

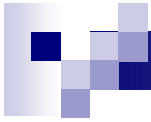
# Performance Analysis of Java Message-Passing Libraries on Fast Ethernet, Myrinet and SCI Clusters

Guillermo L. Taboada, Juan Touriño and Ramon Doallo  
{taboada,juan,doallo}@udc.es



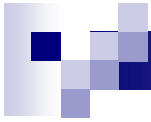
Dept. of Electronics and Systems  
University of A Coruña, Spain

IEEE Cluster 2003 – Hong Kong, DECEMBER 2003



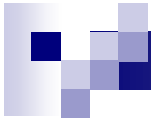
# Overview

- Our analysis combines:
  - State of the Art in Java Message-Passing Libraries
  - Modeling Message-Passing Primitives
  - Performance Analysis on Fast Ethernet, Myrinet and SCI clusters
- Results:
  - Evaluation of the most outstanding Java Message-Passing Libraries and performance implications

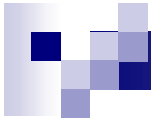


# Outline

- **Introduction**
- **Java Message-Passing Libraries**
- **Modeling Message-Passing Primitives**
- **Experimental Conditions**
- **Experimental Results**
- **Analysis of Performance Results**
- **Conclusions**

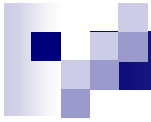






# Introduction

- Related work: Papers by Stankovic and Zhang [1] and by Getov, Gray and Sunderam [2]. Both works evaluate out-of-date libraries and do not derive performance analytical models
- Main contributions:
  - Survey of the state of the art in Java message-passing libraries
  - An updated performance evaluation of these libraries on Fast Ethernet, Myrinet and SCI clusters
  - Their performance analytical models



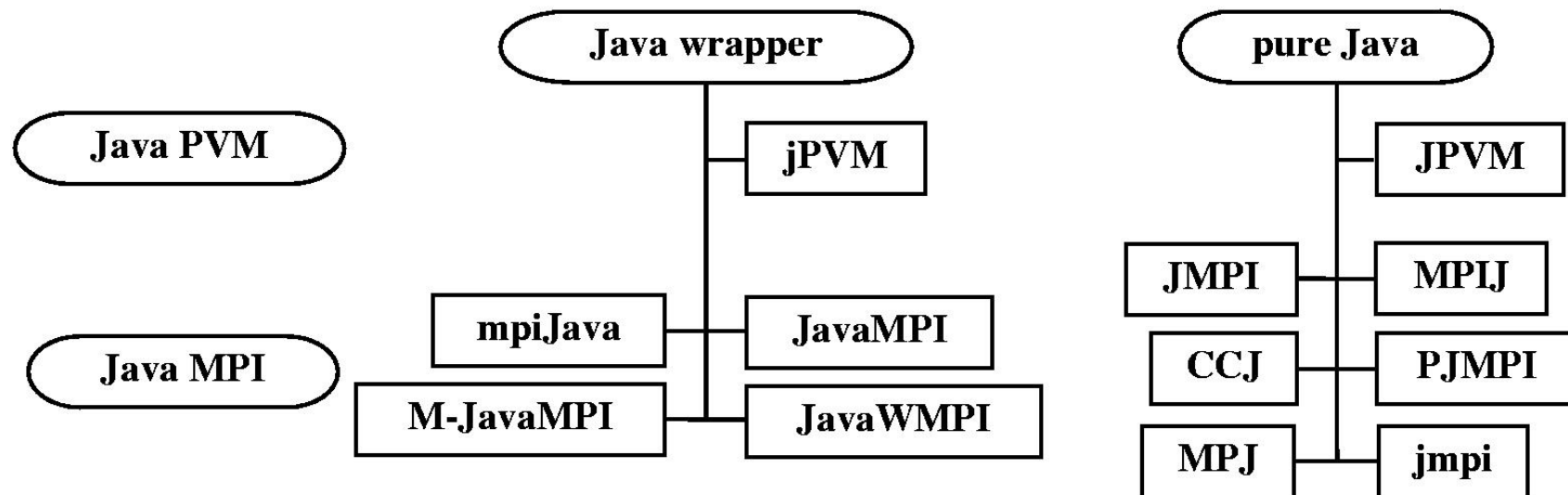
# Java Message-Passing Libraries

- Two main types of implementations of Message-Passing for Java:
  - **Java wrapper** provides efficient MPI communication through calling native methods using JNI. The major drawback is lack of portability
  - **Pure Java** provides a portable message-passing implementation since the whole library is developed in Java, although the communication is less efficient

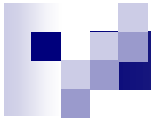


# Java Message-Passing Libraries

## ■ Taxonomy of Existing Libraries:

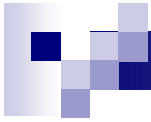






# Java Message-Passing Libraries

- Selected Java Message-Passing libraries:
  - **mpiJava** the most active Java wrapper project
  - **CCJ** pure Java implementation of the Manta Team (Netherlands)
  - **JMPI** pure Java implementation of the Univ. of Massachusetts



# Modeling Message-Passing Primitives

## ■ Point-to-point communications:

$$T(n) = t_s + t_b n$$

$$Bw(n) = n/T(n)$$

$T$ : message latency

$n$ : message size (bytes)

$t_s$ : startup time

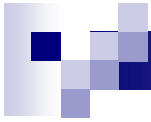
$t_b$ : transfer time (per byte)

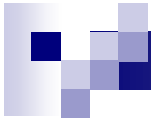
$p$ : processors.

## ■ Collective communications:

$$T(n, p) = t_s(p) + t_b(p) n$$

$$Bw(n, p) = n/T(n, p)$$

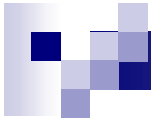






# Experimental Conditions

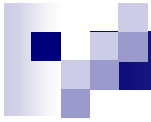
- JVM:
  - IBM JVM 1.4 JITC
  - Sun JVM Sun 1.4.1 HotSpot
- Libraries
  - MPICH 1.2.4, MPICH-GM and SCI-MPICH
  - ScaMPI
  - mpiJava 1.2.5 over native libraries
  - CCJ 0.1
  - JMPI



# Analytical models

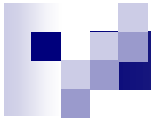
- Send metrics on Fast Ethernet:

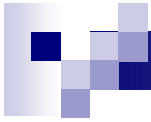
<i>Library</i>	<i>Analytical Model</i>		<i>Experimental Results</i>	
	$t_s \{ \mu s \}$	$t_b \{ ns/B \}$	$T(16B) \{ \mu s \}$	$Bw(1MB) \{ MB/s \}$
<i>MPICH</i>	69	90.3	72	11.061
<i>mpiJava</i>	101	100.7	105	9.923
<i>CCJ</i>	800	138.2	800	7.217
<i>JMPI</i>	4750	154.4	4750	6.281





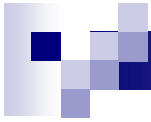


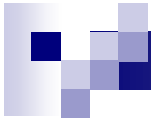




# Analytical models

<i>Library</i>	<i>Analytical Model</i>		<i>Experimental Results</i>	
	$t_s(p) \{ \mu s \}$	$t_b(p) \{ ns/B \}$	$T(16B,8)$	$Bw(1MB,8)$
<i>MPICH</i>	$7+117 \lceil 1p \rceil$	$-0.3+90.4 \lceil 1p \rceil$	364	3.686
<i>mpiJava</i>				





# Analytical models

- Broadcast metrics on SCI ( $\lg = \log_2 p$ ) :

<i>Library</i>	<i>Analytical Model</i>		<i>Experimental Results</i>	
	$t_s(p) \{ \mu s \}$	$t_b(p) \{ ns/B \}$	$T(16B,8)$	$Bw(1MB,8)$
<i>ScaMPI</i>	$6 \lceil \lg \rceil$	$-0.105 + 4.128 \lceil \lg \rceil$	18	81.71
<i>mpiJava</i>	$33 + 7 \lceil \lg \rceil$	$-0.197 + 4.458 \lceil \lg \rceil$	48	76.71
<i>CCJ</i>	$-800 + 1400 \lceil \lg \rceil$	$-4.242 + 93.01 \lceil \lg \rceil$	3400	3.63
<i>JMPI</i>	$-2800 + 2900p$	$-89.71 + 95.52p$	20600	1.45



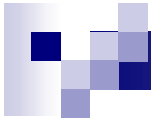
# Analytical models

- Broadcast metrics on Fast Ethernet, Myrinet and SCI :

<i>Library</i>	<i>Analytical Model</i>		<i>Experimental Results</i>	
	$t_s(p) \{ \mu s \}$	$t_b(p) \{ ns/B \}$	$T(16B,8)$	$Bw(1MB,8)$
<i>MPICH</i>	$7+117 \lceil 1p \rceil$	$-0.3+90.4 \lceil 1p \rceil$	364	3.686
<i>mpiJava</i>	$19+124 \lceil 1p \rceil$	$10.1+90.5 \lceil 1p \rceil$	406	3.546
<i>CCJ</i>	$-430+1430 \lceil 1p \rceil$	$6.4+130.4 \lceil 1p \rceil$	3800	2.506
<i>JMPI</i>	$-9302+7151p$	$-123.2+175.7p$	41600	0.752

<i>Library</i>	<i>Analytical Model</i>		<i>Experimental Results</i>	
	$t_s(p) \{ \mu s \}$	$t_b(p) \{ ns/B \}$	$T(16B,8)$	$Bw(1MB,8)$
<i>MPICH-GM</i>	$3+8 \lceil 1p \rceil$	$0.012+5.741 \lceil 1p \rceil$	28	57.77
<i>mpiJava</i>	$20+17 \lceil 1p \rceil$	$0.036+5.263 \lceil 1p \rceil$	101	62.78
<i>CCJ</i>	$-800+1600 \lceil 1p \rceil$	$-10.64+40.69 \lceil 1p \rceil$	4000	8.72
<i>JMPI</i>	$-8617+5356p$	$-61.57+66.7p$	32400	1.80

<i>Library</i>	<i>Analytical Model</i>		<i>Experimental Results</i>	
	$t_s(p) \{ \mu s \}$	$t_b(p) \{ ns/B \}$	$T(16B,8)$	$Bw(1MB,8)$
<i>ScaMPI</i>	$6 \lceil 1p \rceil$	$-0.105+4.128 \lceil 1p \rceil$	18	81.71
<i>mpiJava</i>	$33+7 \lceil 1p \rceil$	$-0.197+4.458 \lceil 1p \rceil$	48	76.71
<i>CCJ</i>	$-800+1400 \lceil 1p \rceil$	$-4.242+93.01 \lceil 1p \rceil$	3400	3.63
<i>JMPI</i>	$-2800+2900p$	$-89.71+95.52p$	20600	1.45



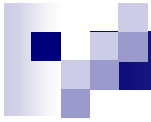
# Analytical models

- Collective metrics on Myrinet (several design issues):

<i>Library</i>	<i>Broadcast Analytical Model</i>		<i>Experimental Results</i>	
	$t_s(p) \{ \mu s \}$	$t_b(p) \{ ns/B \}$	$T(16B,8)$	$Bw(1MB,8)$
<i>MPICH-GM</i>	$3+8 \lceil 1p \rceil$	$0.012+5.741 \lceil 1p \rceil$	28	57.77
<i>mpiJava</i>	$20+17 \lceil 1p \rceil$	$0.036+5.263 \lceil 1p \rceil$	101	62.78
<i>CCJ</i>	$-800+1600 \lceil 1p \rceil$	$-10.64+40.69 \lceil 1p \rceil$	4000	8.72
<i>JMPI</i>	$-8617+5356p$	$-61.57+66.7p$	32400	1.80

<i>Library</i>	<i>Scatter Analytical Model</i>		<i>Experimental Results</i>	
	$t_s(p) \{ \mu s \}$	$t_b(p) \{ ns/B \}$	$T(16B,8)$	$Bw(1MB,8)$
<i>MPICH-GM</i>	$-7+9p$	$4.321+0.414 \lceil 1p \rceil$	65	190.26
<i>mpiJava</i>	$42+10p$	$7.223-0.358 \lceil 1p \rceil$	131	167.95
<i>CCJ</i>	$217+604p$	$19.11+10.38 \lceil 1p \rceil$	5000	18.26
<i>JMPI</i>	$-8287+6438p$	$46.04+11.66 \lceil 1p \rceil$	40400	9.23

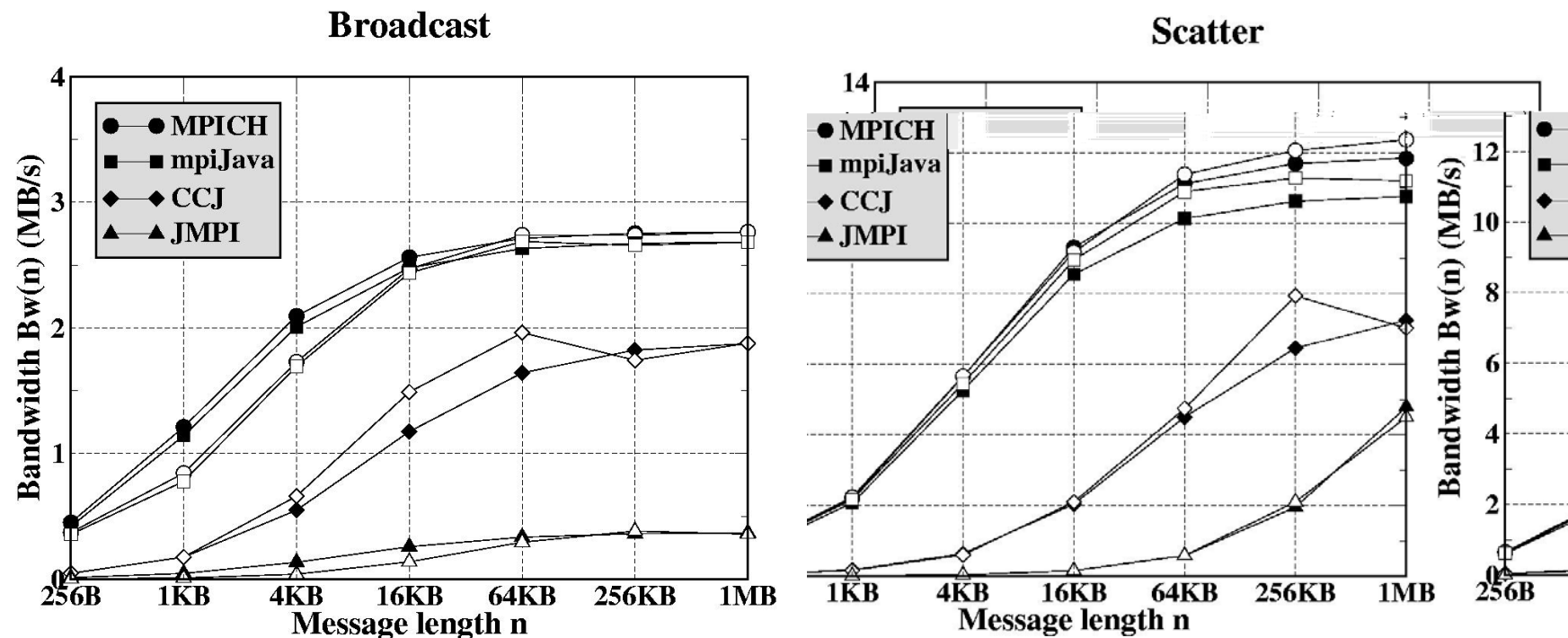
<i>Library</i>	<i>Allgather Analytical Model</i>		<i>Experimental Results</i>	
	$t_s(p) \{ \mu s \}$	$t_b(p) \{ ns/B \}$	$T(16B,8)$	$Bw(1MB,8)$
<i>MPICH-GM</i>	$-10+15p$	$5.308+1.098 \lceil 1p \rceil$	104	116.48
<i>mpiJava</i>	$30+17p$	$8.560+0.483 \lceil 1p \rceil$	173	100.63
<i>CCJ</i>	$817+944p$	$13.60+20.79p$	9400	5.25
<i>JMPI</i>	$-8296+7506p$	$-31.16+41.36p$	42400	2.85



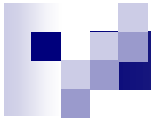


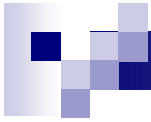
# Experimental Results

- Fast Ethernet: measured and estimated metrics on 16 processors



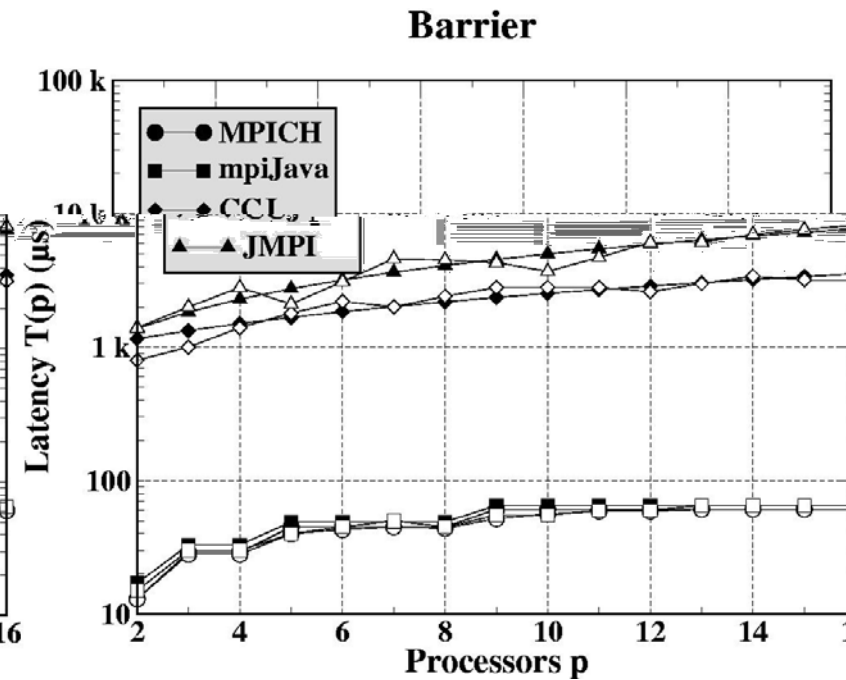
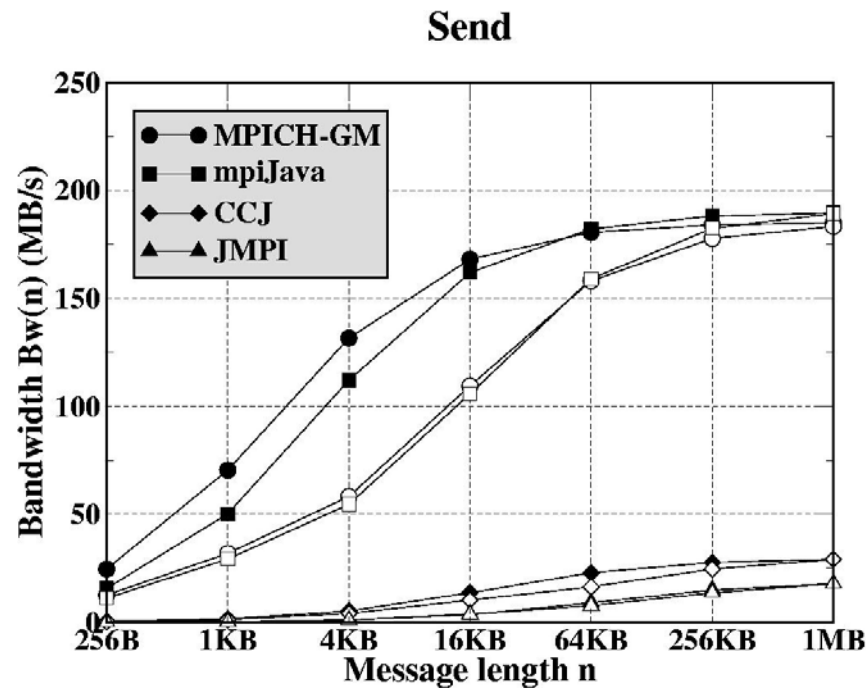


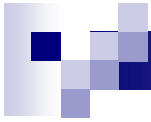


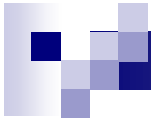


# Experimental Results

- Myrinet: measured and estimated metrics on 2 and 16 processors



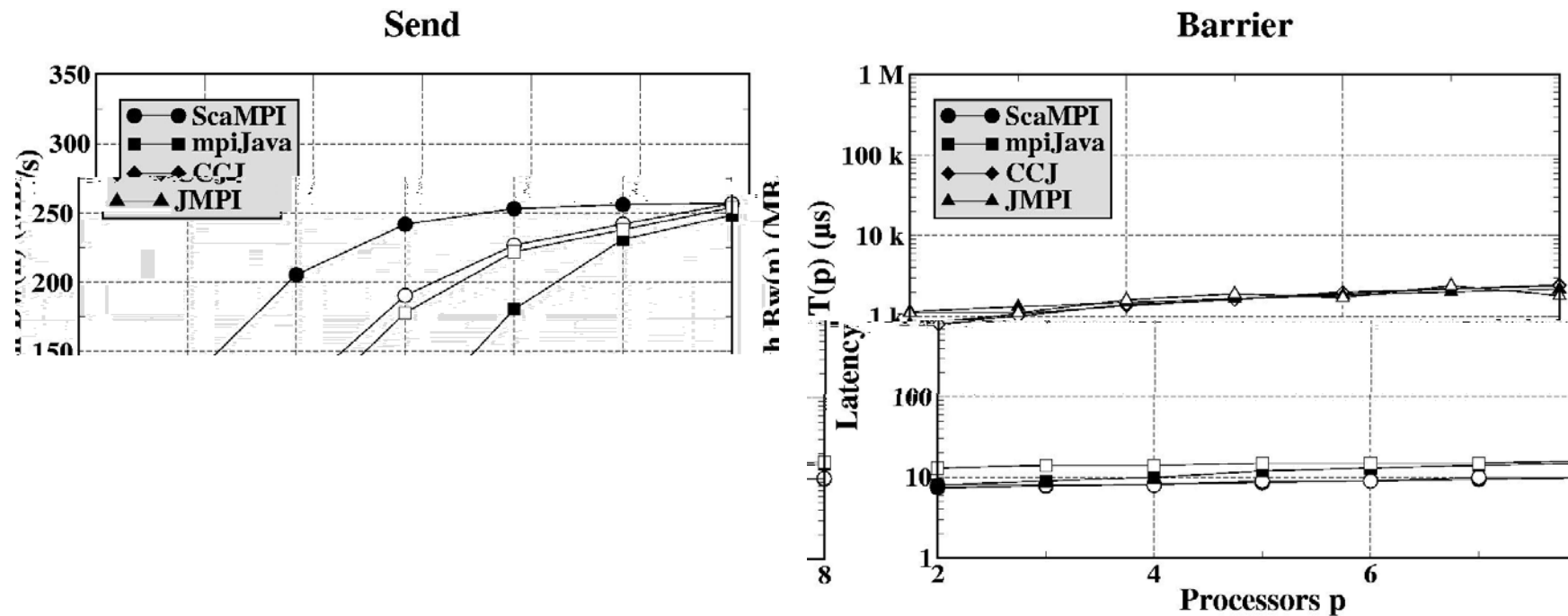


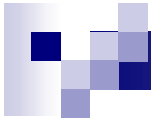




# Experimental Results

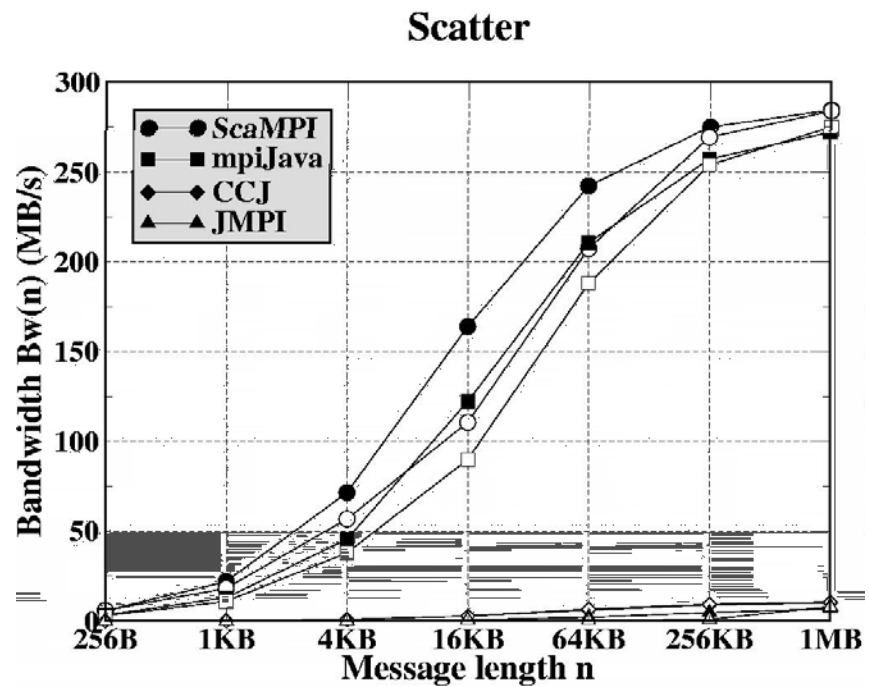
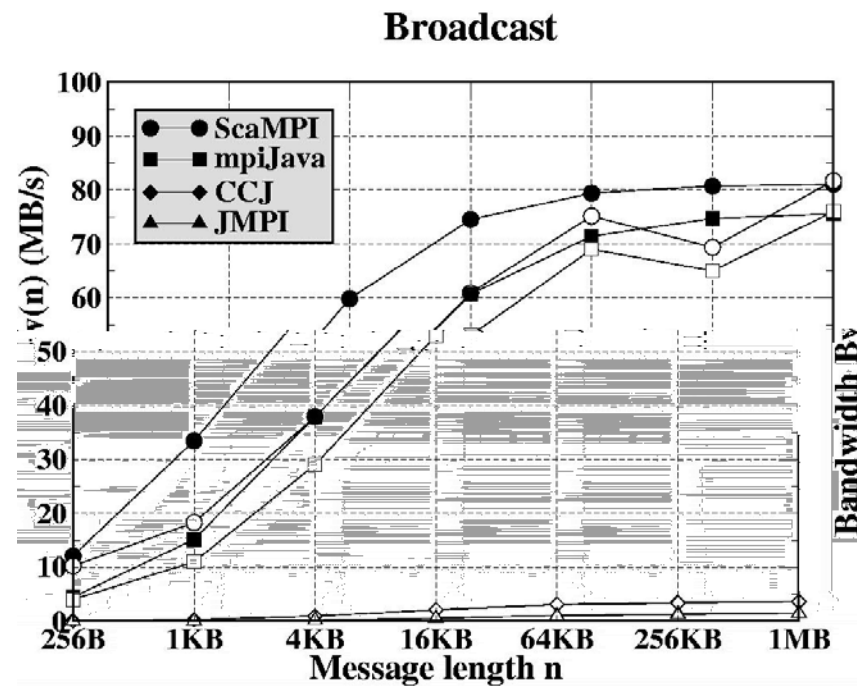
- SCL: measured and estimated metrics on 2 and 8 processors

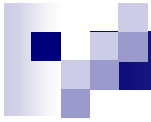




# Experimental Results

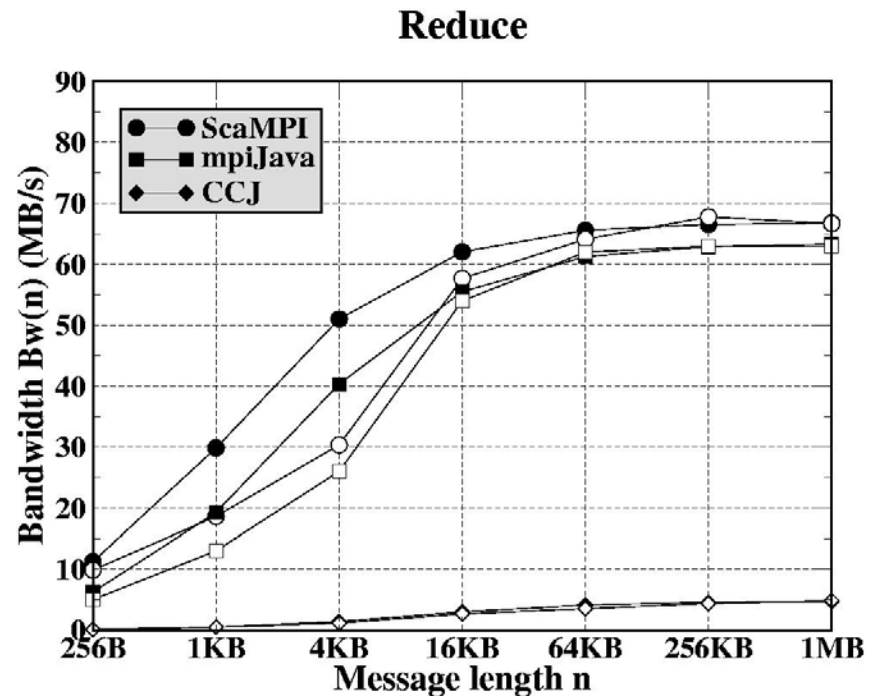
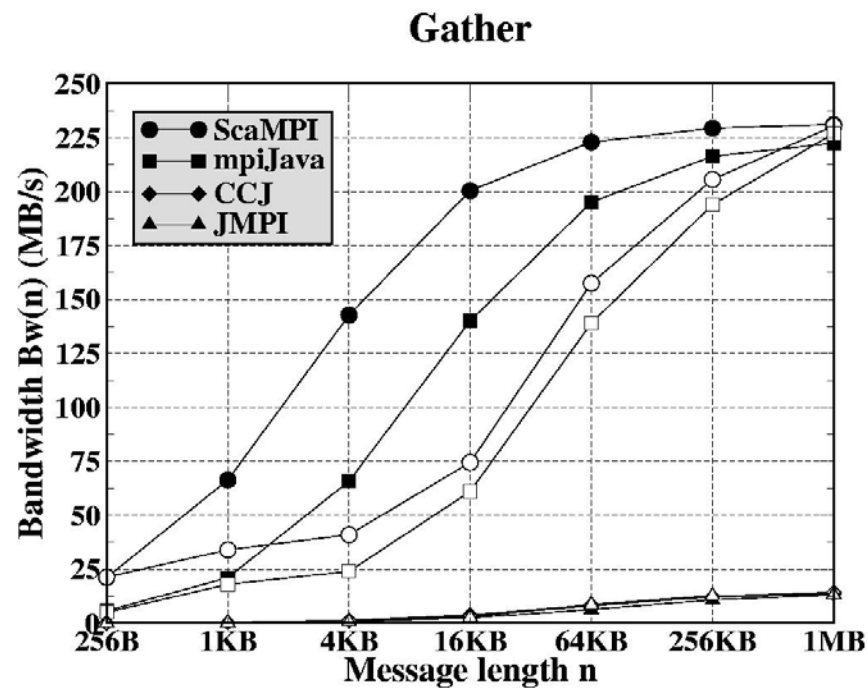
- SCl: measured and estimated metrics on 8 processors





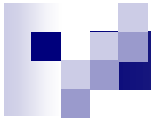
# Experimental Results

- SCI: measured and estimated metrics on 8 processors



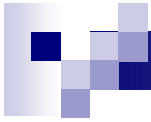






# Performance Issues

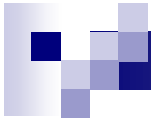
- Design issues. Study of equivalences:
  - Broadcast = Scatter + Allgather
  - Allgather = Gather + Broadcast
  - Allreduce = Reduce + Broadcast
  - Reducescatter = Reduce + Scatter





# Conclusions

- Native and wrapper libraries obtain good results in low latency networks
- Pure Java libraries need IP emulation to work on low latency clusters. The startup is about 15% better on IP over SCI whereas the transfer time is three times faster on IP over Myrinet than over SCI
- Performance results achieved from IP emulation are similar of those achieved on a Fast Ethernet cluster



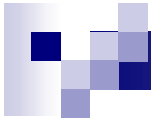
# Future Work

- Research efforts should concentrate on minimizing RMI overhead to consolidate and enhance the use of pure Java message-passing codes
- This is the topic of our future work, the optimization of RMI protocol, specially on low latency networks (Myrinet and SCI)



# References

- [1] Stankovic, N., Zhang, K. An Evaluation of Java Implementations of Message-Passing. Software-Practice and Experience 30(7) (2000) 741-763
- [2] W. Getov, P. Gray, V. Sunderam. MPI and Java-MPI: Contrasts and Comparisons of Low-Level Communications Performance. In Proc. of Supercomputing Conference, SC'99, Portland, OR (1999)



# Contact

- Guillermo L. Taboada  
[taboada@udc.es](mailto:taboada@udc.es)
- Computer Architecture Group  
Department of Electronics and Systems  
University of A Coruña
- <http://www.des.udc.es/~gltaboada>

# Performance Analysis of Java Message-Passing Libraries on Fast Ethernet, Myrinet and SCI Clusters

Guillermo L. Taboada, Juan Touriño and Ramon Doallo  
{taboada,juan,doallo}@udc.es



Dept. of Electronics and Systems  
University of A Coruña, Spain

IEEE Cluster 2003 – Hong Kong, DECEMBER 2003