

Implementation and Analysis of Nonblocking Collective Operations on SCI Networks

Christian Kaiser

Torsten Hoefler



Boris Bierbaum, Thomas Bemmerl

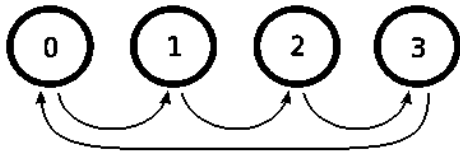


LEHRSTUHL FÜR BETRIEBSSYSTEME

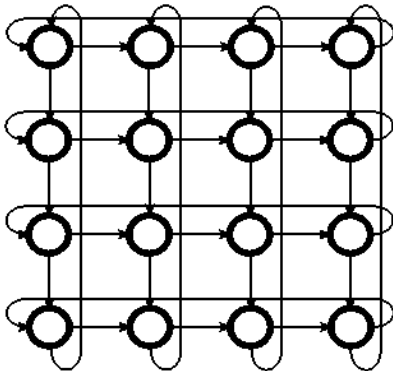
Univ.-Prof. Dr. habil. Thomas Bemmerl

RWTHAACHEN
UNIVERSITY

Ringlet:



2D Torus:



- **IEEE Std 1596-1992**
- **Memory Coupled Clusters**
- **Data Transfer: PIO and DMA**
- **SISCI User-Level Interface**
- **16 x Intel Pentium D, 2.8 GHz**
- **SCI: D352 (IB: Mellanox DDR x4)**

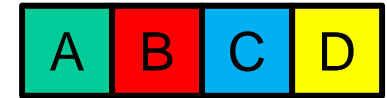
source

destination

Process 0



Process 1 (root)



Process 2



Process 3



source

destination

Process 0



Process 1 (root)



Process 2



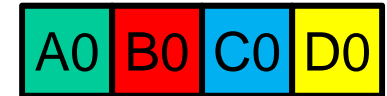
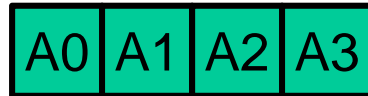
Process 3



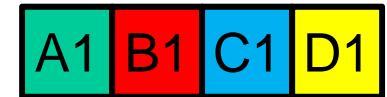
source

destination

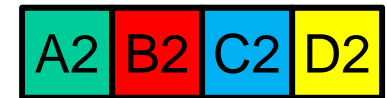
Process 0



Process 1 (root)



Process 2



Process 3



source

destination

Process 0



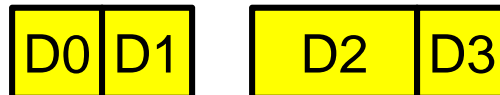
Process 1 (root)

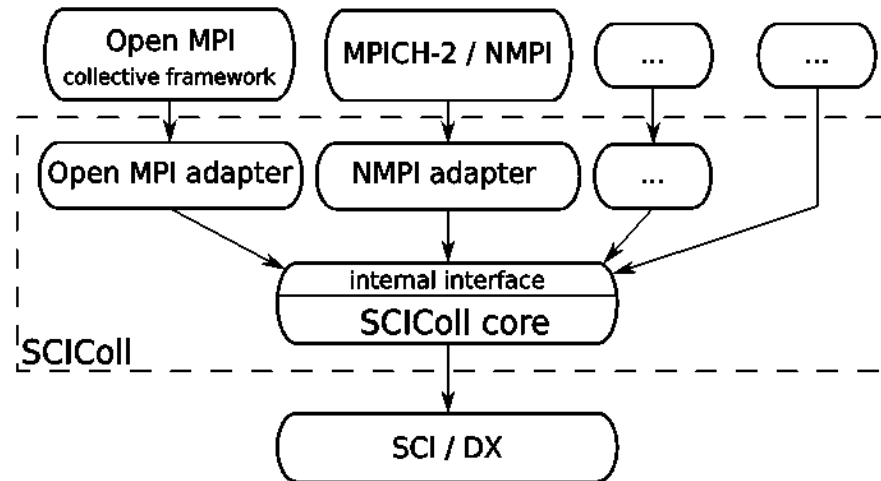


Process 2



Process 3

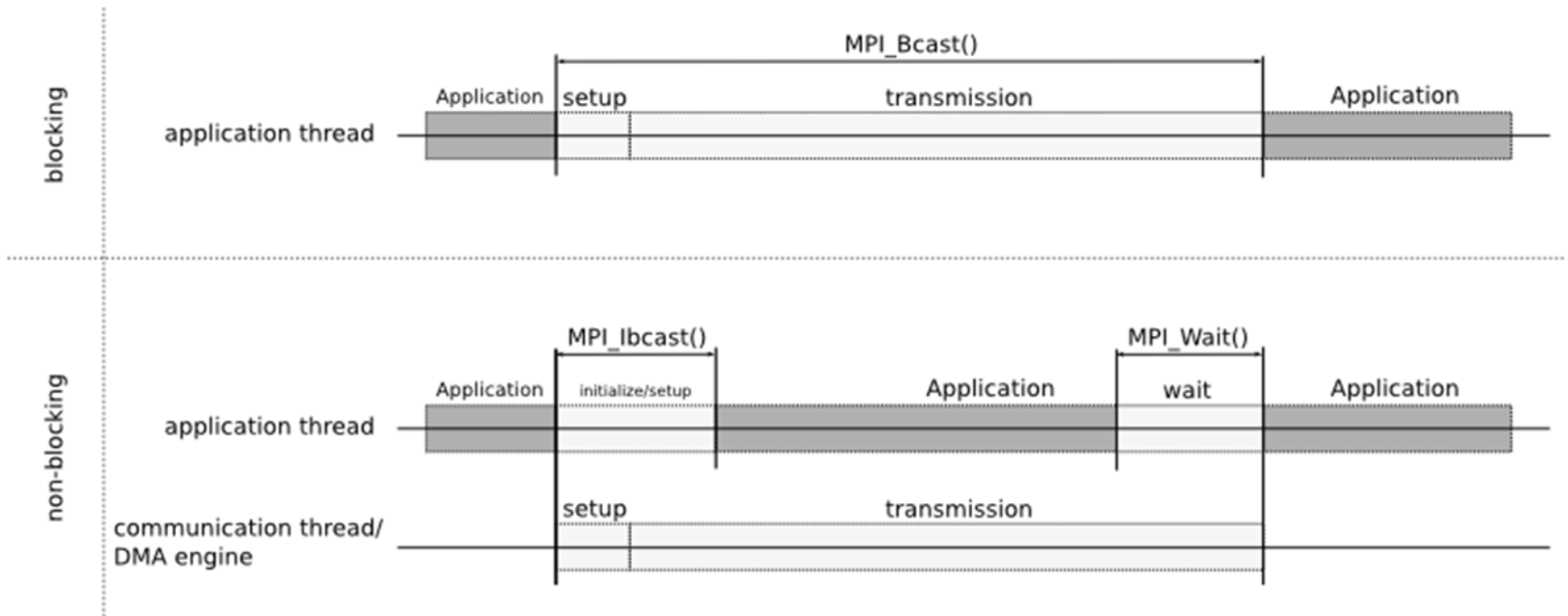




Purpose:

- Study collective communication algorithms for SCI clusters
- Support multiple MPI libraries: Open MPI, NMPI
- Support arbitrary communication libraries: LibNBC

Purpose: Overlap of Computation and Communication



MPI-2.0 JoD: Split Collectives

MPI_BCAST_BEGIN(buffer, count, datatype, root, comm)

MPI_BCAST_END(buffer, comm)

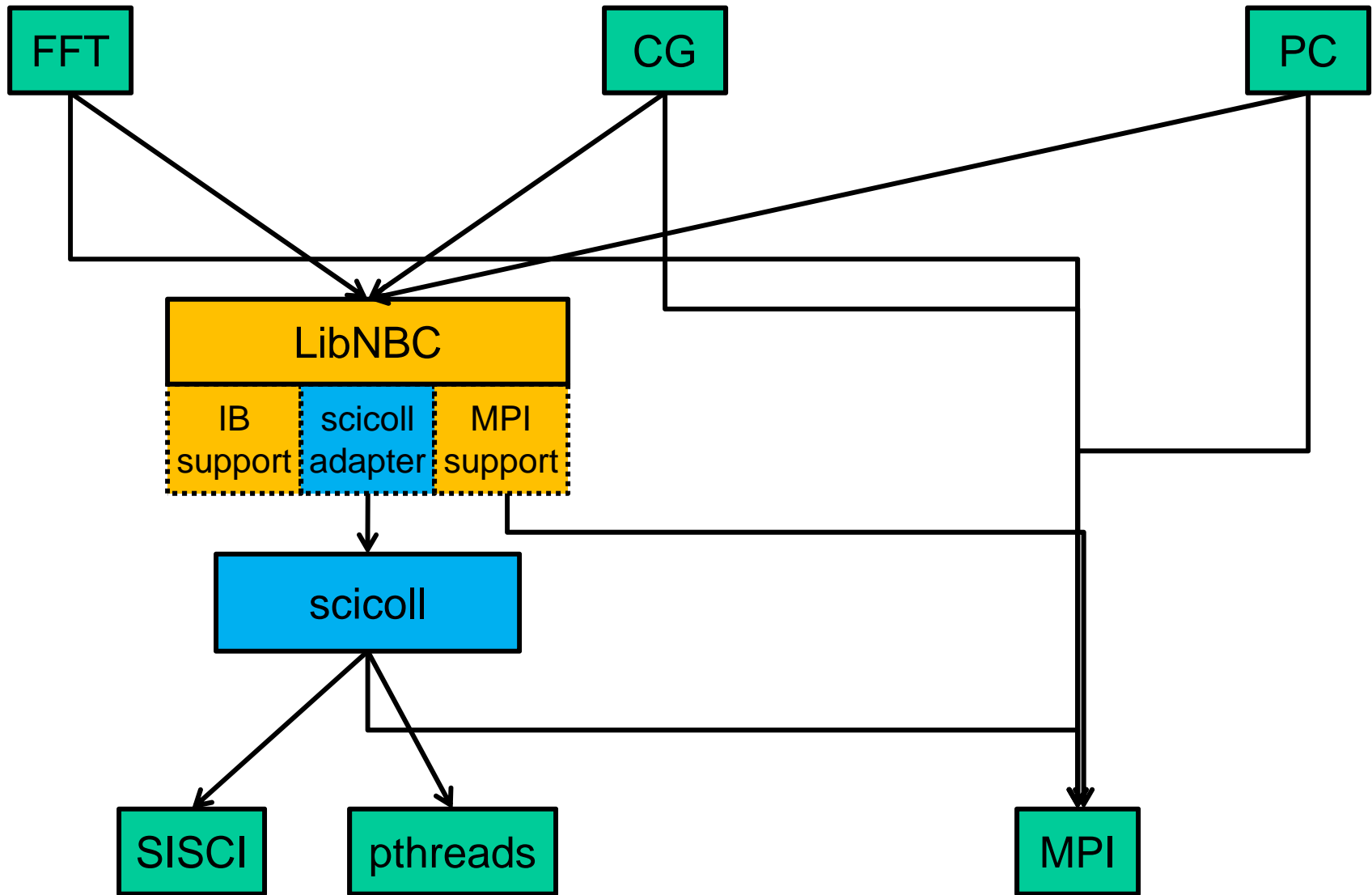
MPI-2.1:

- Implement with nonblocking Point-to-Point operations
- Blocking collectives in separate thread

MPI-3 Draft:

MPI_IBCAST(buffer, count, datatype, root, comm, request)

MPI_WAIT(request, status)



So far:

- Promising results with NBC via LibNBC
- Research done on InfiniBand clusters

Therefore:

What about a very different network architecture?

Implementation considerations:

- Use algorithms different from blocking version?
- PIO vs DMA
- Use background thread?

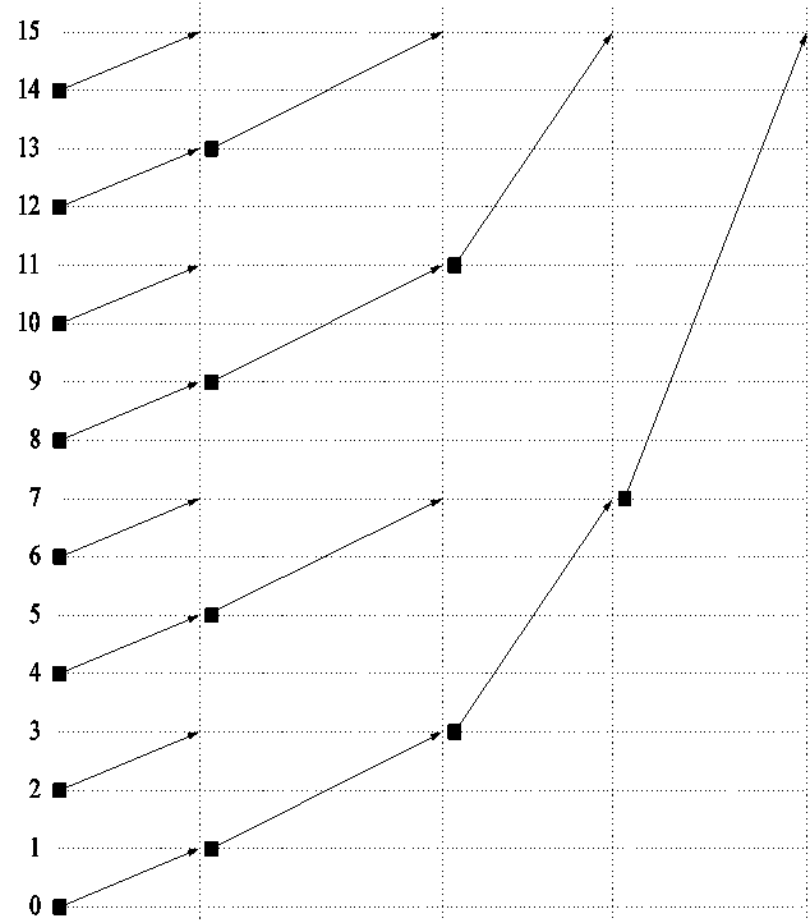
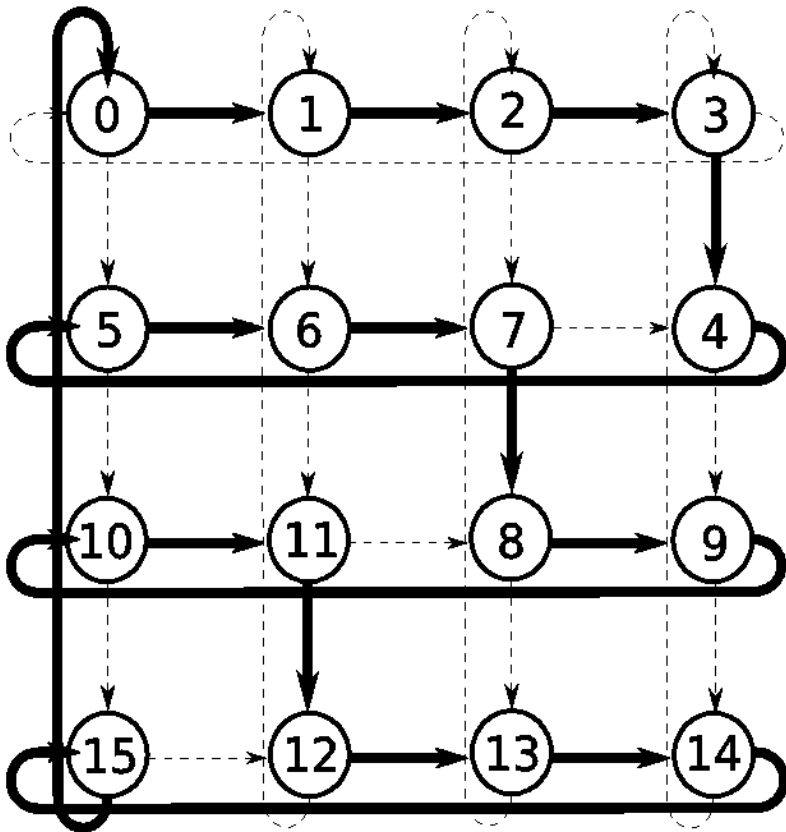
Synthetic:

NBCBench: measures the communication overhead / overlap potential

Application Kernels:

- **CG (Alltoallv)**: 3D Grid, overlaps computation with halo zone exchange
- **PC (Gatherv)**: overlaps compression with gathering of previous results
- **FFT (Alltoall)**: parallel matrix transpose, overlaps data exchange for z transpose with computation for x and y

- **Underlying concept:** Hamiltonian Path in a 2D torus
- **Algorithms:** Binary Tree, Binomial Tree, Flat Tree, Sequential Transmission



Gather(v):

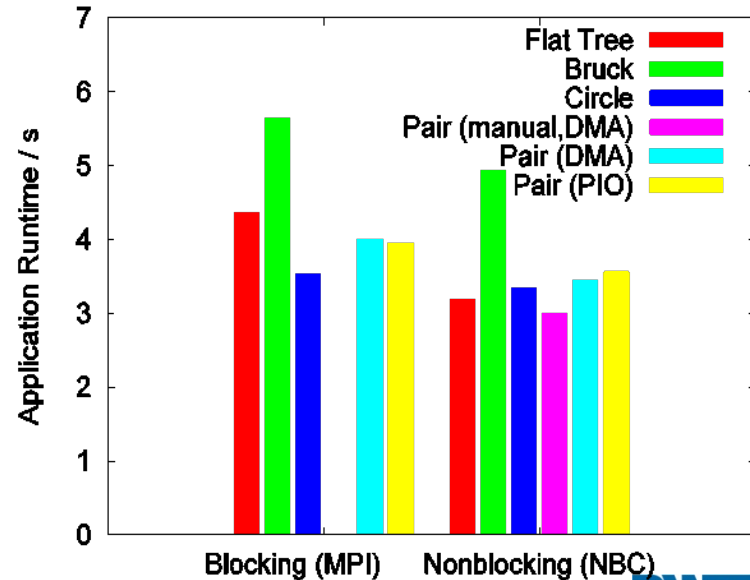
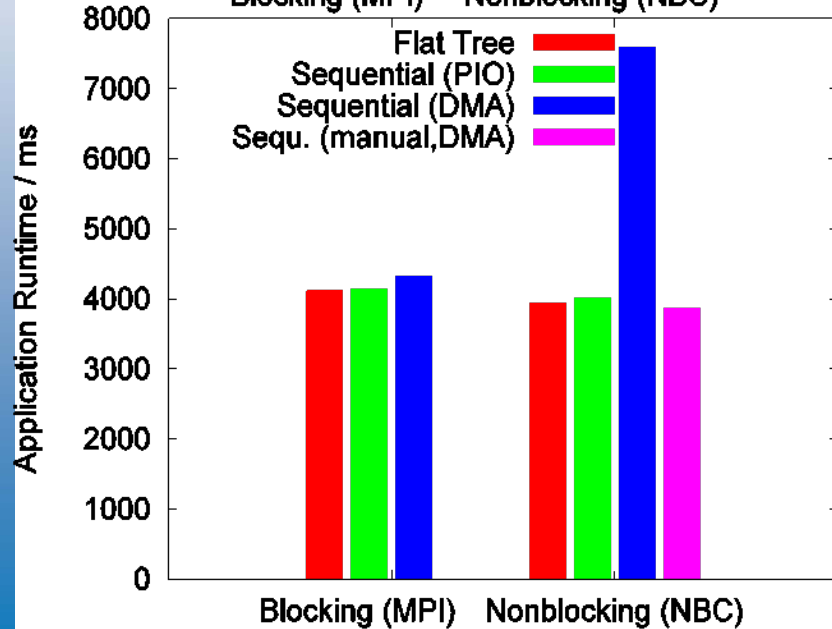
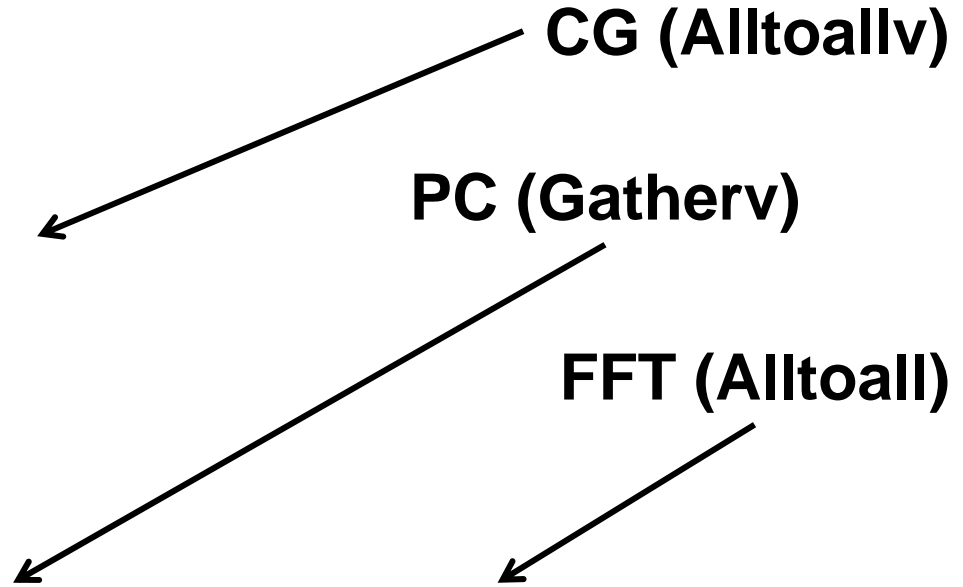
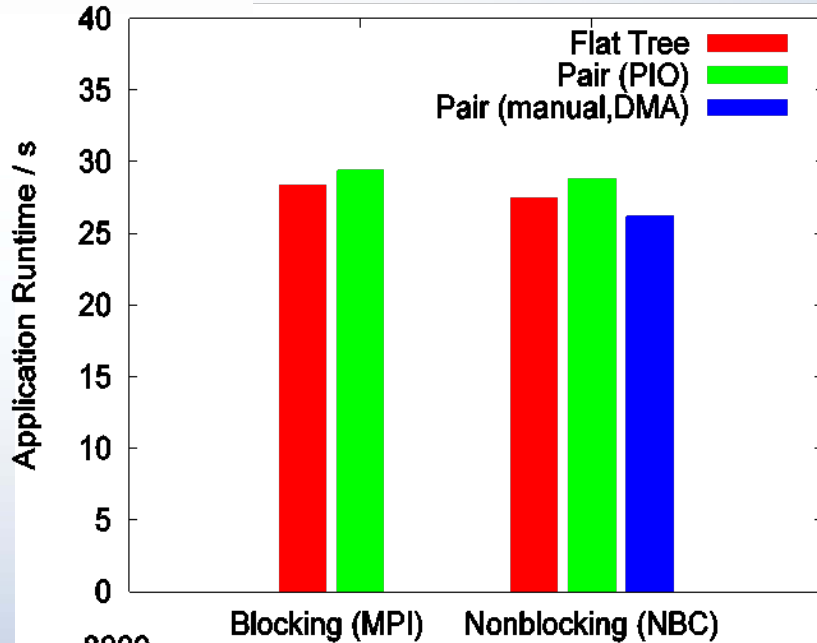
- Additional progress thread: Binary Tree (PIO), Binomial Tree (PIO), Flat Tree (PIO), Sequential Transmission (PIO, DMA)
- Single Thread with manual progress: Sequential Transmission
- Vector Variant: Flat Tree and Sequential Transmission



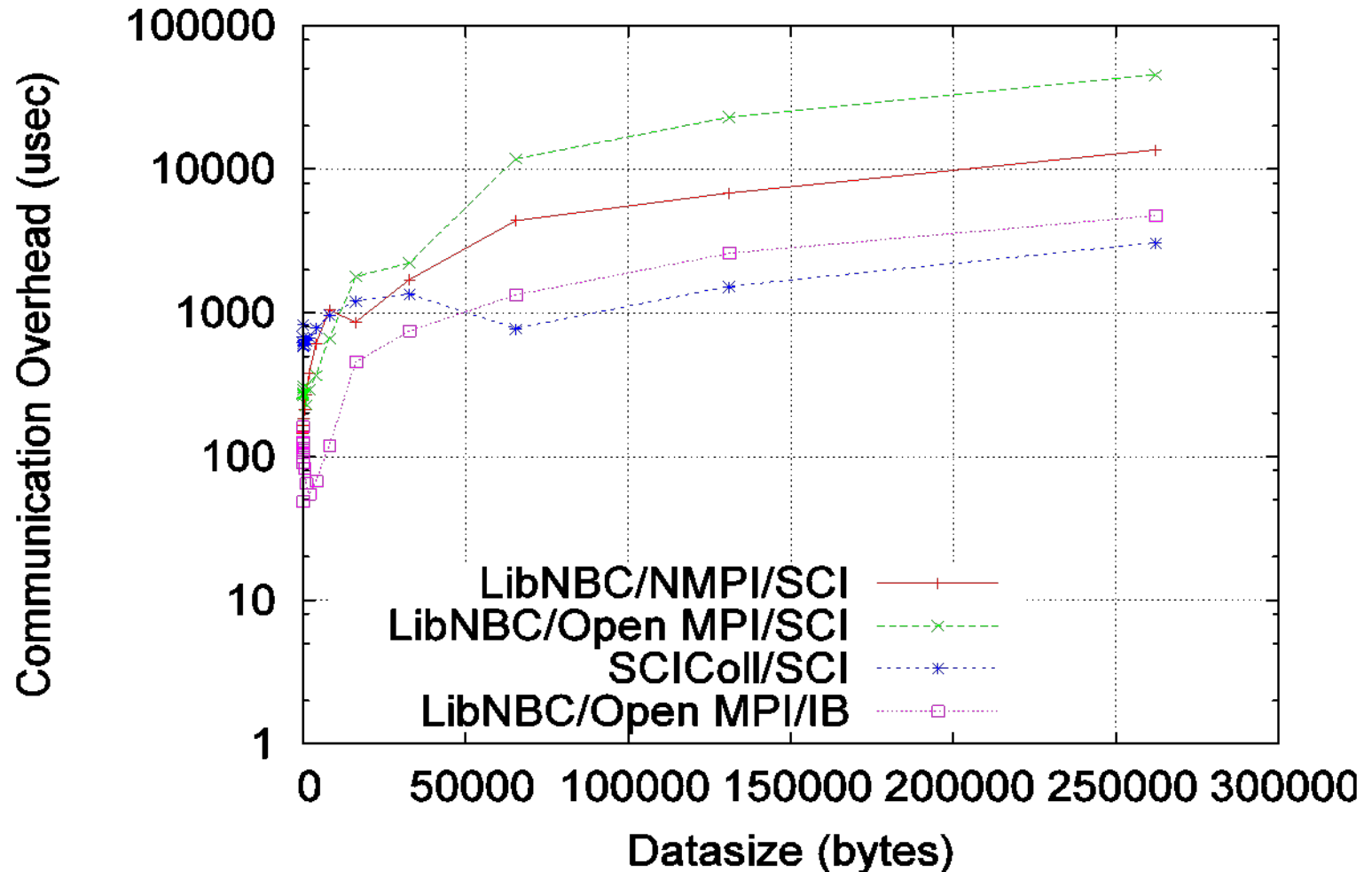
Alltoall(v):

- Additional progress thread: Bruck (PIO), Pairwise Exchange (PIO), Ring (PIO), Flat Tree (PIO)
- Single Thread with manual progress: Pairwise Exchange (DMA)
- Vector Variant: Pairwise Exchange, Flat Tree

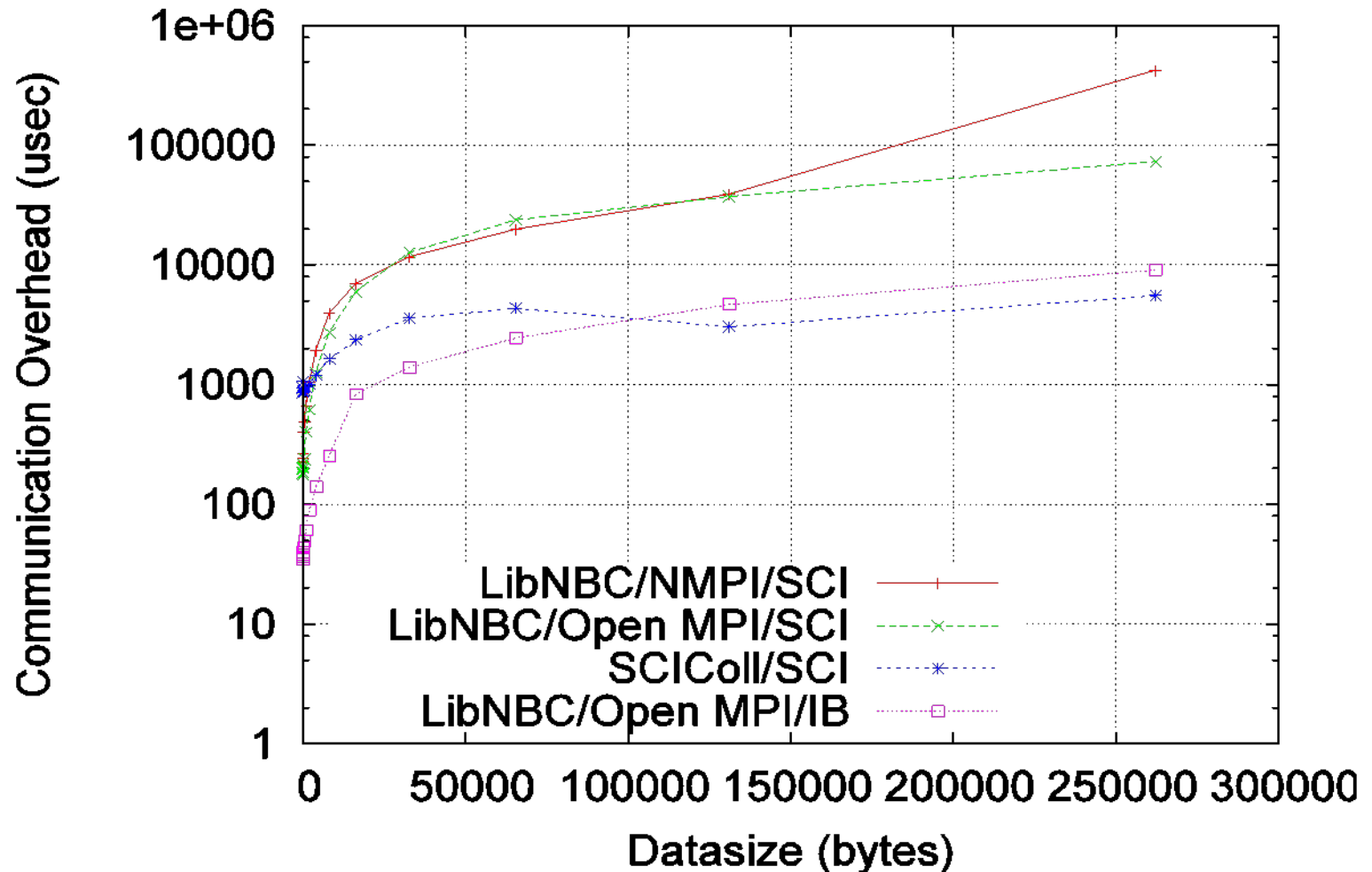
Application Kernels: Algorithms

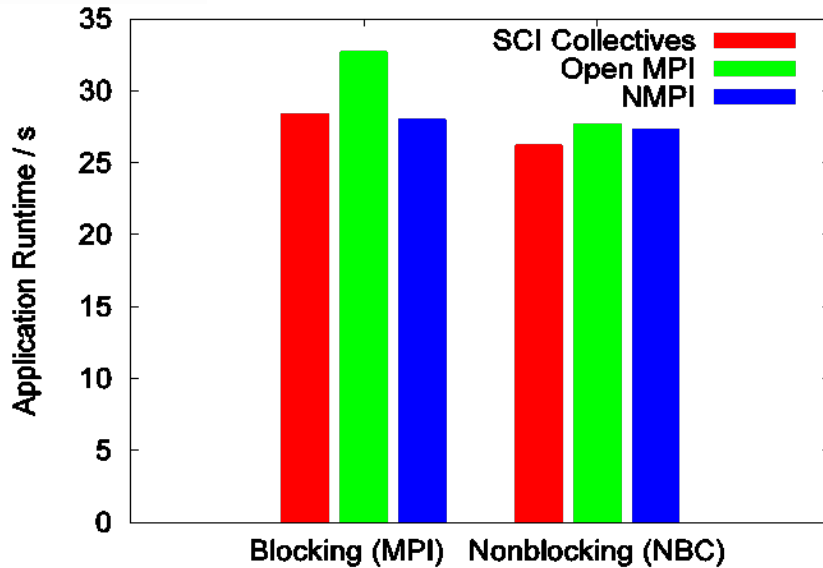


Gather



Alltoall

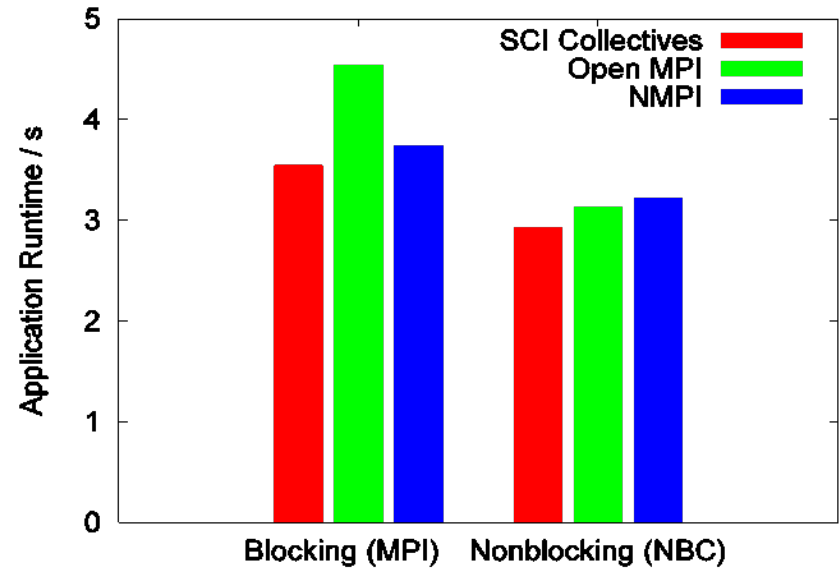
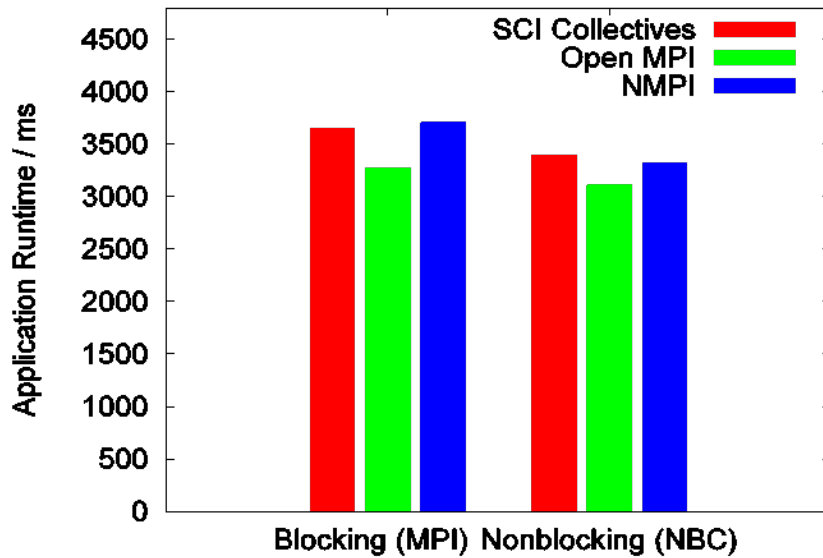




CG (Alltoallv)

PC (Gatherv)

FFT (Alltoall)



What we've done:

Implement nonblocking Gather(v) and Alltoall(v) collective operations on SCI clusters with different algorithms and implementation alternatives

What we found out:

- Applications can benefit from nonblocking collectives on SCI clusters in spite of inferior DMA performance
- Best implementation method: DMA in a single thread, PIO is usually used for blocking collectives
- Issues with multiple threads

**Thank you for
your attention!**