



Performance Insights into Device-initiated RMA using Kokkos Remote Spaces

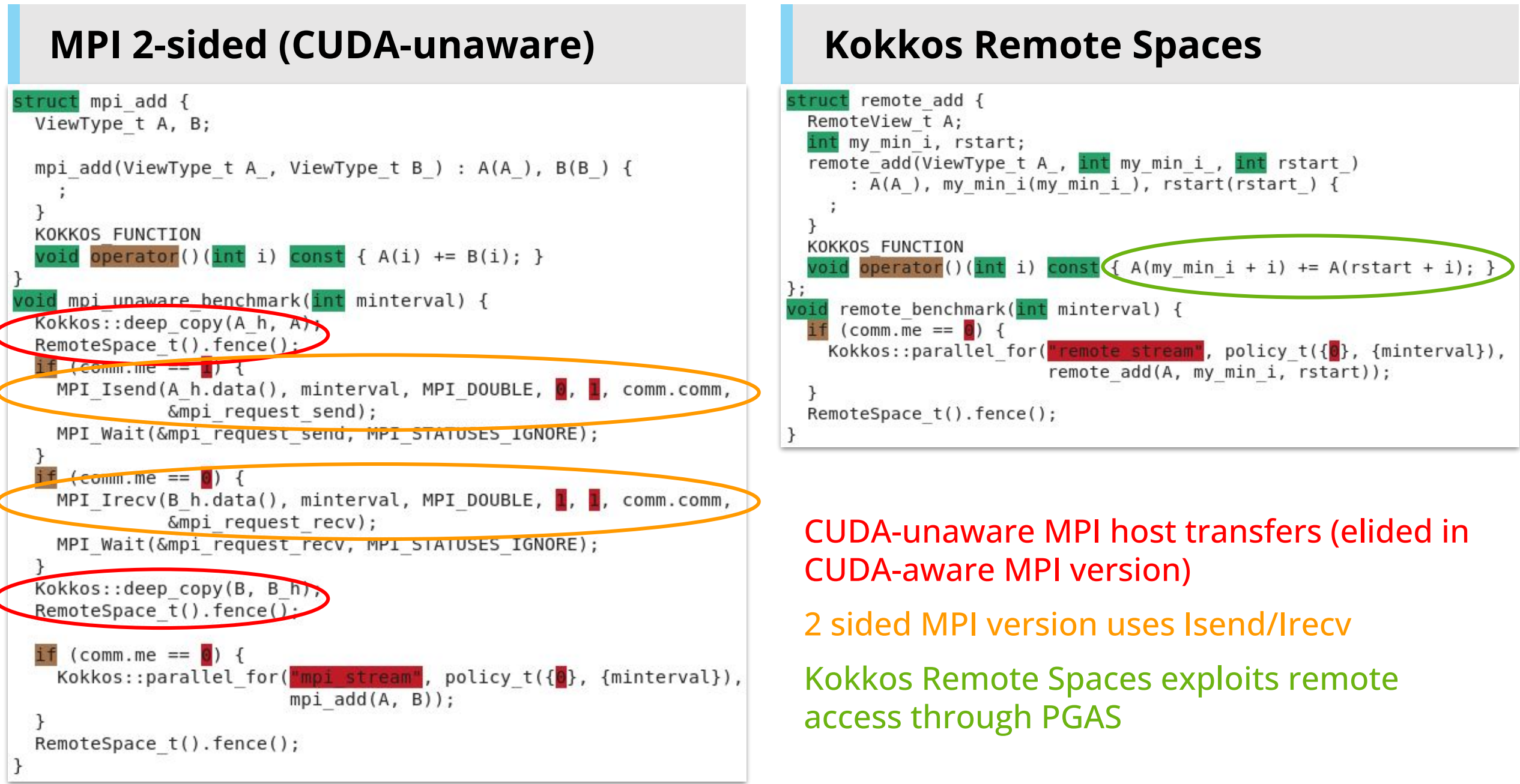
Daniel Mishler^{1,2}, Jan Ciesko¹, Stephen L. Olivier¹, George Bosilca²

¹Sandia National Laboratories, ²University of Tennessee, Knoxville

Achieving scalable performance on supercomputers requires careful coordination of communication and computation. Often, MPI applications rely on buffering, packing, and sorting techniques to accommodate a two-sided API, minimize communication overhead, and achieve performance goals. As interconnects between accelerators become more performant and scalable, programming models such as SHMEM may have the opportunity to enable bandwidth maximization along with ease of programming. In this work, we take a closer look at device-initiated PGAS programming models using NVIDIA Corp’s NVSHMEM communication library and our interface through the Kokkos Remote Spaces project. We show that benchmarks can benefit from this programming model in terms of performance and programmability. We anticipate similar results for mini-applications.

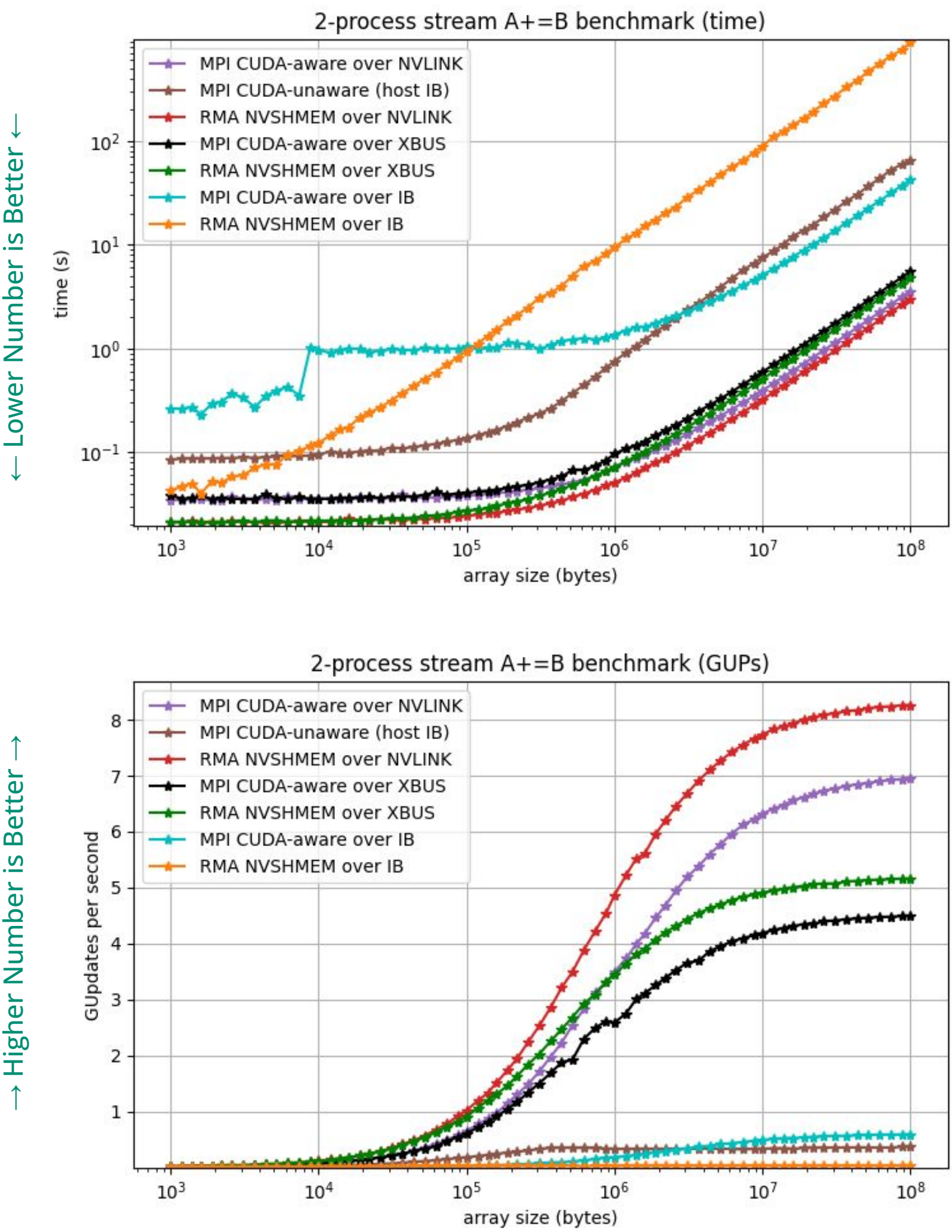
Kokkos Remote Spaces

- Kokkos: <https://github.com/kokkos/kokkos/>
 - Cross-platform framework which exposes parallelism.
 - Used in production applications at Sandia: reliable and well maintained.
- Kokkos Remote Spaces: <https://github.com/kokkos/kokkos-remote-spaces>
 - Also cross-platform framework on top of Kokkos: exposes PGAS & SHMEM.
 - Objective: provide a reliable and usable software layer for programmers.
- Potential application benefits
 - 3x less LOC for conjugate gradient compared to MPI two sided version.
 - Enables more flexible ad hoc data exchange for irregular applications (compared to send-lists in MPI).



Above: traces visualized with NVIDIA Nsight that show time taken for host-device transfers that can be skipped with device-device communication exposed through Kokkos Remote Spaces.

Results



Above are two representations of results from streaming $A(i) += B(i)$ element-wise on 2xV100 NVIDIA GPUs. The benchmark provides 2 arrays, A belonging to process 0, and B belonging to process 1.

MPI: Uses 2 sided MPI to exchange data. B must be packed into a buffer and sent to process 0.

RMA: A and B both accessed through PGAS. Kokkos Remote Spaces exposes NVSHMEM to perform the operations.

We performed the runs with CUDA-aware and CUDA-unaware MPI, the latter being a host transfer. “XBUS” represents non-NVLINK on-node GPU to GPU connections.

Poor performance is observed in Kokkos Remote Spaces over IB because at the time of this poster’s publishing, data transfers only occur element-wise in Kokkos Remote Spaces.

Kokkos Remote Spaces Current Status

Feature	Current Capability
Supported Vendors	NVIDIA Corp, AMD
RMA Operations	Put, Get
Parallel Patterns	All routines Kokkos supports
Debugging support	All routines Kokkos supports

Future Challenges

- Add support for using the shmem calls for more than a single element (e.g. send 128 bytes in one call) to enable OSU latency benchmark.
- Include MPI one-sided codesize and performance in benchmarks.
- Investigate remaining performance anomalies.
- Evaluate on miniapplications (e.g. heat3d, conjugate gradient).