

# Experimental Evaluation of BSP Programming Libraries

Peter Krusche

Department of Computer Science  
Centre for Scientific Computing  
University of Warwick

[peter@dcs.warwick.ac.uk](mailto:peter@dcs.warwick.ac.uk)

# Outline

## Motivation...

Study and compare the communication characteristics and performance predictability of 'BSP-style' communication libraries.

Introduction

• Outline

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

Conclusion

## Outline

1. The BSP Model
2. BSP Programming Libraries
3. Benchmarking
4. Performance and Predictability

# The BSP Model

The BSP model for parallel algorithms was used with some slight adaptations.

The model has been adapted here to use seconds instead of flops as a base unit for the running time:

- $p$  identical processor/memory pairs (computing nodes), computation speed  $f$
- Arbitrary interconnection network, latency  $l$ , bandwidth  $g$

Introduction

The BSP Model

• The BSP Model

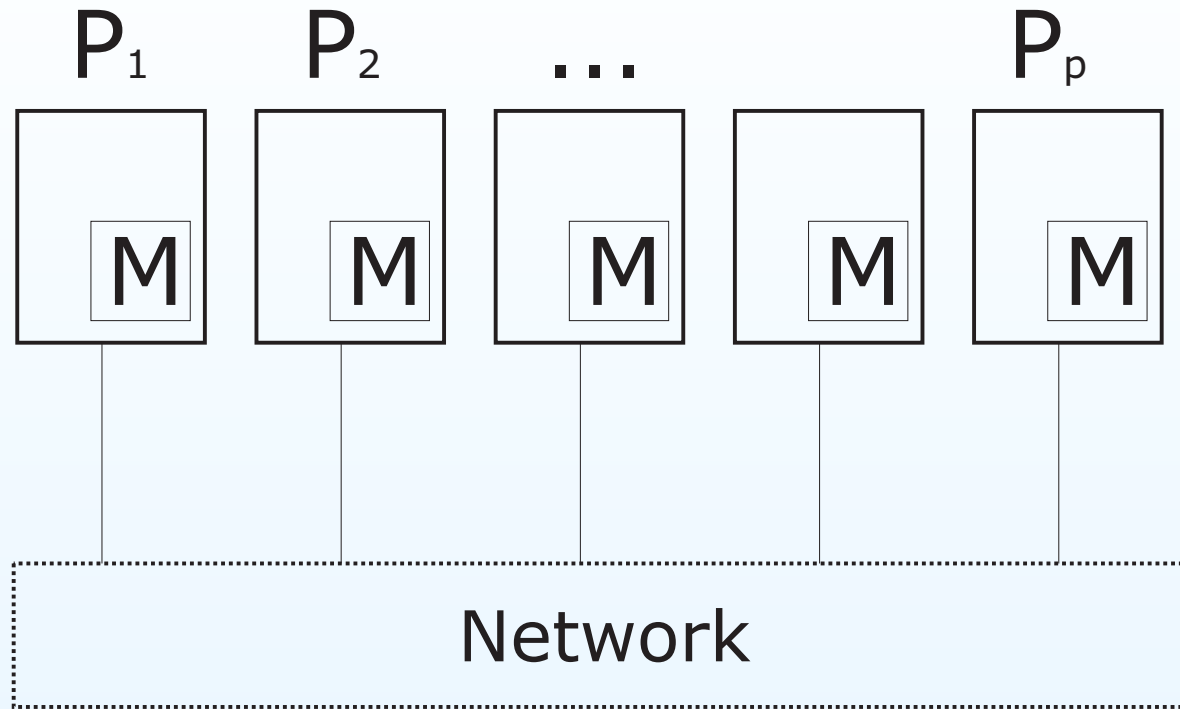
BSP Libraries

Benchmarking

Performance/Predictions

Conclusion

# The BSP Model



- Programs are SPMD
- Execution takes place in *supersteps*
- Cost Formula :  $T = f \cdot W + g \cdot H + l \cdot S$
- As a base unit for communications, 8-byte doubles will be used

Introduction

The BSP Model

• The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

Conclusion

# BSP Programming

Introduction

The BSP Model

BSP Libraries

- **BSP Programming**
- The BSPlib Standard
- BSPlib Implementations
- Other Libraries
- 'BSP-style' Programming in MPI

Benchmarking

Performance/Predictions

Conclusion

## 'BSP-style' programming using a conventional communications library (MPI/Cray shmem/...)

- Barrier synchronizations for creating superstep structure
- Many libraries already provide functionality for one sided communication/direct remote memory access (DRMA)

## Using a specialized library (The Oxford BSP Toolset/PUB/CGMlib/...)

- Specialized communication primitives (bulk synchronous message passing/DRMA)
- Some libraries (Oxford Toolset, PUB) include optimized barrier synchronization functions and routing
- Higher level of abstraction

# The BSPlib Standard

## Communication primitives:

- DRMA: buffered and unbuffered `put`, `get`
- BSMP: `send` and `move`
- Synchronization
- Combining and Broadcasting

For the experiments, a BSPlib-style wrapper library was created.

Introduction

The BSP Model

BSP Libraries

- BSP Programming
- **The BSPlib Standard**
- BSPlib Implementations
- Other Libraries
- 'BSP-style' Programming in MPI

Benchmarking

Performance/Predictions

Conclusion

# BSPLib Implementations

## The Oxford BSP Toolset

- Supports 3 kinds of base architecture: message passing, shared memory, DRMA
- Experiments used message passing MPI interface
- Last release from '98, compatibility issues on more modern systems

## PUB

- Support for message passing and shared memory architectures
- Experiments used message passing MPI interface
- Additional support for oblivious synchronization, processor groups
- Less trouble with setup on all systems
- Advanced functionality e.g. for process migration

Introduction

The BSP Model

BSP Libraries

- BSP Programming
- The BSPLib Standard
- **BSPLib Implementations**
- Other Libraries
- 'BSP-style' Programming in MPI

Benchmarking

Performance/Predictions

Conclusion

# Other Libraries

## CGMlib

- Runs on top of message passing MPI
- Includes set of algorithms for sorting, list ranking, etc.
- Abstract C++ interface
- Lists of abstract datatypes with constant size are used for data exchange

## SSCRAP

- Uses MPI (message passing) or Posix (SHMEM) for data exchange
- Support for DRMA, BSMP, conventional message passing, collective operations, etc.
- 'Soft' synchronization (send or receive)
- C++ interface

Introduction

The BSP Model

BSP Libraries

- BSP Programming
- The BSPLib Standard
- BSPLib Implementations
- **Other Libraries**
- 'BSP-style' Programming in MPI

Benchmarking

Performance/Predictions

Conclusion



# 'BSP-style' Programming in MPI

Approach here: a BSPlib style MPI-1 library was implemented naively (without message combining, etc.)

- Isend/Recv for data exchange
- Barrier synchronization
- Emulated DRMA on top

Advantage: no overhead for `send/put`

Drawback: high latency, presumably overhead for `get` operations

Introduction

The BSP Model

BSP Libraries

- BSP Programming
- The BSPlib Standard
- BSPlib Implementations
- Other Libraries
- 'BSP-style' Programming in MPI

Benchmarking

Performance/Predictions

Conclusion

# Systems used

Measurements on parallel machines at the Centre for Scientific Computing:

`aracari`: IBM cluster,  $64 \times$  2-way SMP Pentium3  
1.4 GHz/128 GB of memory  
(Interconnection Network: Myrinet 2000,  
MPI: `mpich-gm`)

`argus`: Linux cluster,  $31 \times$  2-way SMP Pentium4  
Xeon 2.6 GHz processors/62 GB of memory  
(Interconnection Network: 100Mbit Ethernet,  
MPI: `mpich-p4`)

Introduction

The BSP Model

BSP Libraries

Benchmarking

- **Systems used**
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface  
(`aracari`)
- Bandwidth Surface  
(`argus`)
- Latency
- Benchmark Summary

Performance/Predictions

Conclusion

# Measuring $f$

Introduction

The BSP Model

BSP Libraries

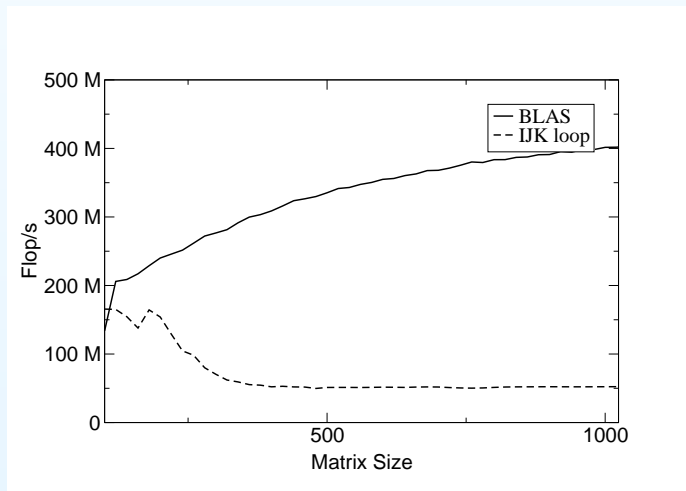
Benchmarking

- Systems used
- **Measuring  $f$**
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- Bandwidth Surface (argus)
- Latency
- Benchmark Summary

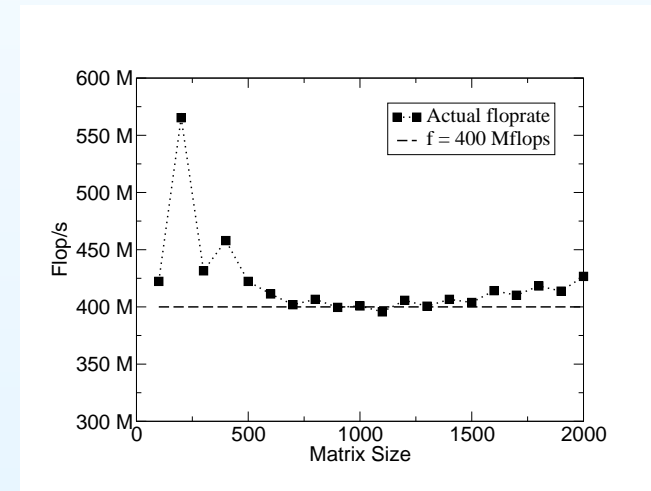
Performance/Predictions

Conclusion

Measuring algorithm performance on one node:



Measuring computation time separately in one run:



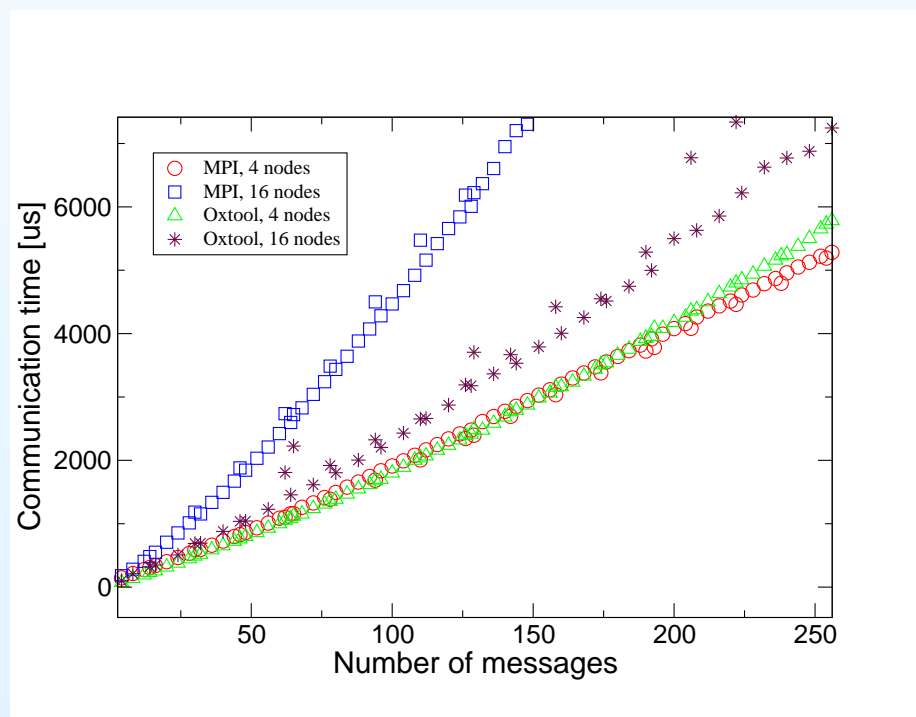
(Example for Matrix-Matrix multiplication)

# Measuring $g$ and $l$

Problems encountered: realistic values of  $g$  and  $l$  depend on

- The number of processors that are used
- The communications pattern
- The communication volume

E.g. for all-to-all communication on `aracari`



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (`aracari`)
- Bandwidth Surface (`argus`)
- Latency
- Benchmark Summary

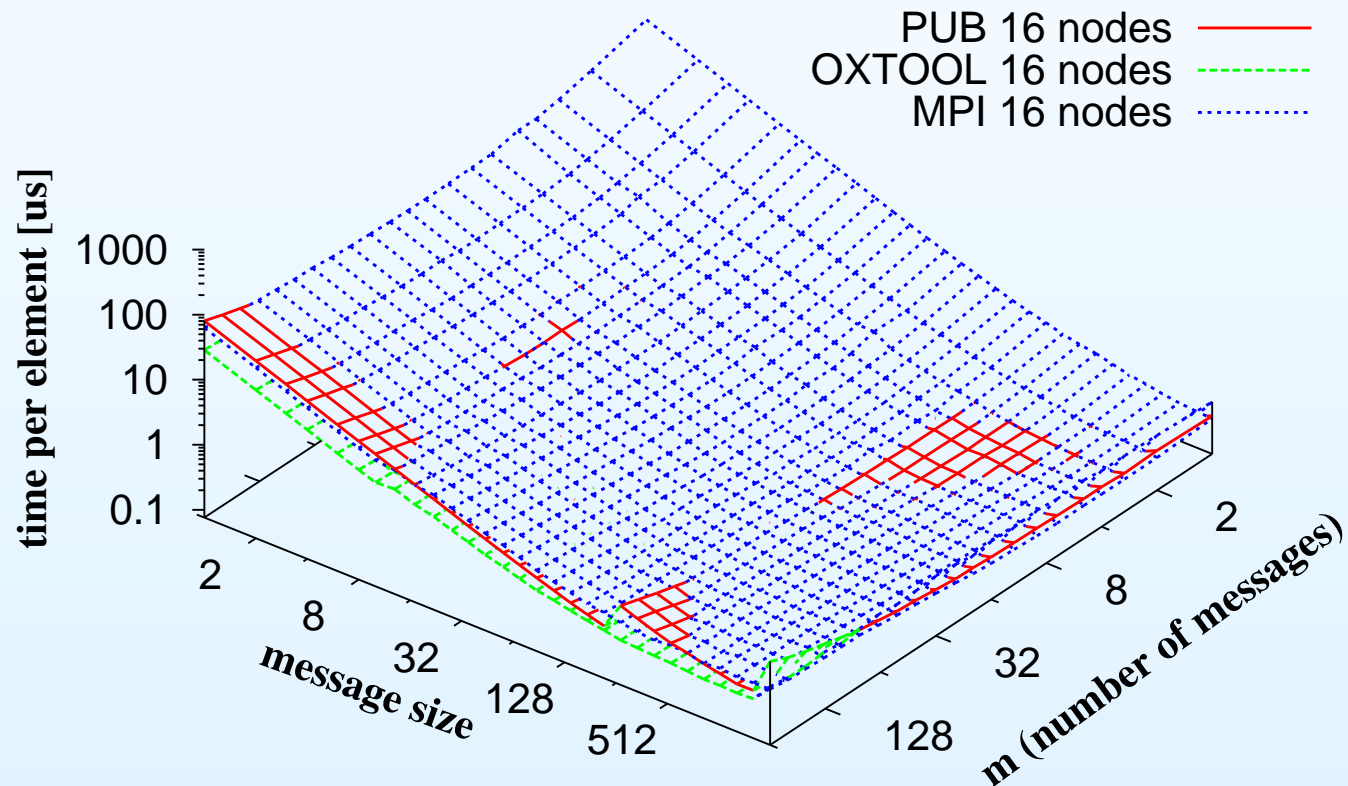
Performance/Predictions

Conclusion

# Bandwidth Surface (aracari)

For a better picture, the effective bandwidth can be sampled depending on message size and count.

All-to-all communication on `aracari`:



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- **Bandwidth Surface (aracari)**
- Bandwidth Surface (argus)
- Latency
- Benchmark Summary

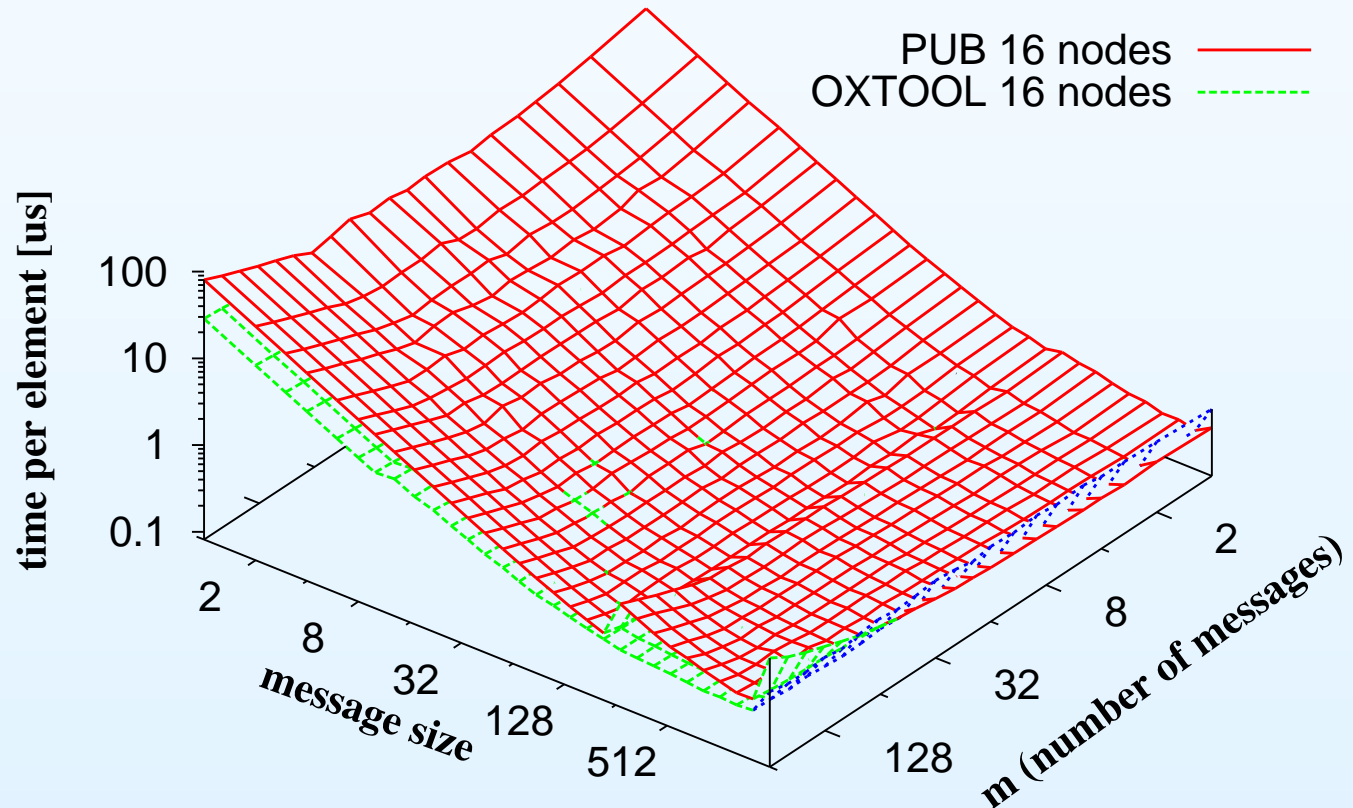
Performance/Predictions

Conclusion

# Bandwidth Surface (aracari)

For a better picture, the effective bandwidth can be sampled depending on message size and count.

All-to-all communication on `aracari`:



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- Bandwidth Surface (argus)
- Latency
- Benchmark Summary

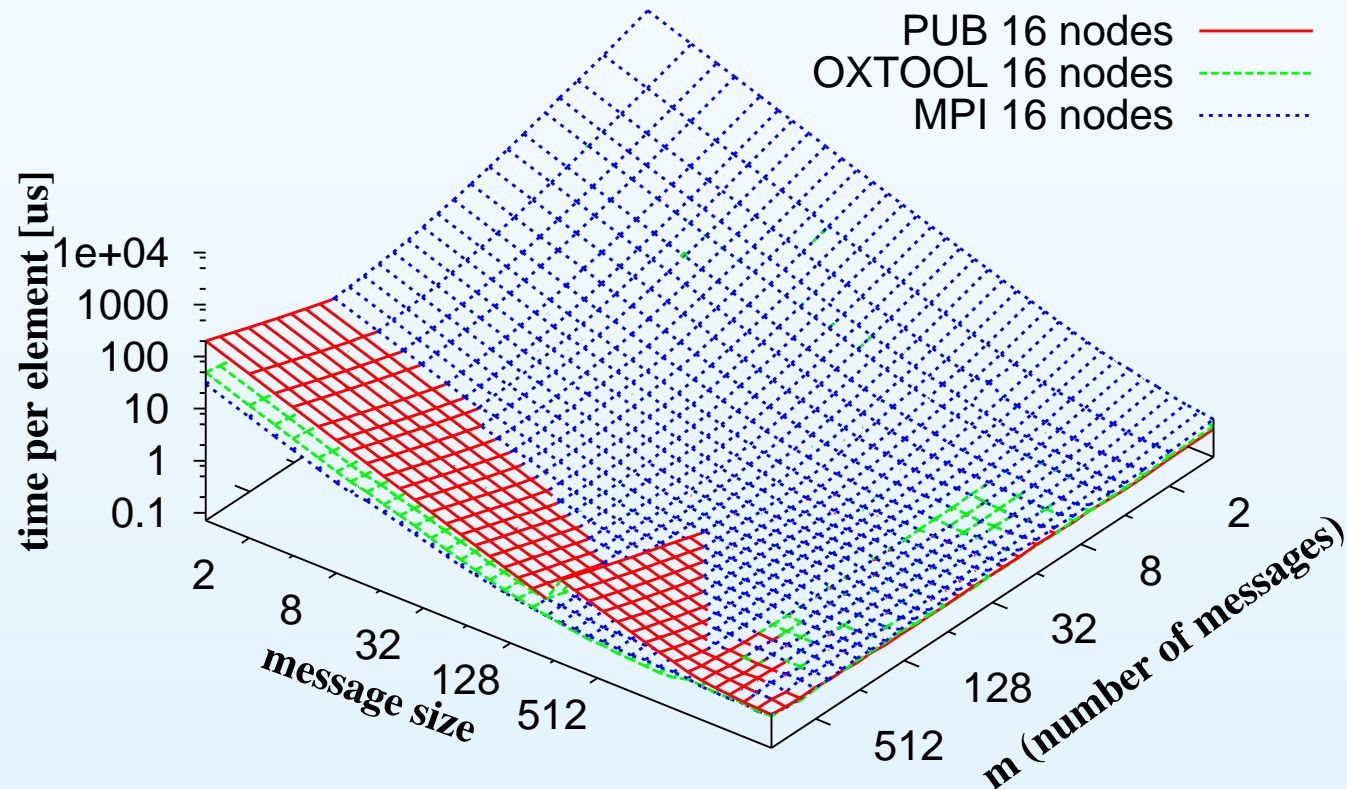
Performance/Predictions

Conclusion

# Bandwidth Surface (aracari)

For a better picture, the effective bandwidth can be sampled depending on message size and count.

Random permutation on `aracari`:



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- **Bandwidth Surface (aracari)**
- Bandwidth Surface (argus)
- Latency
- Benchmark Summary

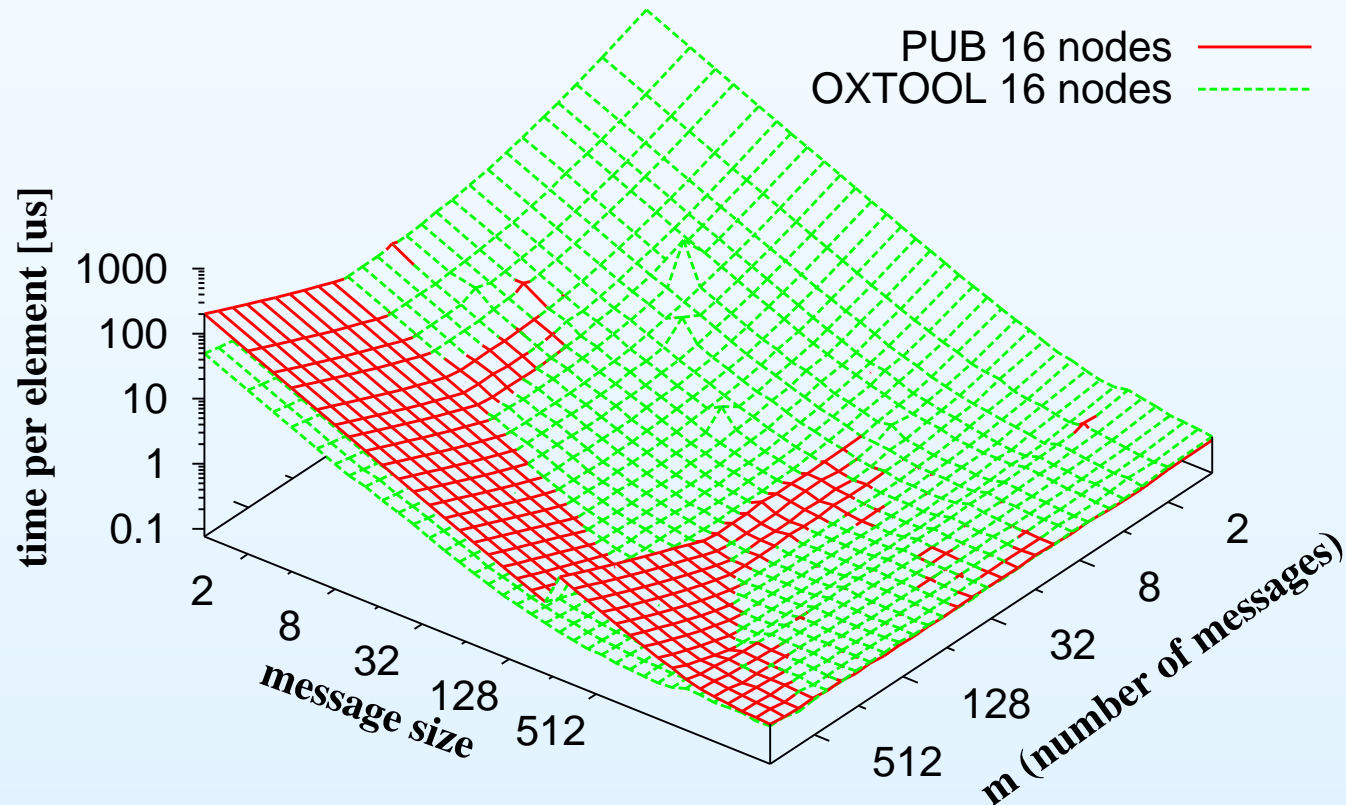
Performance/Predictions

Conclusion

# Bandwidth Surface (aracari)

For a better picture, the effective bandwidth can be sampled depending on message size and count.

Random permutation on aracari:



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- Bandwidth Surface (argus)
- Latency
- Benchmark Summary

Performance/Predictions

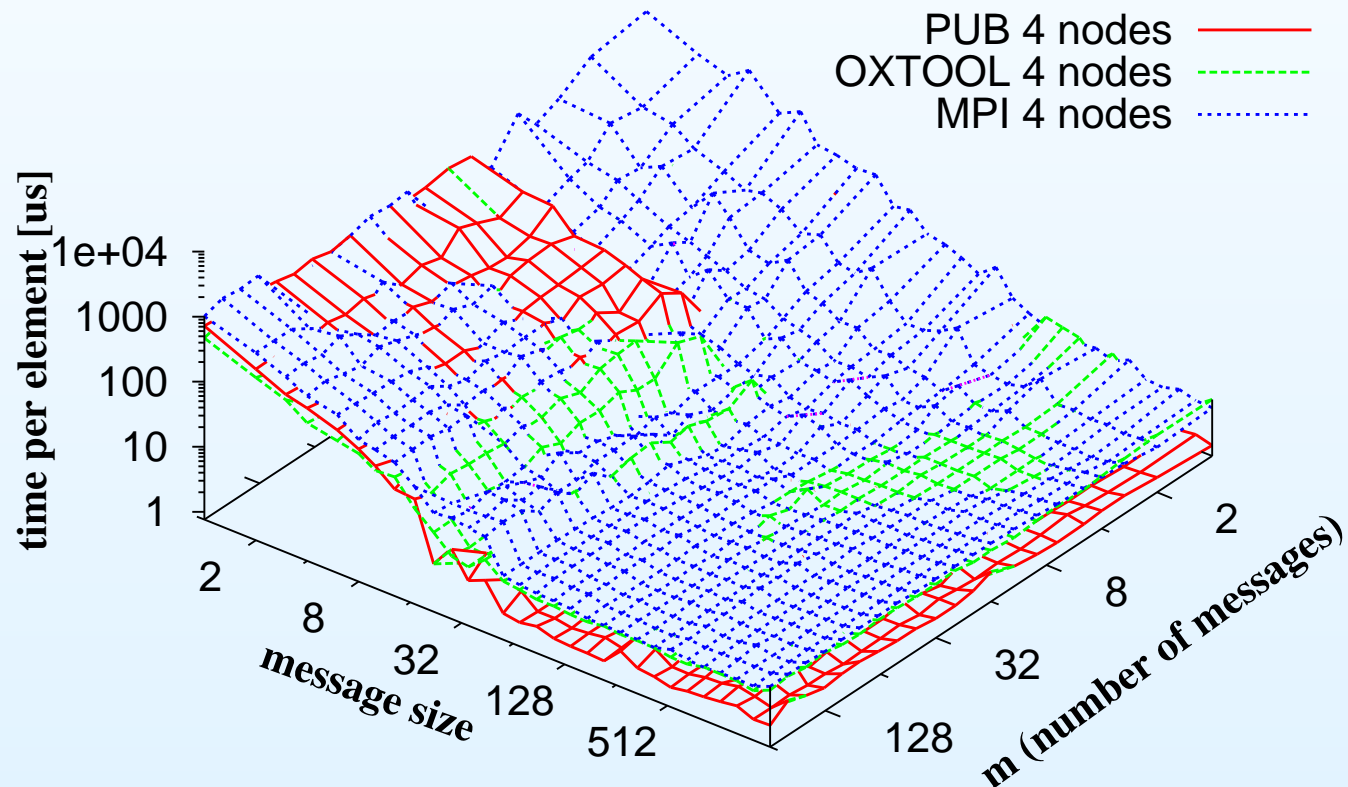
Conclusion



# Bandwidth Surface (argus)

The picture looks different on the slower communications network

All-to-all communication on `argus`:



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- **Bandwidth Surface (argus)**
- Latency
- Benchmark Summary

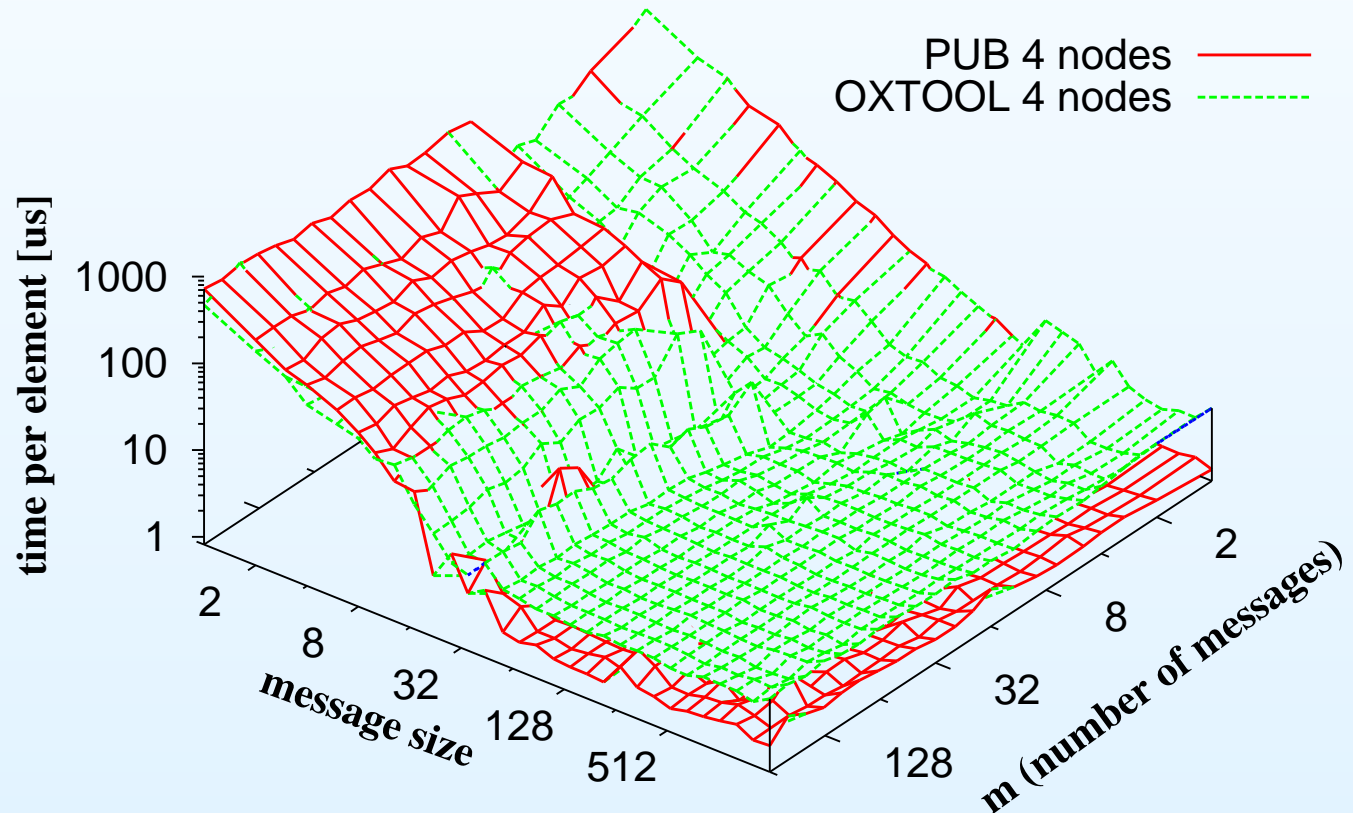
Performance/Predictions

Conclusion

# Bandwidth Surface (argus)

The picture looks different on the slower communications network

All-to-all communication on `argus`:



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- **Bandwidth Surface (argus)**
- Latency
- Benchmark Summary

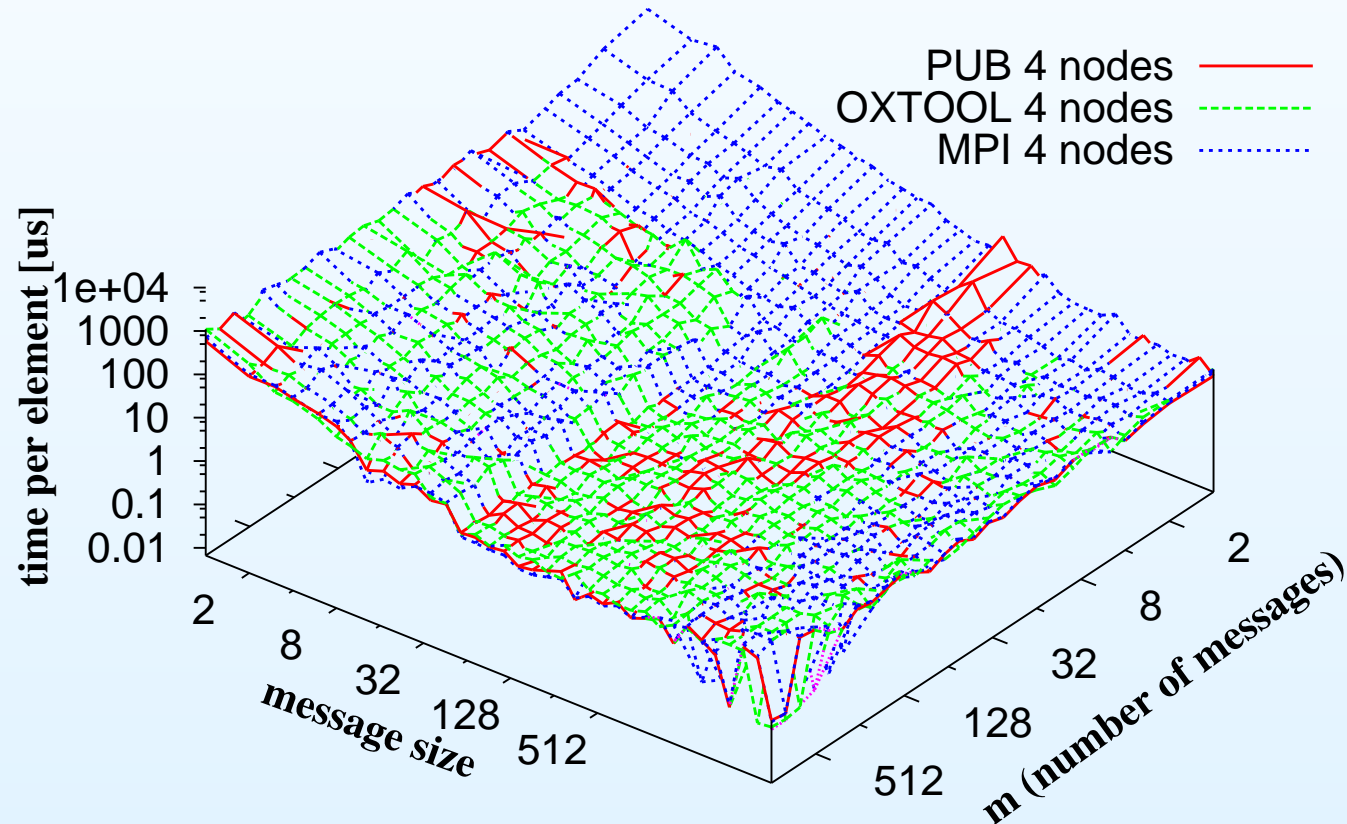
Performance/Predictions

Conclusion

# Bandwidth Surface (argus)

The picture looks different on the slower communications network

Random permutation on argus:



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- **Bandwidth Surface (argus)**
- Latency
- Benchmark Summary

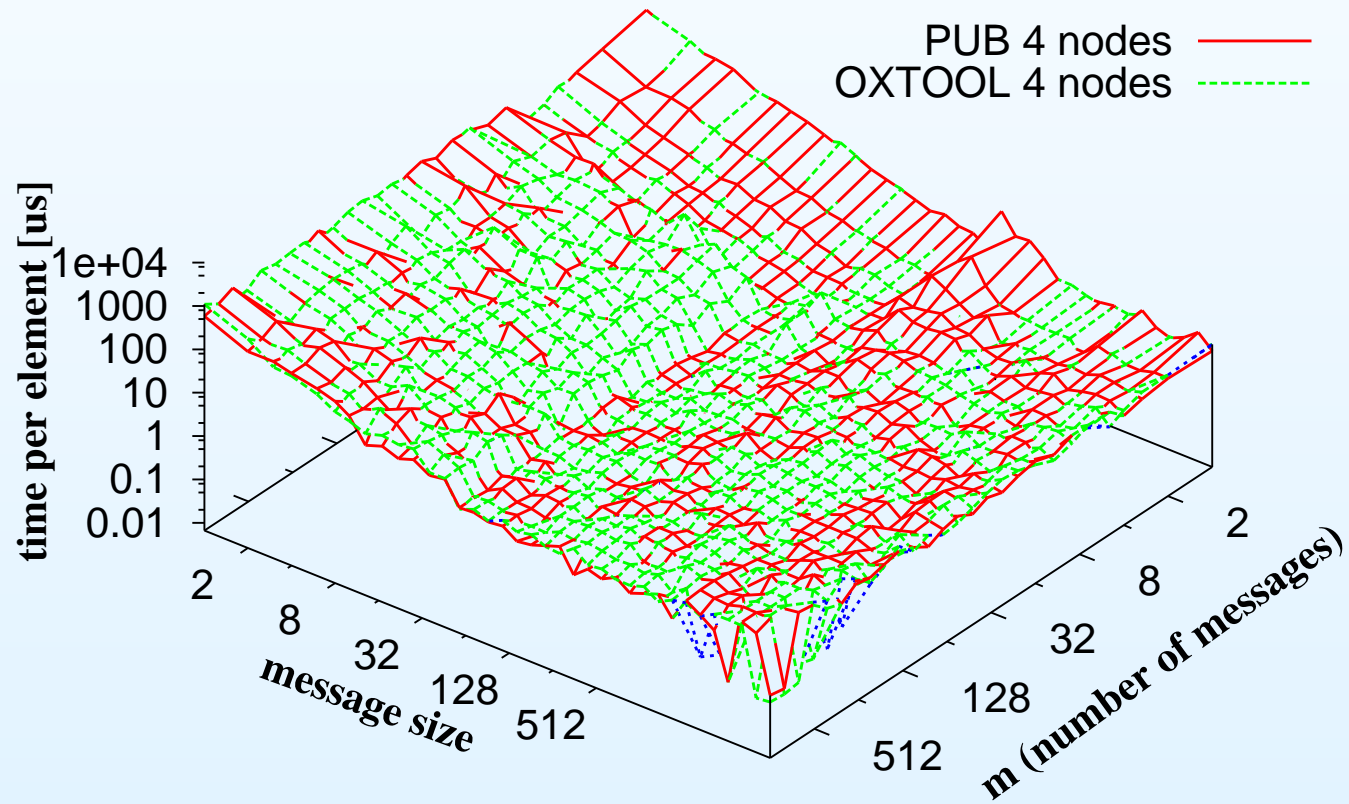
Performance/Predictions

Conclusion

# Bandwidth Surface (argus)

The picture looks different on the slower communications network

Random permutation on argus:



Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- **Bandwidth Surface (argus)**
- Latency
- Benchmark Summary

Performance/Predictions

Conclusion

# Latency

The latency can be measured for synchronizations preceded by different types of communication:

	<b>MPI</b>	<b>Oxtool</b>	<b>PUB</b>
<hr/>			
<i>aracari 4 processors</i>			
<i>l</i> (low)	210 $\mu\text{s}$	43 $\mu\text{s}$	39 $\mu\text{s}$
<i>l</i> (high)	230 $\mu\text{s}$	67 $\mu\text{s}$	55 $\mu\text{s}$
<i>l</i> (all-to-all)	252 $\mu\text{s}$	89 $\mu\text{s}$	72 $\mu\text{s}$
<hr/>			
<i>aracari 32 processors</i>			
<i>l</i> (low)	2203 $\mu\text{s}$	621 $\mu\text{s}$	142 $\mu\text{s}$
<i>l</i> (high)	2242 $\mu\text{s}$	638 $\mu\text{s}$	163 $\mu\text{s}$
<i>l</i> (all-to-all)	2881 $\mu\text{s}$	1250 $\mu\text{s}$	750 $\mu\text{s}$
<hr/>			
<i>argus 4 processors</i>			
<i>l</i> (low)	5642 $\mu\text{s}$	796 $\mu\text{s}$	975 $\mu\text{s}$
<i>l</i> (high)	5789 $\mu\text{s}$	1442 $\mu\text{s}$	1176 $\mu\text{s}$
<i>l</i> (all-to-all)	5086 $\mu\text{s}$	1613 $\mu\text{s}$	871 $\mu\text{s}$
<hr/>			

Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (*aracari*)
- Bandwidth Surface (*argus*)
- **Latency**
- Benchmark Summary

Performance/Predictions

Conclusion

# Benchmark Summary

Bandwidth depends on message count, message size and the communications pattern

On `aracari`:

- Best all-to-all performance: Oxtool
- Best random permutation performance (few messages): PUB,  $> 64$  messages: Oxtool
- Best self communication performance (few messages): PUB,  $> 32$  messages: Oxtool
- MPI: good performance when message size is large

Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (`aracari`)
- Bandwidth Surface (`argus`)
- Latency
- **Benchmark Summary**

Performance/Predictions

Conclusion

# Benchmark Summary

Bandwidth depends on message count, message size and the communications pattern

On argus:

- Best all-to-all performance: PUB
- Random permutation: little difference between PUB and Oxtool
- Best self communication performance (few messages): Oxtool
- MPI: good performance when message size and count are larger than 16/32 doubles

Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- Bandwidth Surface (argus)
- Latency
- **Benchmark Summary**

Performance/Predictions

Conclusion

# Benchmark Summary

## Latency:

- PUB consistently has best latency (without using the faster ‘oblivious’ synchronization)
- As expected, ‘naive’ MPI library has highest latency

Introduction

The BSP Model

BSP Libraries

Benchmarking

- Systems used
- Measuring  $f$
- Measuring  $g$  and  $l$
- Bandwidth Surface (aracari)
- Bandwidth Surface (argus)
- Latency
- **Benchmark Summary**

Performance/Predictions

Conclusion

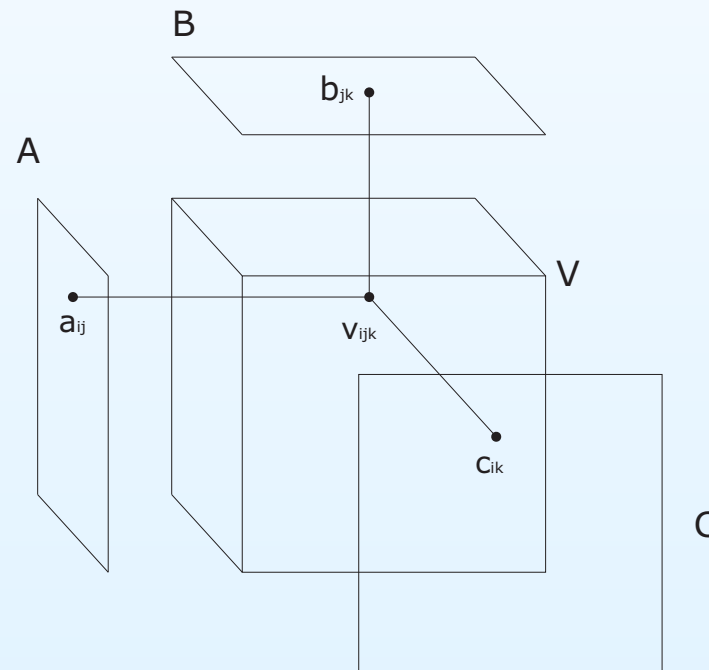


# BSP Matrix-Matrix Multiplication

We want to compute the product of two dense  $n \times n$  matrices  $A$  and  $B$

Simple formula:

$$c_{ik} = \sum_{j=1}^n a_{ij} b_{jk}, \text{ having } A = [a_{ij}], \quad B = [b_{ij}], \quad C = [c_{ij}]$$



Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

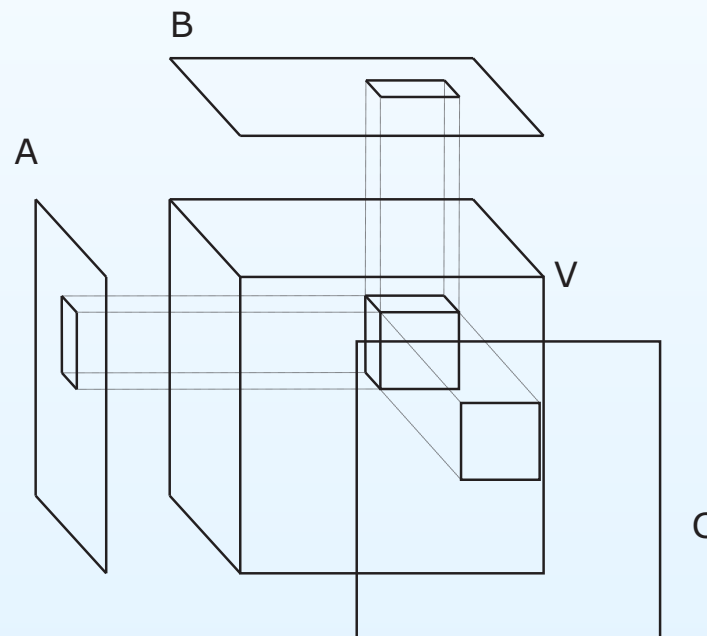
- **BSP Matrix-Matrix Multiplication**
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on `aracari`
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion

# BSP Matrix-Matrix Multiplication (2)

Block decomposition into  $q$  blocks for memory efficient parallel algorithm:

$$C_{IK} = \sum_{J=1}^{q^{1/3}} V_{IJK} \quad \text{with } I, K = 1, 2, \dots, q^{1/3}$$



Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

- BSP Matrix-Matrix Multiplication
- **BSP Matrix-Matrix Multiplication (2)**
- Why this algorithm?
- Prediction model
- Prediction results on aracari
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion

# Why this algorithm?

- Communication block size can be controlled by parameter  $q$
- Message combining has to be used when using fixed initial data distribution (block-cyclic with block width  $n/\sqrt{p}$ )
- Can be compared e.g. to PBLAS
- ‘Nice’ version can predistribute the blocks before the computation to avoid spikes because of data distribution

Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- **Why this algorithm?**
- Prediction model
- Prediction results on `aracari`
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

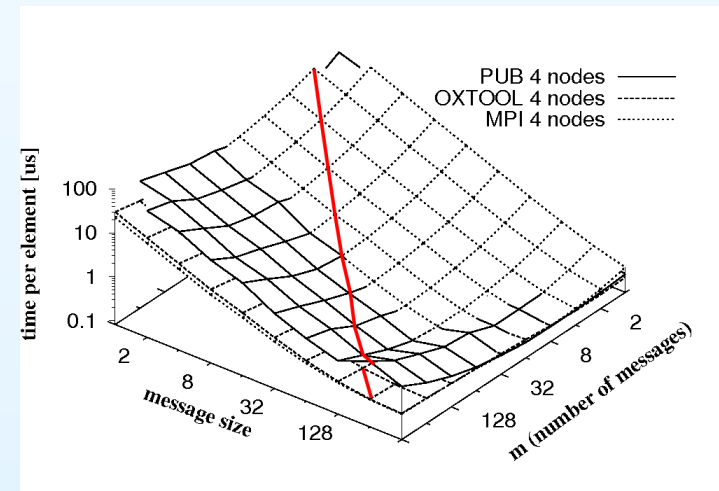
Conclusion

# Prediction model

BSP running time:

$$T = f \cdot \left[ \frac{q}{p} \right] \cdot \frac{n^3}{q} + g \cdot \left[ \frac{q}{p} \right] \cdot \frac{n^2}{q^{2/3}} \cdot \left( 2 + 1/q^{1/3} \right) + l \cdot 2 \left[ \frac{q}{p} \right]$$

Two matrices are transferred row by row  
→ value of  $g$  is taken from the red line as value for maximum matrix size



Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

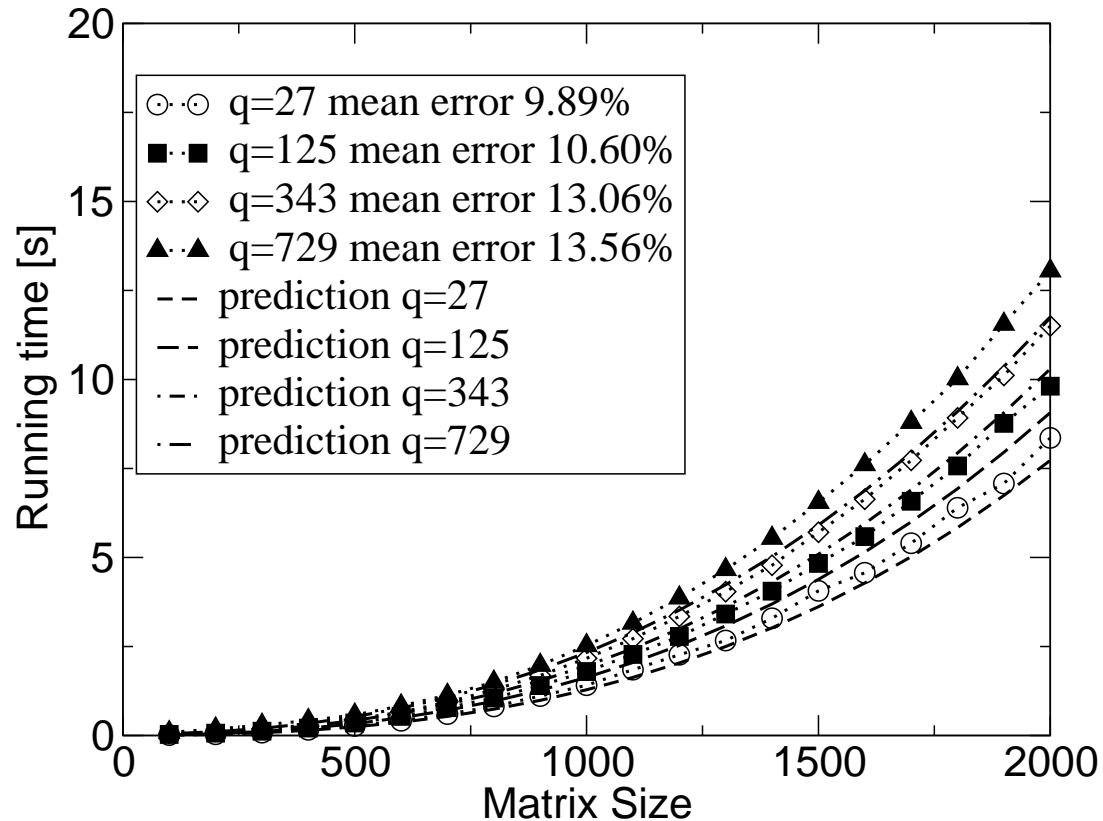
- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- **Prediction model**
- Prediction results on `aracari`
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion

# Prediction results on aracari

## Oxtool, using 4 processors

$$p = 4, 1/f = 4e+08 \quad g = 3.5e-07, l = 8.9e-05$$



Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

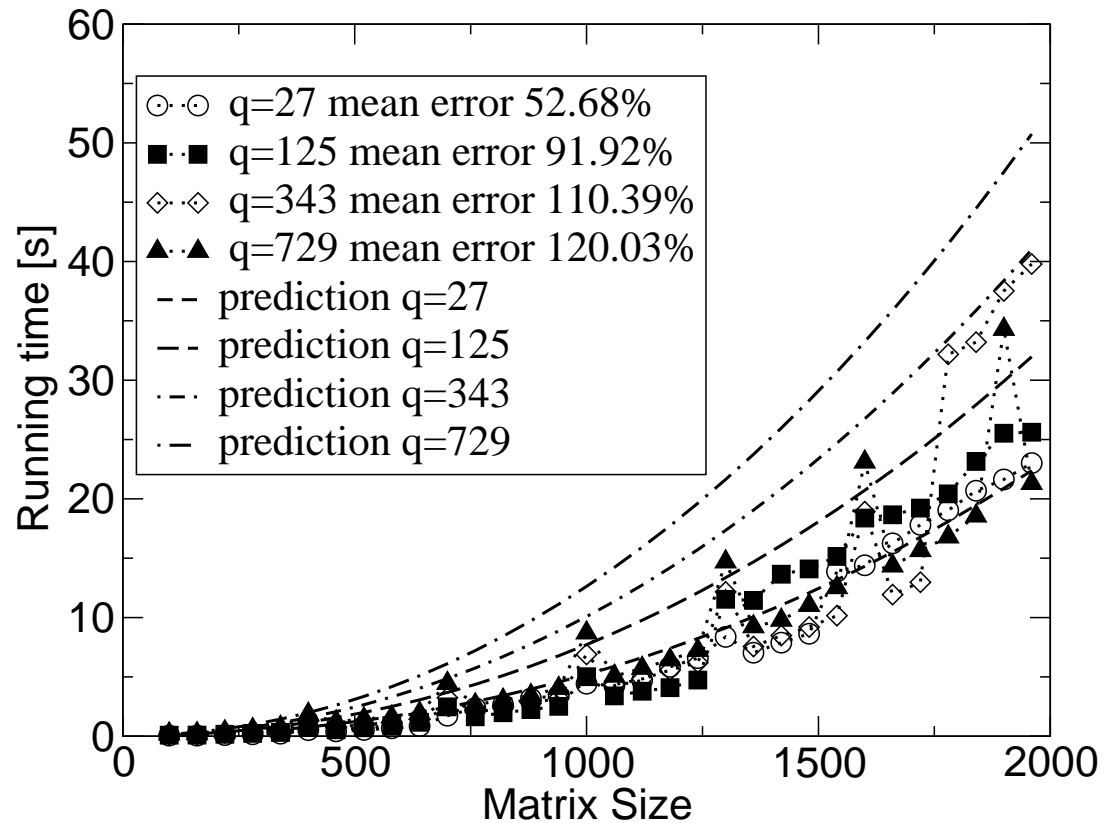
- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on *aracari*
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion

# Prediction results on aracari

## PUB, using 4 processors

$$p = 4, 1/f = 4e+08 \quad g = 2.5e-06, l = 7.2e-05$$



Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

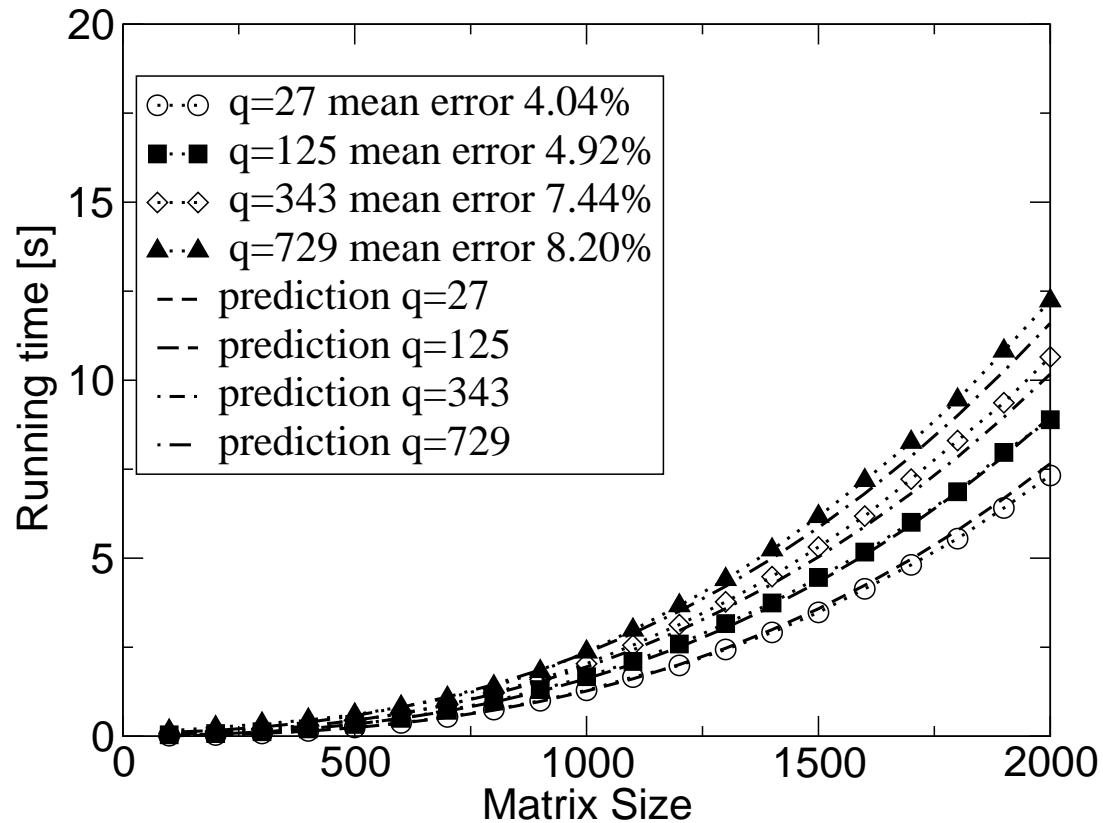
- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on *aracari*
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion

# Prediction results on aracari

## MPI, using 4 processors

$p = 4$ ,  $1/f = 4e+08$   $g = 3.4e-07$ ,  $l = 0.00025$



Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

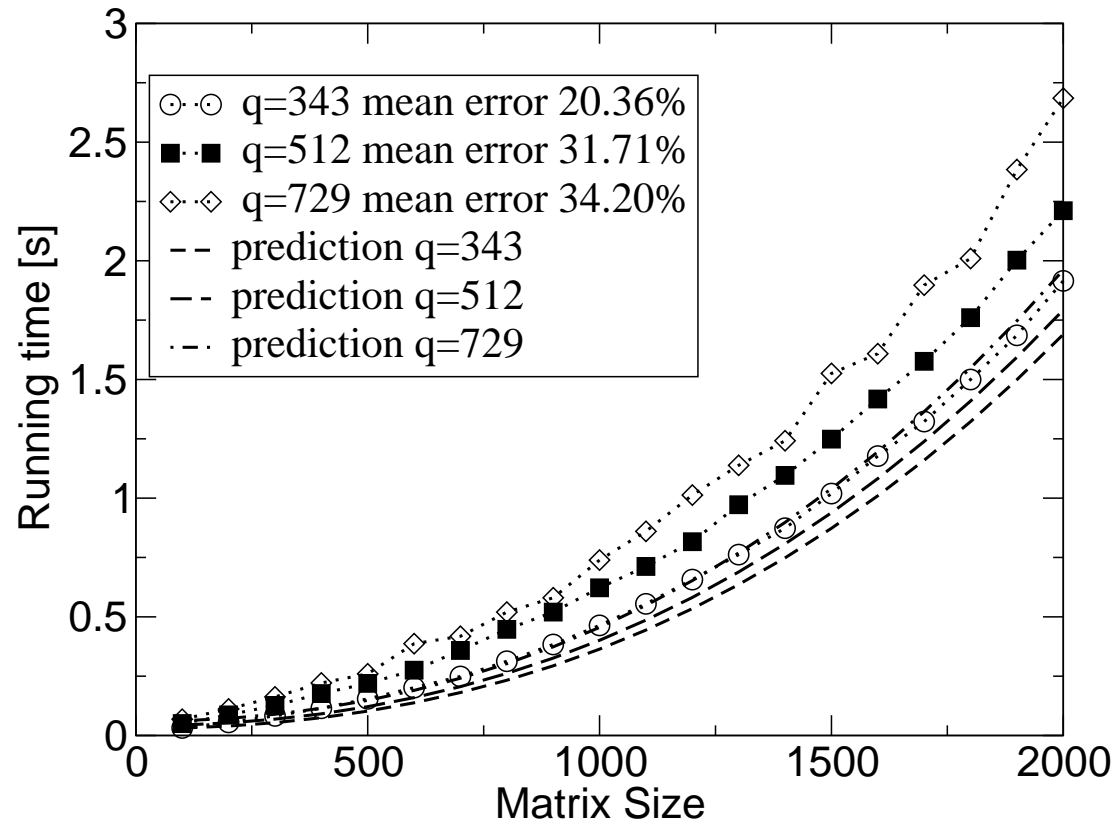
- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on *aracari*
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion

# Prediction results, more processors

## Oxtool, using 32 processors

$p = 32$ ,  $1/f = 4e+08$   $g = 5.3e-07$ ,  $l = 0.0013$



Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

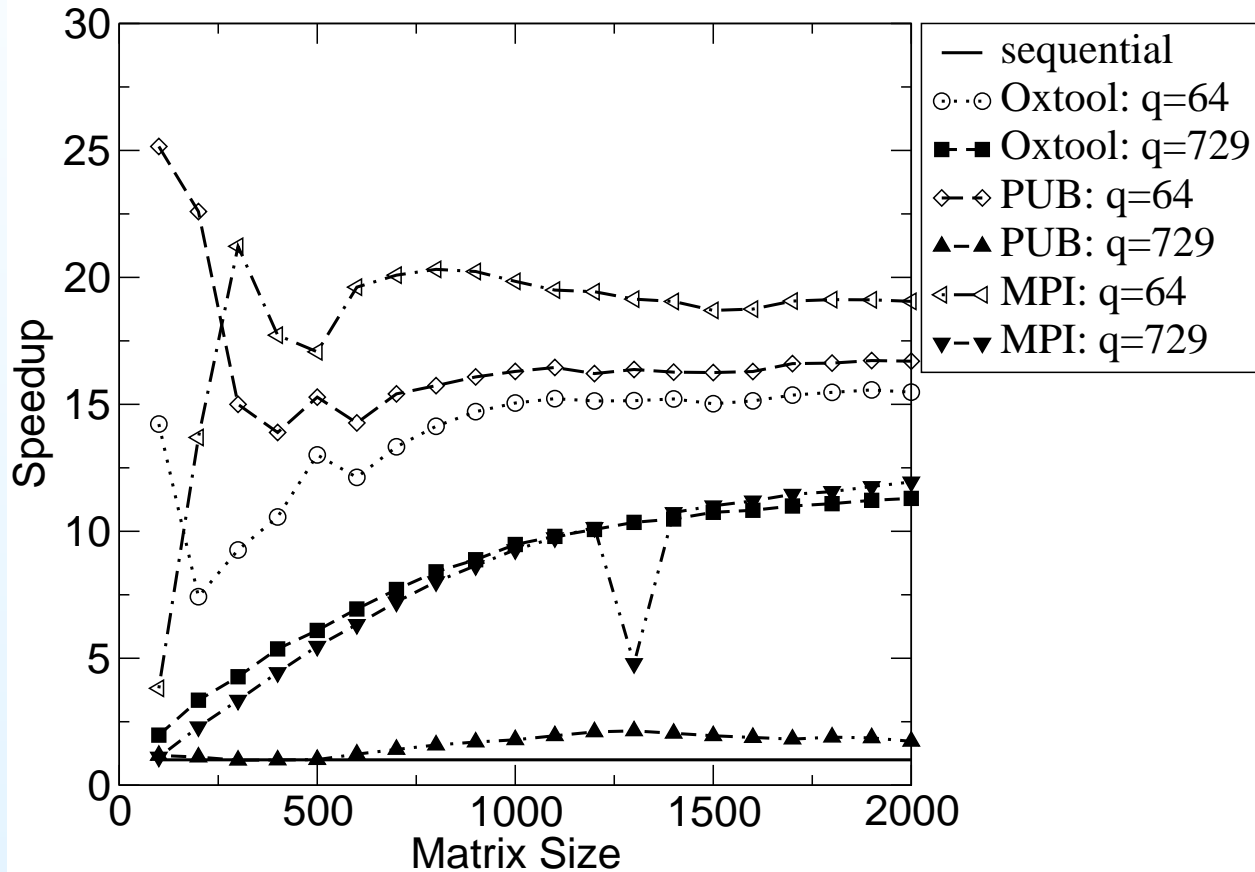
- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on aracari
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion



# Speedup Results (1)

aracari, using 16 processors



Introduction

The BSP Model

BSP Libraries

Benchmarking

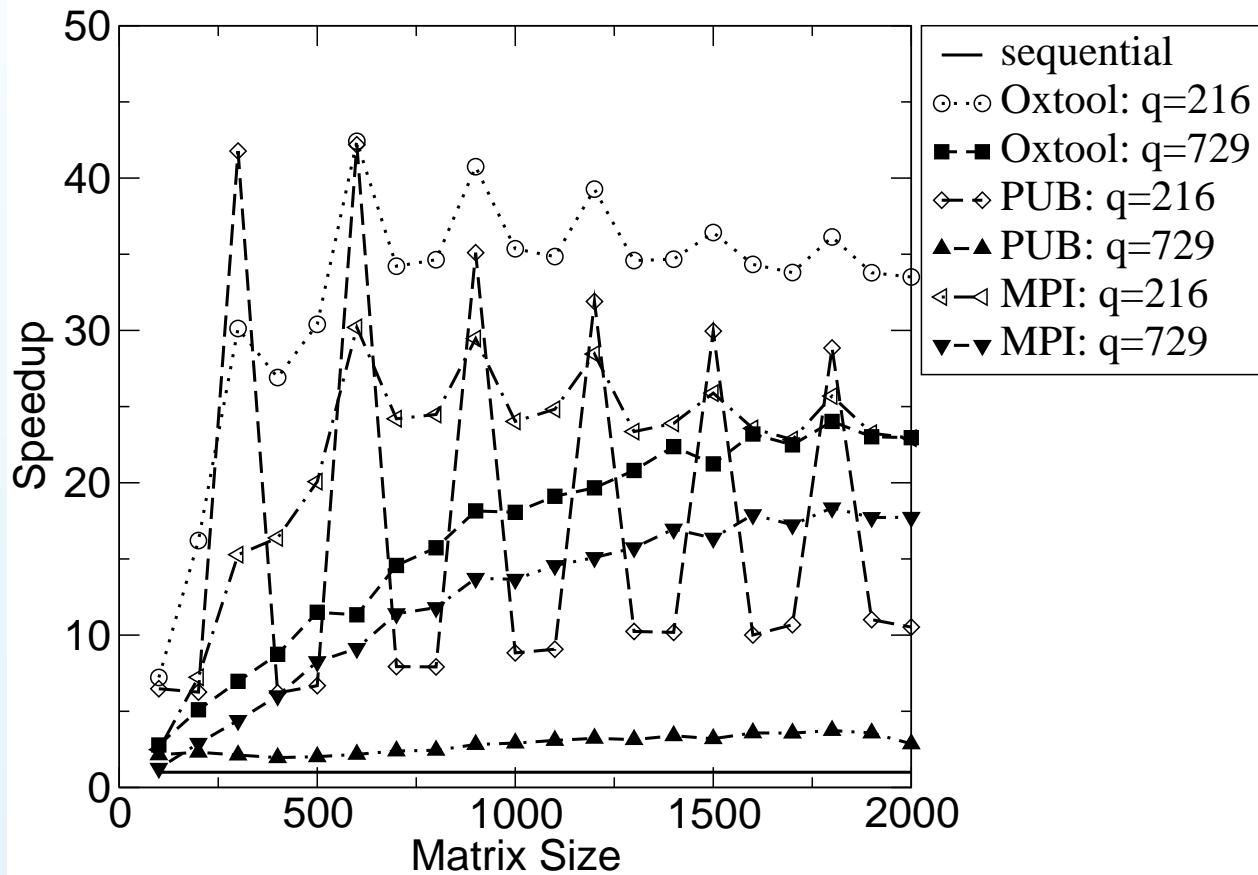
Performance/Predictions

- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on aracari
- Prediction results, more processors
- **Speedup Results (1)**
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion

# Speedup Results (2)

aracari, using 32 processors (spikes when matrix size mod 6 == 0)



Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on aracari
- Prediction results, more processors
- Speedup Results (1)
- **Speedup Results (2)**
- Performance Comparison with PBLAS

Conclusion

# Performance Comparison with PBLAS

Introduction

The BSP Model

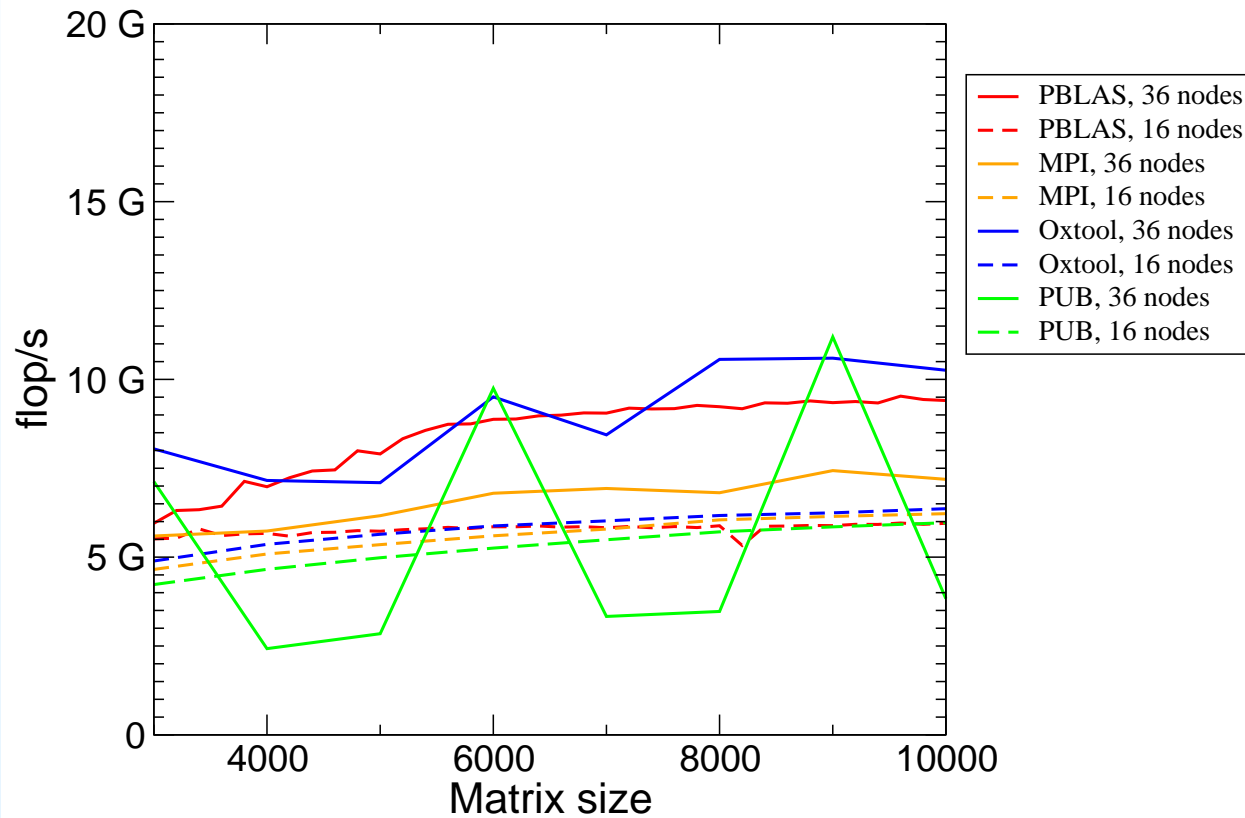
BSP Libraries

Benchmarking

Performance/Predictions

- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on `aracari`
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion



# Performance Comparison with PBLAS

Introduction

The BSP Model

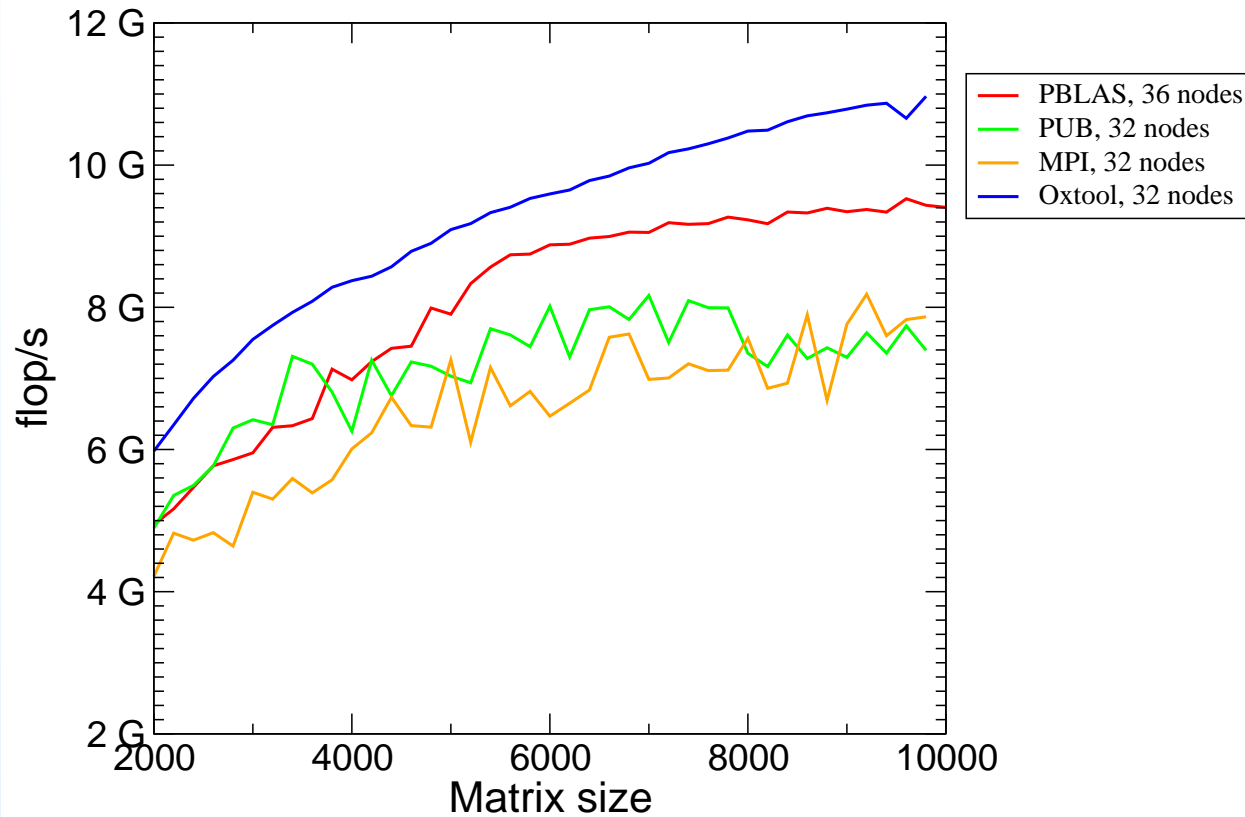
BSP Libraries

Benchmarking

Performance/Predictions

- BSP Matrix-Matrix Multiplication
- BSP Matrix-Matrix Multiplication (2)
- Why this algorithm?
- Prediction model
- Prediction results on `aracari`
- Prediction results, more processors
- Speedup Results (1)
- Speedup Results (2)
- Performance Comparison with PBLAS

Conclusion



# Summary

Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

Conclusion

- Summary
- Further Work

- Despite restrictions due to BSP model, all implementations reach good speedup on `aracari` when number of blocks is low
- Overall benchmark results look better for PUB
- Oxtool has best matrix multiplication performance on `aracari` (Myrinet)
- PUB has best matrix multiplication performance on `argus` (Ethernet)
- Predictably no real speedup on `argus`, due to slow communications network and fast nodes
- Performance of simple BSP algorithm is comparable with PBLAS

## Further Work

- Run experiments on shared memory machine
- Use more different communication patterns for benchmarking
- Study other algorithms with different communication patterns
- Keeping simplicity, extend prediction model for more accuracy

Introduction

The BSP Model

BSP Libraries

Benchmarking

Performance/Predictions

Conclusion

- Summary

- **Further Work**