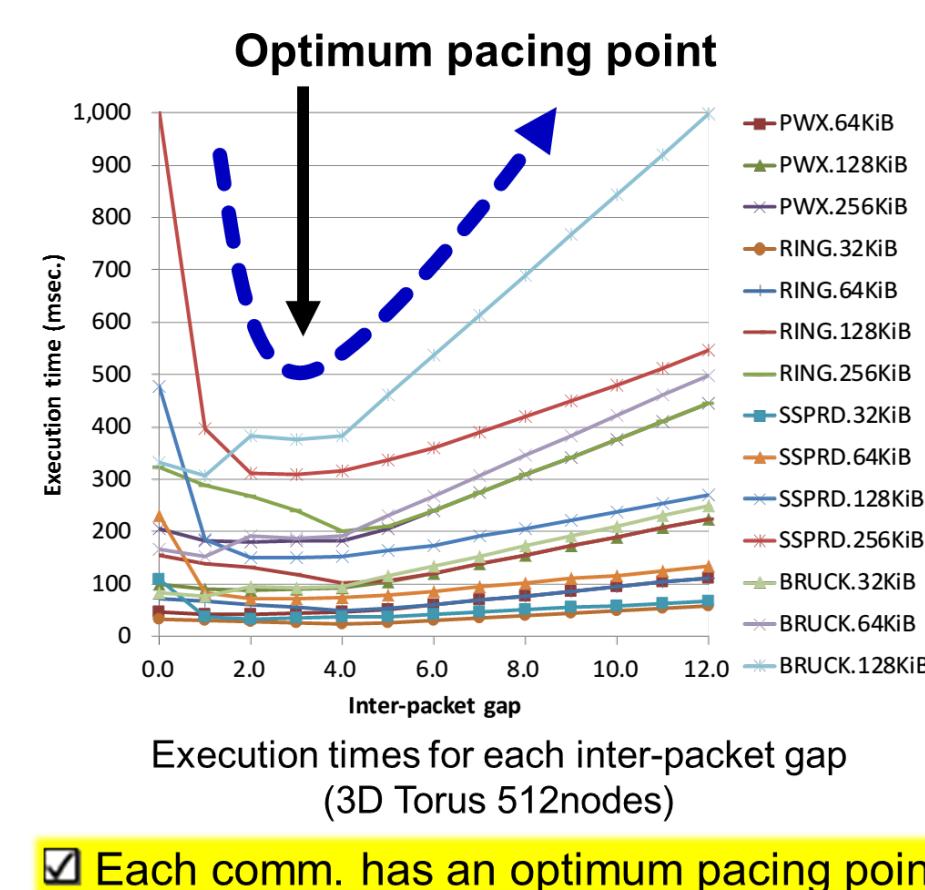
“PACKET PACING” MAKES YOUR COMMUNICATION FAST !

INTRODUCTION

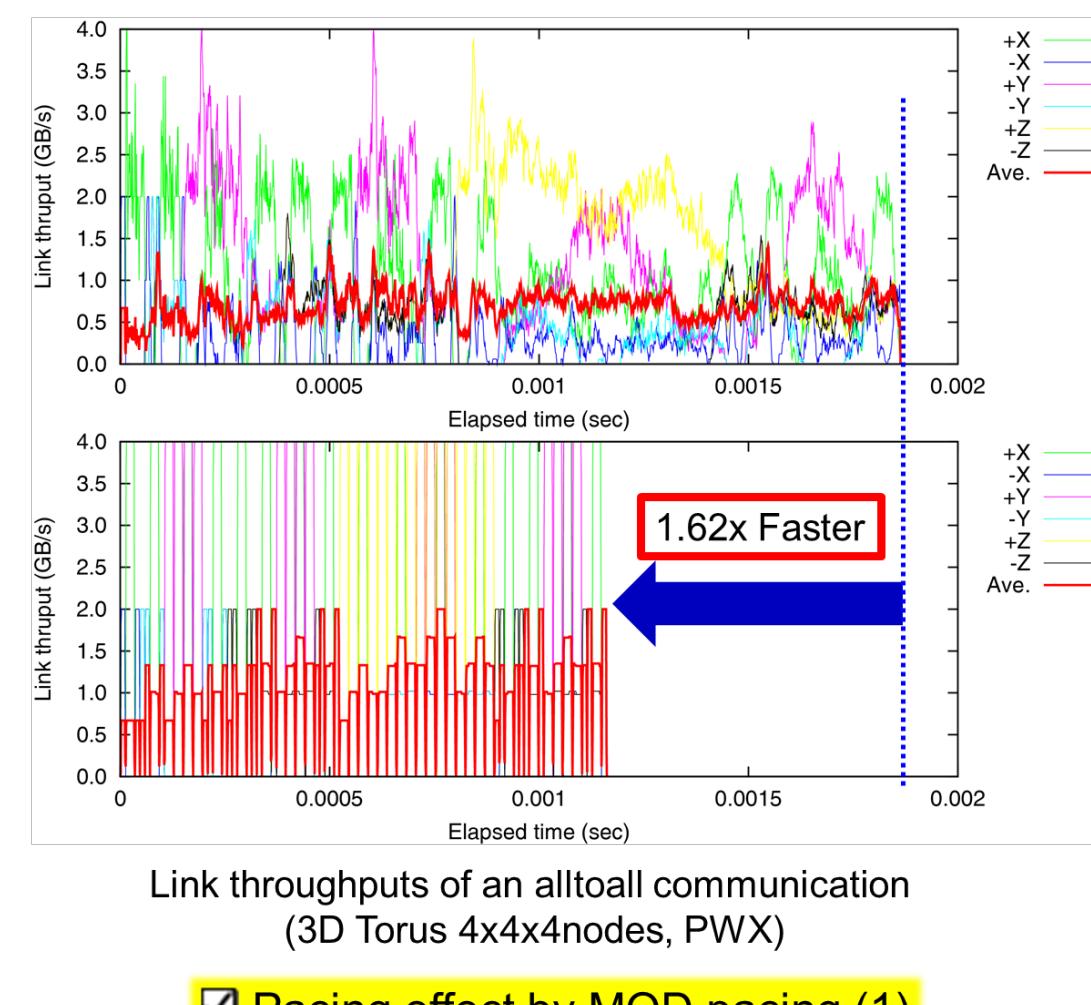
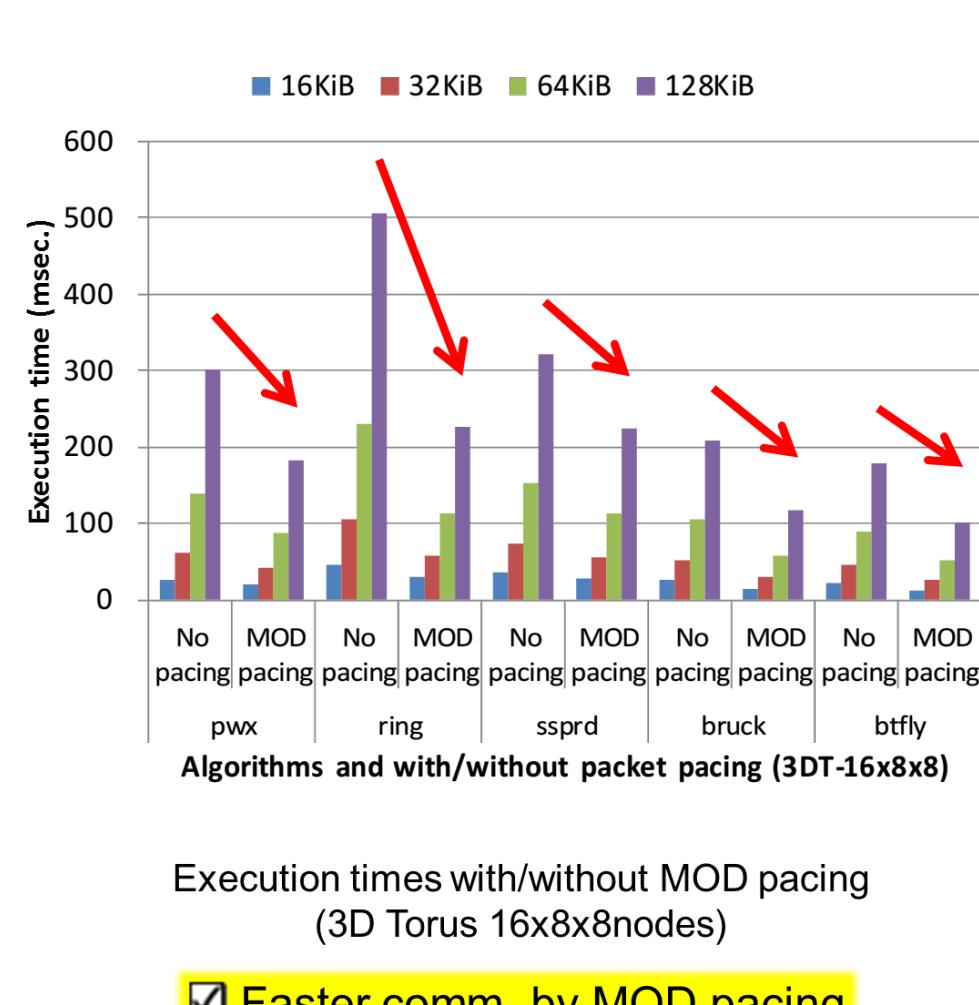
- Reducing communication latency is still important in HPC
 - Some collective comm. may not be able to use in the near future
 - Time-consuming comm. is caused by *network congestion* in heavy traffic.
 - Packet pacing avoids the critical network congestion

WHAT IS PACKET PACING?

- Insertion of non-sending period (**inter-packet gap**) between sending packets
 - Controls packet injection rate by interleaving packets
 - Gives an optimum inter-packet gap for each message.
 - Reduces network congestion
 - Decreases stop & go latency.
 - Maximizes throughput & minimizes network latency
 - Applications with heavy traffic comm. get faster.
- Inter-packet gap
 - An interval between sending packets.
 - Gap = 0: No packet pacing.
 - Gap = N: Transfer time for N packets.



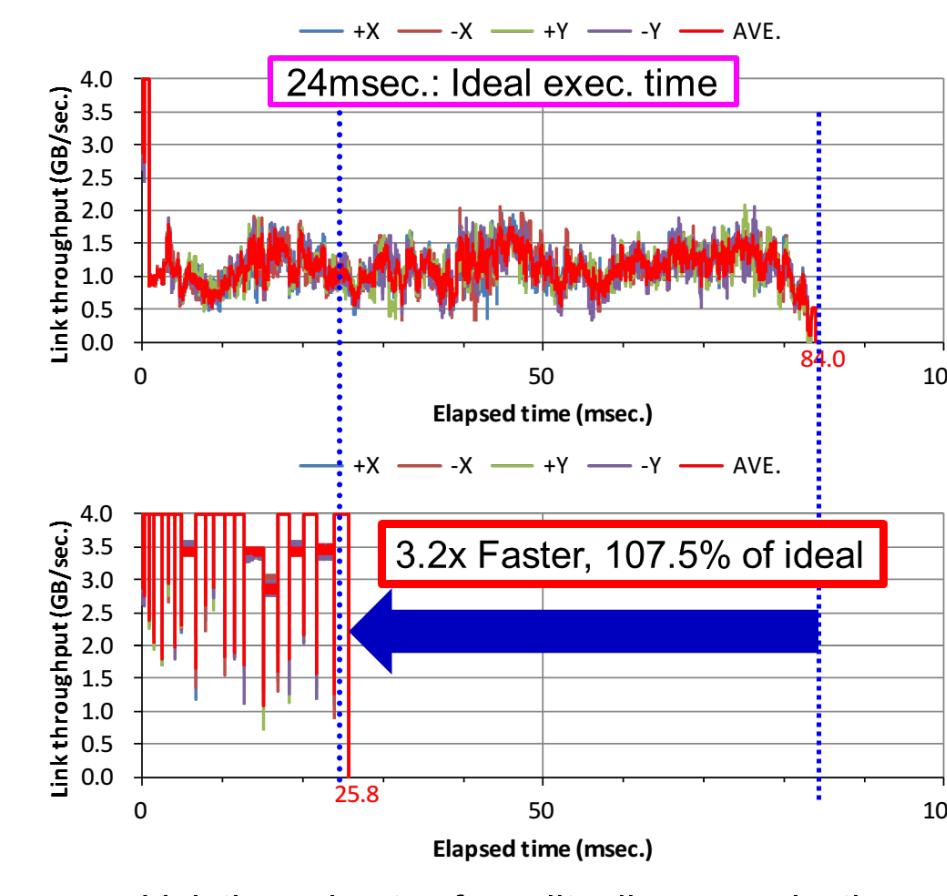
EXPERIMENTAL RESULTS ON SIMULATION (NSIM)



Execution times with/without MOD pacing (3D Torus 16x8x8nodes)

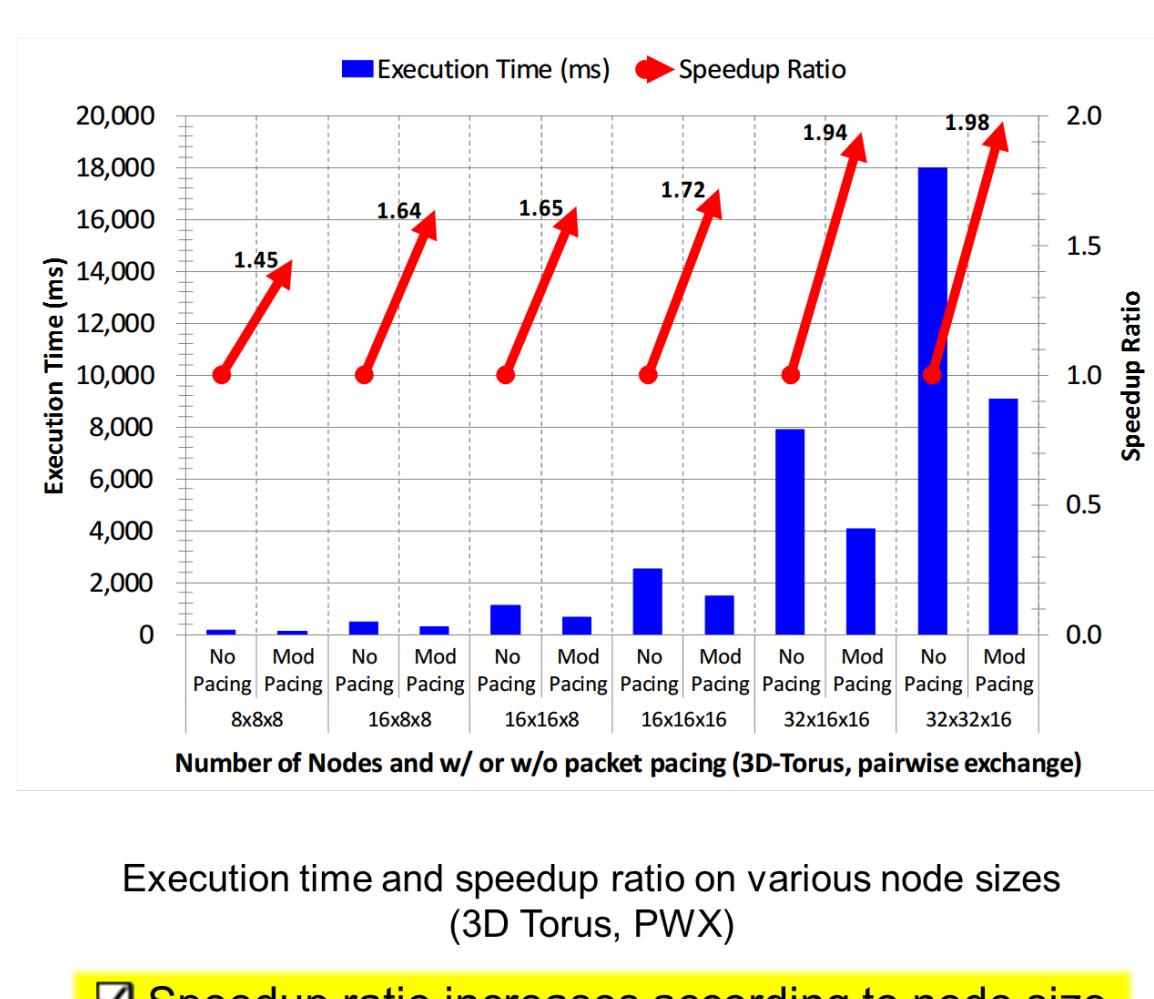
Faster comm. by MOD pacing

- Alltoall algorithm: PWX(Pairwise exchange), RING, SSPRD(Simple spread), BRUCK, BFLY(Butterfly), and A2AT(by Ishihata et al. at Tokyo Univ. of Tech.).
- Link bandwidth: 4GB/s. Routing algorithm: DOR(Dimension Ordered Routing).



Link throughputs of an alloverall communication (2DT-9x9 nodes, A2AT)

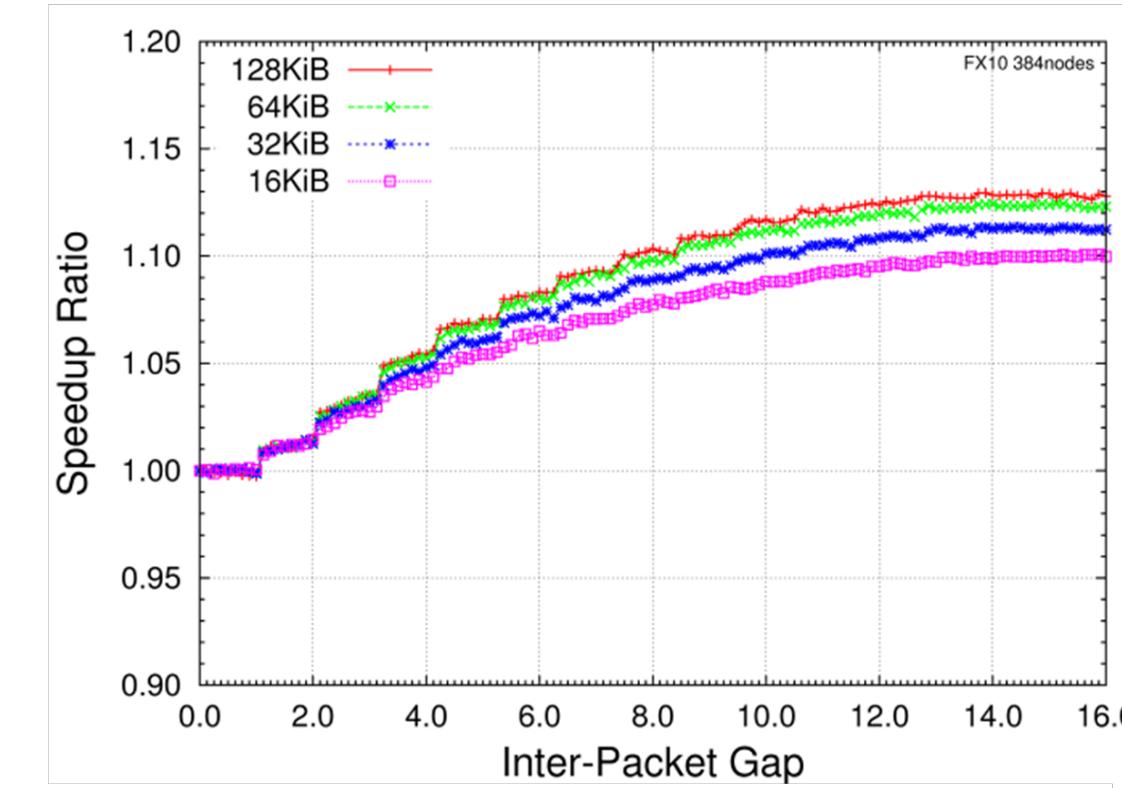
Pacing effect by MOD pacing (2)



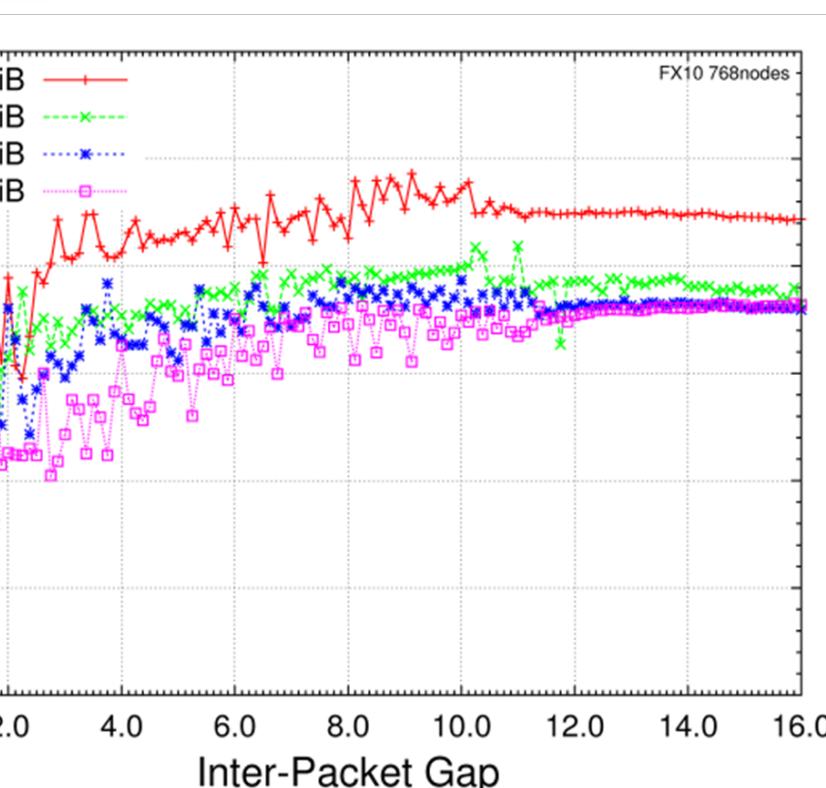
Execution time and speedup ratio on various node sizes (3D Torus, PWX)

Speedup ratio increases according to node size

- Random ring communications (from HPCC) with packet pacing.
- Measured on Fujitsu FX10 at Kyushu Univ., JAPAN.



384 nodes (6,144 procs.)



768 nodes (12,288 procs.)

Optimum pacing point and effect of packet pacing grows as node size increases

SUMMARY

Packet pacing : controls packet injection aggressively to accelerate HPC comm.

- Each heavy traffic communication has its own **optimum pacing point**.
- The **MOD pacing strategy** finds the best pacing point.
- The **Effectiveness improves** by message size and/or node size.



Hidetomo Shibamura <shibamura@isit.or.jp>

Institute of Systems, Information Technologies and Nanotechnologies
Japan Science and Technology Agency, CREST