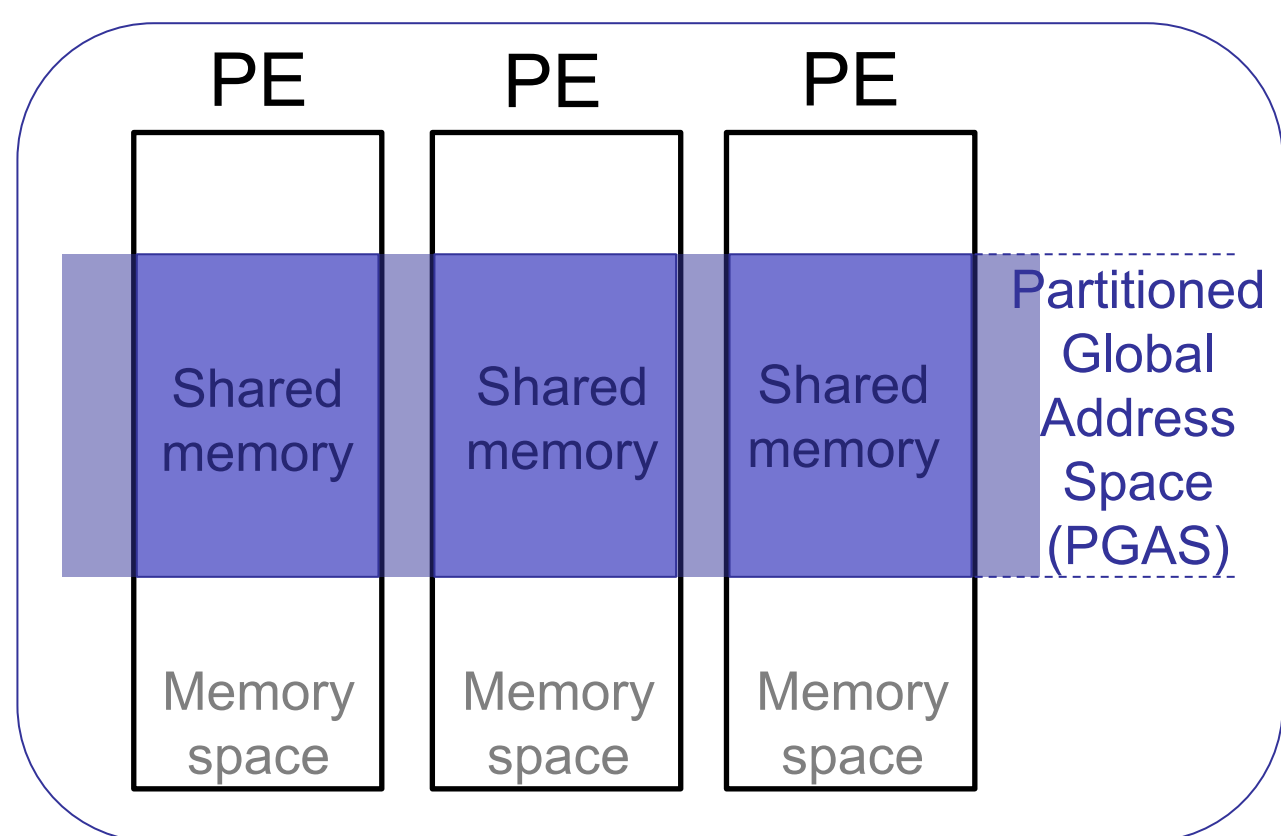


## Introduction and Motivation

- OpenSHMEM is a widely used PGAS programming model in the HPC area
- The *team* concept and team-based collective communication in the latest OpenSHMEM Specification v1.5 are similar to the communicator and collective communication in MPI



- The different design approaches (the native one-sided and the MPI-based two-sided communication) can lead to different performance characteristics on HPC clusters (even for the same collective routine)

- We characterize two aspects that can influence the performance
  - Synchronization methods
  - Collective algorithms

- We compare the native one-sided design and the MPI-based two-sided design in collectives using benchmarks, and show how big the performance difference between them

## Experiment Setup

- Library Selection:**
  - One-sided design: Sandia OpenSHMEM (SOS) [1]
  - Two-sided design: OpenSHMEM over MPI (OSHMPI) [2][3]
- Platform:**
  - Bebop HPC cluster [4]. Each node has 36-core Intel Broadwell CPU. The nodes are connected with Omni-Path Fabric
- MPI Implementation:** MPICH [5][6]
- Benchmark:** OSU Micro-Benchmarks [7]

## Conclusion

- The performance of OpenSHMEM varies a lot (up to 10X) with several factors
- The comparison exposes the performance characteristics of different collective communication designs in OpenSHMEM
- We believe the performance characterizations can give the community some insights for future research avenues

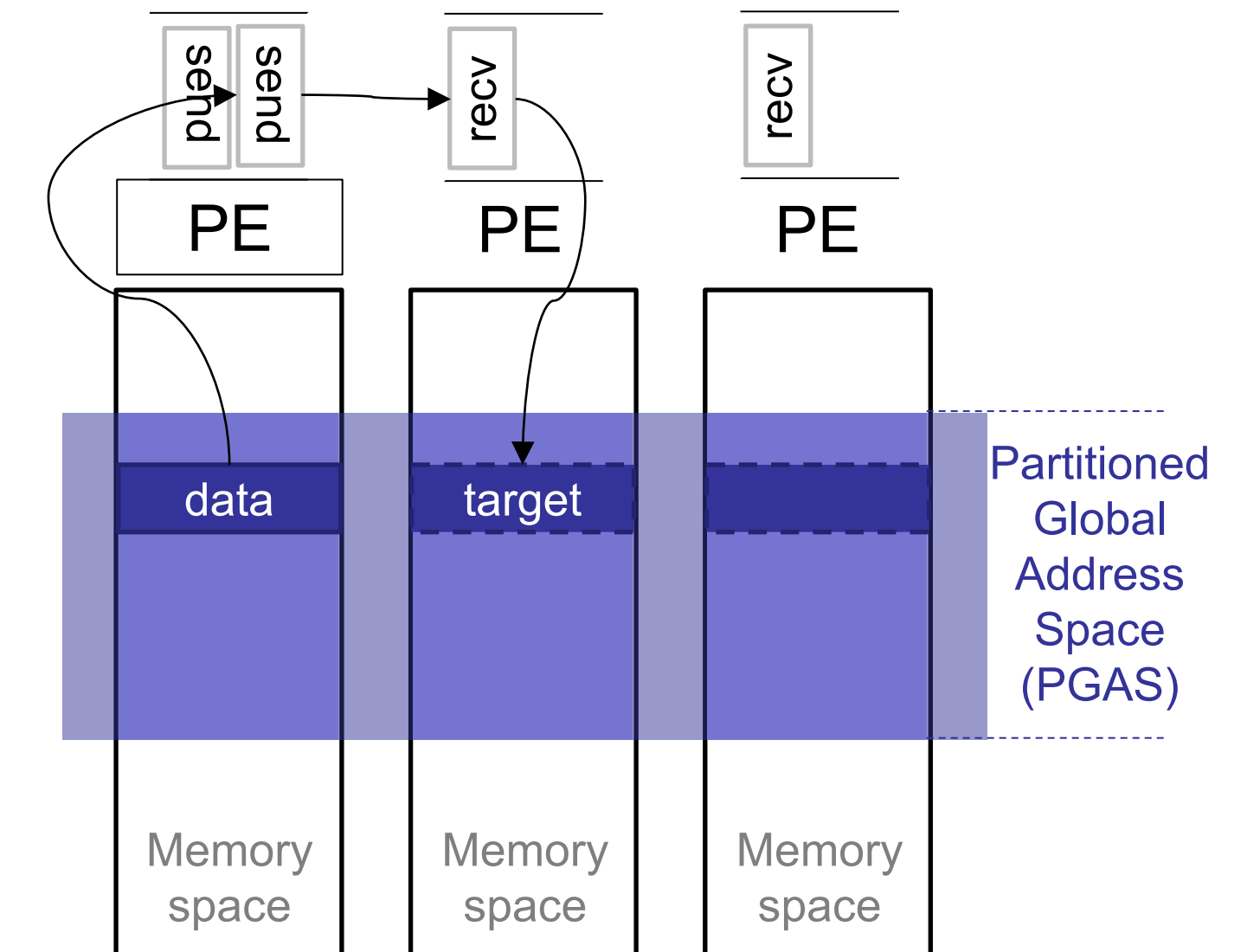
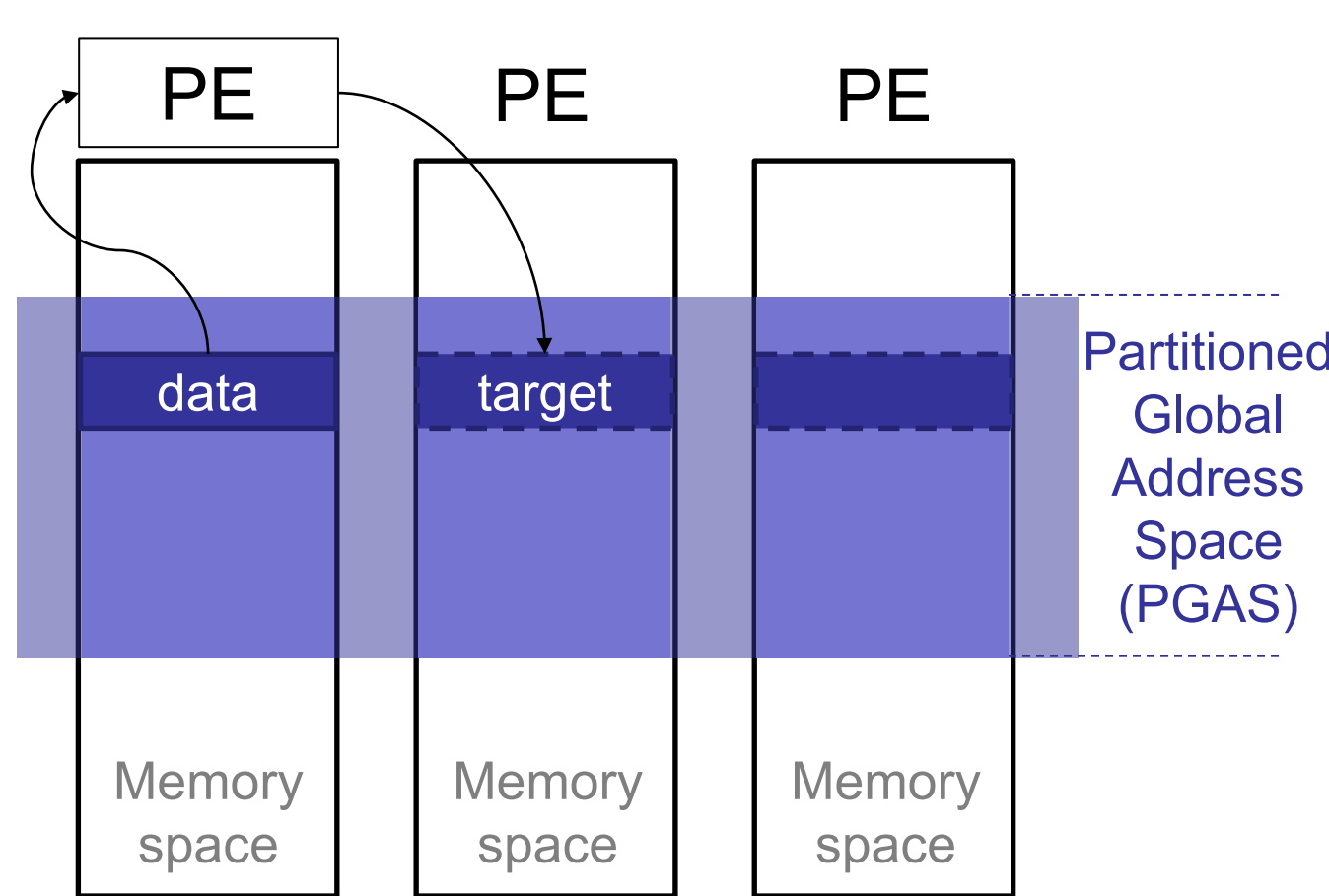
## References

- S. Corporation, "Sandia-openshmem/sos." [Online]. Available: <https://github.com/Sandia-OpenSHMEM/SOS>
- M. Si, H. Fu, J. R. Hammond, and P. Balaji, "OpenShmem over mpi as a performance contender: Thorough analysis and optimizations," in OpenSHMEM and Related Technologies. OpenSHMEM in the Era of Exascale and Smart Networks, S. Poole, O. Hernandez, M. Baker, and T. Curtis, Eds. Cham: Springer International Publishing, 2022, pp. 39–60.
- A. N. Laboratory, "Oshmpi: Openshmem implementation over mpi." [Online]. Available: <https://github.com/pmodels/oshmpi>
- "Bebop." [Online]. Available: <https://www.lrcr.anl.gov/systems/resources/bebop/>
- A. N. Laboratory, "Mpich: High-performance portable mpi." [Online]. Available: <https://www.mpich.org/>
- "Official mpich repository." [Online]. Available: <https://github.com/pmodels/mpich>
- D. K. Panda, H. Subramoni, C.-H. Chu, and M. Bayatpour, "The mvapich project: Transforming research into high-performance mpi library for hpc community," *Journal of Computational Science*, vol. 52, p. 101208, 2021, case Studies in Translational Computer Science. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S187750320305093>

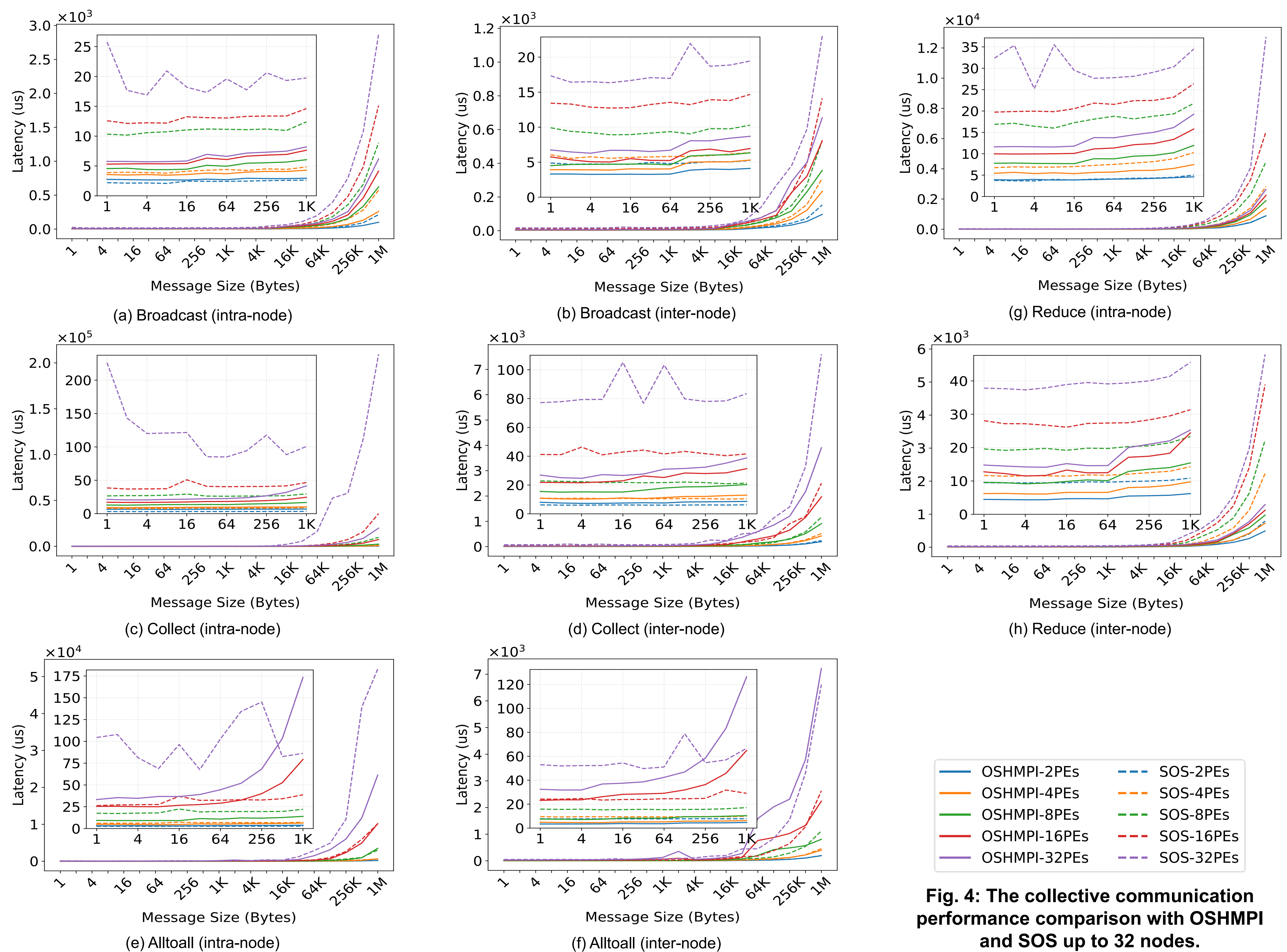
This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration, and by the U.S. Department of Energy, Office of Science, under Contract DE-AC02-06CH11357. We gratefully acknowledge the computing resources provided on Bebop (and/or Swing and/or Blues), a high-performance computing cluster operated by the Laboratory Computing Resource Center at Argonne National Laboratory. This work was supported in part by the NSF research grant OAC #2321123 and a DOE research grant DE-SC0024207.

## Designs of OpenSHMEM Libraries

- Sandia OpenSHMEM (SOS)**
  - Native one-sided design -- a source PE *puts* data into (or gets data from) the shared global memory of the target PE
  - The synchronization of the one-sided design does not require the target PE to actively acknowledge data receiving to the source PE
  - SOS collectives natively implement some algorithms – Tree, Ring, Round-Robin algorithms
- OpenSHMEM over MPI (OSHMPI)**
  - MPI-based two-sided design -- a source PE *sends* data, which matches with a *recv* called by a target PE
  - Less flexible synchronization. It needs *send* and *recv* coordination between the source and the target PEs
  - OSHMPI collectives essentially call MPI collectives. Taking MPICH as an example, it provides rich collective algorithms – Tree, Ring, Bruck, Recursive Doubling, etc. The algorithm can be chosen at runtime by scheduling



## Performance Characterization



Key Findings	Reason Explanations
OSHMPI (two-sided) is faster than SOS (one-sided) in most cases.	The existing two-sided design (e.g., OSHMPI) can inherit the advantages in well-optimized MPI engine, while the one-sided design (e.g., SOS) still needs optimized implementations.
SOS shows lower or comparable latency in certain cases: <ul style="list-style-type: none"> <li>Communication with a small number of PEs like 2, for <i>collect</i>, <i>intra-node broadcast</i>, <i>alltoall</i>, and <i>reduce</i> collectives.</li> <li>The <i>alltoall</i> collective for medium messages in <i>intra-node</i> communication and medium-large messages in <i>inter-node</i> communication with more PEs, like 16 and 32.</li> </ul>	With a small number of PEs, basic point-to-point primitive performance and synchronization method dominate the performance. SOS's native one-sided design can achieve better overlapping with simpler synchronizations. Besides, SOS also has some specific optimizations on Omni-Path Fabric. For the <i>alltoall</i> collective, MPI internal routine needs further optimization or tuning for these particular settings.
SOS shows unstable performance in the case of <i>intra-node</i> communication with 32 PEs.	SOS uses extra helper threads to progress the internal communication tasks and thus it is oversubscribed.