

IEEE CCGRID 2025

The 25th IEEE International Symposium on Cluster, Cloud and Internet Computing
Tromsø, Norway

ActorISx: Exploiting Asynchrony for Scalable High-Performance Integer Sort

Youssef Elmougy, Shubhendra Singhal, **Akihiro Hayashi***, and Vivek Sarkar

Habanero Extreme Scale Software Research Lab

Georgia Institute of **Technology**

* Presenting Author: ahayashi@gatech.edu

Integer and Data Movement Performance in HPC

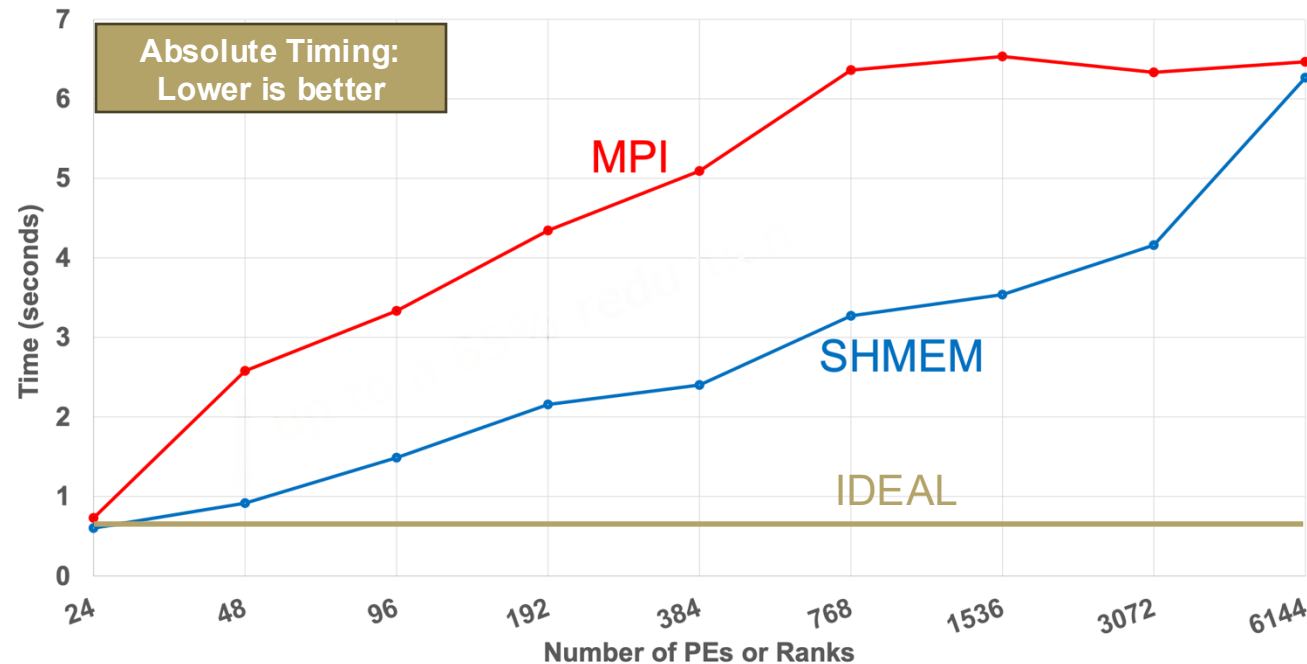
- Integer sorting remains a critical operation in HPC, playing a vital role in numerous applications
 - Graph Analytics
 - Distributed Sparse Linear Algebra
 - Scientific Simulations
 - and more!
- The Integer Sort (IS) kernel from the NAS Parallel Benchmark (NPB) has served as a standard for assessing integer and data movement performance in HPC systems since 90s
 - Two-sided MPI-based Bucket Sort Algorithm with `MPI_Alltoallv`

IS — The NPB 2 implementation of the IS kernel benchmark is based on a bucket sort. The number of keys ranked, number of processors used, and number of buckets employed are all presumed to be powers of two. This simplifies the coding effort and leads to a compact program. The number of buckets is a tuning parameter. On the systems tested, best performance was obtained when the number of buckets was half that which gives best load balancing. Communication costs are dominated by an `MPI_Alltoallv`, wherein each processor sends to all others those keys which fall in the key range of the recipient.

ISx: Addressing limitations of the original NPB IS

- The development of **ISx** addressed several limitations of the original NPB IS benchmark by introducing OpenSHMEM's one-sided communication

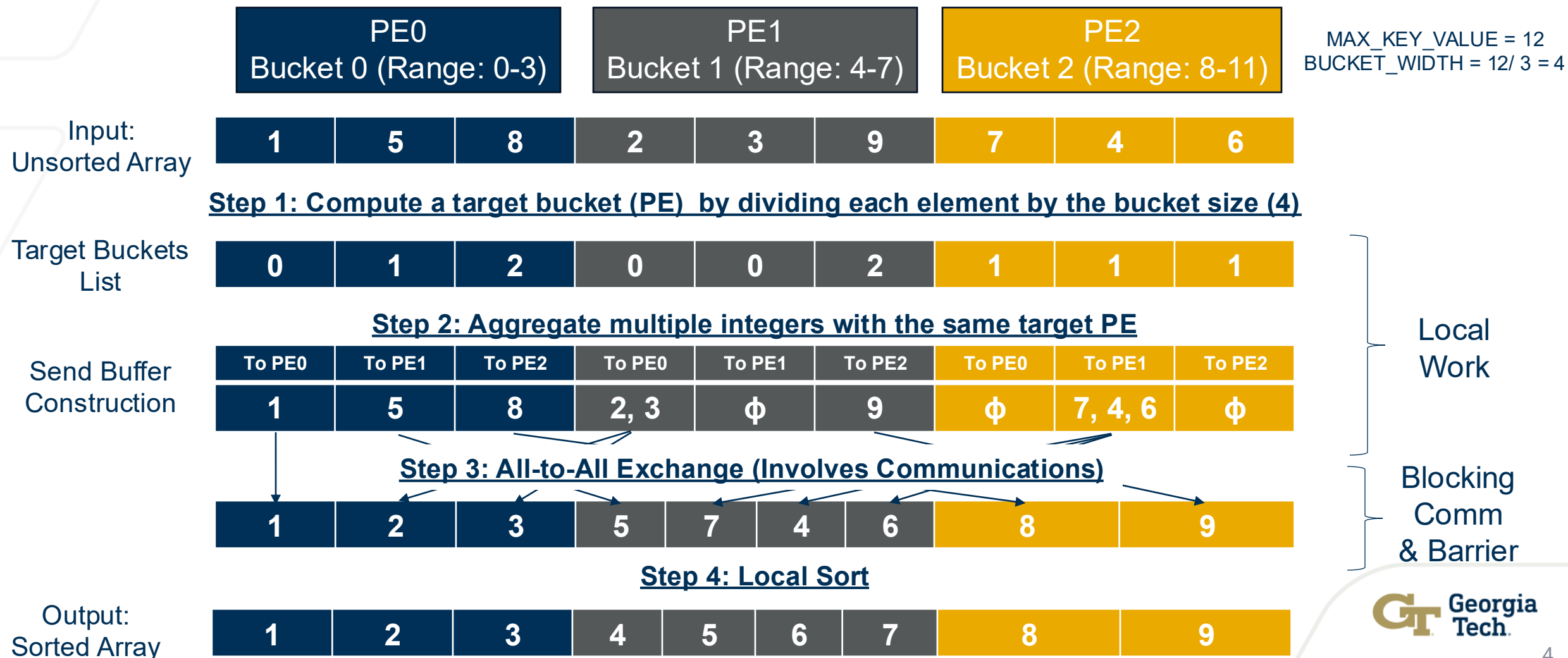
Weak-Scaling Performance of the Key Exchange part of ISx (2^{27} keys per PE)



Research Question: Can we come closer to the IDEAL performance on modern distributed HPC systems?

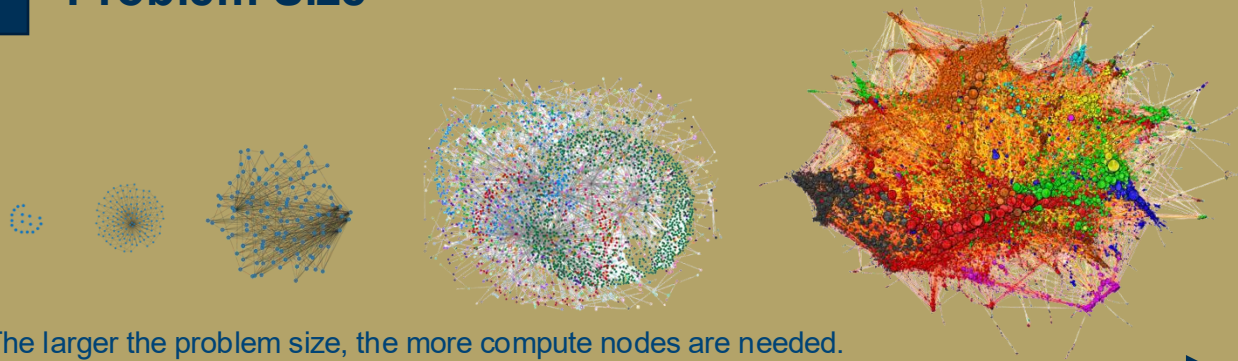
ISx: High-Level Overview

Goal: Sort a distributed array on 3 PEs (Ranks) using the Bucket Sort Algorithm
Input: an unsorted array, Output: a sorted array



The 3 Challenges in Scalable Integer Sort:

1 Problem Size

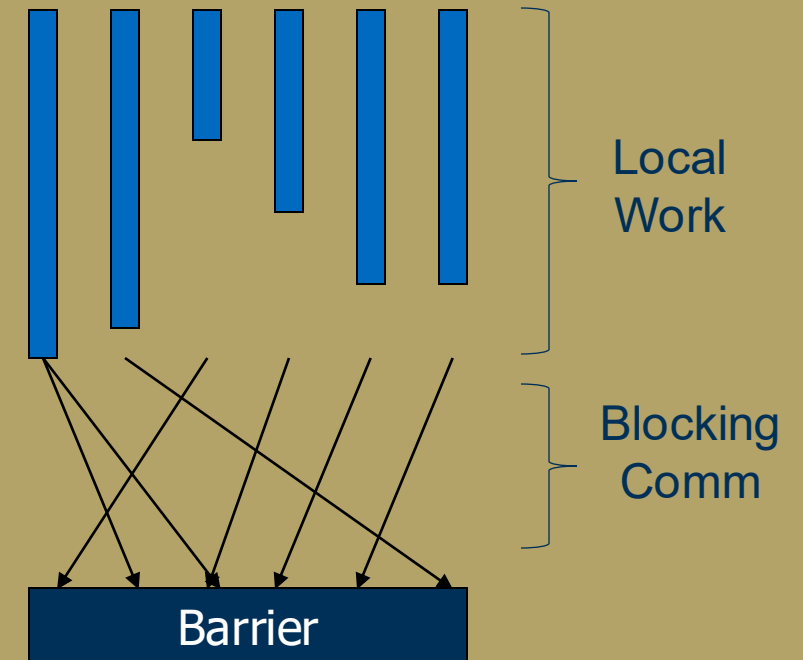


2 Irregular Messaging for Arbitrary Destinations

From	To PE0	To PE1	To PE2	...
PE0	1	5	8	
PE1	2, 3	ϕ	9	
PE2	ϕ	7, 4, 6	ϕ	
...				

3

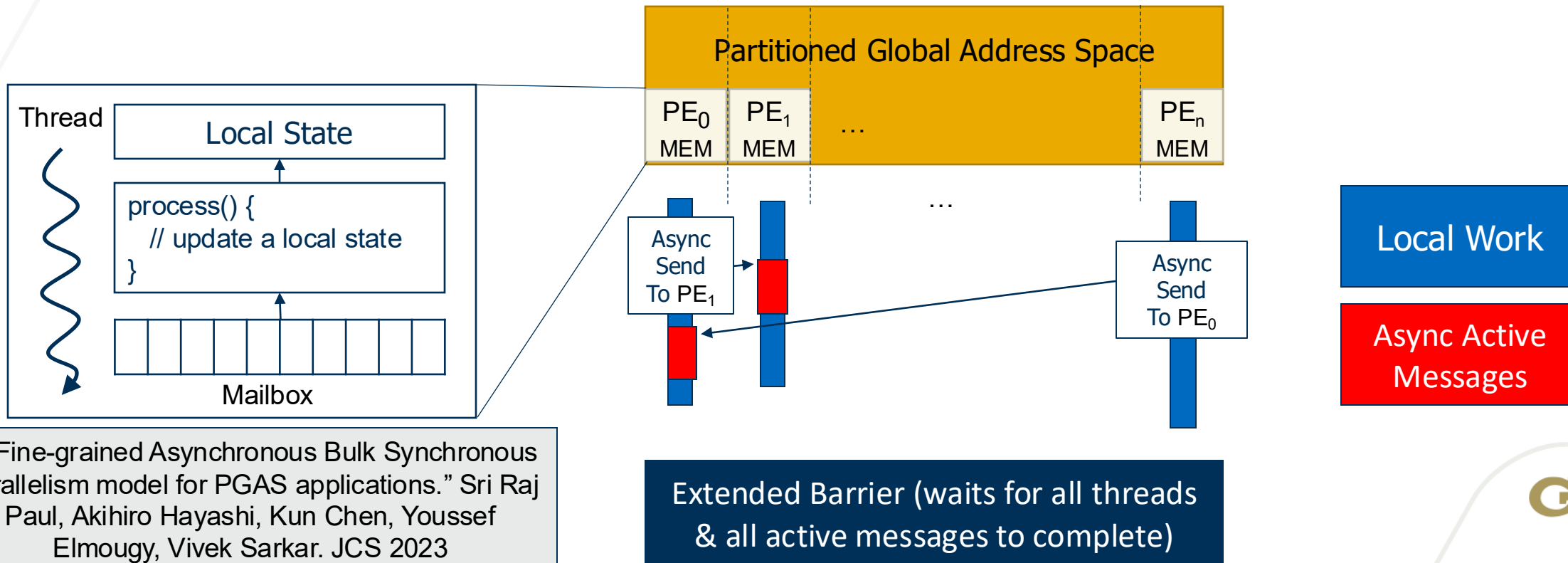
Bulk Synchronous Parallelism (BSP)



This paper studies the scalability of our actor-based approach to overcome the inherent challenges of traditional programming models by performing large-scale integer sorting.

Our Vision: Fine-grained-Asynchronous Bulk-Synchronous Model (FA-BSP)

- FA-BSP = PGAS + Actor-based asynchronous messaging + BSP
 - Actor = PE = a thread/rank that owns a slice of the global address space
 - Communications between PEs are achieved via asynchronous active messages
 - The extended barrier waits for all active messages to complete in a superstep



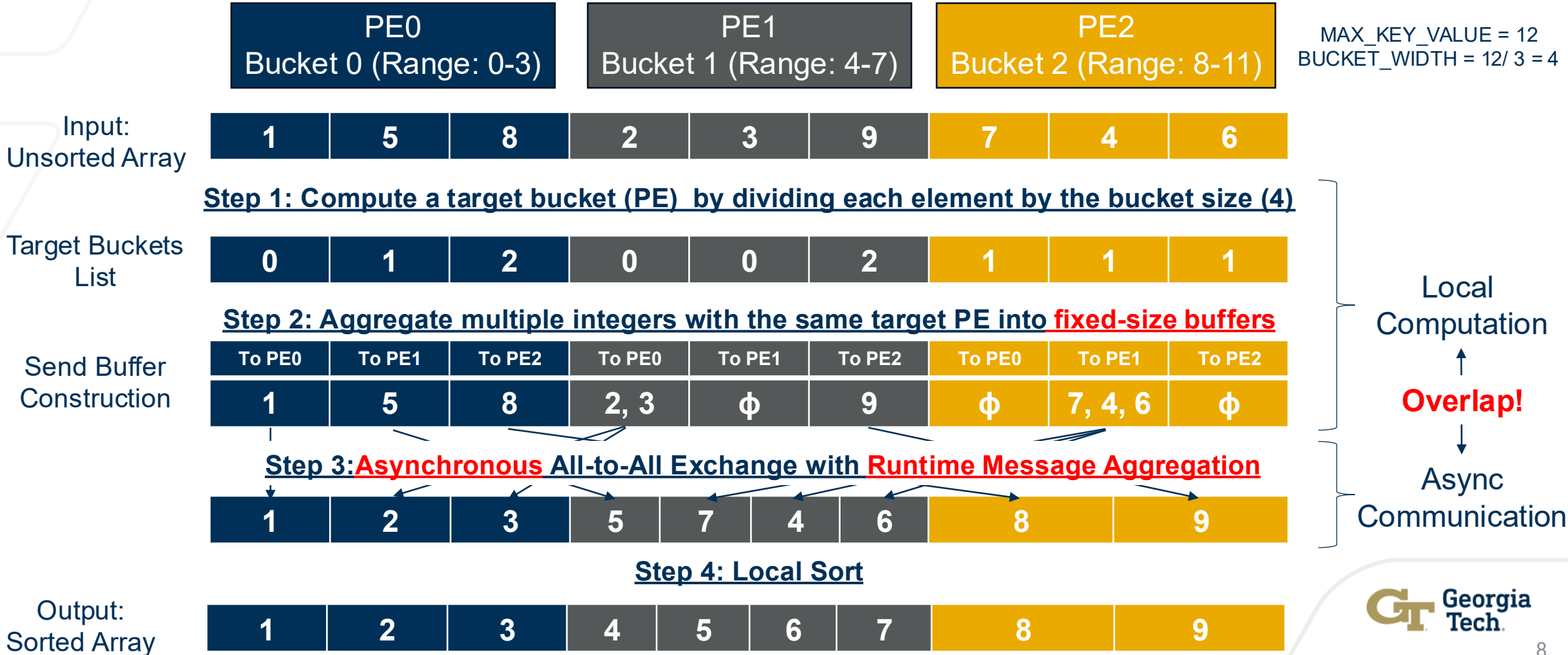
Many distributed graph/non-graph algorithms can be implemented using our FA-BSP model

- Our FA-BSP model is well-suited for graph and non-graph algorithms
- Bale Kernels (JCS'23)
 - Histogram
 - Index Gather
 - Permute Matrix
 - Random Permutation
 - Transpose Matrix
 - Triangle Counting
 - Toposort
- Other Graph kernels
 - Triangle Centrality (SCALE Challenge at CCGRID'24)
 - Page Rank (SCALE Challenge at CCGRID'23)
 - Jaccard Coefficients (ISC'24)
 - Triangle Counting (SC'23 Poster)
 - BFS
- IARPA AGILE Workflows
 - Graph Neural Networks
 - Pattern Matching (IPDPSW'25)
 - K-mer Counting (IPDPS'25)
 - Influence Maximization (SC'24)



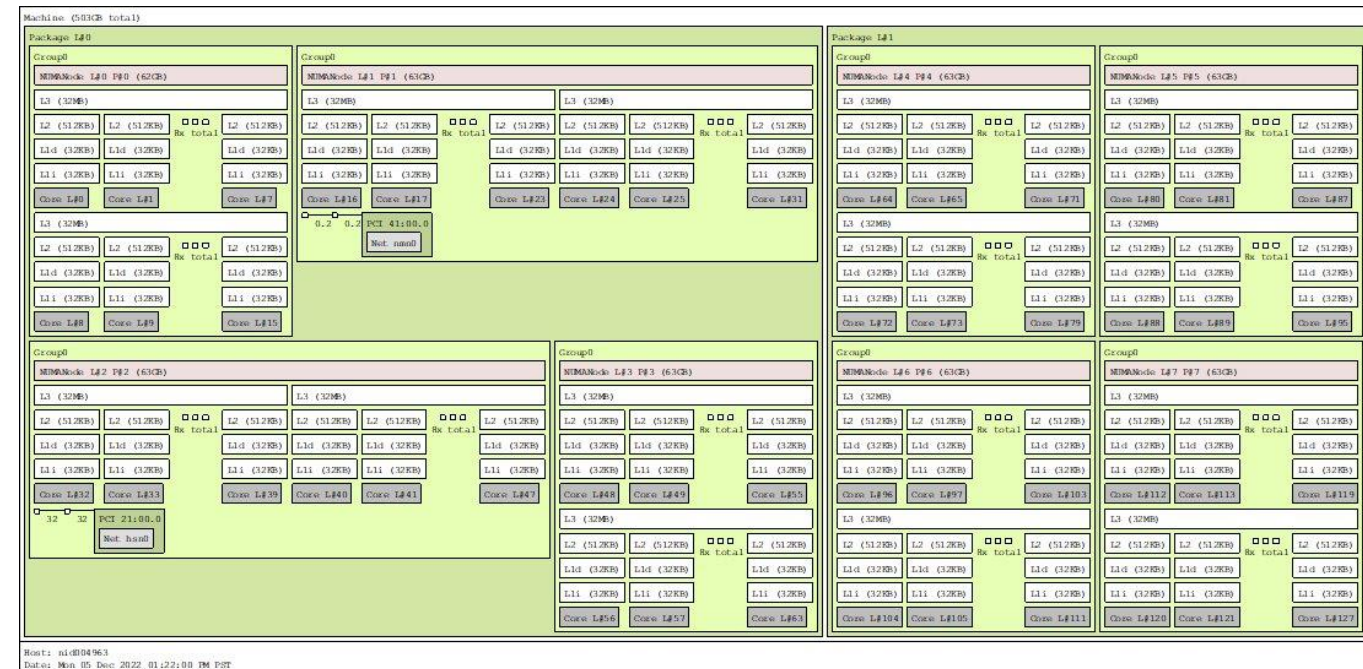
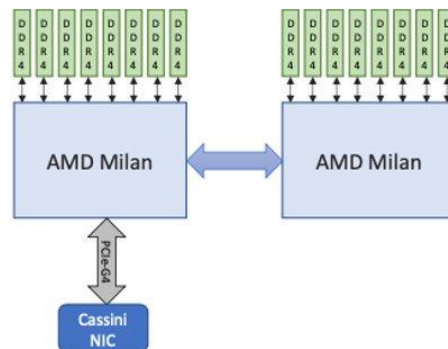
ActorISx: Exploiting Asynchrony for Scalable High-Performance Integer Sort

ActorISx enables asynchronous point-to-point communication with application-level and runtime-level message aggregation



Experimental Setup and Architecture

- Experiments conducted on the CPU nodes of the **Perlmutter supercomputer** at the National Energy Research Scientific Computing Center (NERSC)
 - 2x AMD EPYC 7763 (Milan) CPUs
 - 64 physical cores per CPU
 - 512 GB memory
 - 1x HPE Cray Slingshot Interconnect
- Results for **different dimensions of scalability** are presented



Picture borrowed from: <https://docs.nersc.gov/systems/perlmutter/architecture/>

Dimensions of Scalability

The paper studies two weak-scaling scenarios from the original ISx:

	WEAK SCALING	ISO WEAK SCALING
# of cores	Variable ($2^6 - 2^{14}$)	
Keys per core	Constant (2^{27})	
MAX_KEY_VALUE	Constant (2^{28})	Variable ($2^{19} - 2^{27}$)
BUCKET_WIDTH = MAX_KEY_VALUE / # of cores	Variable (2^{22} to 2^{14})	Constant (2^{13})
Stability of Communication	Unstable	Stable

Example
(Ultimate Case)

MAX_KEY_VALUE = 4
BUCKET_WIDTH = 4 / # of cores

2 cores

1	0	3	2
---	---	---	---

No communication

4 cores

1	0	3	2
---	---	---	---



Unstable (depends on # of cores)

MAX_KEY_VALUE = 4 - 8
BUCKET_WIDTH = 2

2 cores

1	0	3	2
---	---	---	---

No communication

4 cores

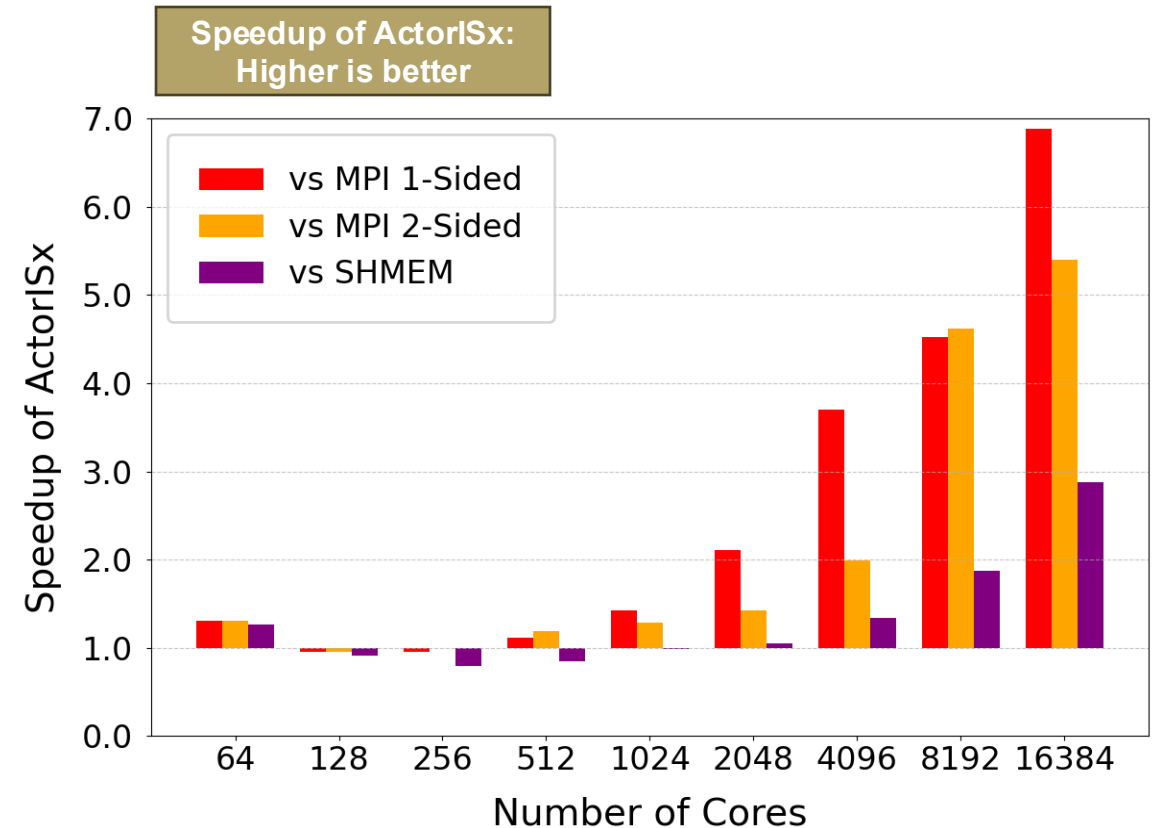
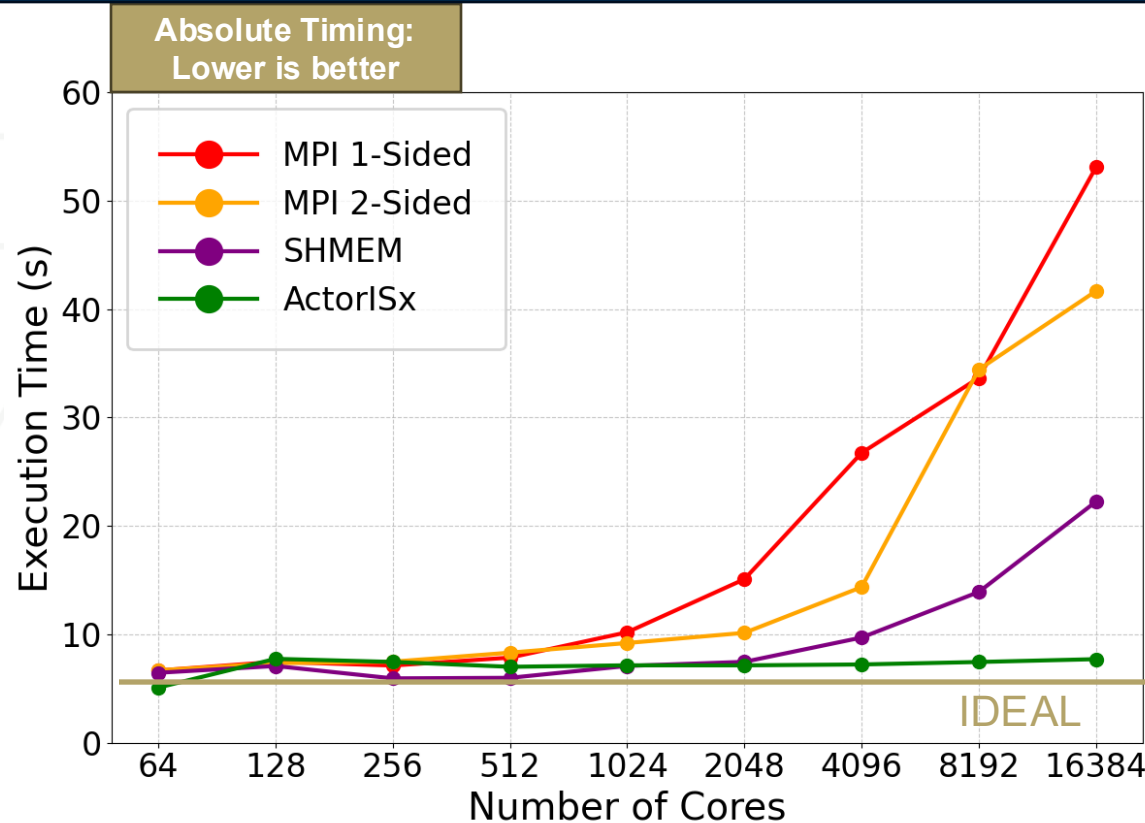
1	0	3	2	4	5	6	7
---	---	---	---	---	---	---	---

No communication

Stable (independent of core count)

Weak Scaling Results on Perlmutter (up to 16k cores): CONSTANT max key (2^{28}): Unstable communication

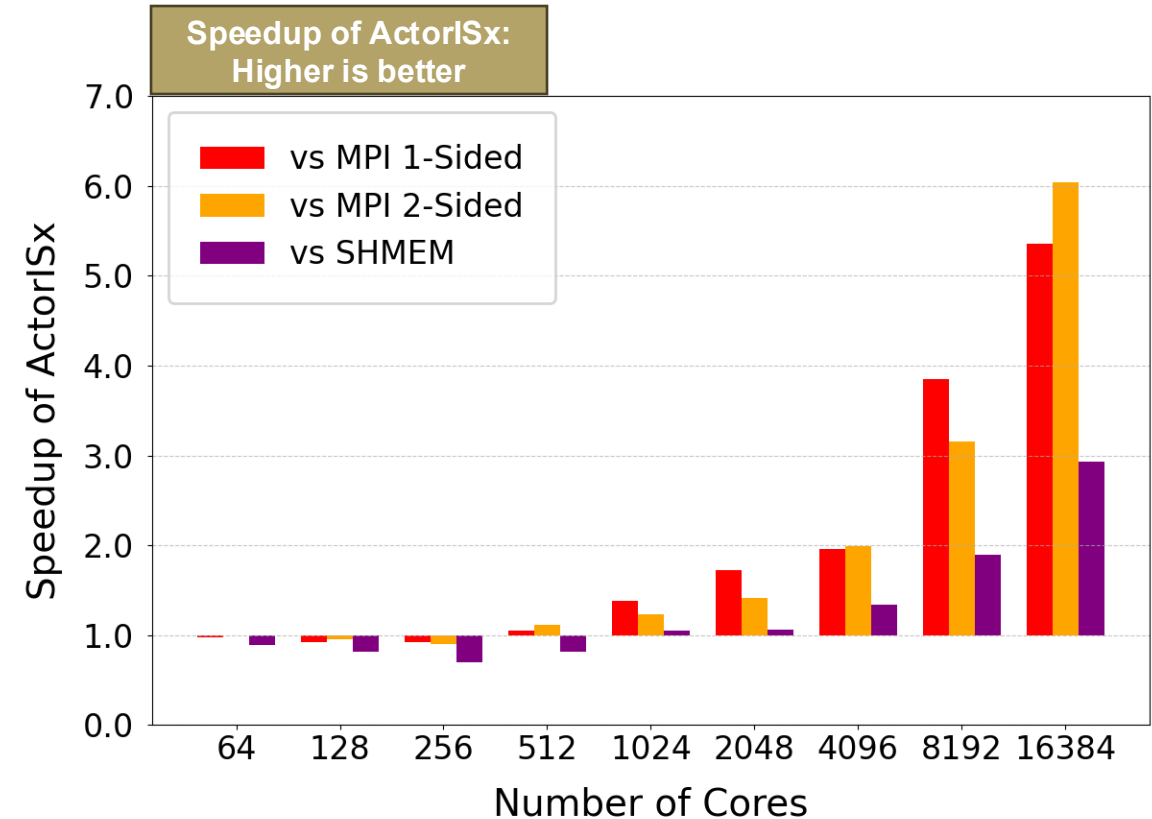
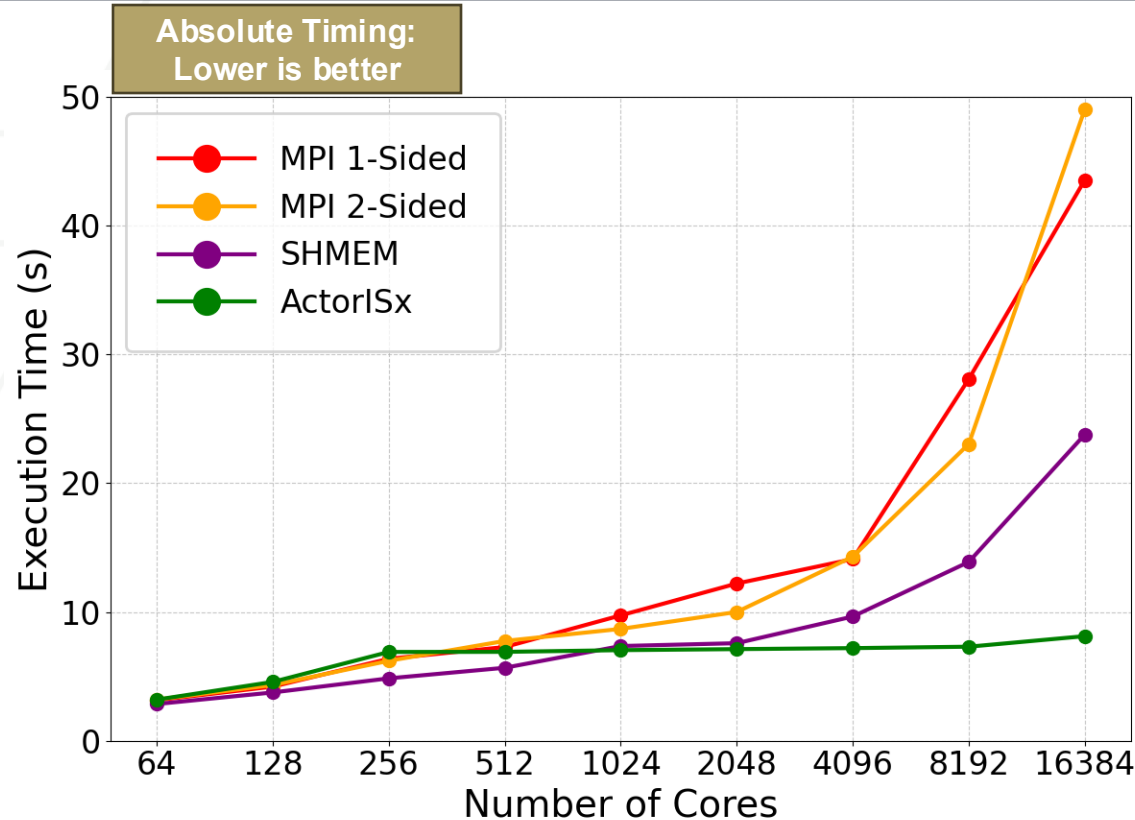
2^{27} keys per core \Rightarrow 2^{33} keys (64 cores), 2^{41} keys (16k cores)



- **ActorISx** achieves almost an ideal weak-scaling result thanks to
 - Asynchronous communication
 - Two-level message aggregation
 - Multi-hop routing

ISO Weak Scaling Results on Perlmutter (up to 16k cores): INCREASE max key ($2^{19} - 2^{27}$): Stable communication

2^{27} keys per core $\Rightarrow 2^{33}$ keys (64 cores), 2^{41} keys (16k cores)



- **ActorISx** achieves almost an ideal weak-scaling result thanks to
 - Asynchronous communication
 - Two-level message aggregation
 - Multi-hop routing

Impact of the Solution

- Our algorithm has shown efficient **scalability** and **performance**
- The extensibility of this algorithm has four major impacts:

1

By leveraging the actor model, we have addressed the fundamental limitations of NPB IS while maintaining high performance and scalability

2

Our algorithm can be applied to other data-intensive and communication-intensive applications

3

Our actor-based approach is a viable alternative to traditional MPI-based and SHMEM-based approaches

4

Our solution shows that actor-based approaches will play an increasingly important role in future HPC systems

DEMO

ActorISx: Exploiting Asynchrony for Scalable High-Performance Integer Sort

You can try this at home... Just visit hclib-actor.com !

The screenshot shows a web browser with the URL `hclib-actor.com/background/bsp/`. The page title is "Bulk Synchronous Parallel". The left sidebar contains a navigation menu with sections: "Hclib-Actor Documentation", "Home", "Background", "Theory" (with a sub-link "Bulk Synchronous Parallel"), "Partitioned Global Address Space", "Actor Model", "Practice", "OpenSHMEM", "Bale", "Summary", "spmat", "libgetput", "Habanero-C Library (HCLib)", "Getting Started", "Containers" (with sub-links "Docker" and "Singularity"), "Clusters/Supercomputers" (with sub-link "NERSC/ORNL/PACE"), and "Writing Hclib-Actor Programs". The main content area has the heading "Bulk Synchronous Parallel" and the sub-heading "What is the bulk synchronous parallel model?". Below this, it states: "The Bulk Synchronous Parallel (BSP) model is one of the most popular parallel computation models." and "The model consists of:" followed by a bulleted list: "• A set of processor-memory pairs.", "• A communication network that delivers messages in a point-to-point manner.", and "• Efficient barrier synchronization for all or a subset of the processes." To the right of the main content is a "Table of contents" with links: "What is the bulk synchronous parallel model?", "Single Program Multiple Data (SPMD) Programming", and "Further Readings". At the bottom of the page is a diagram illustrating the BSP model. It shows a horizontal row of blue vertical bars representing "Virtual Processors" labeled PE_0, PE_1, PE_2, \dots . A bracket on the left labeled "SUPERSTEP" spans the height of these bars. Below the bars, a horizontal orange bar represents "Barrier Synchronization". Arrows labeled "Inter-processor Communications" point from the bottom of each processor bar to the orange barrier bar. To the right of the diagram, the text "Local Computation" is positioned above the processors, and "Barrier Synchronization" is positioned below the orange bar.

ACKNOWLEDGEMENT

This research is based upon work supported by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), through the Advanced Graphical Intelligence Logical Computing Environment (AGILE) research program, under Army Research Office (ARO) contract number W911NF22C0083. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the ODNI, IARPA, or the U.S. Government.

IEEE/ACM CCGRID 2025

The 25th IEEE/ACM International Symposium on Cluster, Cloud and Internet Computing
Tromsø, Norway

Actor!Sx: Exploiting Asynchrony for Scalable High-Performance Integer Sort

Youssef Elmougy, Shubhendra Singhal, **Akihiro Hayashi**, and Vivek Sarkar

Habanero Extreme Scale Software Research Lab
Georgia Institute of Technology

Thank you for your attention!