

GML – Global Matrix Library Overview

by

Juemin Zhang(Jeff) @ IBM

Computational Function

- Cell-wise operations
 - Add, subtract, multiply, division
 - Sum, Euclidean distance, etc

- Matrix multiply

- $C = A^T * B^T + C$

Optional

<code>C.mult(A, B, addOn);</code>
<code>C.transMult(A, B, addOn);</code>
<code>C.multTrans(A, B, addOn);</code>
<code>C.transMultTans(A, B, addOn);</code>

- Result is stored in invoking object, which is also the return object
 - Chained operations

Matrix Type

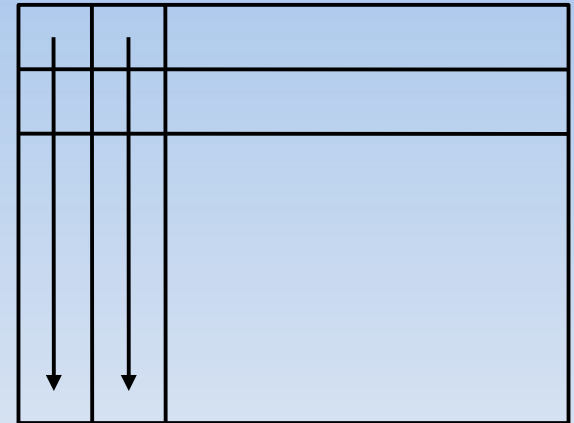
- Dense matrix
- Sparse matrix
- Block matrix
- Distributed dense matrix
- Distributed sparse matrix
- Duplicated dense matrix
- Duplicated sparse matrix

Matrix methods

- Constructor
 - Matrix dimensions, storage, block partitioning, block distribution, sparsity
- Make
 - Memory allocation
- Clone, and CopyTo
- Initialization
- Matrix element access `this(x,y)` and set
- Matrix subset: `column(s)` and `row(s)`
- `toString`, `print`

Dense Matrix

- Memory layout
 - Column-major memory storage
- BLAS driver (double precision)
 - Level 1, 2, and 3
 - Dense-dense operation
- X10 driver
 - Dense-dense
 - Dense-sparse
- Result is stored in dense



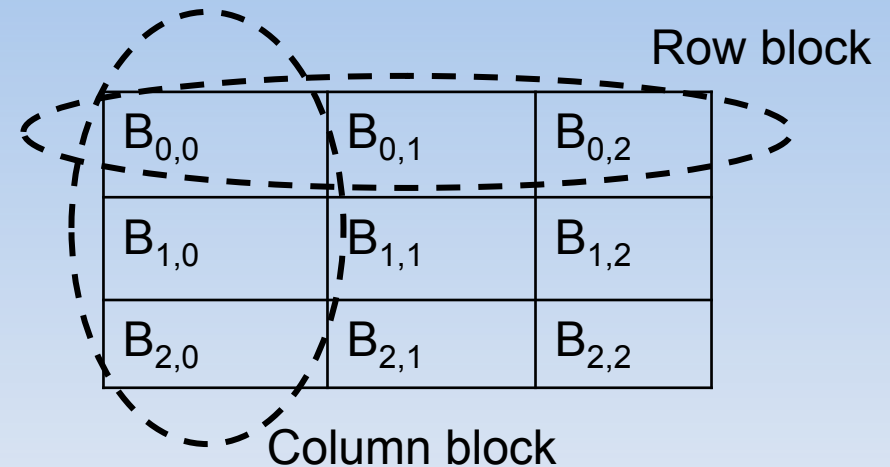
0	for (var c:Int=0; c<B.N; c++) {
1	for (var k:Int=0; k<A.N; k++) {
2	val v2 = B(k, c);
3	for (var r:Int=0; r<A.M; r++) {
4	val v1= A(r, k);
5	C.d(startcol+r) = v1 * v2;
6	}
7	startcol += C.M;
8	}
	}

Sparse Matrix

- Compressed sparse column
 - *Matrix elements, each of which contains index and nonzero value pair, are read first by column*
- 2D compressed array
 - 1D compressed array
- Compress array
 - Index array
 - Value array
- Not good for inserting or modification
 - *Sparse matrix is not used to store matrix computation result*

Matrix Blocks and Block Matrix

- Matrix partitioning - Grid



- Matrix blocks

- Row block ID and column block ID
- Dense/Sparse matrix

- Block matrix

- Number of row blocks, number of column blocks
- List of rows for all row blocks, and list of columns for all column blocks
- List of dense/sparse blocks

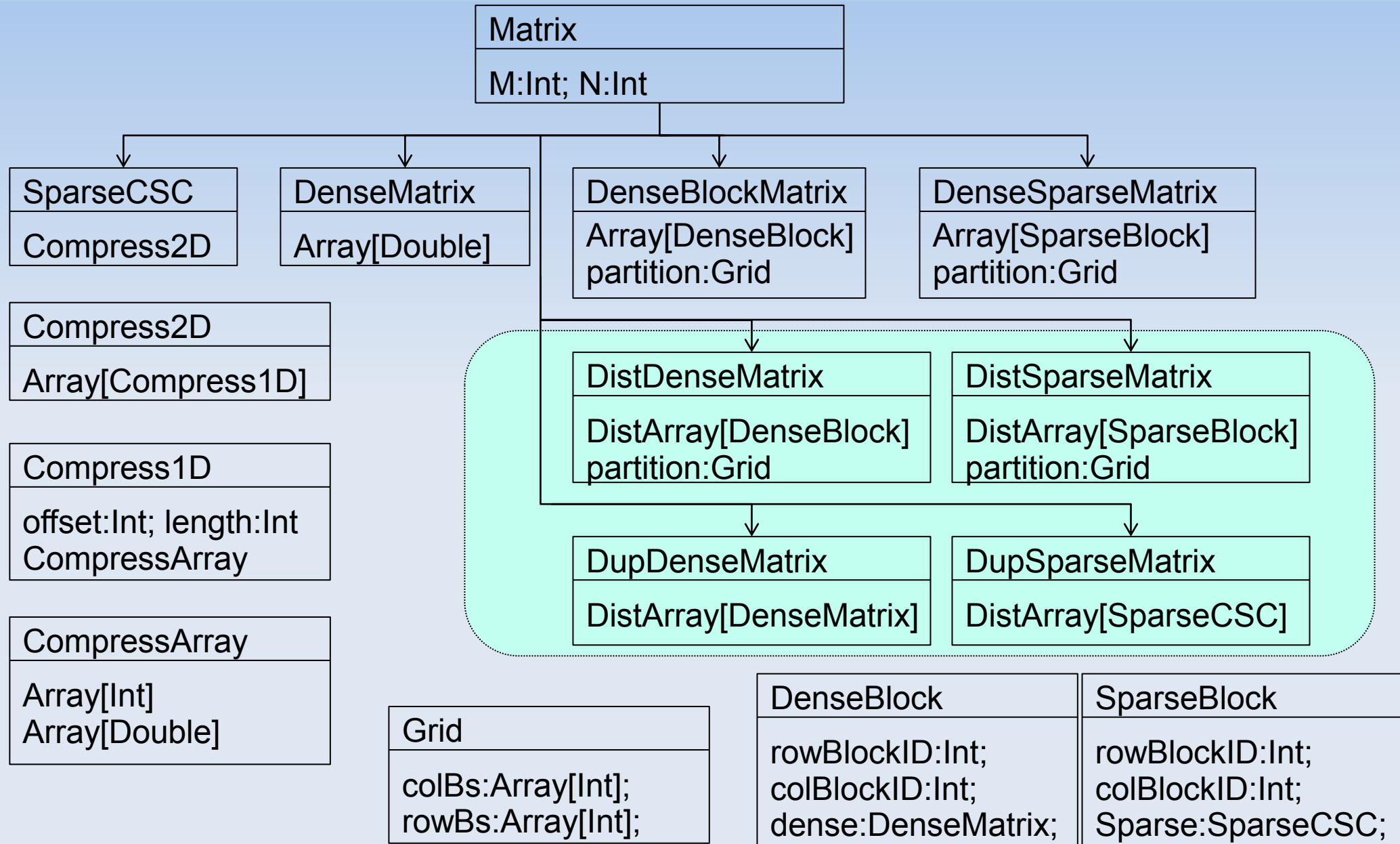
Distributed and Duplicated Matrix

- Distributed dense/sparse matrix
 - Matrix blocks unique distributed on DistArray

Place ₀ B _{0,0}	Place ₃ B _{0,1}	Place ₆ B _{0,2}
Place ₁ B _{1,0}	Place ₄ B _{1,1}	Place ₇ B _{1,2}
Place ₂ B _{2,0}	Place ₅ B _{2,1}	Place ₈ B _{2,2}

- Duplicated dense/sparse matrix
 - Replicated matrices are uniquely distributed at all places on DistArray

Matrix Class Structure



Access data from different place

- X10 implicit data capture
 - At (..) { }
 - Expensive, but easy to use

- Explicit data copy

- Remote array copy
 - LocalPlaceHandle, array on DistArray and dense/sparse on DistArray
 - Portable to different transports: sockets, lapi and pami
- MPI wrapped calls (only for native backend on MPI transport)
 - Array on DistArray and dense/sparse on DistArray
 - Optimized collective operations
- Fast, but difficult in programming

Data declare	val data:.....;

Remote execution block	at (...) {
	...;
Remote data capture	... data ...;
	}

Inter-place data communication

- Place-to-place communication
- Collective communication
 - Bcast:
 - Synchronize copies of duplicated matrix
 - Scatter:
 - Copy data from block matrix to distributed one
 - Copy data from single-column/row matrix to distributed one
 - Copy data from single row block matrix to distributed one
 - Gather
 - Reverse of scatter
 - Reduce
 - Sum of all blocks

Build GML

- Setup “system_setting.mk”
 - X10 compiler path
 - MPI/g++
 - BLAS path
- Build GML library for native/managed backend
 - Native backend: make native
 - Native backend + MPI transport: make native_mpi
 - Managed backend: make managed
- Build application
 - Add: -classpath [gml_path]/[**managed|native|native_mpi**].gml.jar
 - Add: -x10lib [gml_path]/[**managed|native native_mpi**].gml.properties

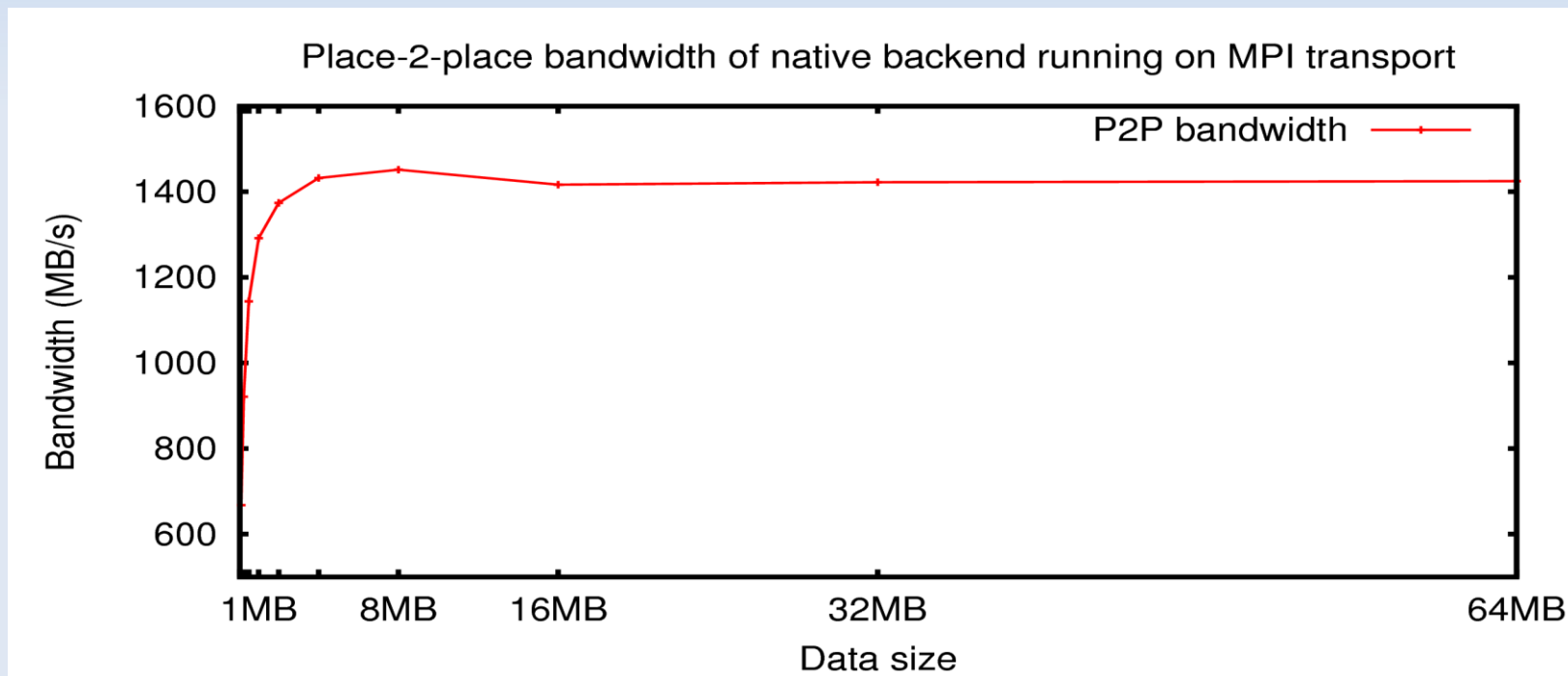
Benchmark platform - Triloka

- 10-node cluster
 - Node: 2 Quad-Core AMD Opteron(tm) 2356
- Memory: 16-GB each node
- Network
 - Infiniband (MPI): 20 Gb/sec (4X DDR)
 - Ethernet (socket): 1Gb/sec

Communication Benchmark

Place-to-Place

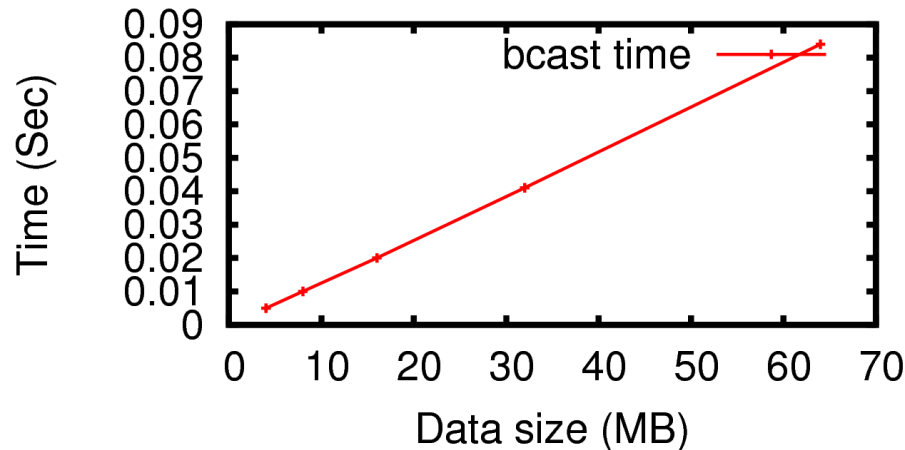
- Native backend running on MPI transport (InfiniBand 20Gb/s)
- X10 Peak performance: 1.45GB/s
 - Intel MPI benchmark peak performance 1.48GB/s



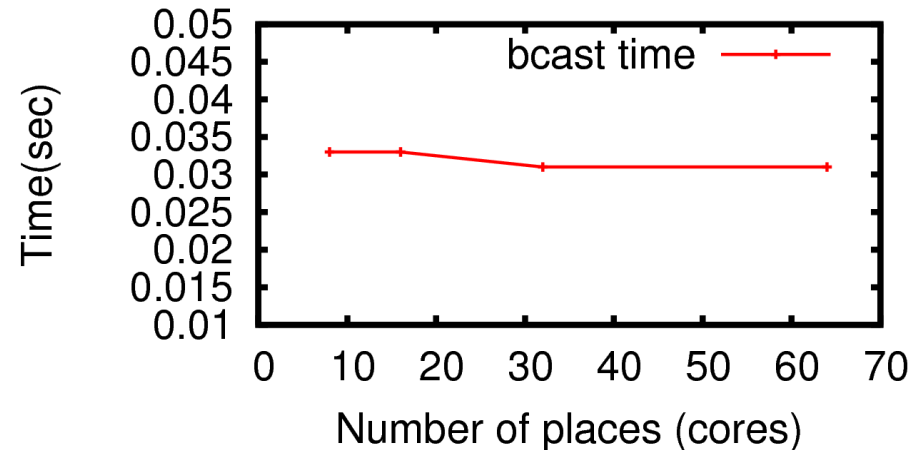
Collective Communication Benchmarks: bcast, gather

Bcast benchmark of native backend running on MPI-transport

Performance on 8-node 8-place

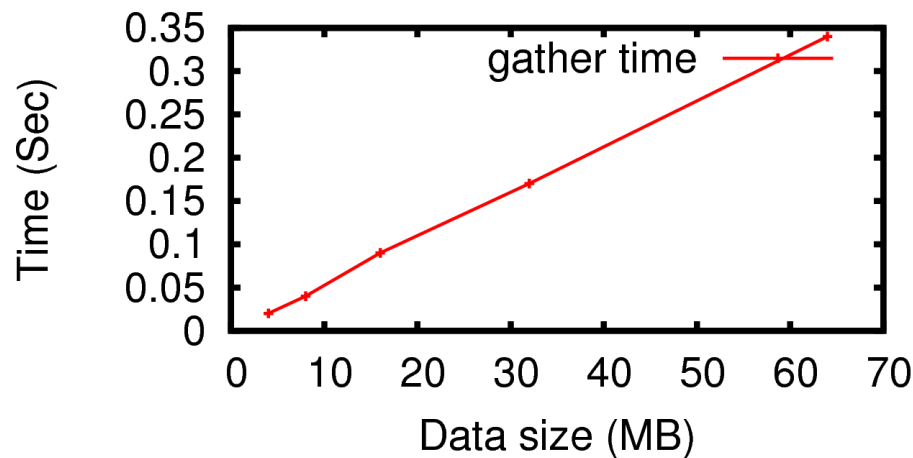


Bcast 8-MB data package

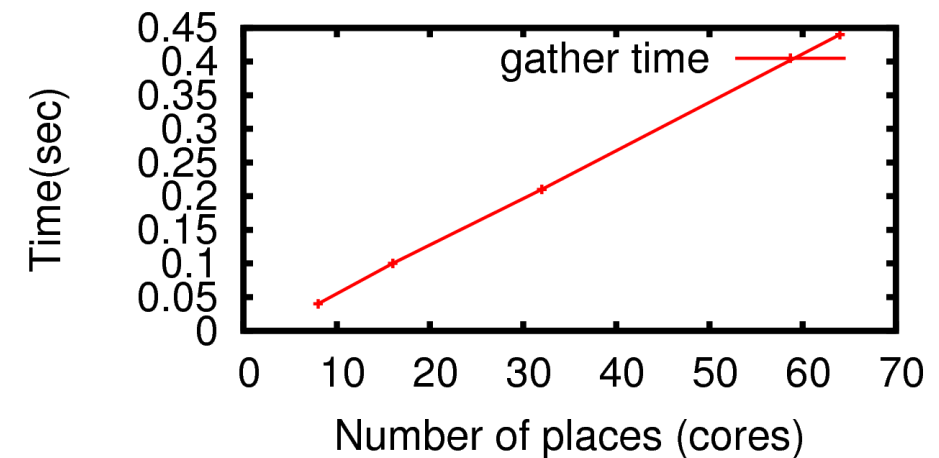


Gather benchmark of native backend running on MPI transport

Performance on 8-node 8-place



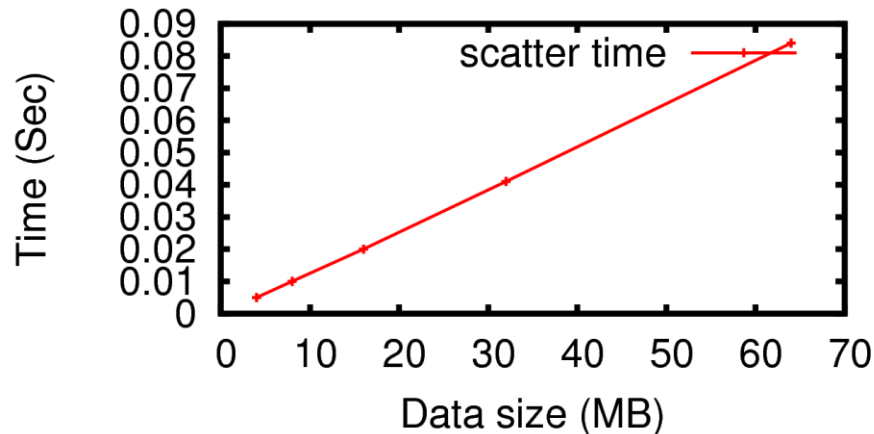
Gather 8-MB data packages



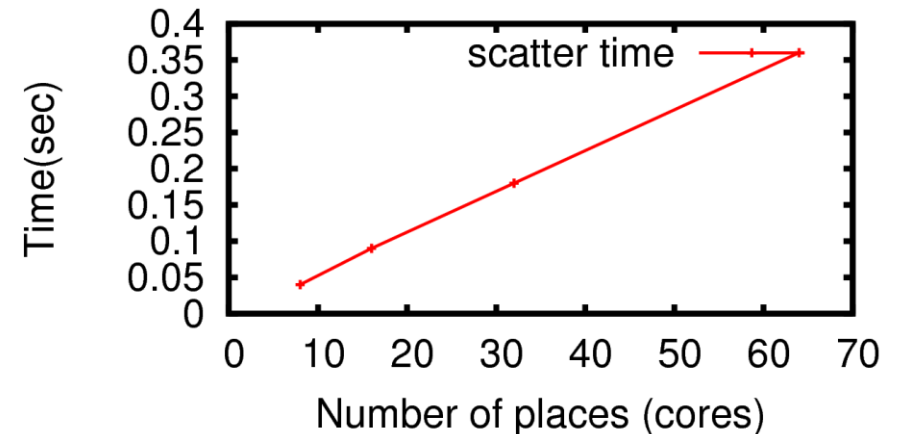
Collective Communication Benchmarks: scatter, reduce

Scatter benchmark of native backend running on MPI transport

Scatter performance on 8-node 8-place

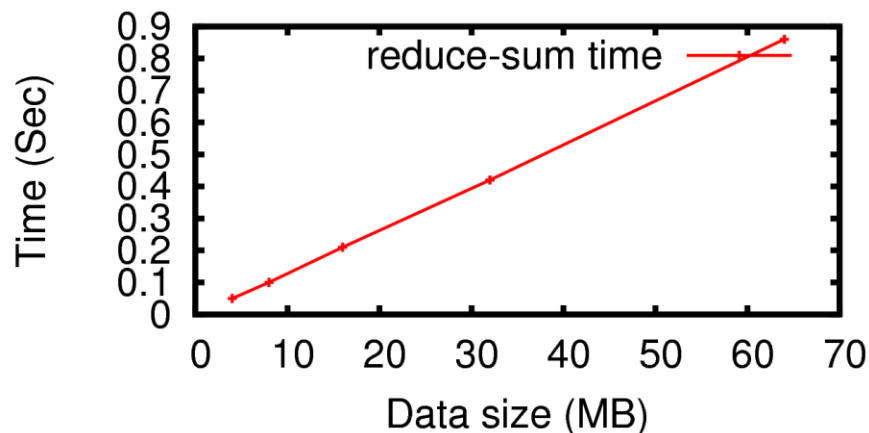


Scatter 8-MB data packages

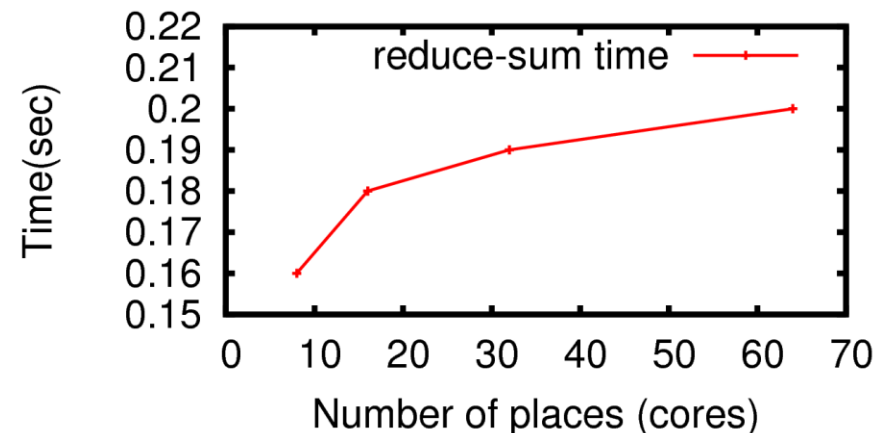


Reduce-Sum benchmark of native backend running on MPI transport

Reduce-sum performance on 8-node 8-place

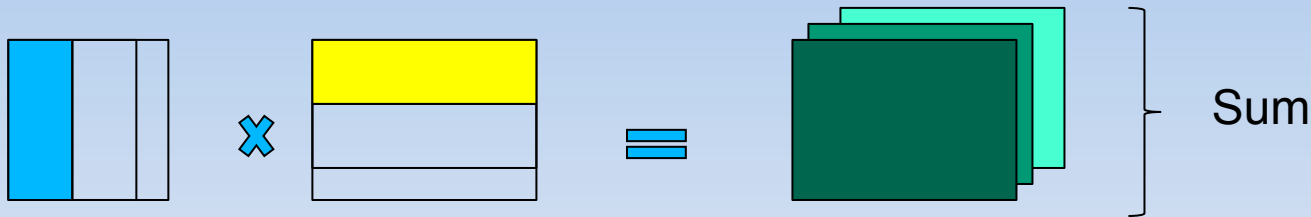


Reduce sum of 8-MB data packages



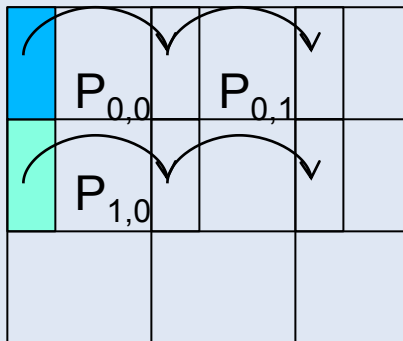
SUMMA $A \times B$

- Local (in place) matrix * matrix

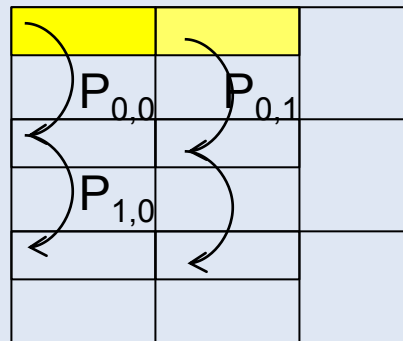


- Matrices are partitioned in the same way

Input A
row block ring-cast



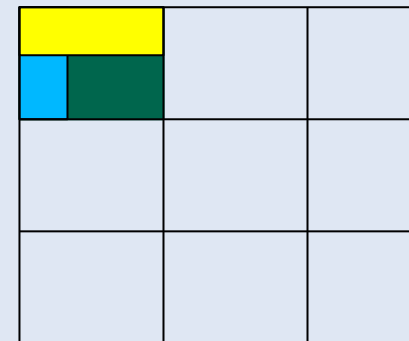
Input B
column block ring-cast



\times

$=$

Result



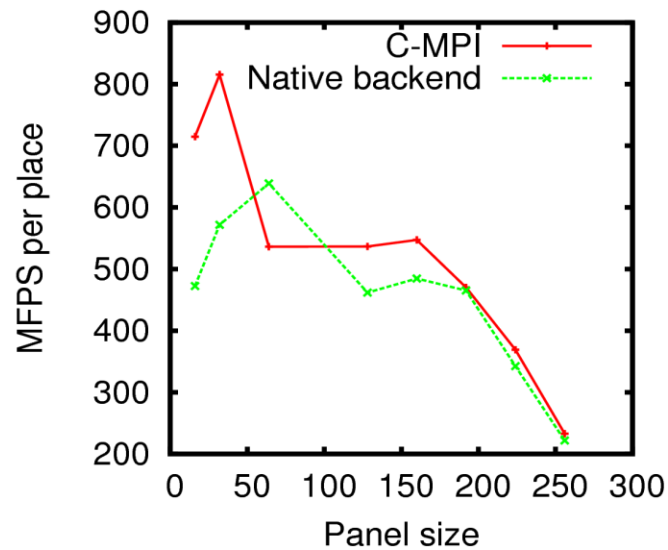
Test Settings

- DistDenseMatrix SUMMA
 - Computation: BLAS library
 - Comparing:
 - C-MPI SUMMA implementation
 - X10 SUMMA implementation native backend running on MPI transport
- DistSparseMatrix SUMMA
 - Computation: X10 sparse*sparse driver

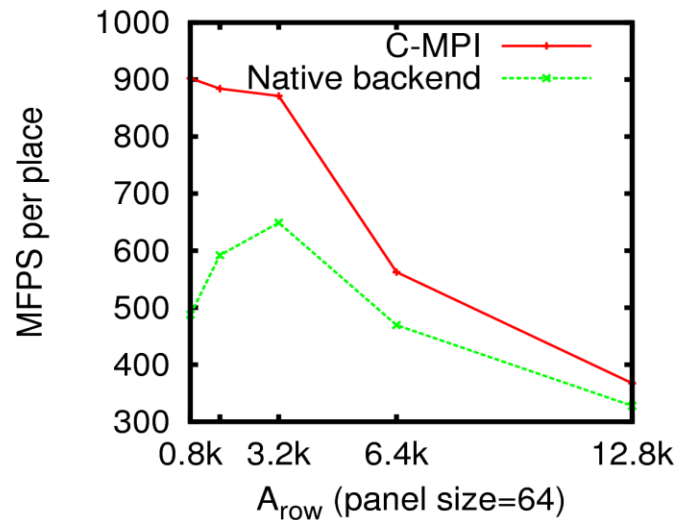
DistDenseMatrix SUMMA

Distributed dense * dense (SUMMA)

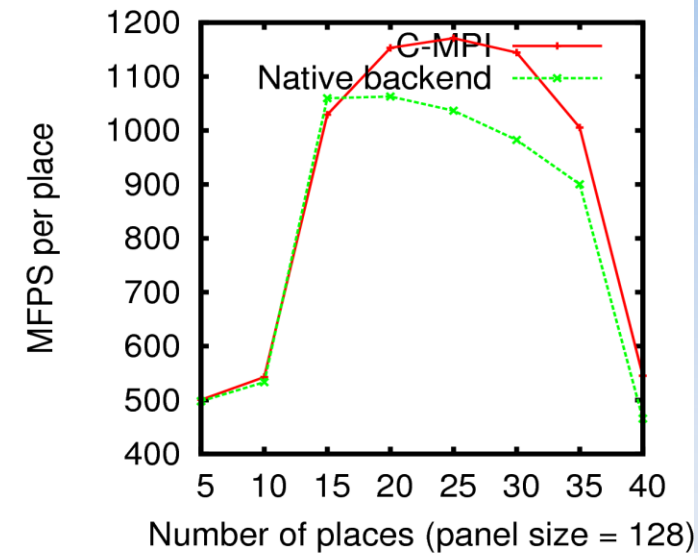
$A(4k,2k) \times B(2k,2k)$ on 8-node 40-place



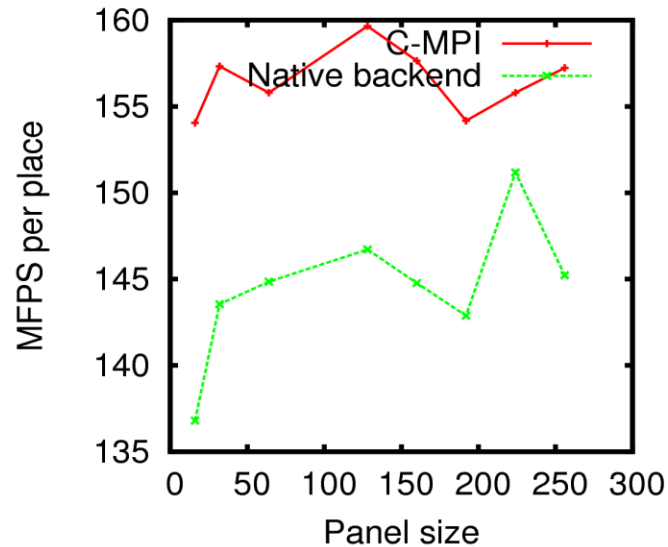
$A(A_{\text{row}},2k) \times B(2k,2k)$



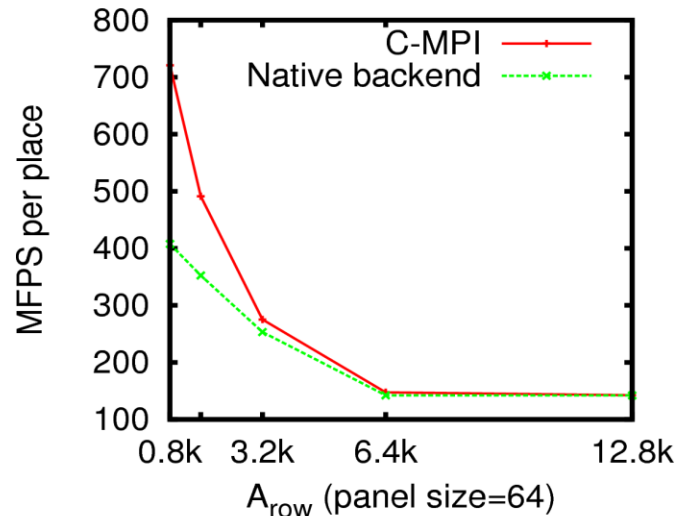
$A(3.2k,2K) \times B(2K,2K)$ on 5 nodes



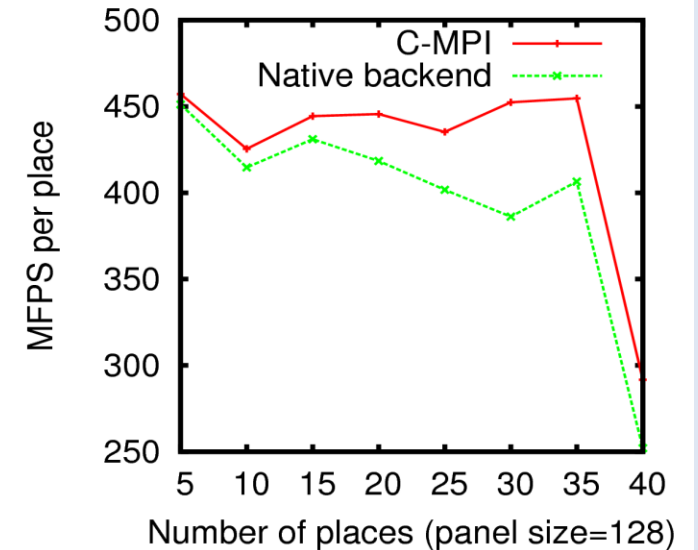
$A(4k,2k) \times B(2k,2k)^T$ on 8-node 40-place



$A(A_{\text{row}},2k) \times B(2k,2k)^T$ on 8-node 40-place



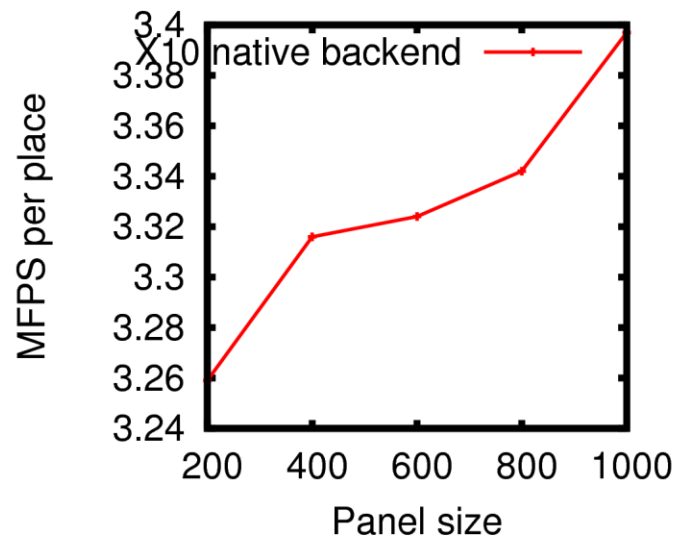
$A(3.2k,2K) \times B(2K,2K)^T$ on 5 nodes



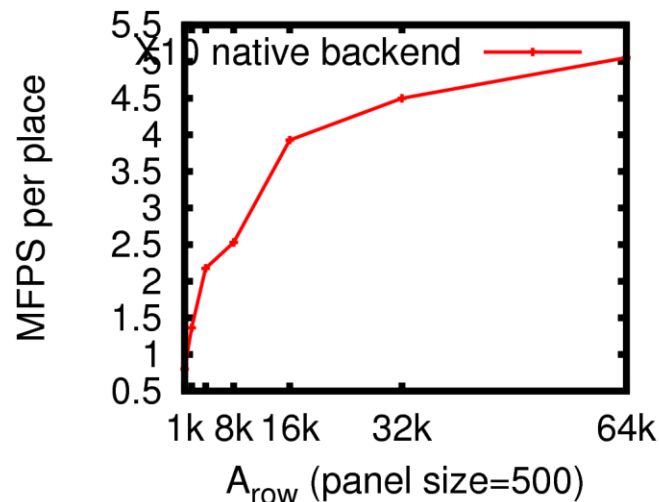
SUMMA (DistSparseMatrix)

Sparse * sparse (density 0.01) SUMMA benchmark of native backend on MPI transport

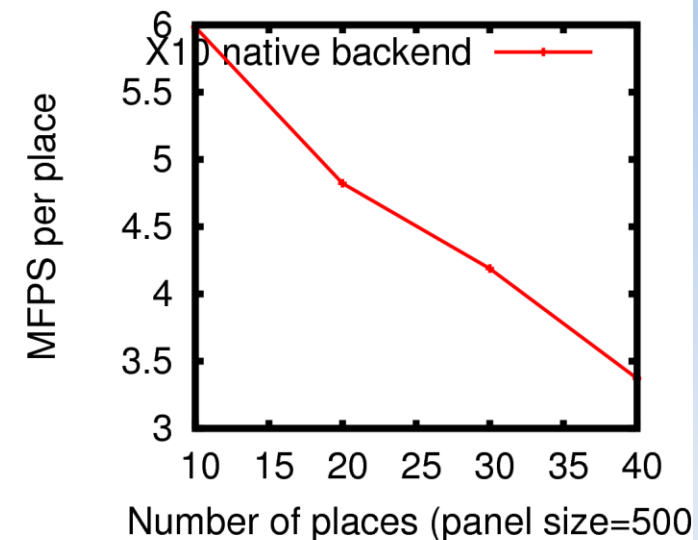
A(10k,10k) x B(10k,10k) on 8-n 40-pl



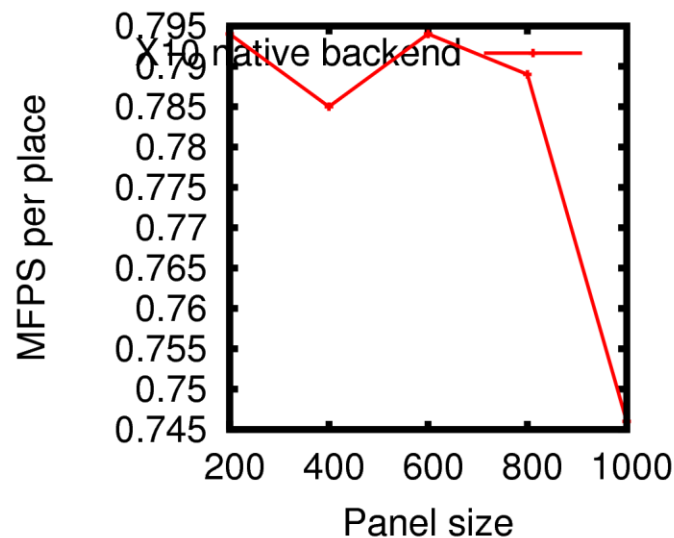
A(A_{row},10k) x B(10k,10k) 8-n 40-pl



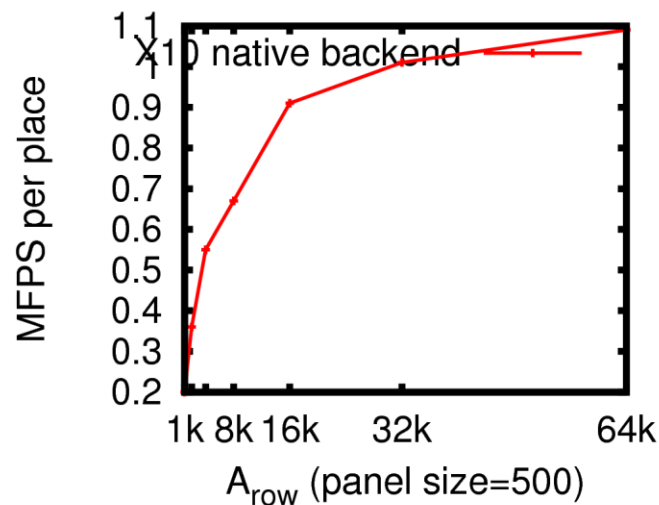
A(10k,10k) x B(10k,10k) on 5 node



A(10k,10k) x B(10k,10k)^T on 8-n 40-pl



A(A_{row},10k) x B(10k,10,k)^T on 8-n 40-pl



A(10k,10k) x B(10k,10k)^T on 5 node

