

# MPJava: High-Performance Message Passing in Java using Java.nio

Bill Pugh  
Jaime Spacco

University of Maryland, College Park

# Message Passing Interface (MPI)

- MPI standardized how we pass data on a cluster
- MPI:
  - Single Processor Multiple Data (SPMD)
  - Provides point-to-point as well as collective communications
  - Is a set of library routines
  - Is an interface with several free and commercial implementations available
  - source code is portable
  - Has C, Fortran and C++ bindings, but not Java

# Previous Java + MPI work:

- mpiJava (Carpenter)
  - Native wrappers to C libraries
  - Worse performance than native MPI
- jmpj
  - Pure-Java implementation of proposed standard for Java/MPI bindings
  - Also bad performance compared to native MPI

# MPJava

- Pure-Java Message Passing framework
- Makes extensive use of java.nio
  - *select()* mechanism
  - direct buffers
  - efficient conversions between primitive types
- Provides point-to-point and collective communications similar to MPI
- We experiment with different broadcast algorithms
- We try to use threads
  - More work needed to get this right
- Performance is pretty good

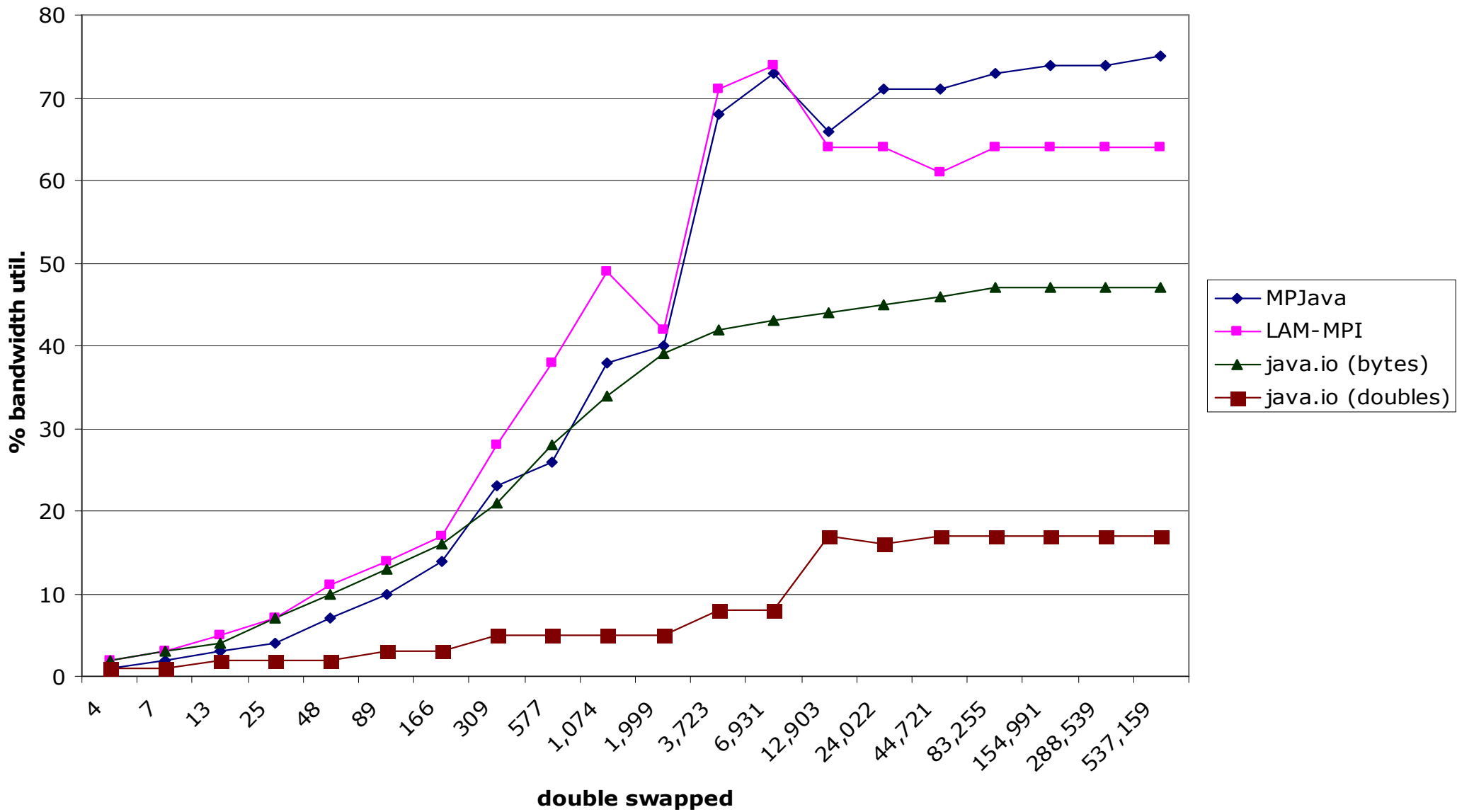
# Benchmarks

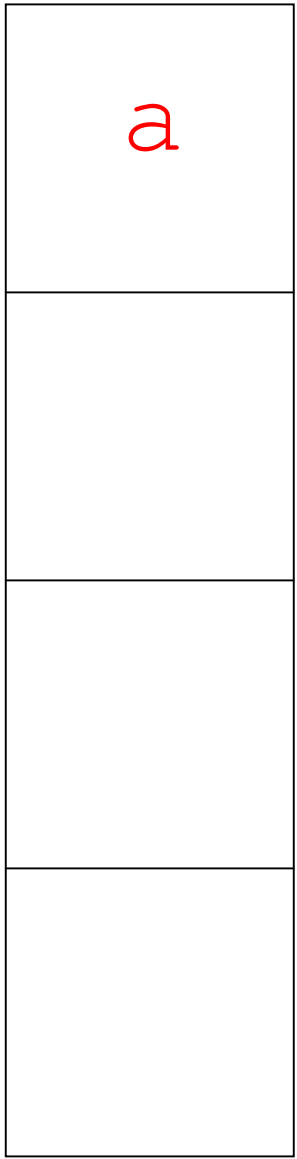
- PingPong
- Alltoall
- NAS Parallel Benchmarks Conjugate Gradient

# Results

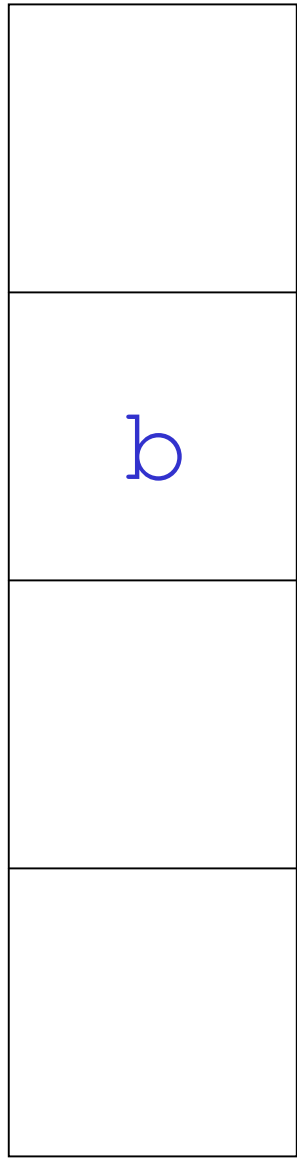
- 50 650 MHz PIII machines
- 768 MB memory
- RedHat 7.3
- two 100 Mbps channel-bonded NICs
- Fortran compiler: g77 v2.96
  - tried a commercial compiler (pgf90) but no difference for these benchmarks
- LAM-MPI 6.5.8
- JDK 1.4.2-b04

# PingPong

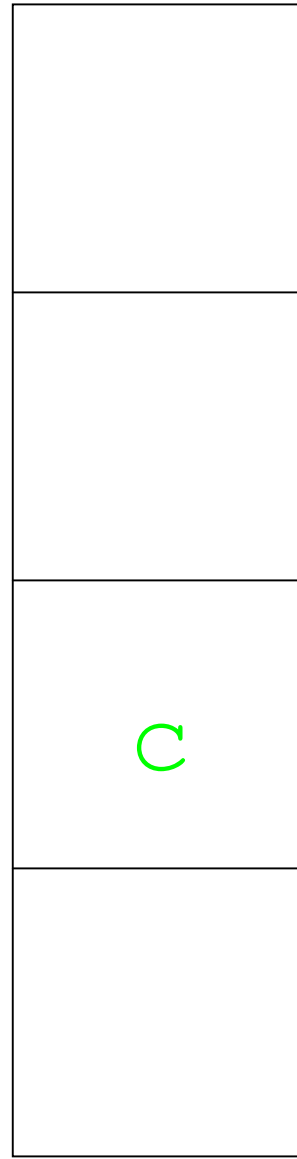




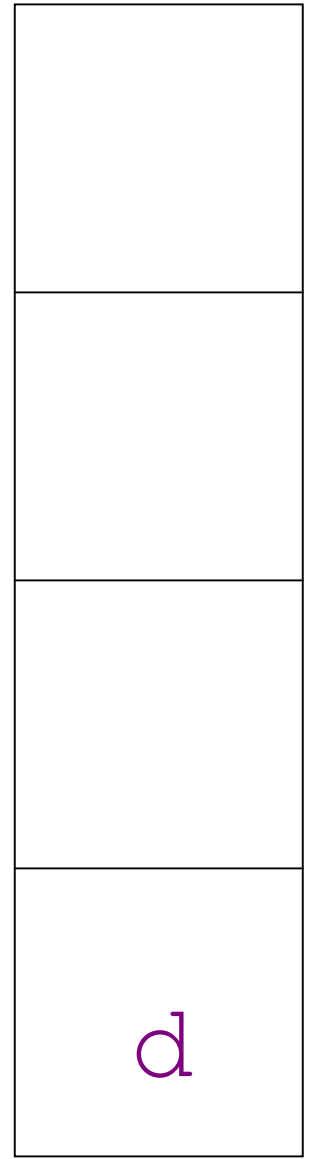
0



1

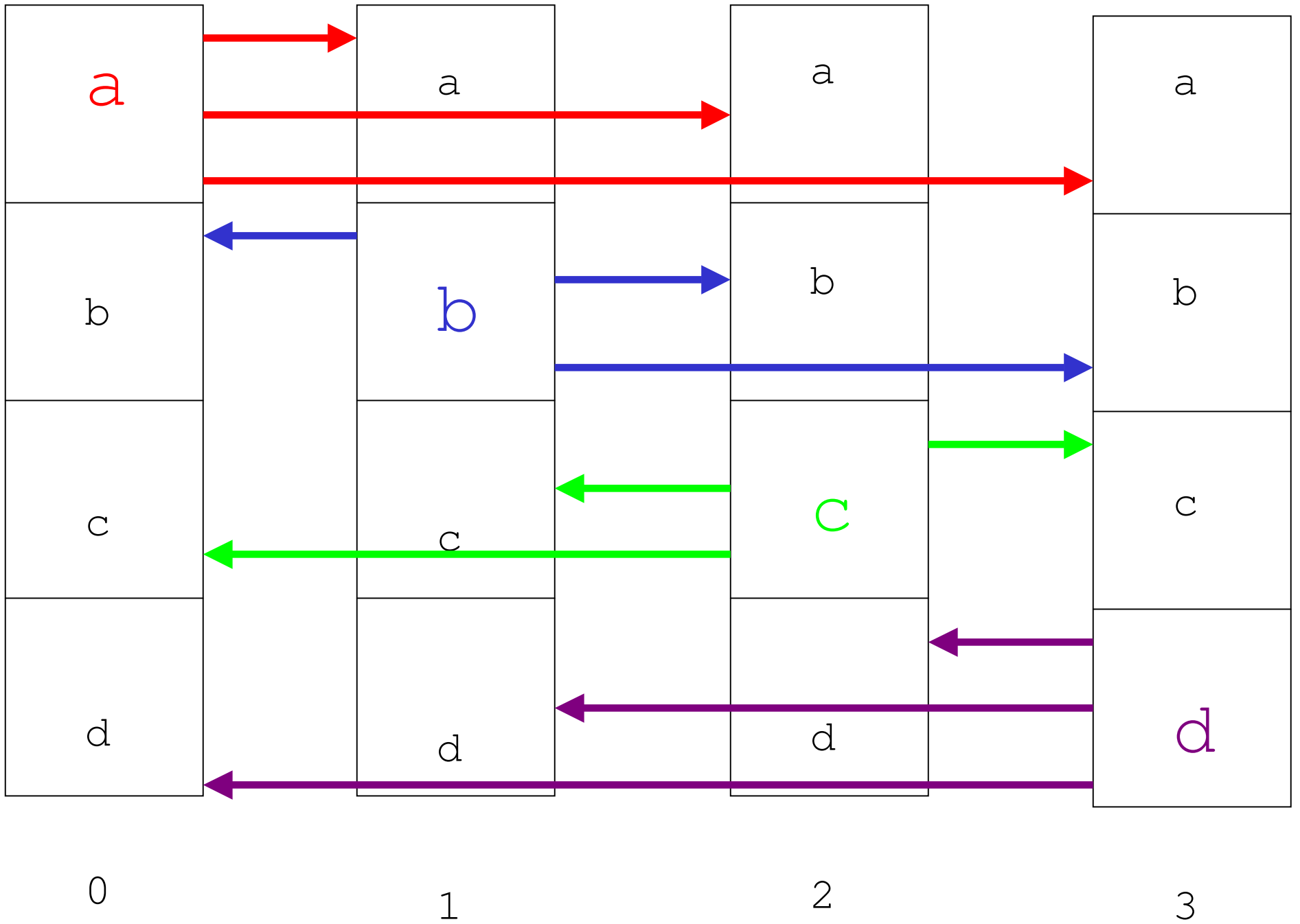


2

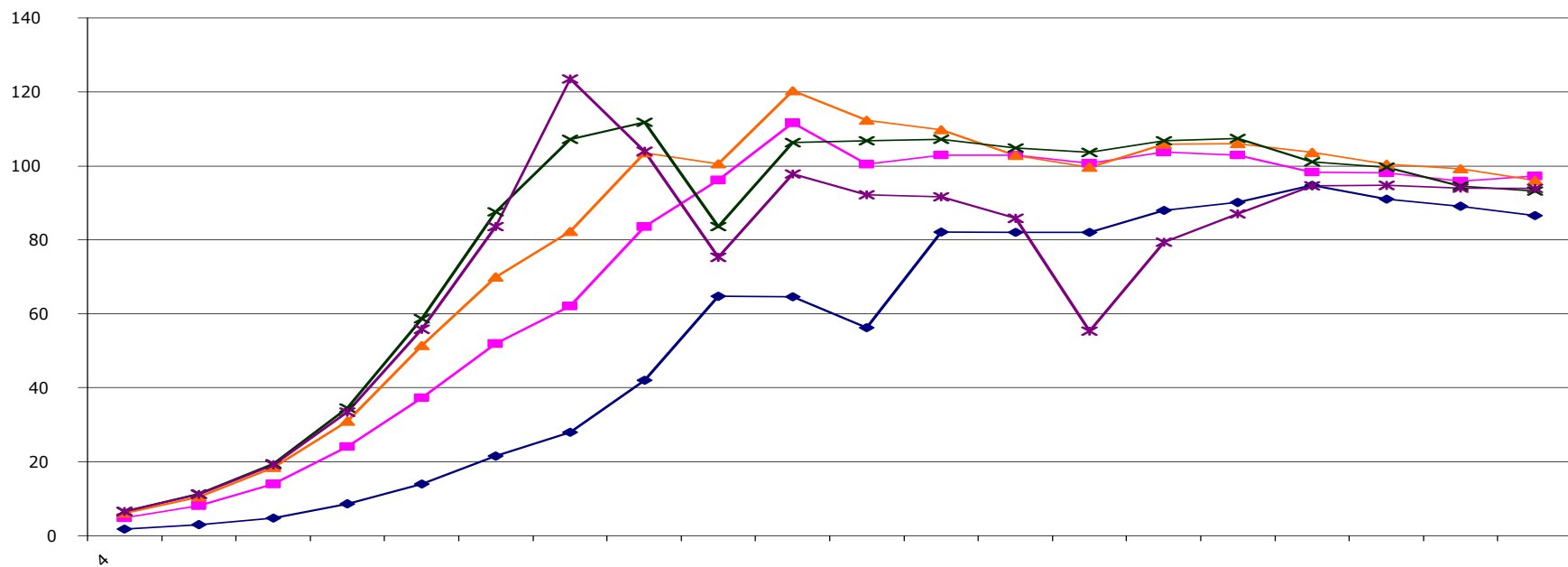


3

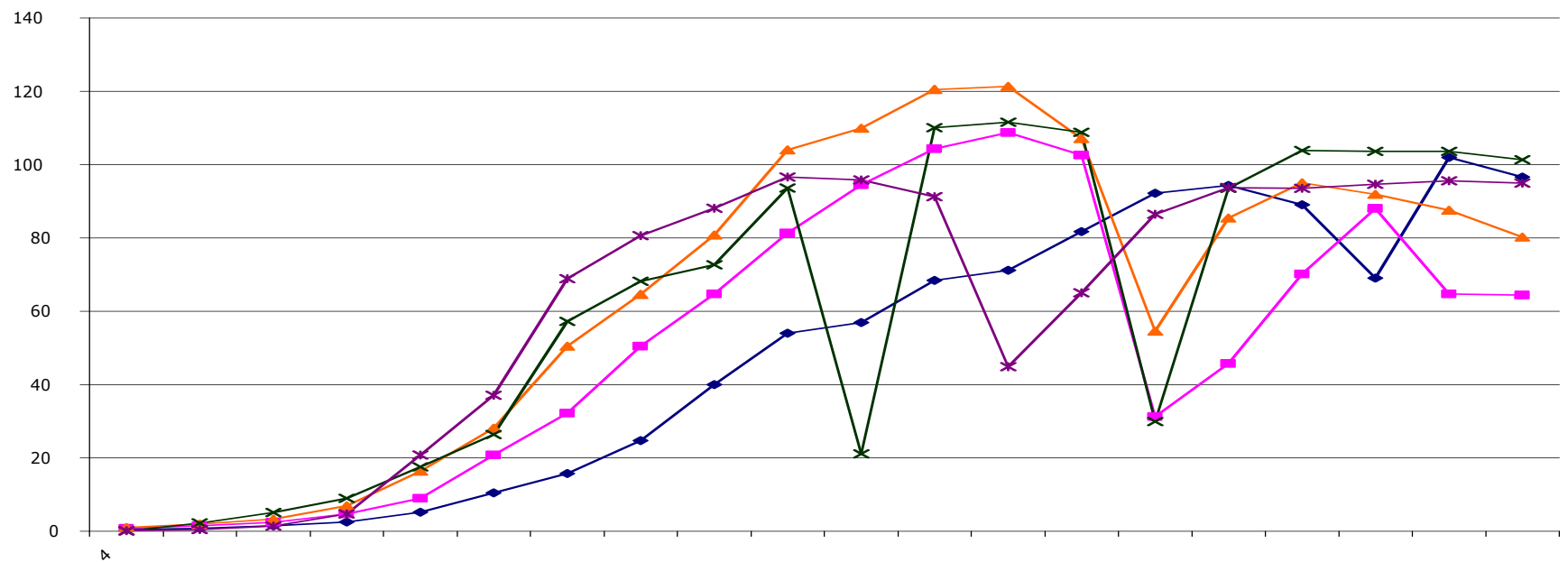


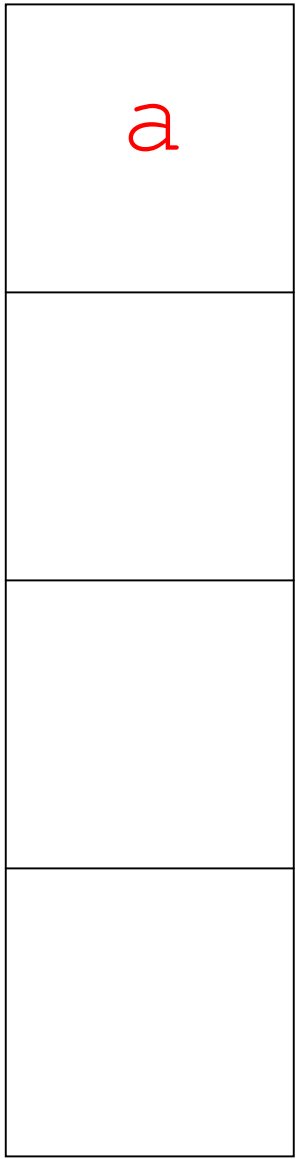


**Alltoall LAM**

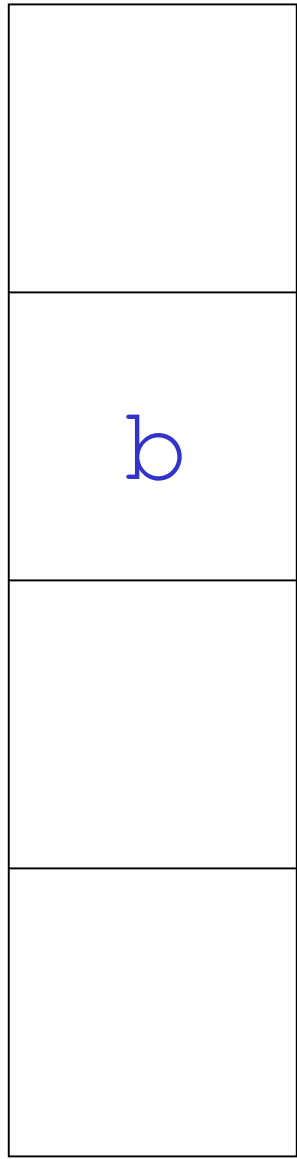


**Alltoall MPJava**

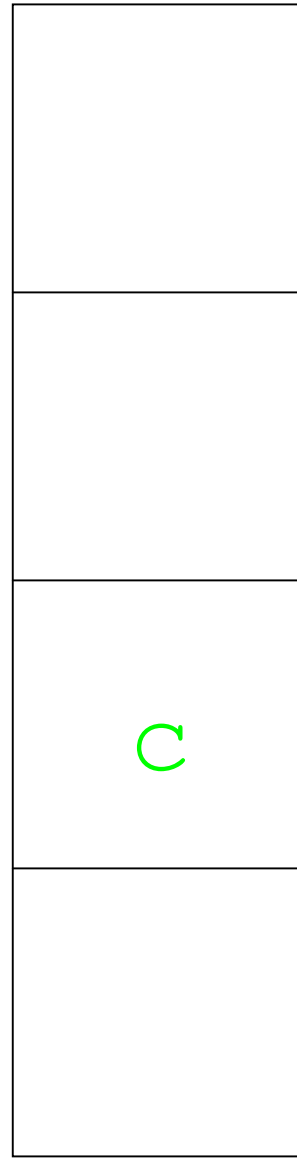




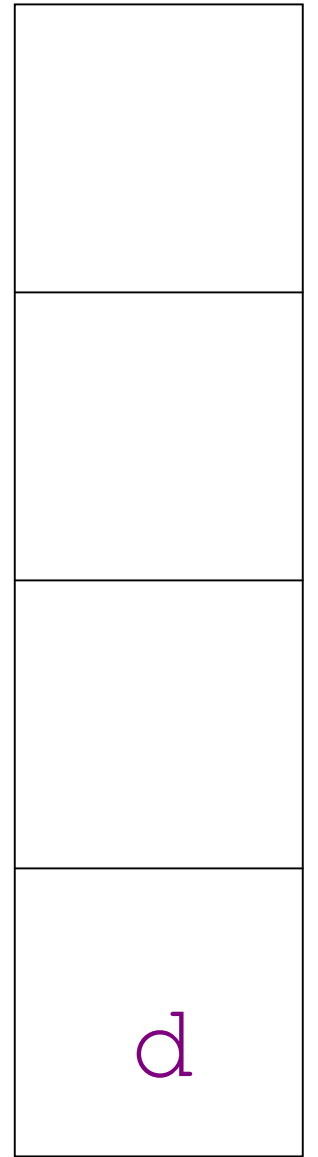
0



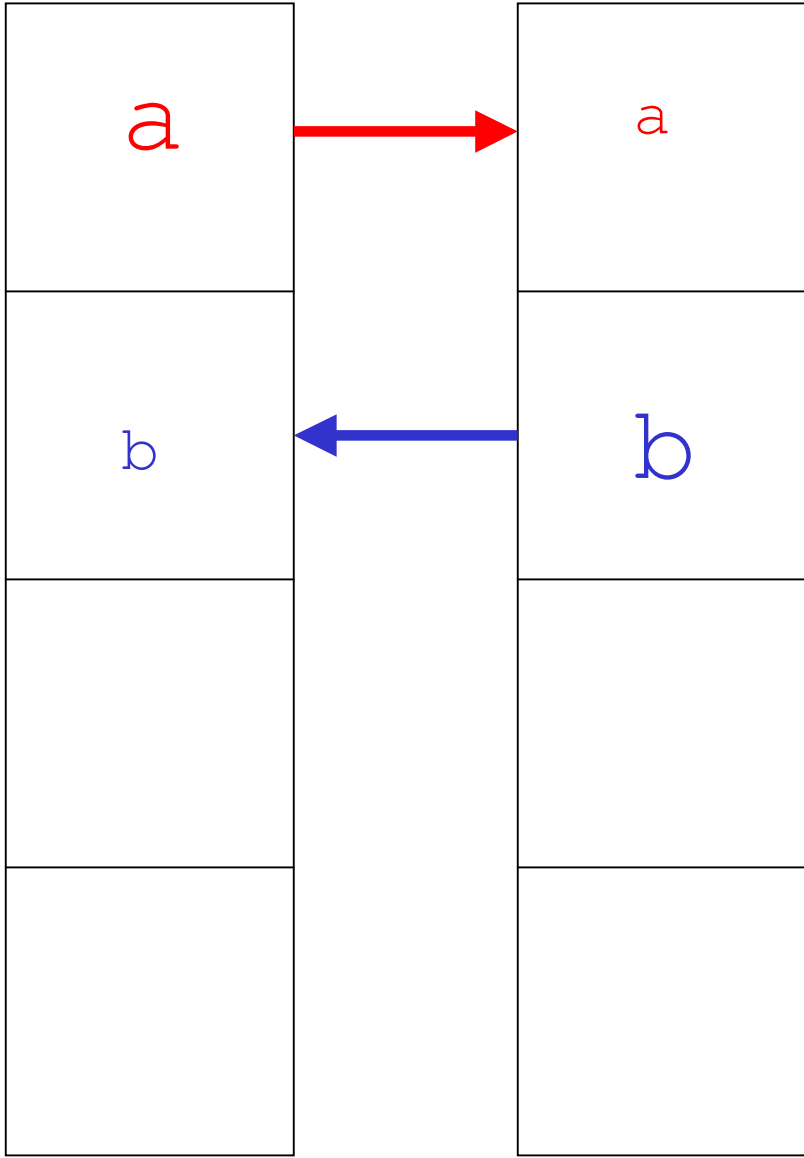
1



2

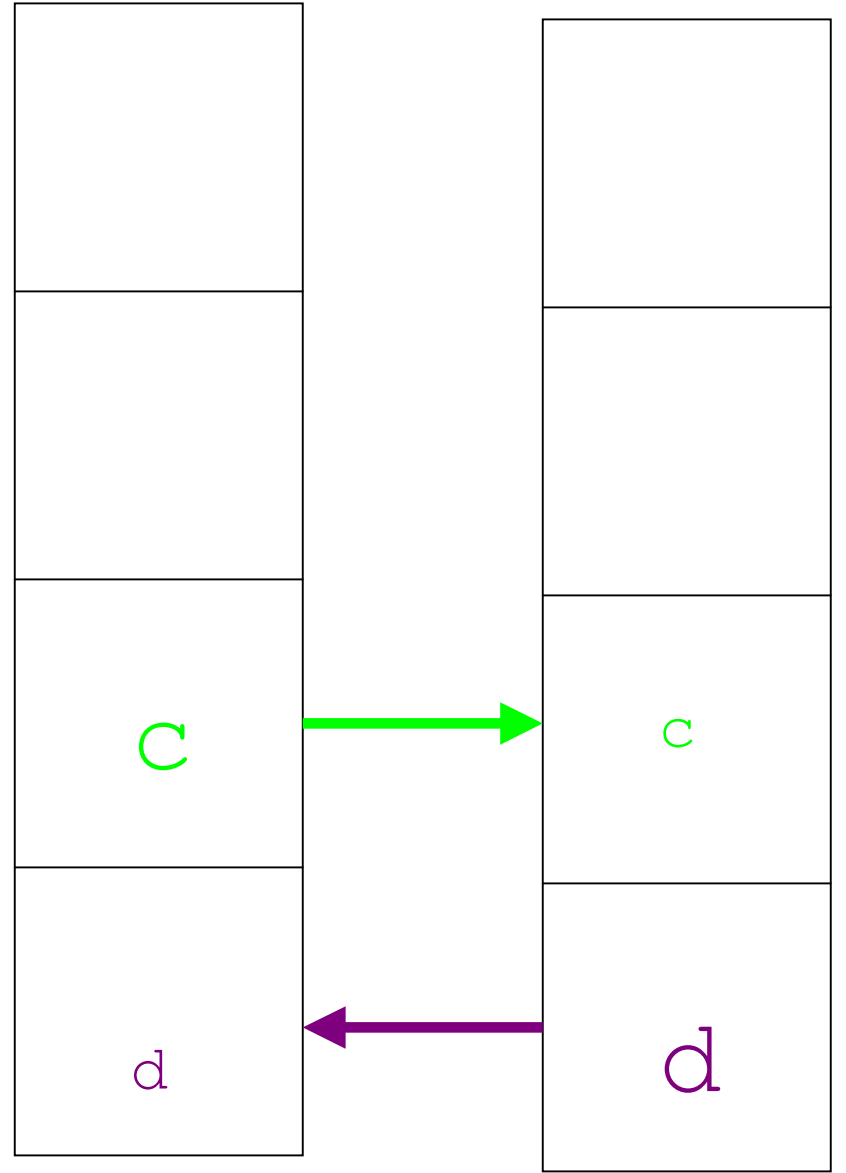


3



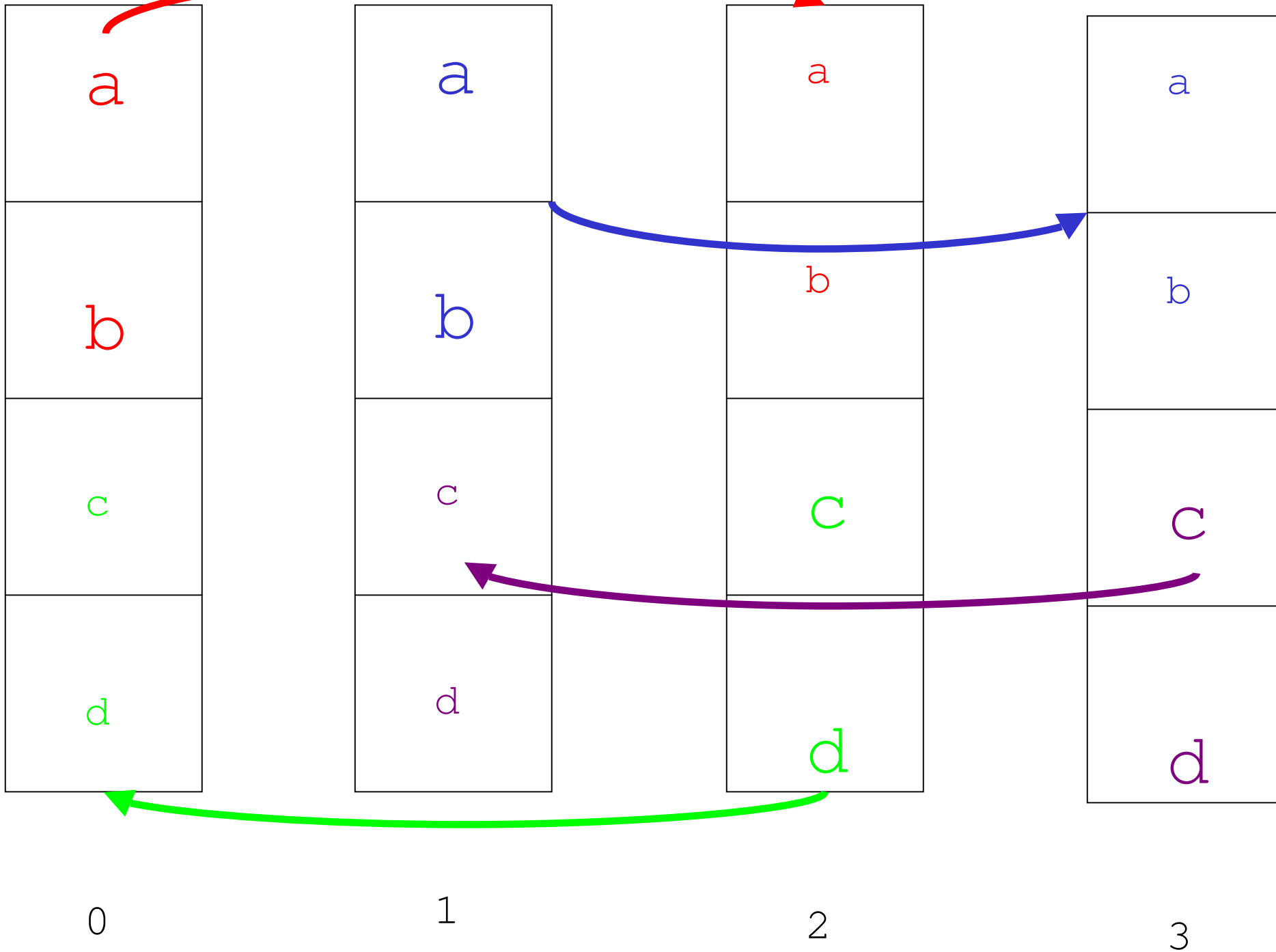
0

1

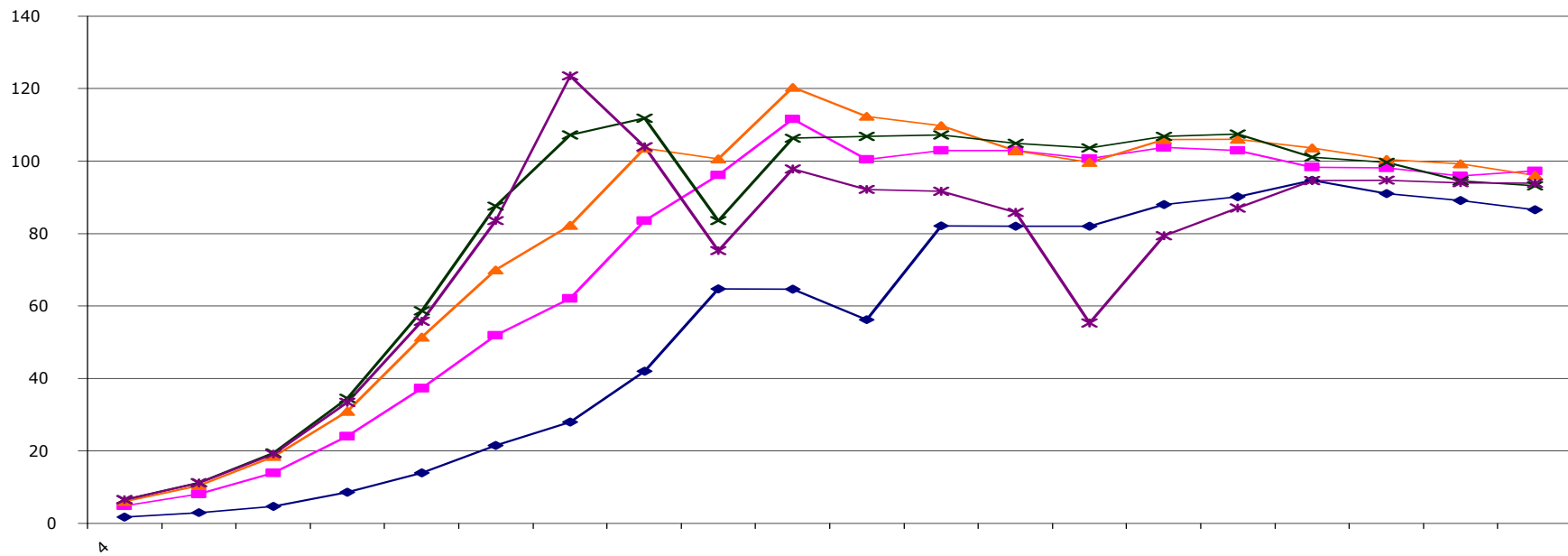


2

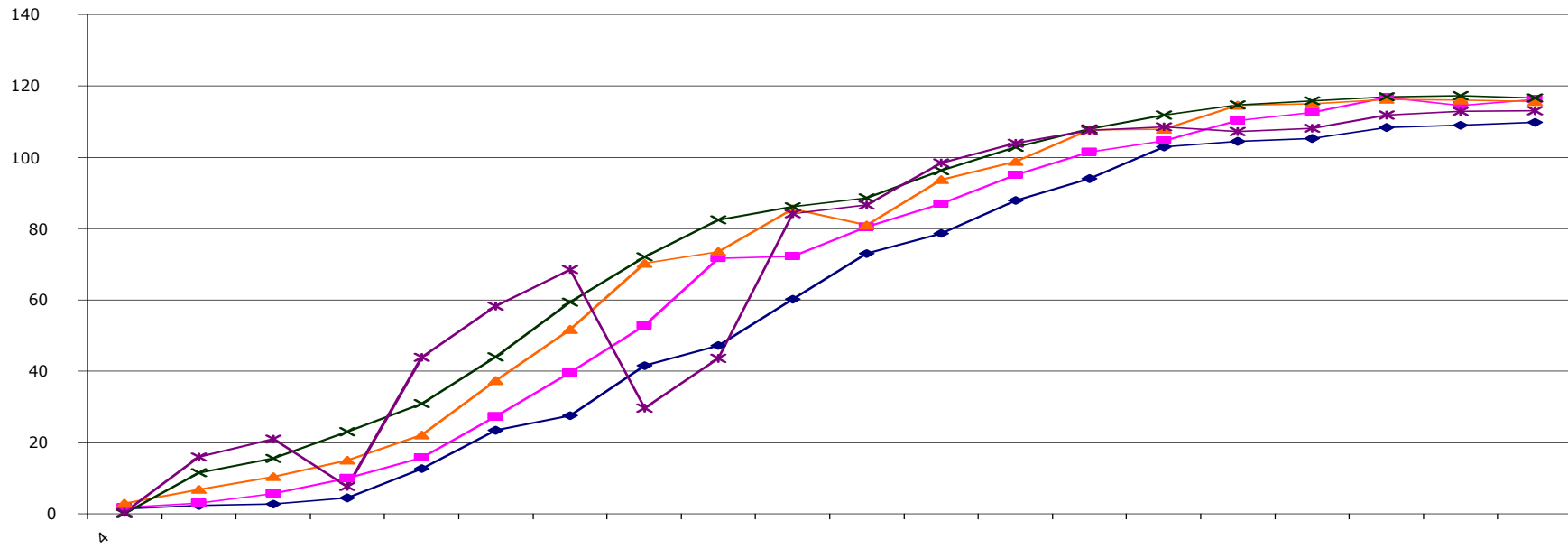
3



Alltoall LAM



Alltoall MPJava (prefix algorithm)



# NAS PB Conjugate Gradient

Class C Sparse Matrix is 150,000 X 150,000

241 nonzero elements per row

36,121,000 total nonzero elements

Class B Sparse Matrix is 75,000 X 75,000

183 nonzero elements per row

13,708,000 total nonzero elements

$$\begin{array}{cccc}
 & \mathbf{A} & & \\
 & & \cdot & \mathbf{p} = & & \mathbf{q} \\
 \begin{array}{|c|c|c|c|}
 \hline
 a & b & c & d \\
 \hline
 e & f & g & h \\
 \hline
 i & j & k & l \\
 \hline
 m & n & o & p \\
 \hline
 \end{array} & & * & \begin{array}{|c|}
 \hline
 w \\
 \hline
 x \\
 \hline
 y \\
 \hline
 z \\
 \hline
 \end{array} & = & \begin{array}{cccc}
 aw + bx + cy + dz \\
 ew + fx + gy + hz \\
 iw + jx + ky + lz \\
 mw + nx + oy + pz
 \end{array}
 \end{array}$$

Simple approach to parallelizing matrix-vector multiple:  
Stripe across rows

$$\begin{array}{l}
 0 \\
 1 \\
 2 \\
 3
 \end{array}
 \begin{array}{|c|c|c|c|}
 \hline
 a & b & c & d \\
 \hline
 e & f & g & h \\
 \hline
 i & j & k & l \\
 \hline
 m & n & o & p \\
 \hline
 \end{array}
 *
 \begin{array}{|c|}
 \hline
 w \\
 \hline
 x \\
 \hline
 y \\
 \hline
 z \\
 \hline
 \end{array}
 =
 \begin{array}{|c|c|c|c|}
 \hline
 aw + bx + cy + dz \\
 \hline
 ew + fx + gy + hz \\
 \hline
 iw + jx + ky + lz \\
 \hline
 mw + nx + oy + pz \\
 \hline
 \end{array}$$

Requires an all-to-all broadcast to reconstruct the vector  $p$



# Multi-Dimensional matrix-vector multiply decomposition

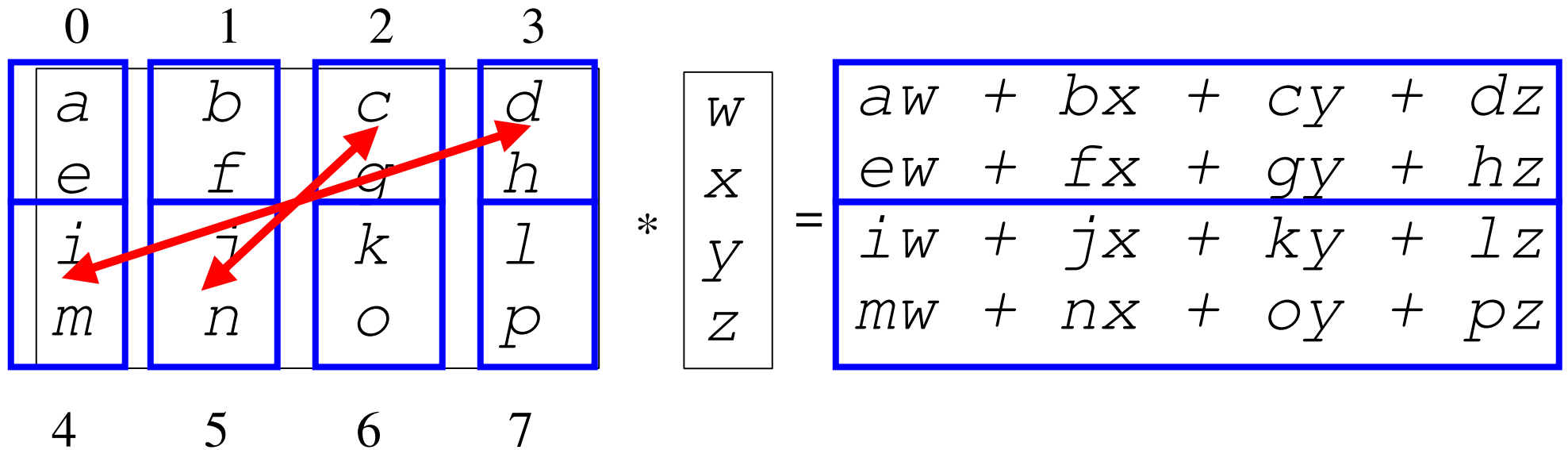
$$\begin{array}{|c|c|c|c|} \hline 0 & 1 & 2 & 3 \\ \hline a & b & c & d \\ \hline e & f & g & h \\ \hline i & j & k & l \\ \hline m & n & o & p \\ \hline 4 & 5 & 6 & 7 \\ \hline \end{array} * \begin{array}{|c|} \hline w \\ \hline x \\ \hline y \\ \hline z \\ \hline \end{array} = \begin{array}{|c|} \hline aw \\ \hline ew \\ \hline iw \\ \hline mw \\ \hline \end{array} + \begin{array}{|c|} \hline bx \\ \hline fx \\ \hline jx \\ \hline nx \\ \hline \end{array} + \begin{array}{|c|} \hline cy \\ \hline gy \\ \hline ky \\ \hline oy \\ \hline \end{array} + \begin{array}{|c|} \hline dz \\ \hline hz \\ \hline lz \\ \hline pz \\ \hline \end{array}$$

# Multi-Dimensional matrix-vector multiply decomposition

$$\begin{array}{|c|c|c|c|} \hline 0 & 1 & 2 & 3 \\ \hline a & b & c & d \\ \hline e & f & g & h \\ \hline i & j & k & l \\ \hline m & n & o & p \\ \hline 4 & 5 & 6 & 7 \\ \hline \end{array} * \begin{array}{|c|} \hline w \\ \hline x \\ \hline y \\ \hline z \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline aw & + & bx & + & cy & + & dz \\ \hline ew & + & fx & + & gy & + & hz \\ \hline iw & + & jx & + & ky & + & lz \\ \hline mw & + & nx & + & oy & + & pz \\ \hline \end{array}$$

Reduction along decomposed rows

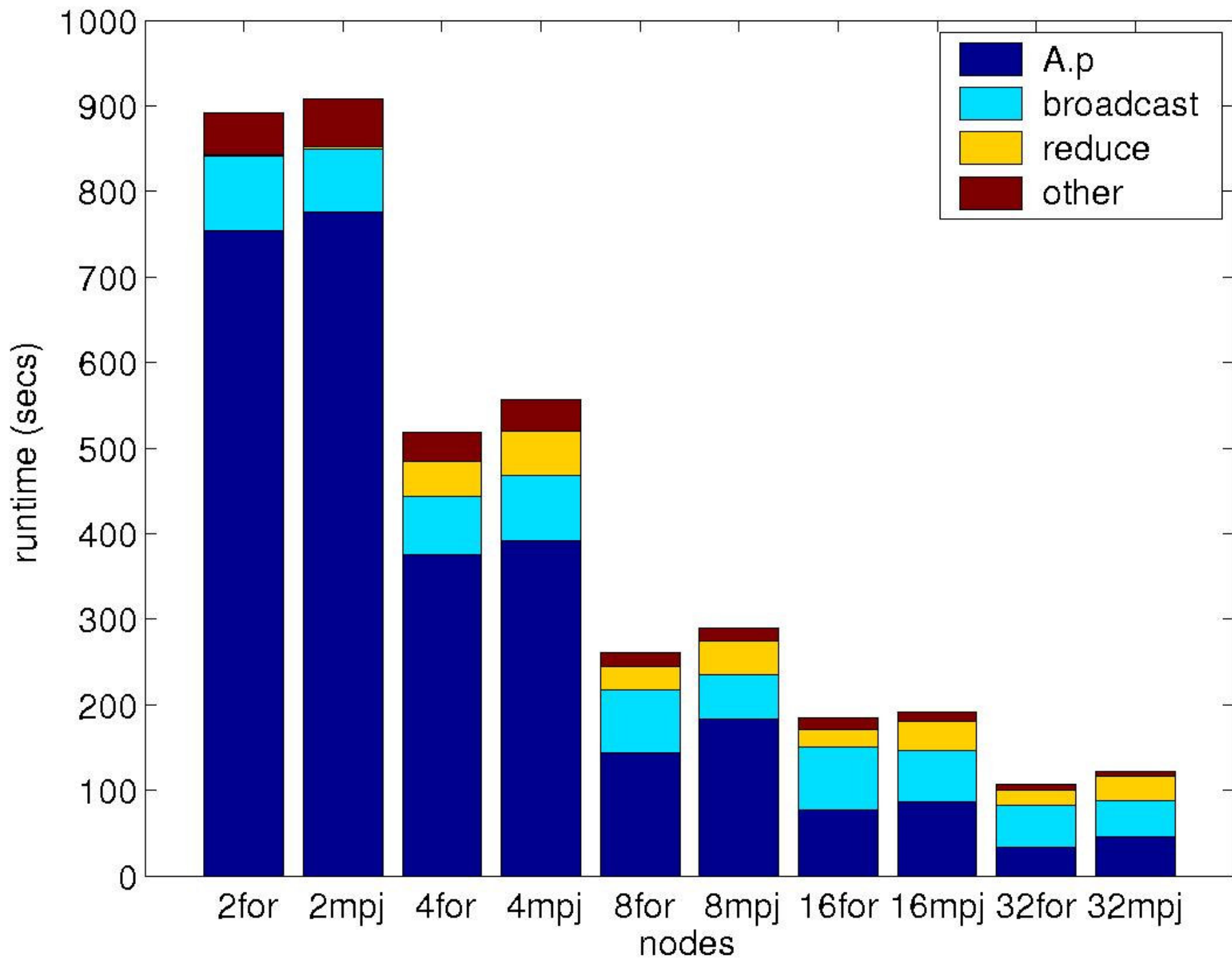
# Multi-Dimensional matrix-vector multiply decomposition



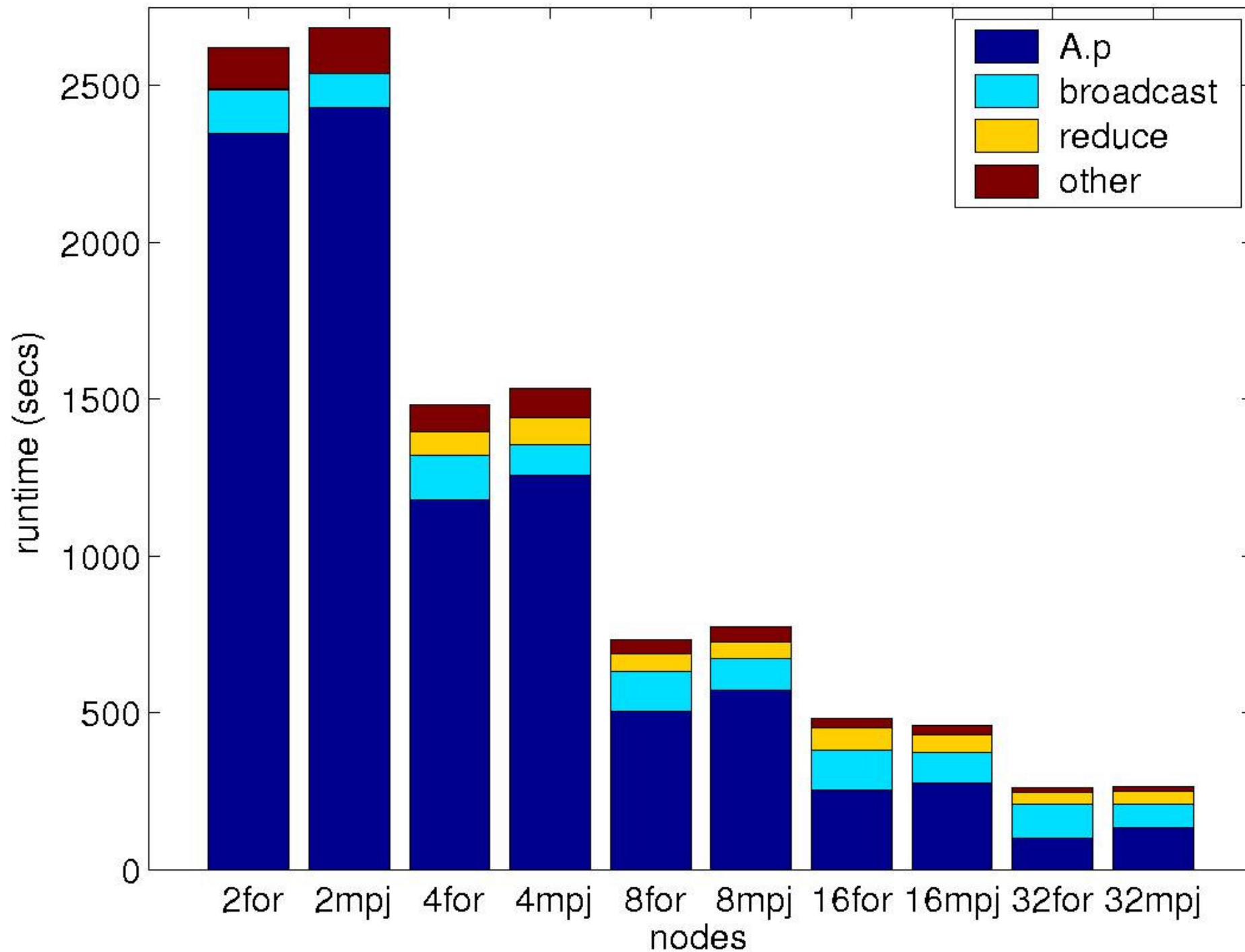
Node 4 needs w, and has y,z  
 Node 3 needs z, has w,x  
 SWAP

Node 2 needs y, and has w,x  
 Node 5 needs x, and has y,z  
 SWAP

Conjugate Gradient, Two-Dimensional algorithm, Class B



Conjugate Gradient, Two-Dimensional algorithm, Class C



# Conclusion

- A pure-Java message passing framework can provide performance competitive with widely available Fortran and MPI implementations
- java.nio is much faster than the older I/O model
- Java Just-In-Time compilers can deliver competitive performance
- Java has many other useful features
  - type safety
  - bounds checks
  - extensive libraries
  - portable
    - Is performance portable too?
  - easy to integrate with databases, webservers, GRID applications

# Future Work

- Exploiting asynchronous pipes
  - Great for work-stealing and work-sharing algorithms, but...
  - subject to Thread scheduling woes
- What about clusters of SMPs?
  - Different bottlenecks
  - More use for multiple threads on a single node
  - Importance of interleaving communication and computation
- Is MPI the right target?
  - BLAS, LAPACK, Netsolver, etc. suggest that programmers will use libraries

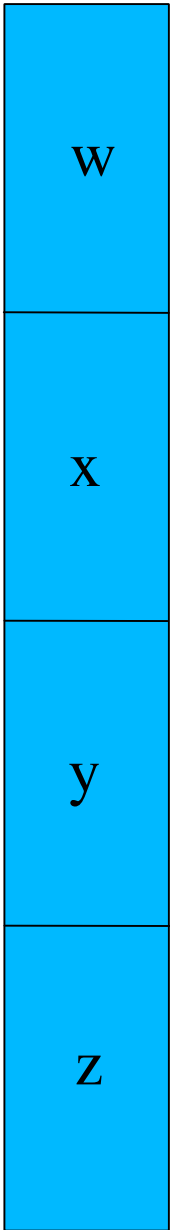
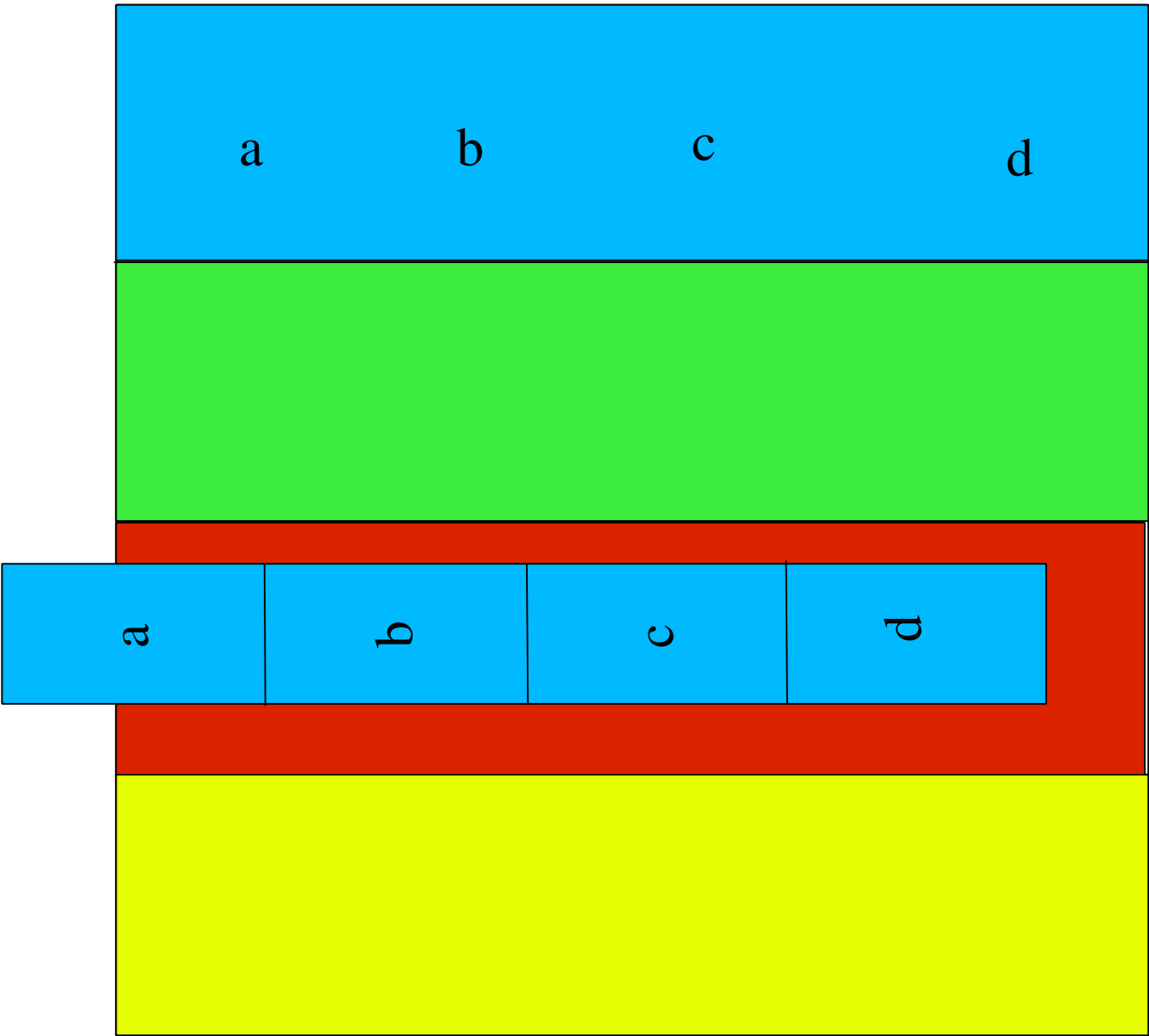
# Where do we go next?

- Java has the reputation that it's too slow for scientific programming!
  - Is this still accurate?
  - Or were we lucky with our benchmarks?
- Interest in message passing for Java was high a couple of years ago, but has waned
  - Because of performance?
- Does anyone care?
  - Are people happy with Fortran?
  - Is there enough interest in Java for scientific computing?



# Java may be fast enough but...

- No operator overloading
- No multiarrays package (yet)
  - Also need syntactic sugar to replace `.get()/.set()` methods with brackets!
- Autoboxing
- Generics (finally available in 1.5)
- Fast, efficient support for a Complex datatype
  - Stack-allocatable objects in general?
- C# provides all/most of these features



NAS PB implementation uses a better algorithm

Multi-Dimensional matrix-vector multiply decomposition

$$\begin{array}{cccc} 0 & 1 & 2 & 3 \\ \boxed{a} & \boxed{b} & \boxed{c} & \boxed{d} \\ \boxed{e} & \boxed{f} & \boxed{g} & \boxed{h} \\ \boxed{i} & \boxed{j} & \boxed{k} & \boxed{l} \\ \boxed{m} & \boxed{n} & \boxed{o} & \boxed{p} \\ 4 & 5 & 6 & 7 \end{array} * \begin{array}{c} w \\ x \\ y \\ z \end{array} = \begin{array}{cccc} \boxed{aw} & + & \boxed{bx} & + & \boxed{cy} & + & \boxed{dz} \\ \boxed{ew} & + & \boxed{fx} & + & \boxed{gy} & + & \boxed{hz} \\ \boxed{iw} & + & \boxed{jx} & + & \boxed{ky} & + & \boxed{lz} \\ \boxed{mw} & + & \boxed{nx} & + & \boxed{oy} & + & \boxed{pz} \end{array}$$

Note the additional swap required for “inner” nodes