# A simple, pipelined all-gather algorithm for large irregular problems

Jesper Larsson Träff, Andreas Ripke,
Christian Siebert,
Pavan Balaji, Rajeev Thakur, William Gropp
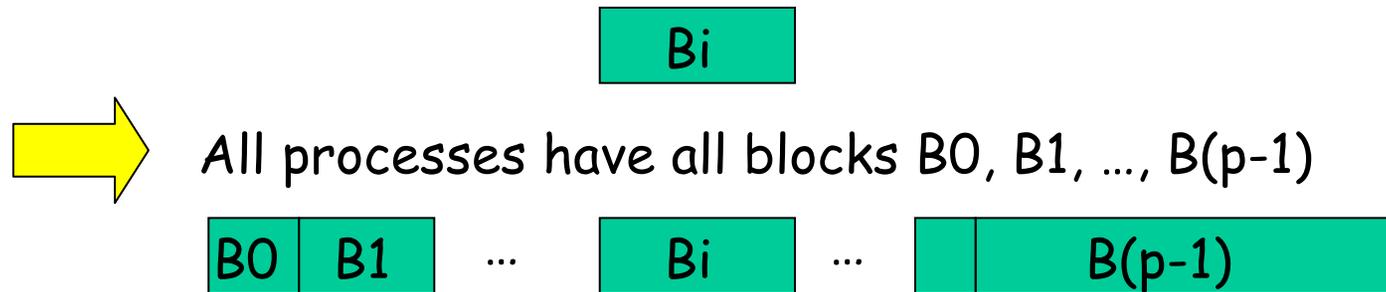NEC Labs Europe, ANL, UIUC

NEC

# An irregular all-gather data-exchange operation

- A set of processes, 0, …, p-1
- Each has a block of data $B_i$ of (possibly) different size



All processes have all blocks B0, B1, …, B(p-1)



MPI_Allgatherv(sendbuf,…,recvbuf,…,counts,…);

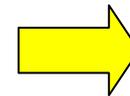counts[i] ≈ |Bi|, all processes know the size of all blocks!

NEC

MPI_Allgatherv used in numerical libraries eg. PETSc
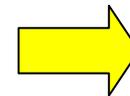
•Other irregular collectives, eg. MPI_Alltoallw

Irregular collectives algorithmically (much!) more difficult – challenging - than regular collectives:

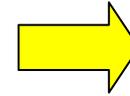1. Different amounts of data between processes (in different rounds)    ⟹ Load imbalance

2. Partial information for the processes (MPI_Alltoallw, MPI_Gatherv, …)    ⟹ Difficult/expensive to compute schedule

3. Schedules    ⟹ Optimality is (NP)-hard!

NEC

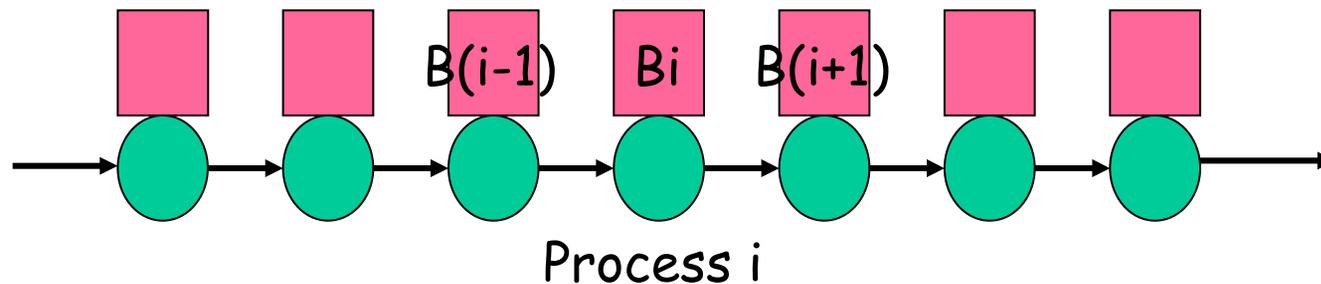**Related work:** **MPI_Allgather** (gossip, broadcast-to-all, …)

- Bruck et al. '97: simultaneous binomial tree algorithm:

- Benson et al. '03: MPI implementation of allgather algorithms on switched networks [non SMP-aware]

- Thakur et al. '04: mpich2 implementations of Bruck and other [non SMP-aware]

- Mamidala et al. '06: SMP implementations

- Träff'06: Graceful degradation from Bruck to linear ring [SMP-aware]

**MPI_Allgatherv**

- Balaji et al '07: optimization in the context of PETSc

# An algorithm for large, regular all-gather problems



Process i

Linear ring: p-1 rounds

- Process i receives block B(i-1-k) and sends block B(i-k) in round k, k=0, …, p-1

NEC

**Analysis:**

Size per process $m_i = |B_i| = m'$ (for regular problem)

Total problem size $m = \sum m_i$

- p-1 (regular) communication rounds
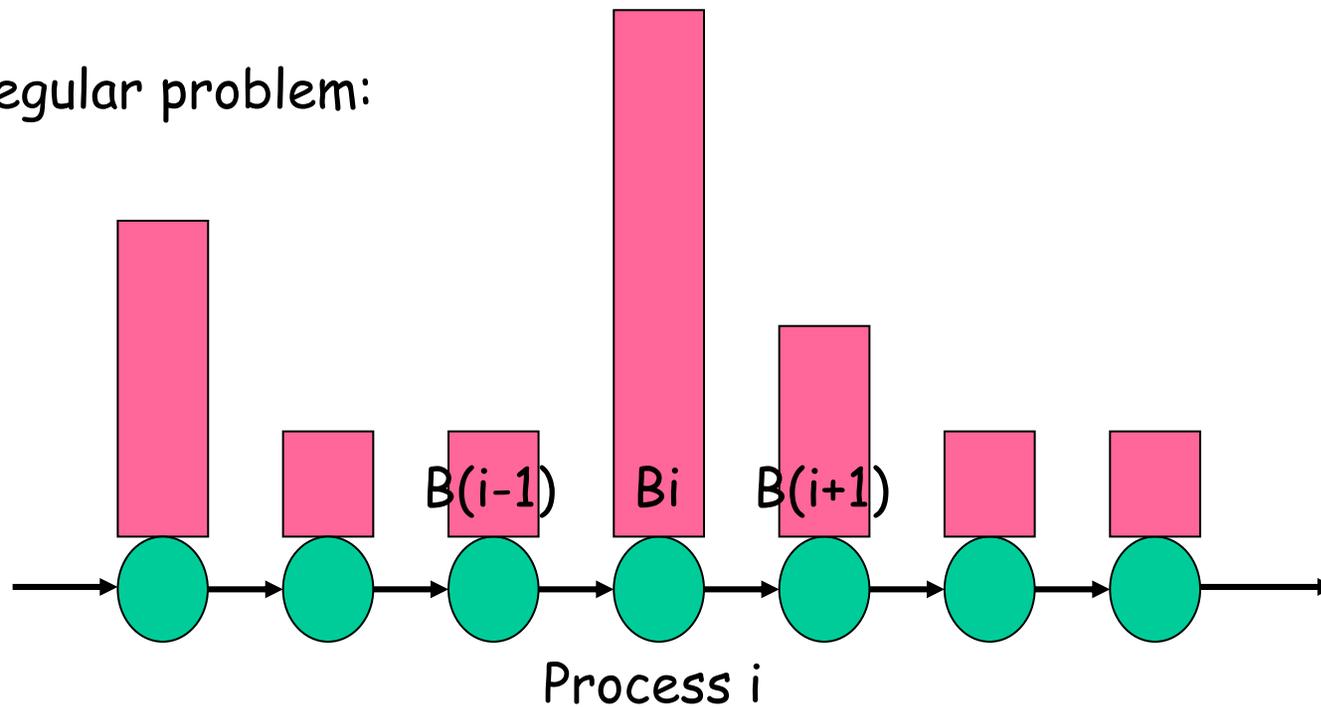- Each round takes time $O(m_i)$, total $O((p-1)m_i) = O(m-m')$

Assumptions:

- Processes can send and receive simultaneously
- Cost of sending/receiving data of size $m'$ is $O(m')$
- Homogeneous communication along ring

NOTE:

For small $m'$ algorithm with log p rounds preferable!

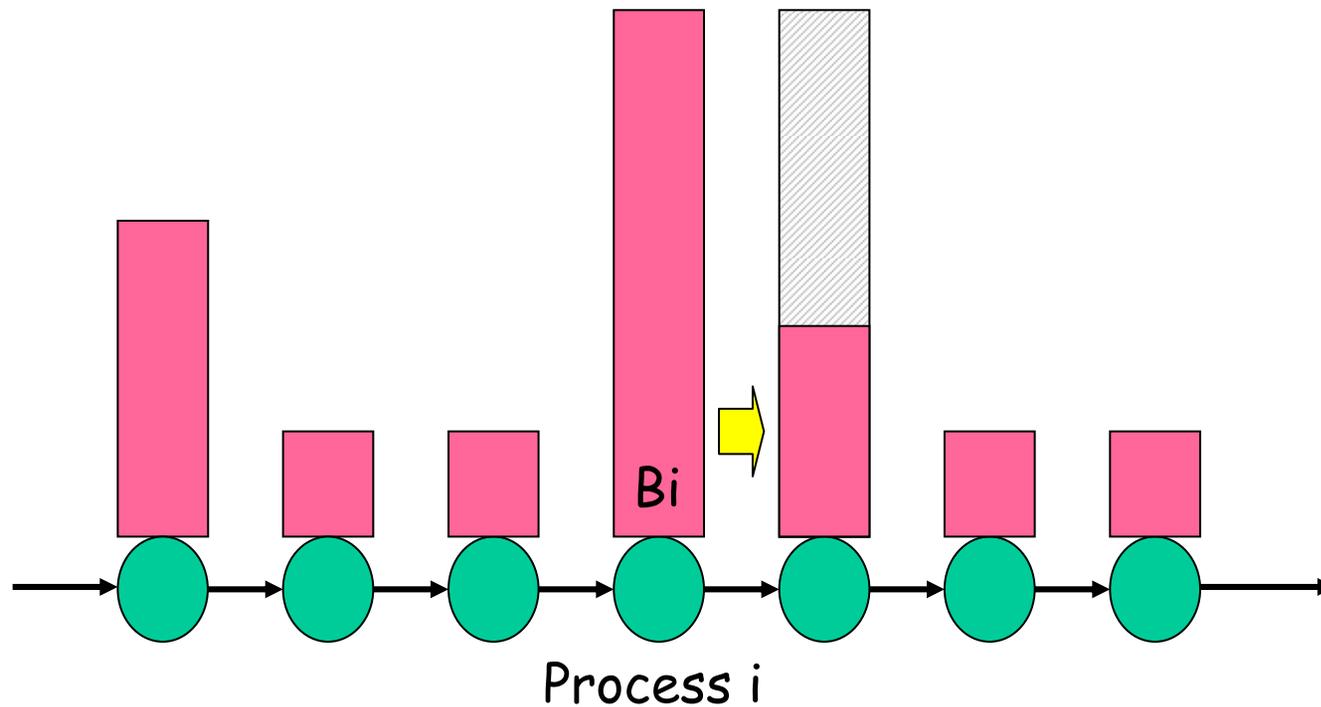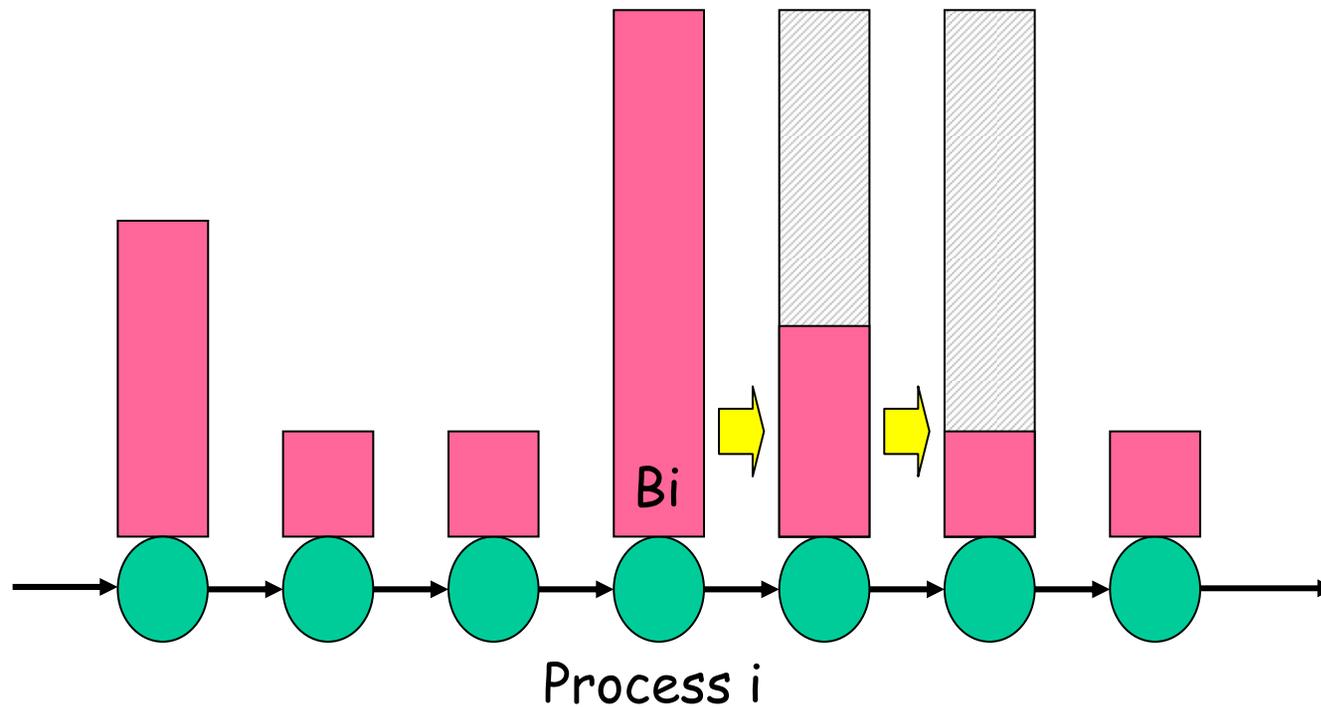Optimal: no idle time, no superfluous data!

NEC

Irregular problem:



B(i-1)    Bi    B(i+1)

Process i
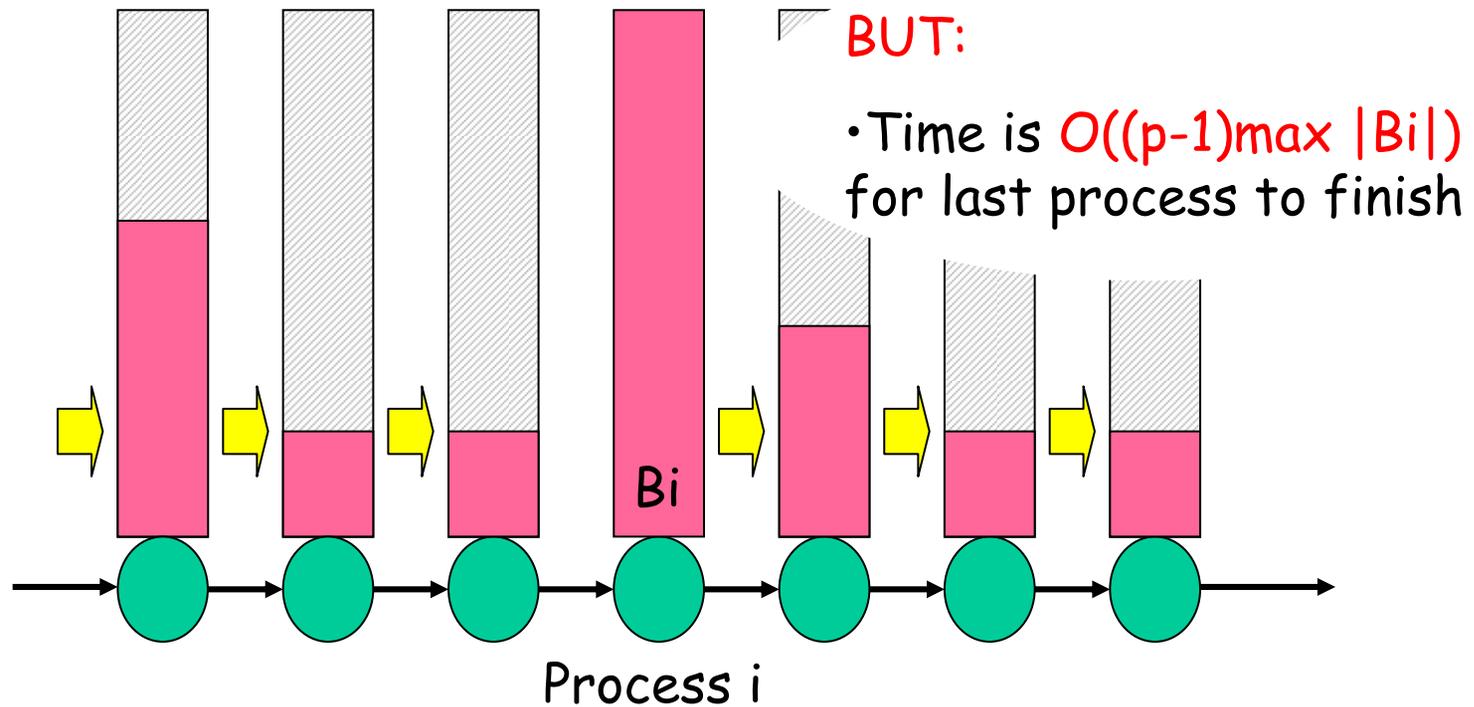
Linear ring: p-1 rounds

•Process i receives block B(i-1-k) from i-1 and sends block B(i-k) to i+1 in round k, k=0, …, p-1

Linear ring: p-1 rounds

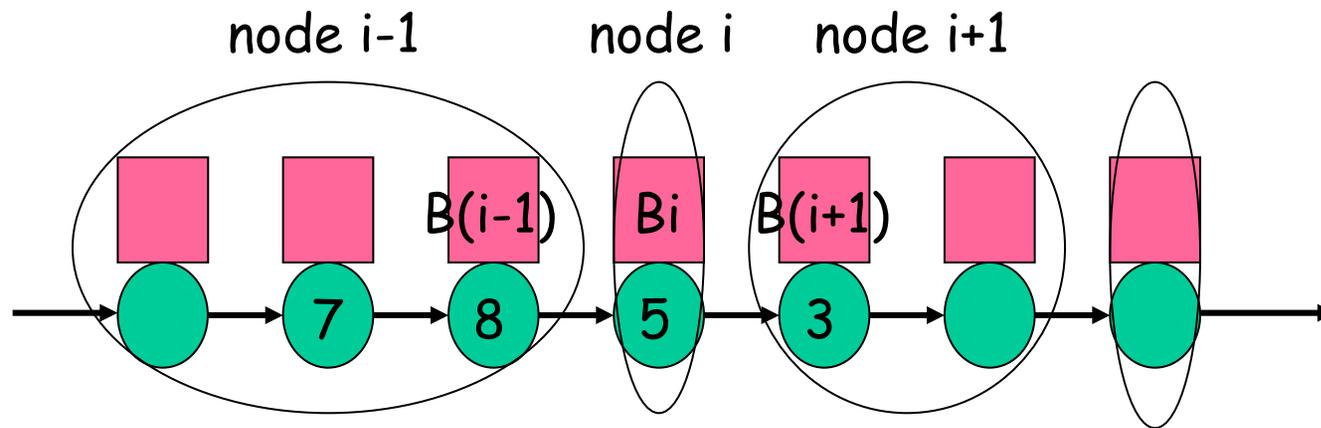•Process i receives block B(i-1-k) from i-1 and sends block B(i-k) to i+1 in round k, k=0, …, p-1

Linear ring: p-1 rounds

- Process i receives block B(i-1-k) from i-1 and sends block
B(i-k) to i+1 in round k, k=0, …, p-1

BUT:

- Time is $O((p-1)\max|B_i|)$ for last process to finish

$B_i$

Process i

Linear ring: p-1 rounds

- Process i receives block $B(i-1-k)$ from i-1 and sends block $B(i-k)$ to i+1 in round k, k=0, …, p-1

NEC

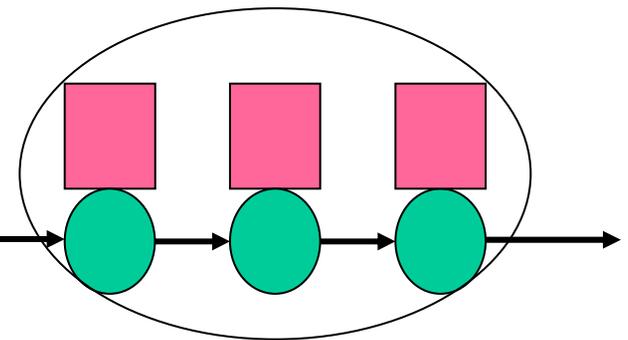**Observation 1:** linear ring works for clustered systems also

node i-1            node i      node i+1

$B(i-1)$   $B_i$   $B(i+1)$

7    8    5    3

Linear ring: p-1 rounds

• Process i receives block $B(i-1-k)$ from i-1 and sends block $B(i-k)$ to i+1 in round k, k=0, …, p-1

**MODIFICATION:** Use virtual ranking, one process per node sends, one process per node receives

**Analysis:**

- p-1 rounds

- Inter-node connections busy in all rounds: one process per node sends, one process per node receives
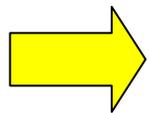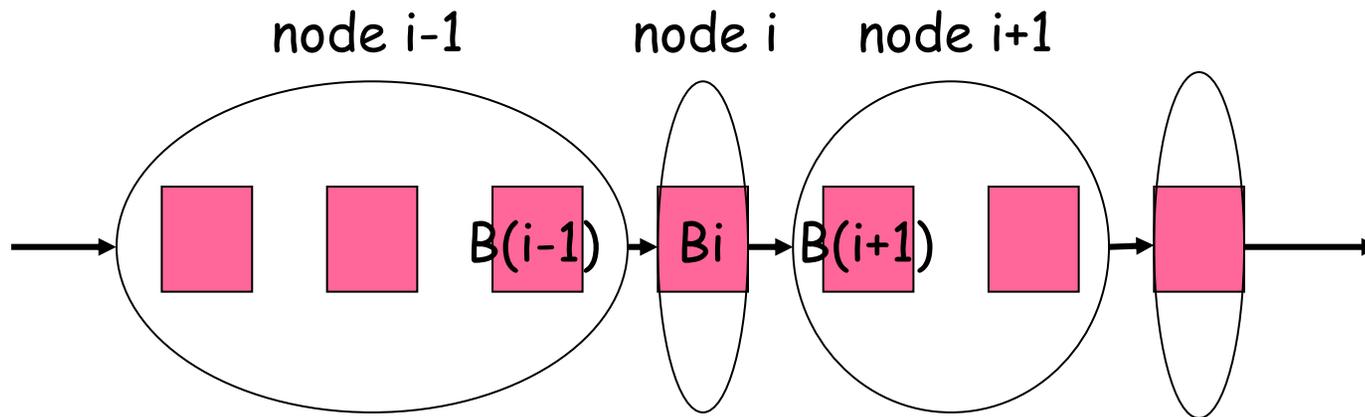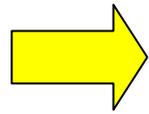
- Each node sends and receives (p-1) blocks

IMPROVEMENT:

a node never receives a block that it already has (replace by intra-node all-gather).
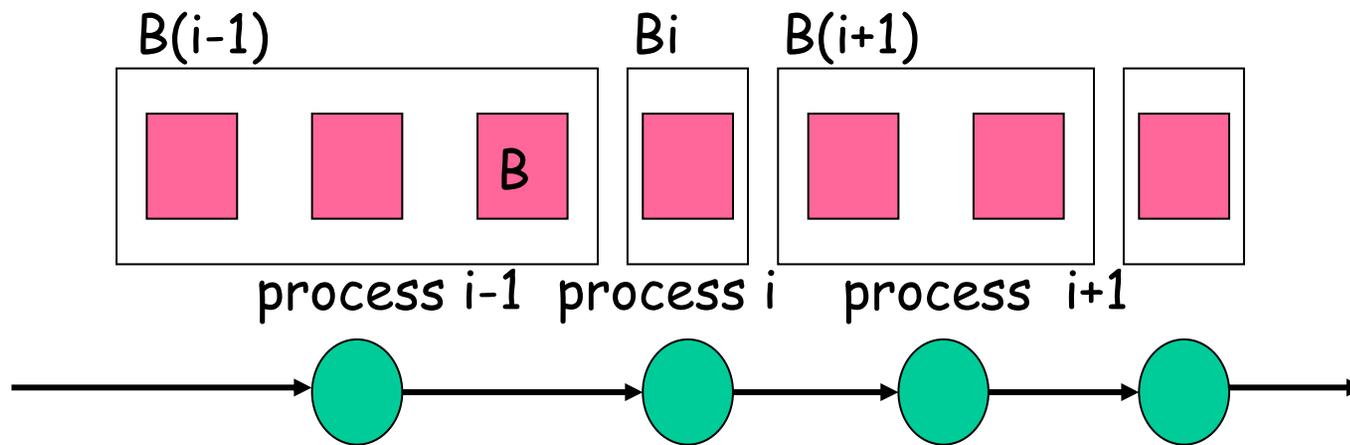
NEC

## Observation 2:

linear ring on cluster solves irregular problem over nodes



Irregular problem can be solved by simulating clustered algorithm

NEC

Handle each block Bi as node of ceil(Bi/B) regular blocks of some size B

B(i-1)          Bi      B(i+1)

B

process i-1  process i   process  i+1

EuroPVM/MPI 2008, Dublin
© NEC Laboratories Europe

NEC

# The algorithm for large, irregular all-gather problems

The blocked/pipelined ring algorithm
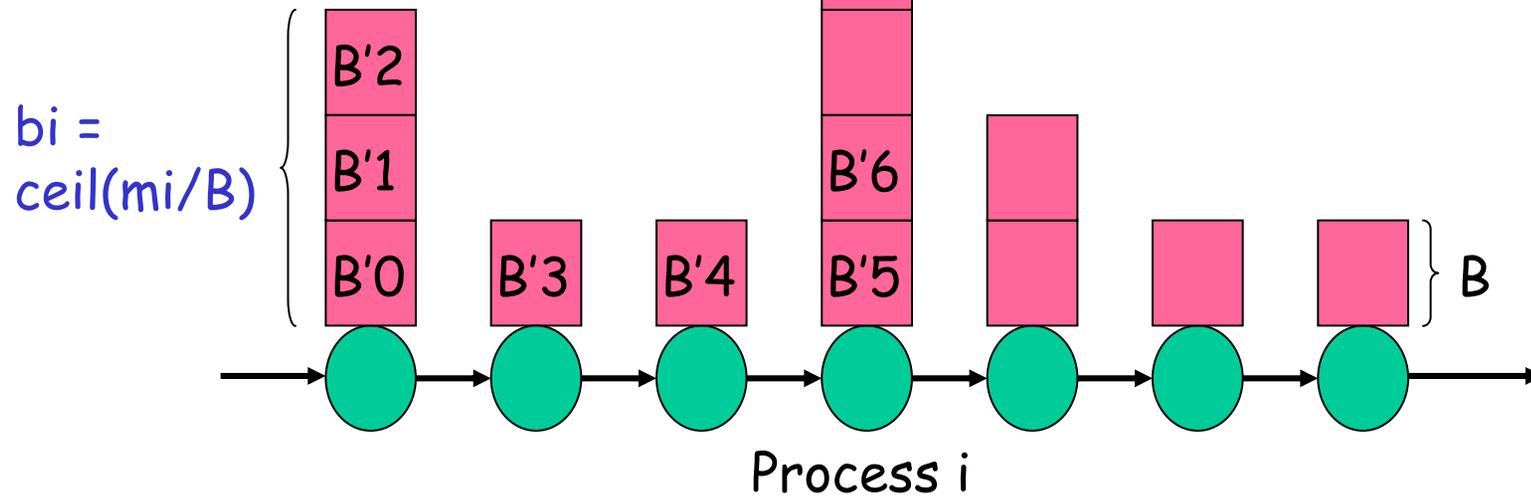
**NEC**

The blocked ring algorithm:



mi

Bi

Process i

1. For each process i, cut data Bi into $b_i = \max(1, \text{ceil}(m_i/B))$ blocks $B'_j$ of some chosen size (at most) B

**The blocked ring algorithm:**
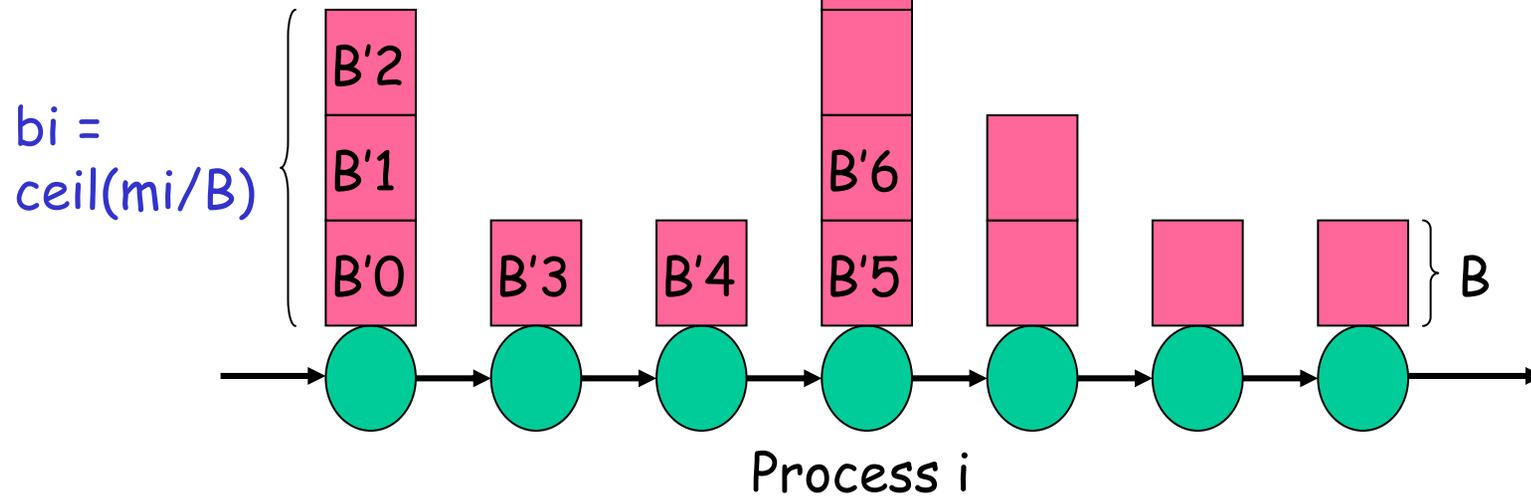
si: first block of process i,
$s_i = \sum_{(j<i)} b_i$

B'9

B'2

$b_i = ceil(m_i/B)$

B'1

B'6

B'0    B'3    B'4    B'5    } B

**Process i**

1. For each process i, cut data Bi into $b_i = max(1, ceil(m_i/B))$ blocks B'j of some chosen size (at most) B – each process has at least one block

Total number of blocks $b = \sum b_i$

**NEC**

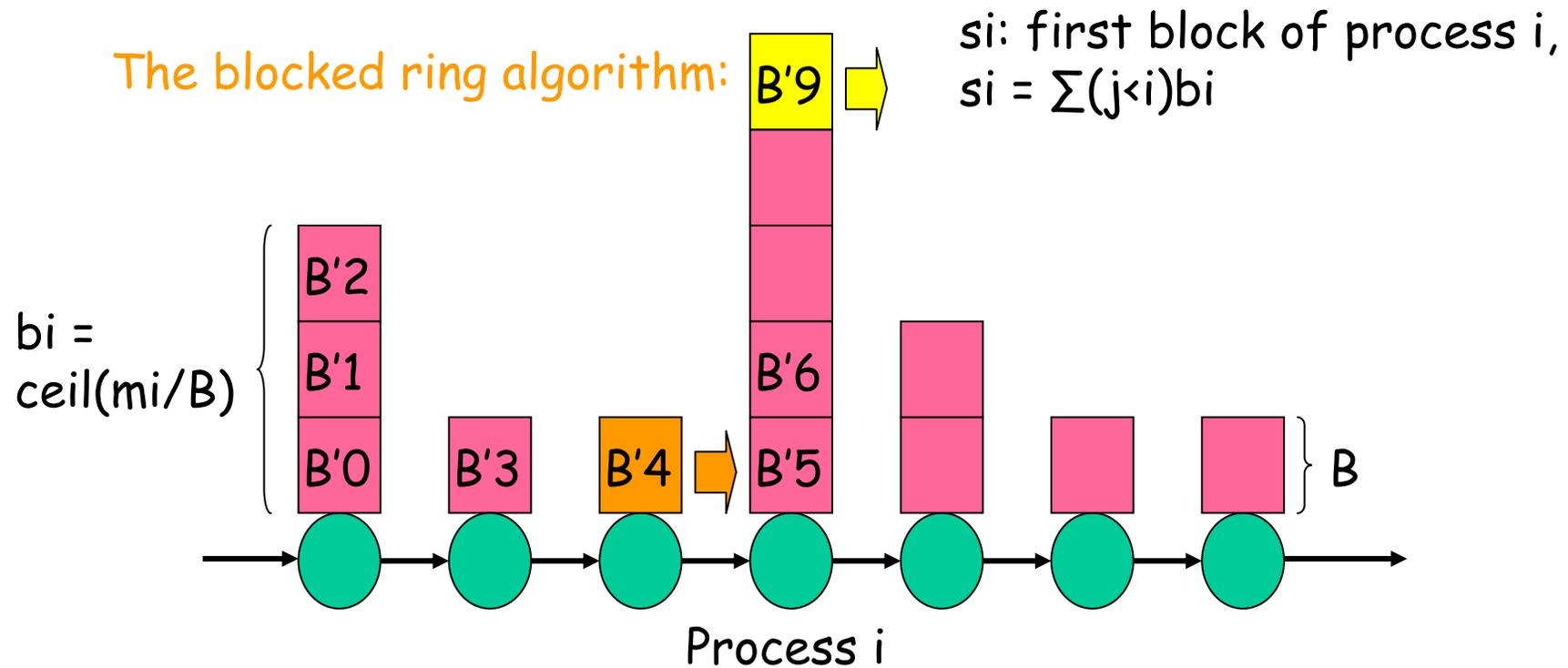The blocked ring algorithm:

si: first block of process i,
$s_i = \sum_{(j<i)} b_i$
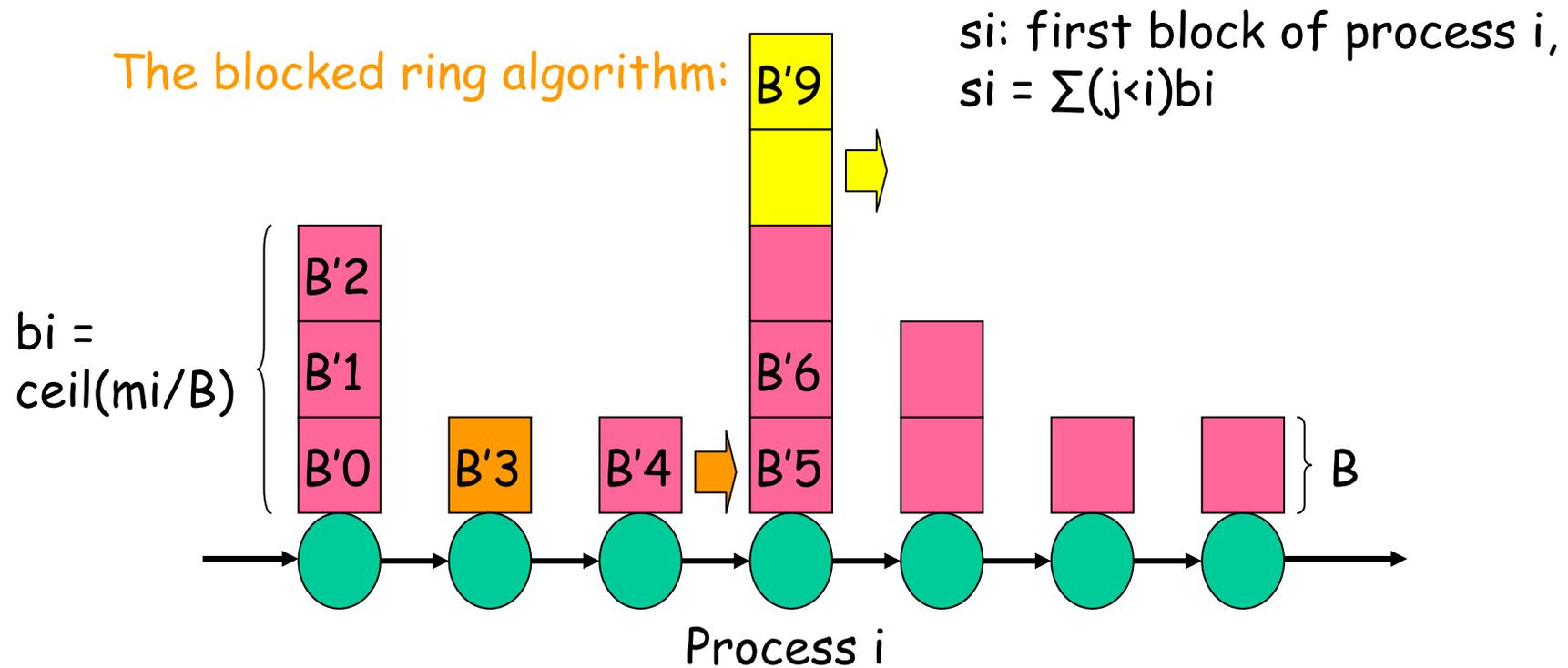
$b_i = ceil(m_i/B)$

B'2
B'1
B'0   B'3   B'4   B'5   B'6   B'9

} B

Process i

2. Run linear pipe on blocks, in round k process i receives
   block B'(si-1-k) from i-1 and sends block B'(si+bi-1-k) to i+1

The blocked ring algorithm:

si: first block of process i,
$s_i = \sum_{(j<i)} b_i$

B'9

$b_i =$
$ceil(m_i/B)$

B'2

B'1

B'0   B'3   B'4   B'5   B'6

} B

Process i

2. Run linear pipe on blocks, in round k process i receives
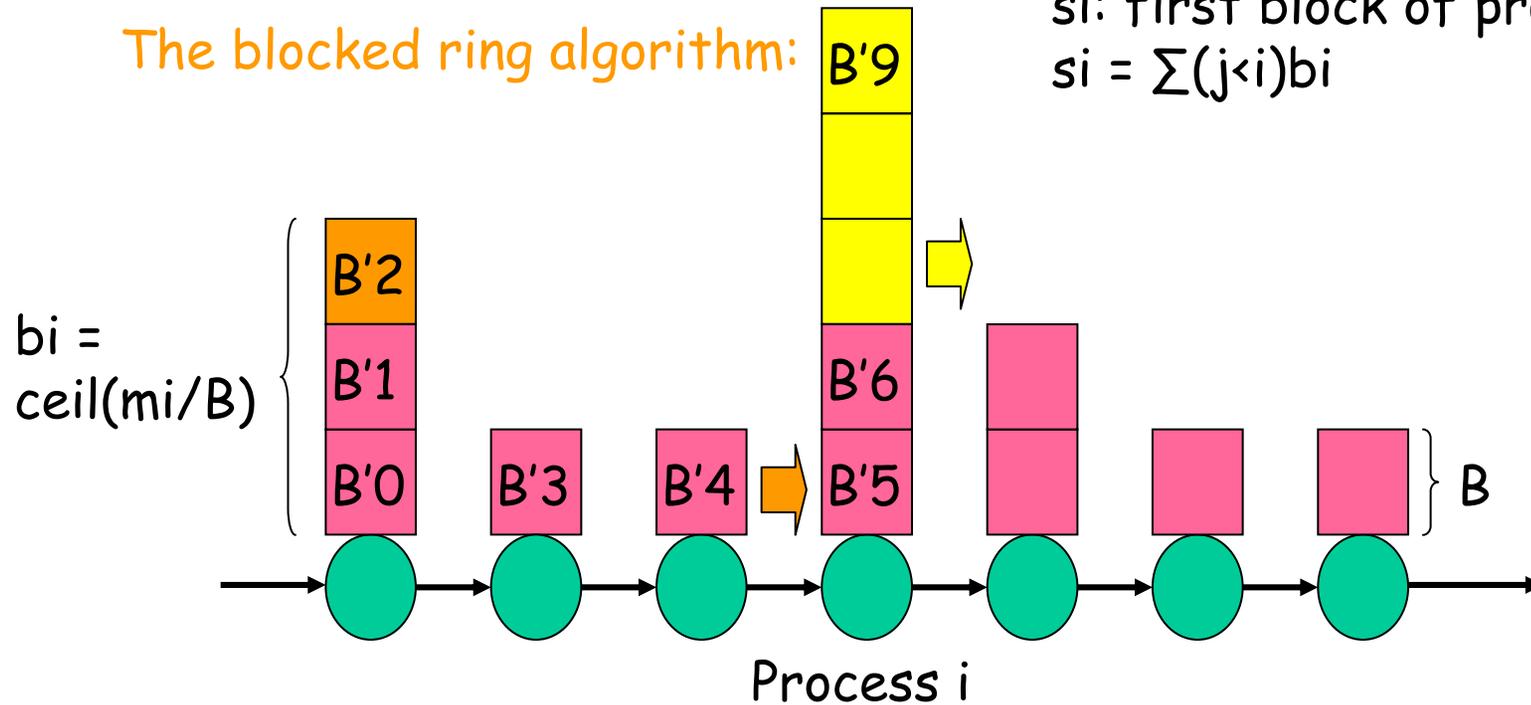   block B'(si-1-k) from i-1 and sends block B'(si+bi-1-k) to i+1

The blocked ring algorithm:

si: first block of process i,
si = ∑(j<i)bi

B'9

$b_i = \text{ceil}(m_i/B)$

B'2

B'1

B'0    B'3    B'4    B'5    B'6

B

Process i

2. Run linear pipe on blocks, in round k process i receives block B'(si-1-k) from i-1 and sends block B'(si+bi-1-k) to i+1
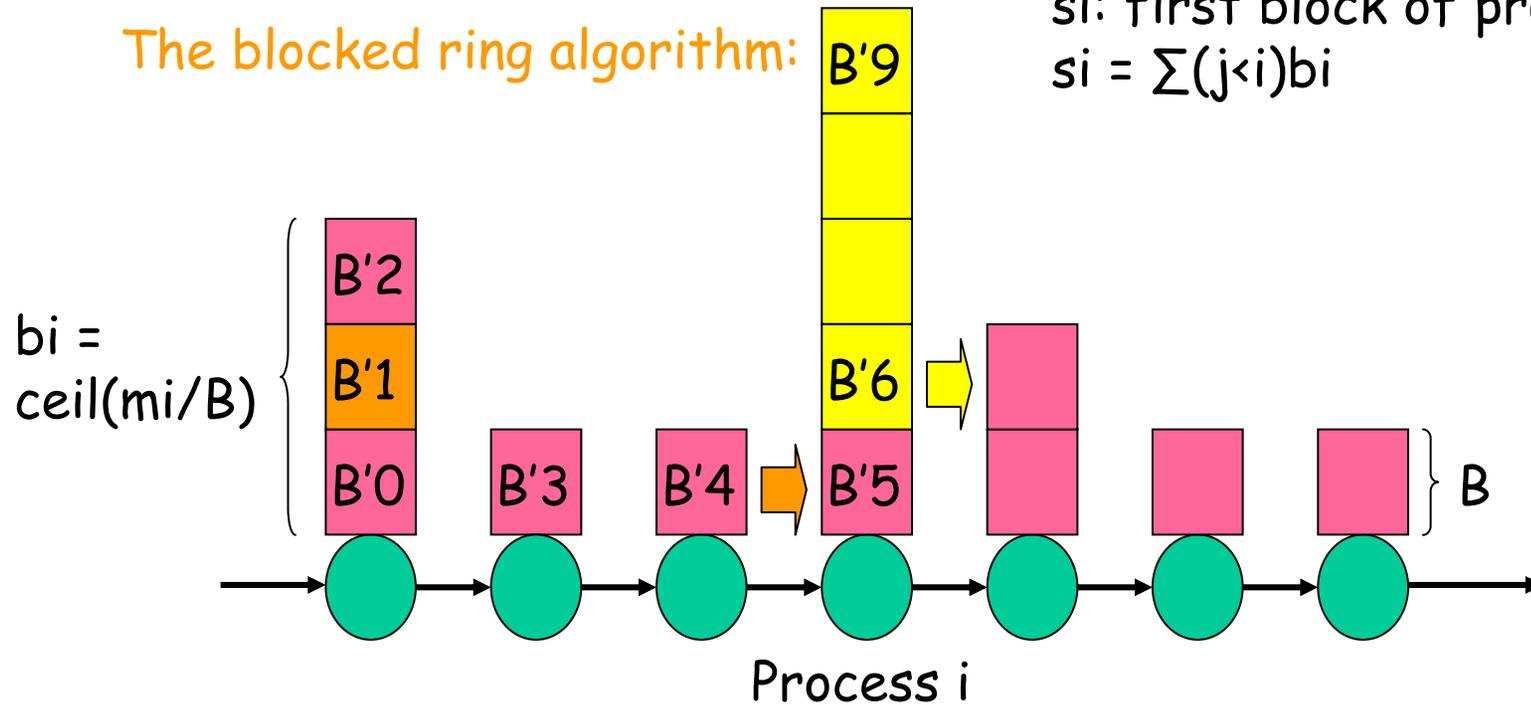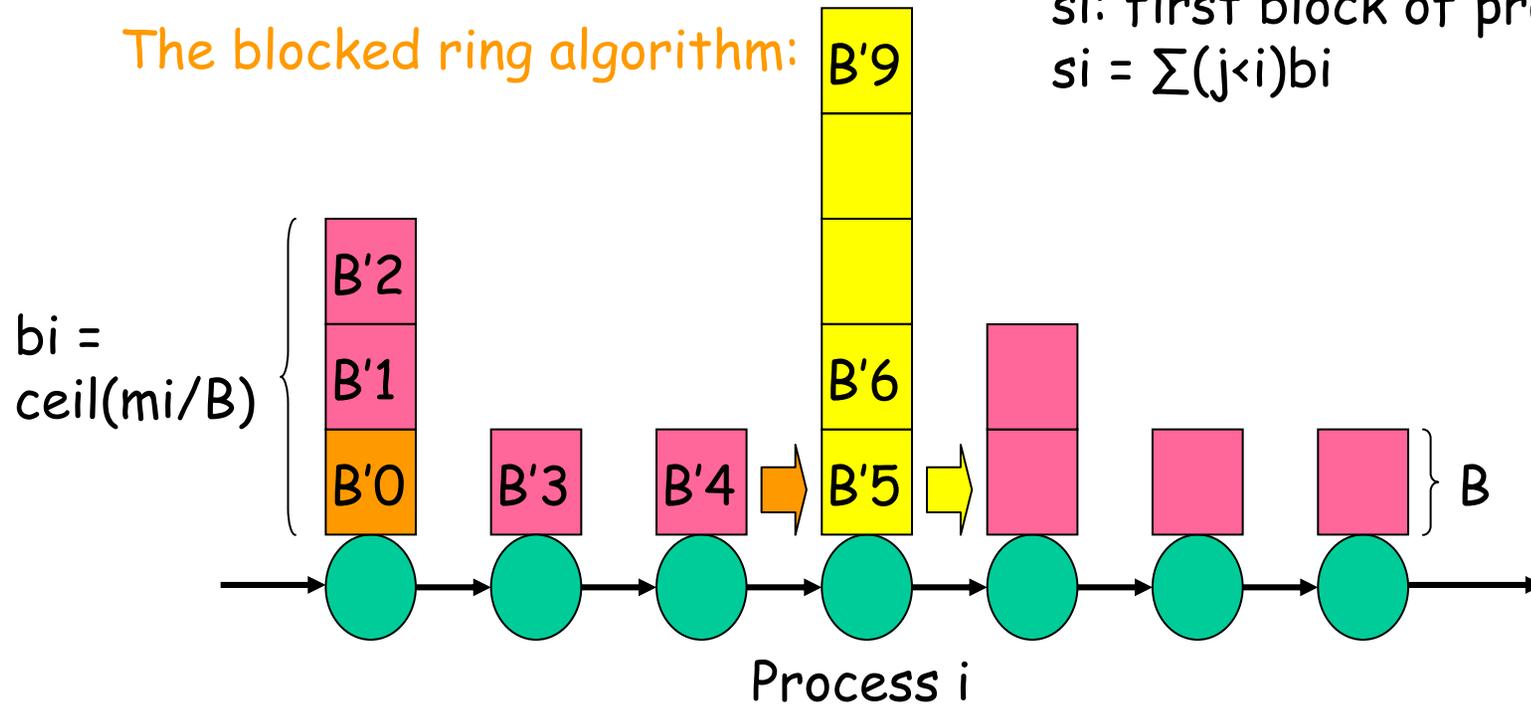
NEC

The blocked ring algorithm:

si: first block of process i,
si = ∑(j<i)bi

B'9

$bi = ceil(mi/B)$

B'2

B'1

B'0

B'3

B'4

B'5

B'6

B

Process i

2. Run linear pipe on blocks, in round k process i receives block B'(si-1-k) from i-1 and sends block B'(si+bi-1-k) to i+1

NEC

The blocked ring algorithm:
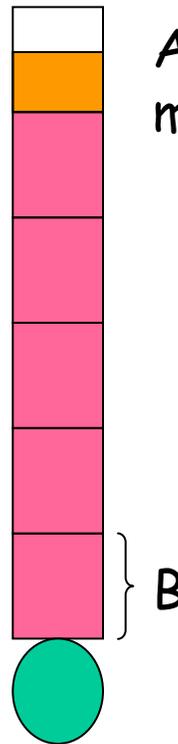
si: first block of process i,
$s_i = \sum_{(j<i)} b_i$

$b_i = \text{ceil}(m_i/B)$

B'2
B'1
B'0
B'3
B'4
B'5
B'6
B'9

B

Process i

2. Run linear pipe on blocks, in round k process i receives block $B'(s_i-1-k)$ from i-1 and sends block $B'(s_i+b_i-1-k)$ to i+1
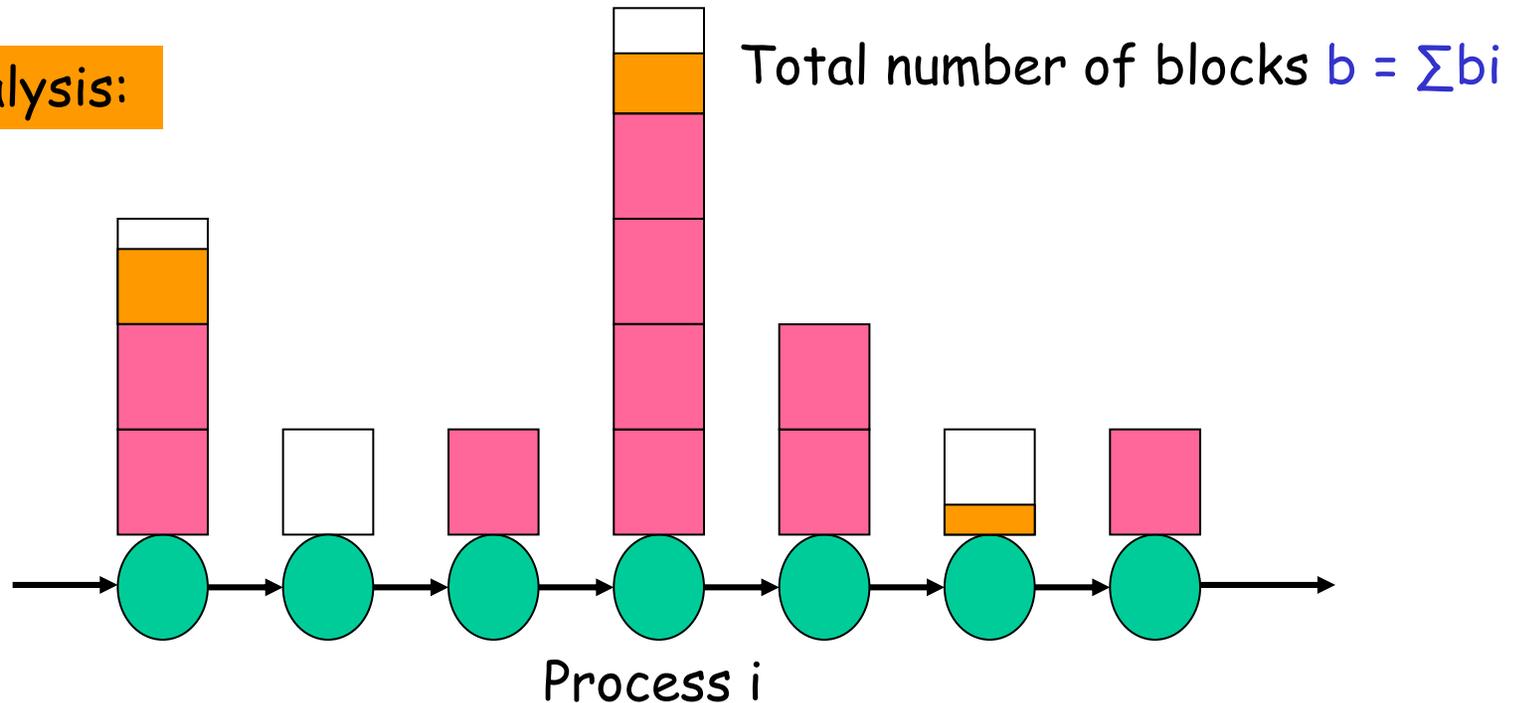
NEC

The blocked ring algorithm:

si: first block of process i,
$s_i = \sum_{(j<i)} b_i$

$b_i = \text{ceil}(m_i/B)$

B'9

B'2

B'1

B'0    B'3    B'4    B'5    B'6

} B

Process i

2. Run linear pipe on blocks, in round k process i receives block $B'(s_i-1-k)$ from i-1 and sends block $B'(s_i+b_i-1-k)$ to i+1

**The blocked ring algorithm:**

At most one block per process
may not be full (may be empty)

B

MODIFICATION:

• Empty blocks are neither sent nor received

• Only actual data of partial blocks is sent/received
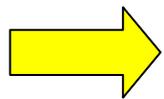
• No process receives a block it already has

**NEC**

**Analysis:**

Total number of blocks $b = \sum b_i$

Process i

- $b-1$ (almost) regular communication rounds

- Each round takes time $O(B)$, total time $O((b-1)B) = O((p+\sum \text{floor}(B_i/B))B) \approx O(m)$

MUCH BETTER than $O((p-1)\max(m_i))$ of linear ring without blocking

**General principle:**

1. Regular collective operation solves similar irregular problem on clustered system

2. By simulation, algorithm for regular problem on clustered system can be converted to algorithm for irregular problem. ASSUMPTION: communication capability of node and processor similar, eg. one ported

3. Irregular operation (on processes) remains irregular problem on clustered system

⟹ Blocked ring algorithm also works for MPI_Allgatherv on clustered system

## Choosing the block size B:

1. All mi>0: choose B=min(mi) – smallest data of some process (not too small, threshold).

2. Fixed block size

3. Some mi=0: number of rounds is m/B+(p-z)/2, z number of mi=0, assuming partial blocks half full. Time per round in linear model is $\alpha+\beta B$, best block size $\approx \sqrt{[2\alpha m/\beta(p+z-2)]}$

1. All processes busy in all rounds, for regular problems algorithm identical to linear ring

2. Simple solution - can lead to load imbalance for some distributions

3. Optimizes pipelining effect, for extreme problems with m0=m, mi=0 similar to pipelined broadcast. Linear model not accurate enough!

## Experimental results

Comparison of blocked ring algorithm to standard ring:

- Performance

- Load balance

- Effect of block size B

## Target systems

- NEC SX-8 - up to 30 nodes used

- Linux clusters with IB and Gig. Ethernet - 16 nodes, 24 nodes

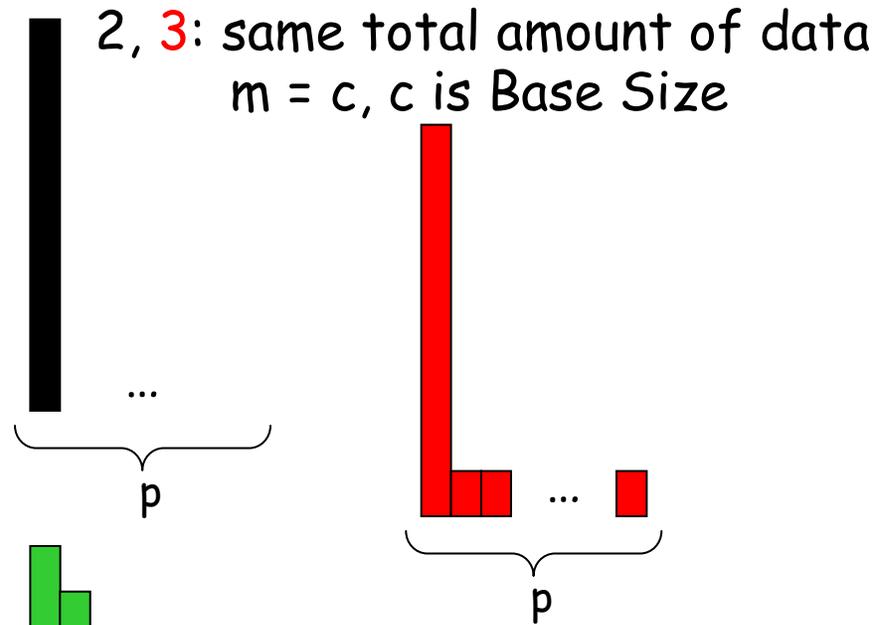- IBM Blue Gene/P - up to 4096 processes

- SiCortex 5832 - 5784 processes

Distributions

1. Regular
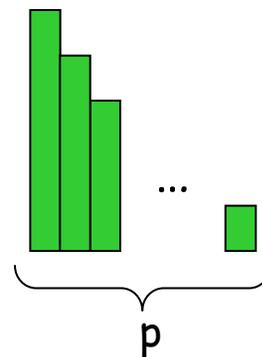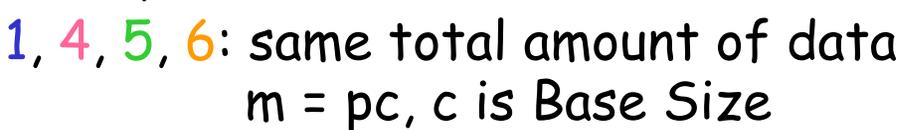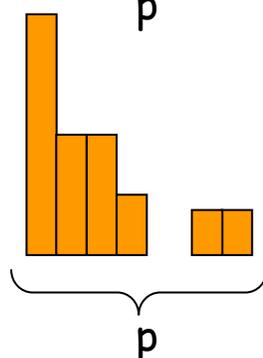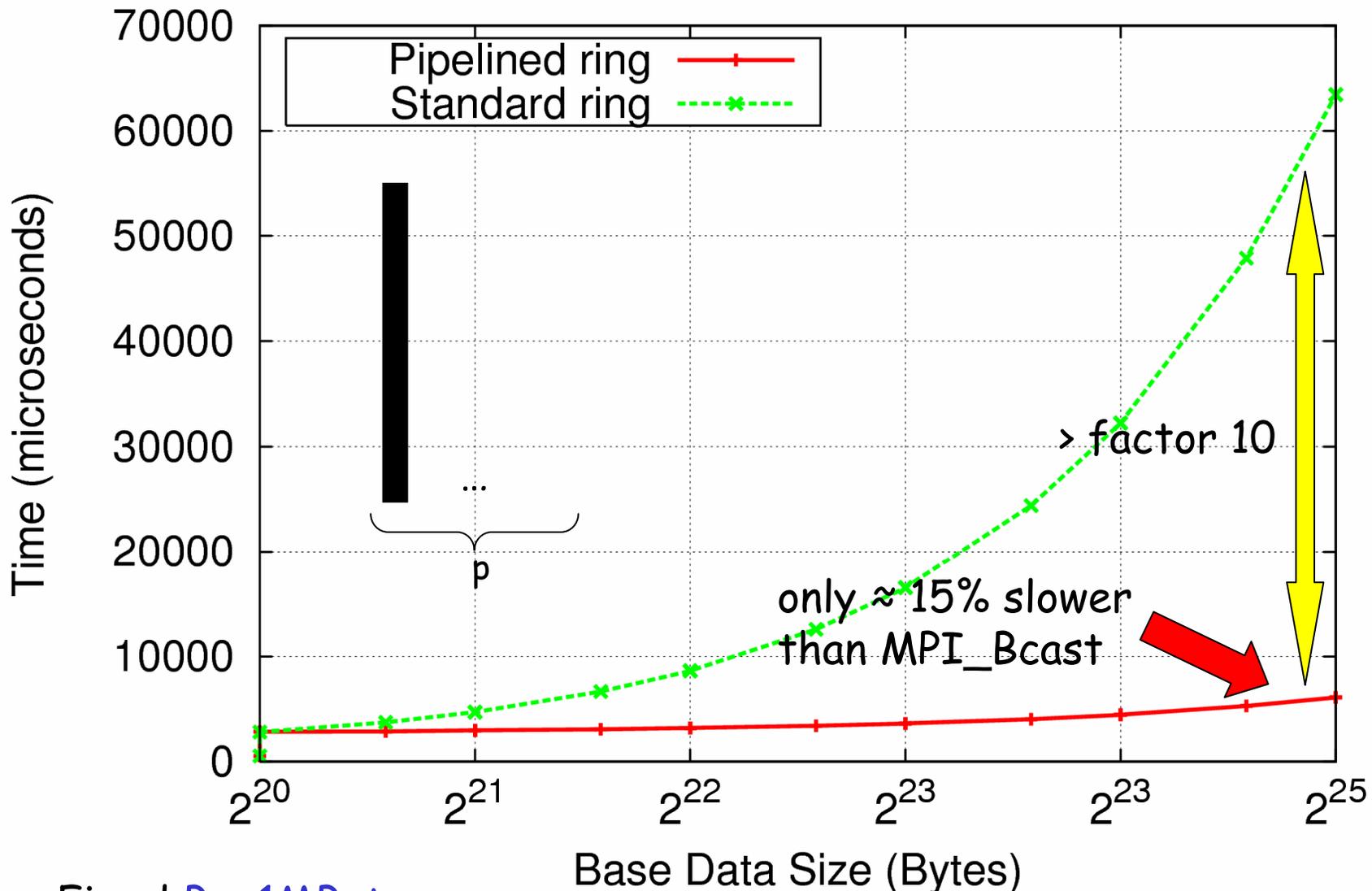2. Broadcast
3. Spike
4. Half
5. Linear
6. Geometric

2, 3: same total amount of data
m = c, c is Base Size

1, 4, 5, 6: same total amount of data
m = pc, c is Base Size

Comparable running times

NEC

**SX-8, 30x1 processes**

MPI_Allgatherv (Bcast)

> factor 10

only ≈ 15% slower than MPI_Bcast

Fixed B = 1MByte

September 8-10, 2008    EuroPVM/MPI 2008, Dublin
© NEC Laboratories Europe
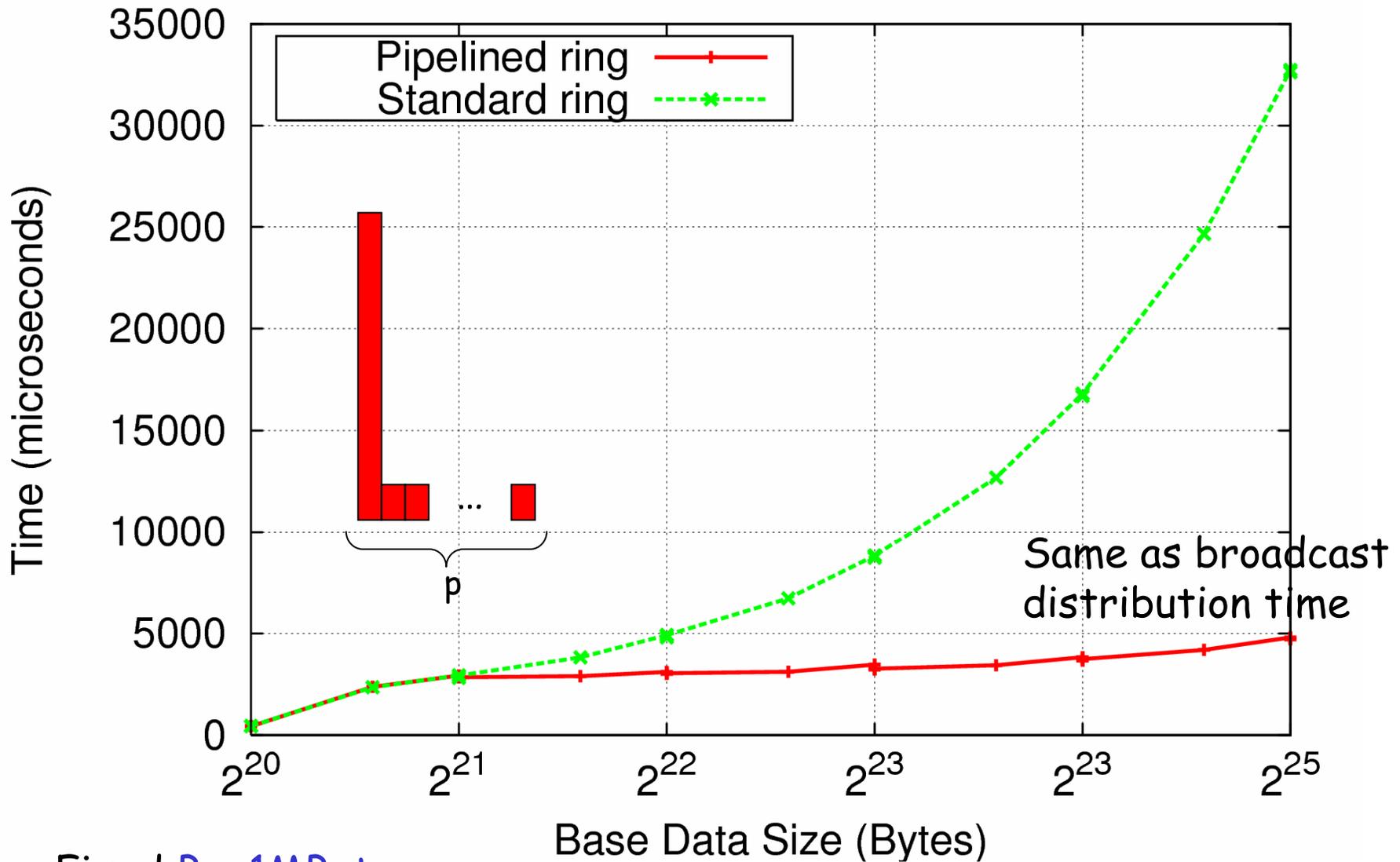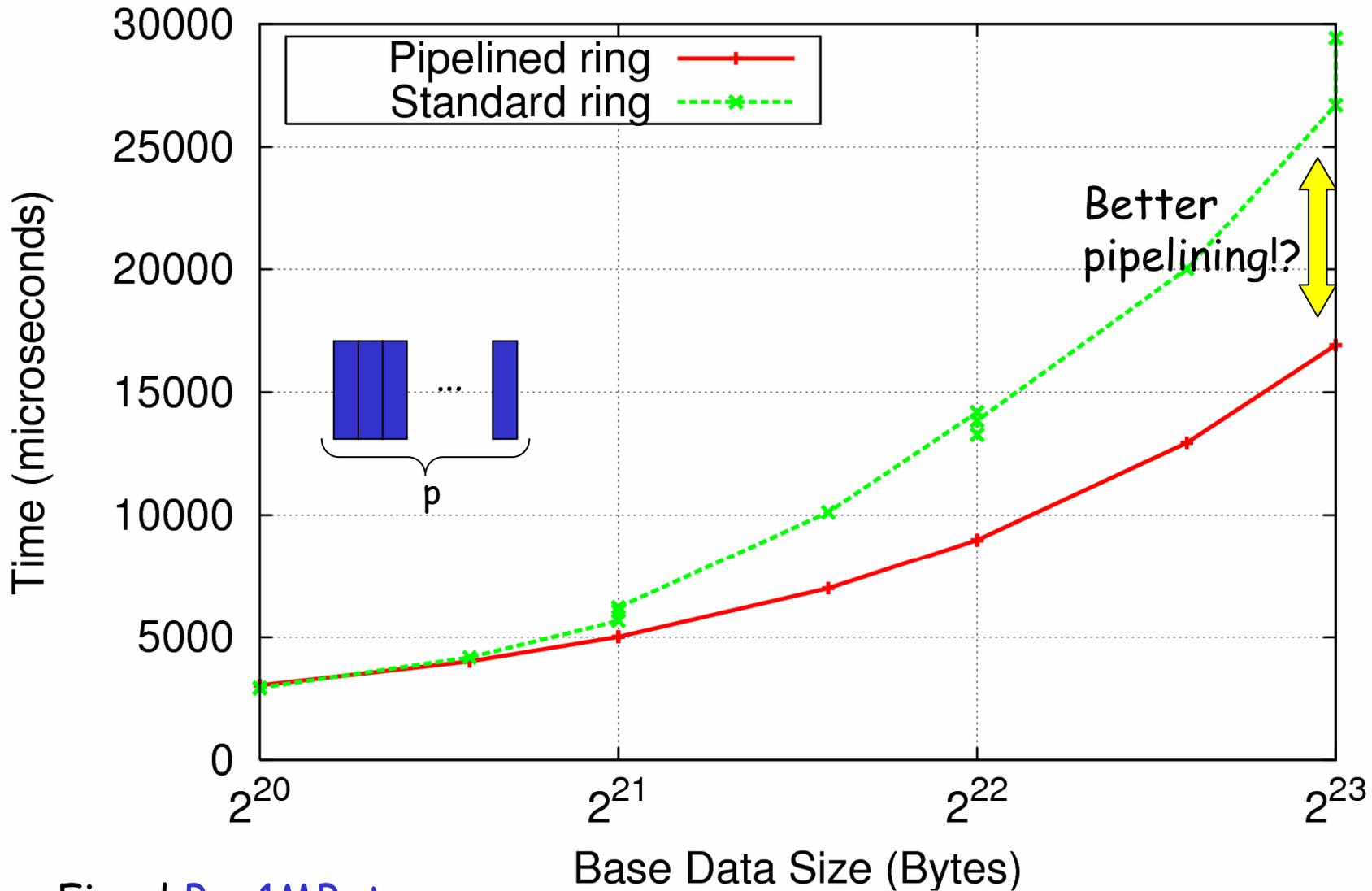
SX-8, 30x8 processes

MPI_Allgatherv (Bcast)

Time (microseconds) vs Base Data Size (Bytes)

- Pipelined ring
- Standard ring

> factor 12

p

Fixed B = 1MByte

September 8-10, 2008    EuroPVM/MPI 2008, Dublin
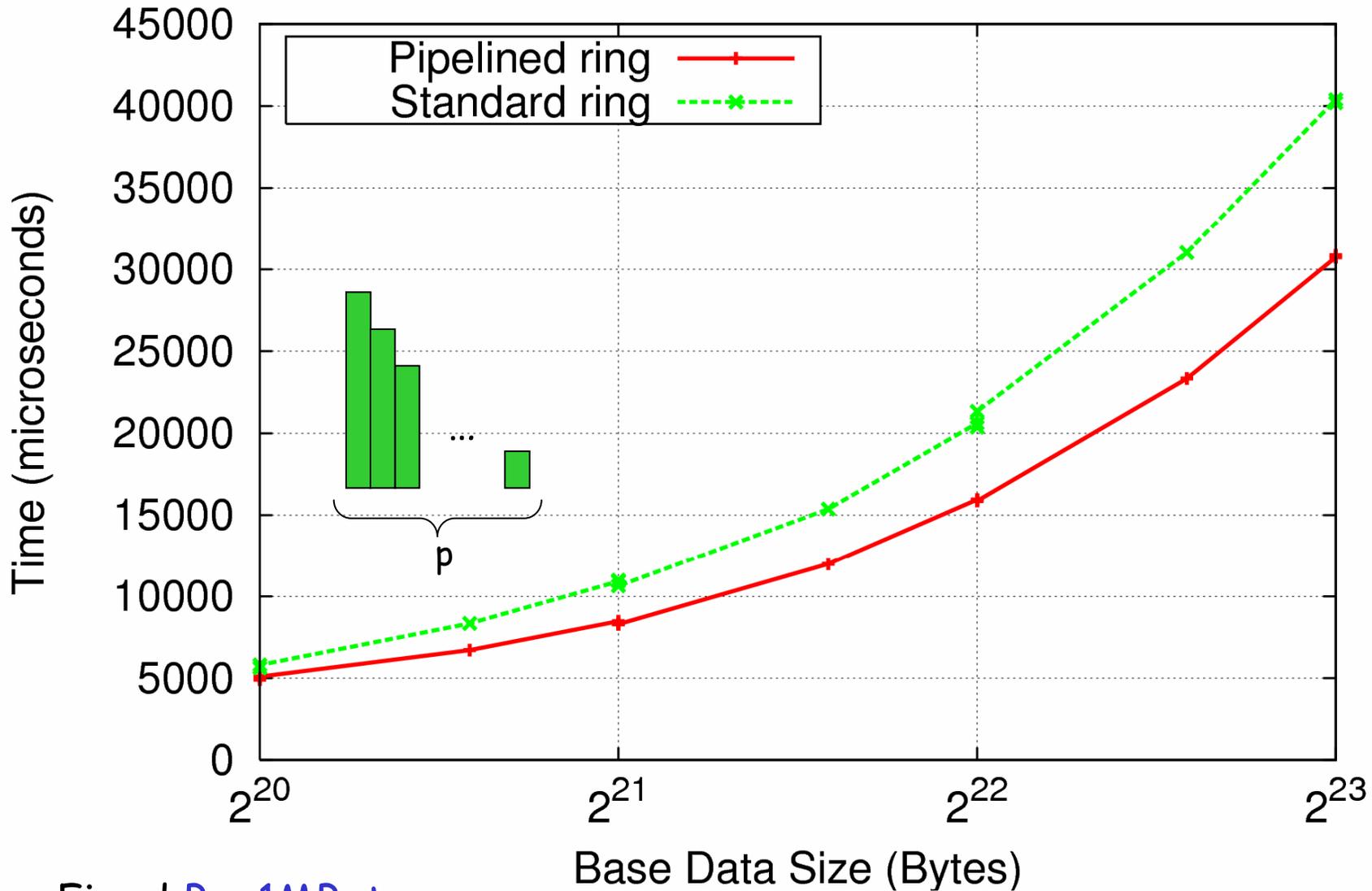© NEC Laboratories Europe

NEC

SX-8, 30x1 processes

MPI_Allgatherv (Spike)

Pipelined ring
Standard ring

Time (microseconds)

35000
30000
25000
20000
15000
10000
5000
0

$2^{20}$   $2^{21}$   $2^{22}$   $2^{23}$   $2^{23}$   $2^{25}$

Base Data Size (Bytes)

p

Same as broadcast distribution time

Fixed B = 1MByte

NEC

SX-8, 30x8 processes

MPI_Allgatherv (Spike)

Fixed B = 1MByte

September 8-10, 2008          EuroPVM/MPI 2008, Dublin
                              © NEC Laboratories Europe
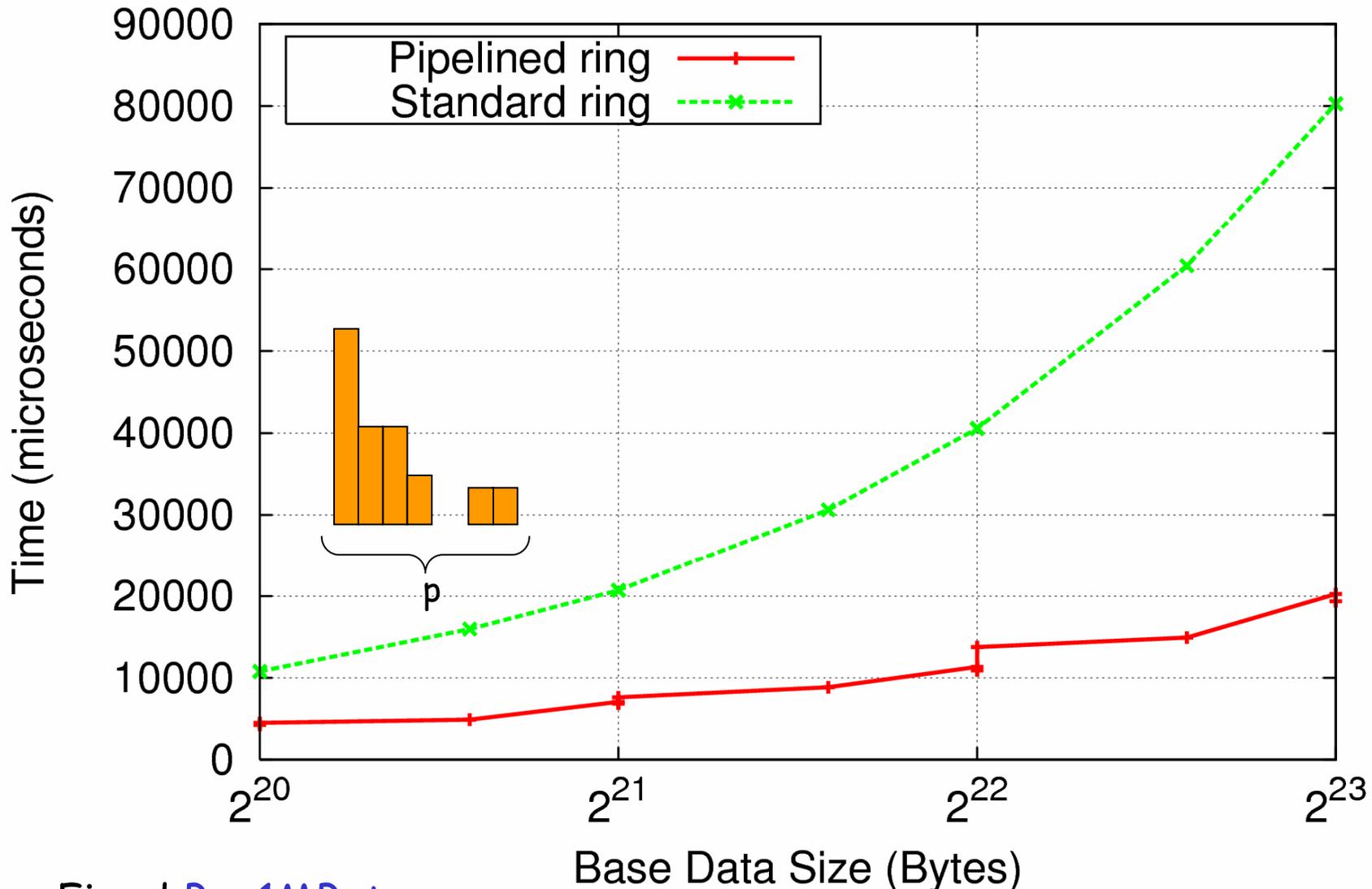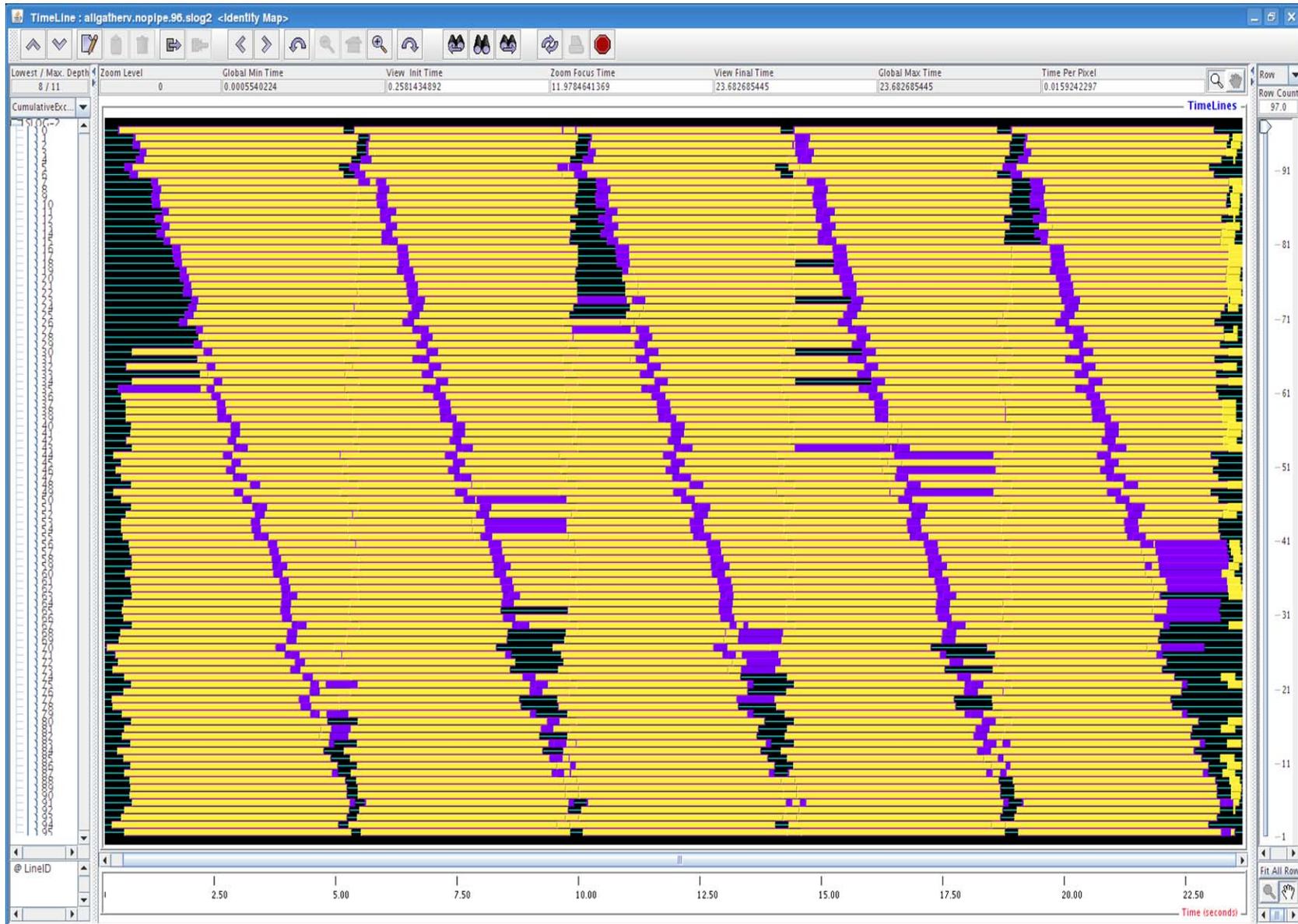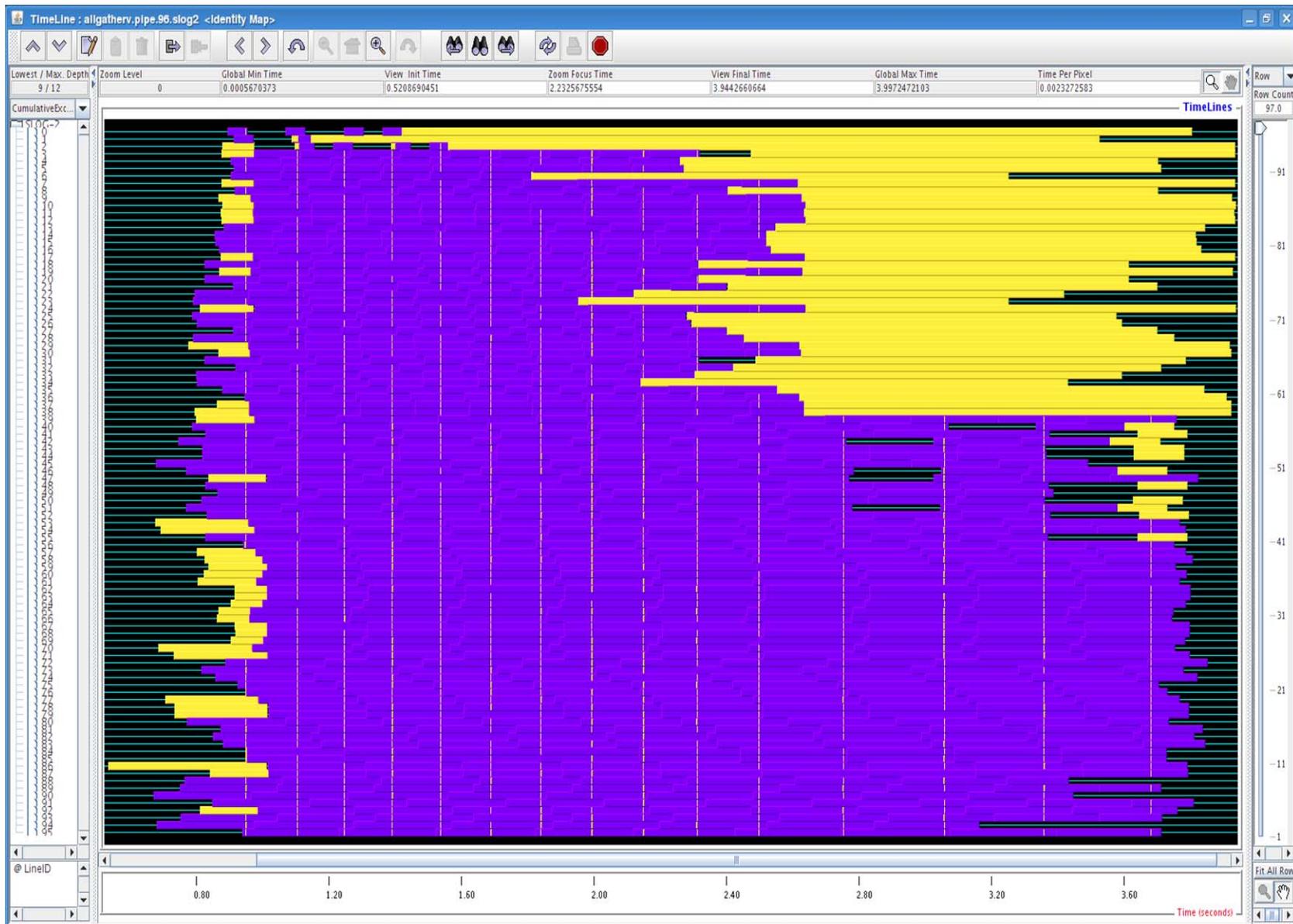
SX-8, 30x1 processes

MPI_Allgatherv (Regular)

Legend:
- Pipelined ring (red solid line with + markers)
- Standard ring (green dashed line with x markers)

Y-axis: Time (microseconds), 0 to 30000
X-axis: Base Data Size (Bytes), $2^{20}$ to $2^{23}$

Better pipelining!?

p

Fixed B = 1MByte

September 8-10, 2008          EuroPVM/MPI 2008, Dublin
                              © NEC Laboratories Europe

NEC

SX-8, 30x1 processes

MPI_Allgatherv (Half)

Fixed B = 1MByte

September 8-10, 2008    EuroPVM/MPI 2008, Dublin
© NEC Laboratories Europe

NEC

**SX-8, 30x1 processes**

MPI_Allgatherv (Geometric curve)

Time (microseconds) vs Base Data Size (Bytes)

Legend:
- Pipelined ring
- Standard ring

Fixed B = 1MByte

September 8-10, 2008

EuroPVM/MPI 2008, Dublin
© NEC Laboratories Europe

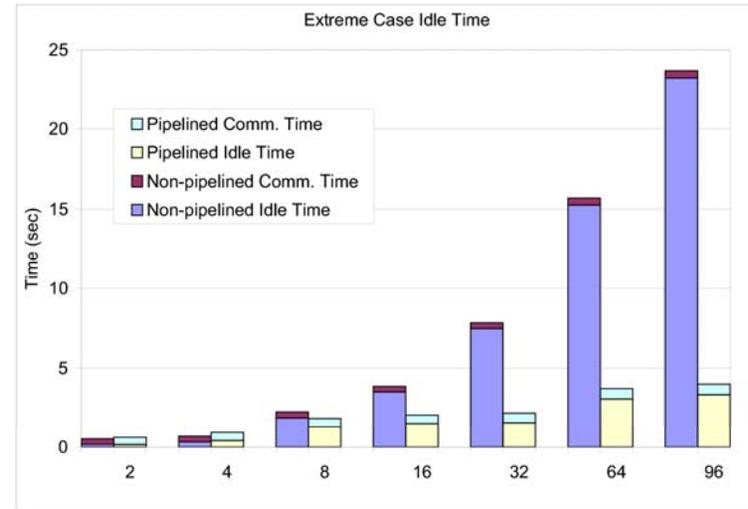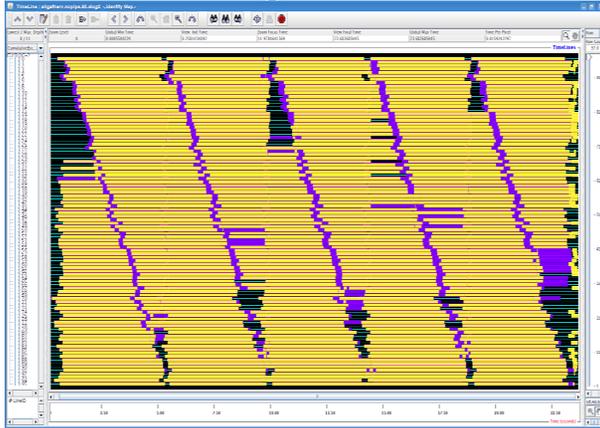NEC

EuroPVM/MPI 2008, Dublin
© NEC Laboratories Europe

NEC

# Linux cluster, 96 processes

### Linear ring





### Blocked ring



Yellow: idle time

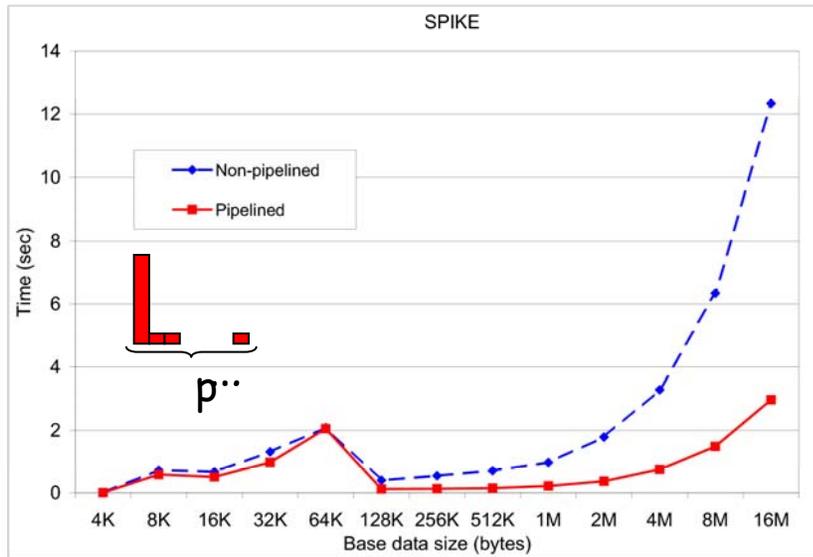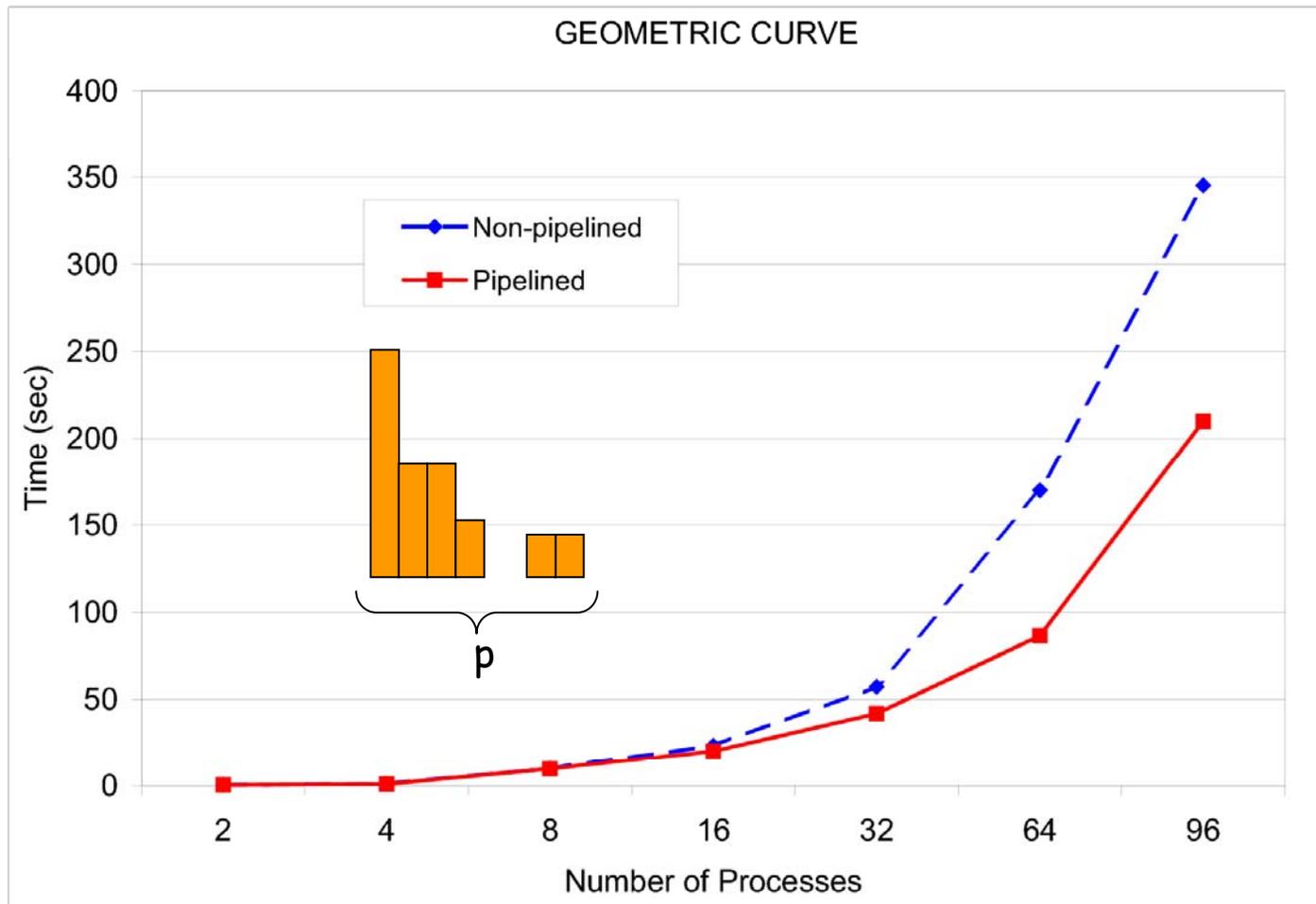Blue: communication time

EuroPVM/MPI 2008, Dublin
© NEC Laboratories Europe

# Linux cluster, 96 processes



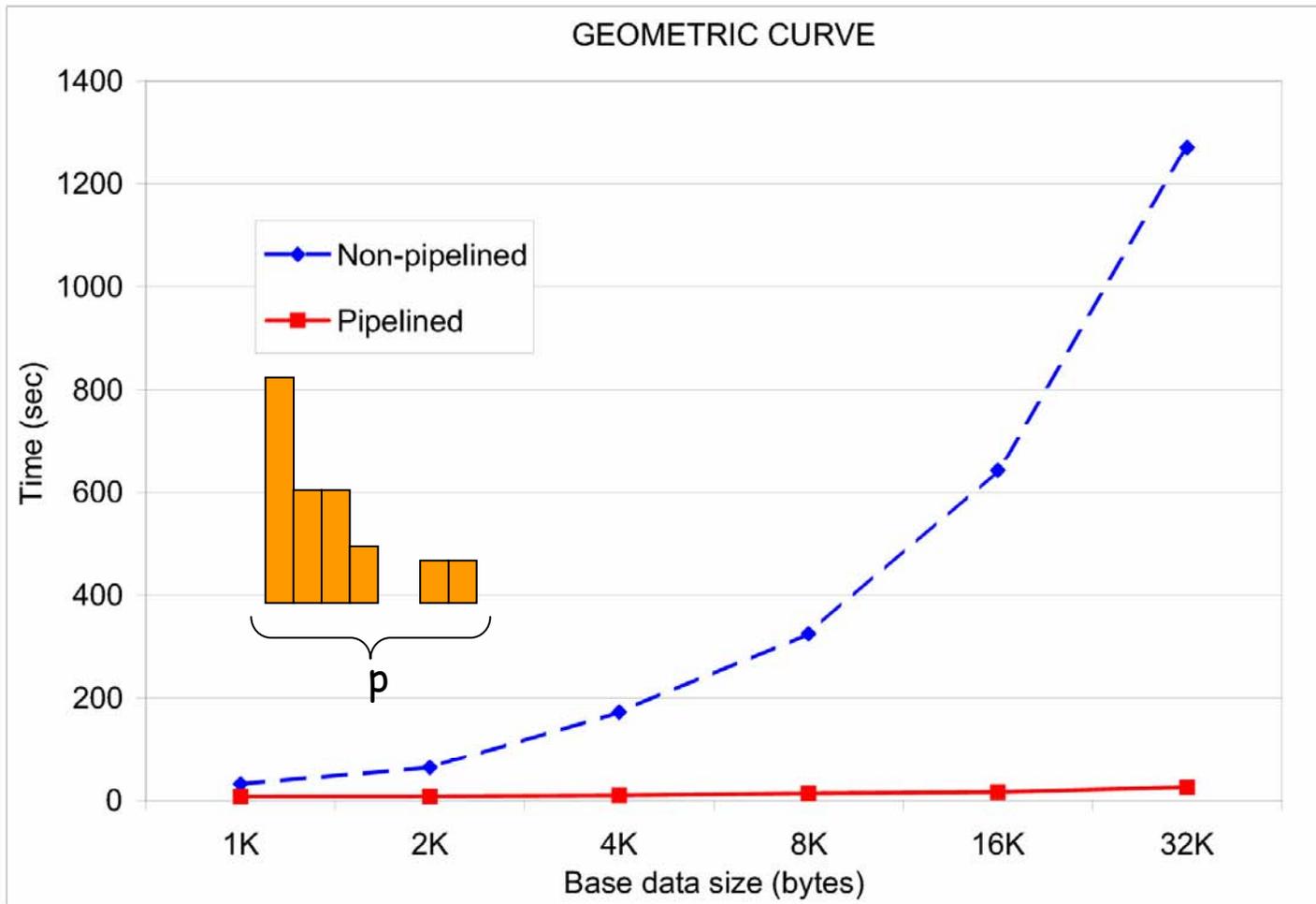Fixed B = 32KByte

# Linux cluster, varying number of processes



Fixed B = 32KByte

EuroPVM/MPI 2008, Dublin
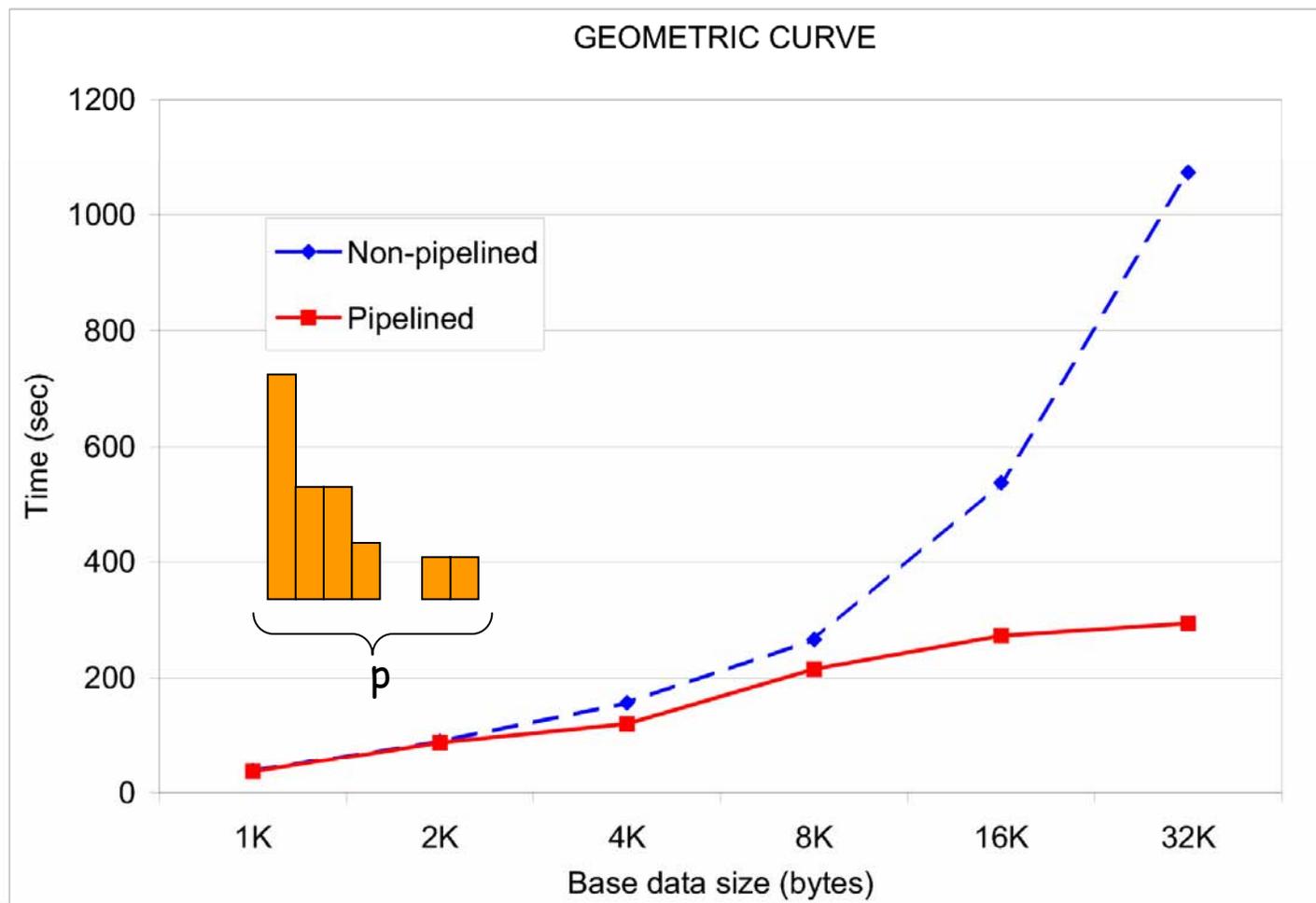© NEC Laboratories Europe

NEC

# Blue Gene/P, 4096 processes



Fixed B = 64KByte

# SiCortex, 5784 processes



GEOMETRIC CURVE

Fixed B = 1MByte

Linux cluster, 16x2 procs

MPI_Allgatherv (Decreasing)

Too small B

Factor 2

Legend: 32k, 64k, 128k, 512k, 1024k, non-pipelined

Time (Seconds) vs Base Data Size (Bytes), from $2^{18}$ to $2^{23}$

p

Linux cluster, 16x2 procs

MPI_Allgatherv (Half full)

September 8-10, 2008          EuroPVM/MPI 2008, Dublin
                             © NEC Laboratories Europe
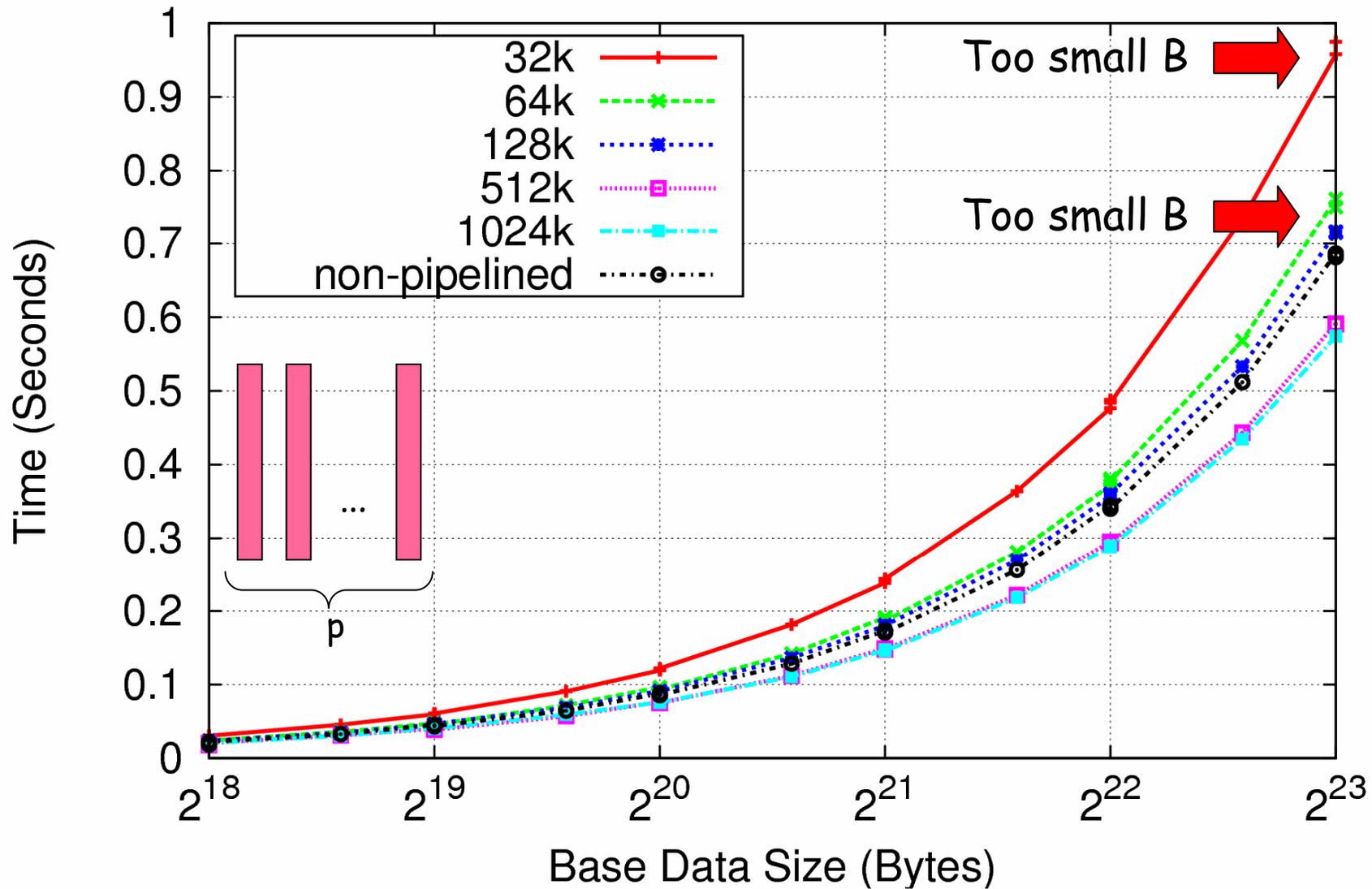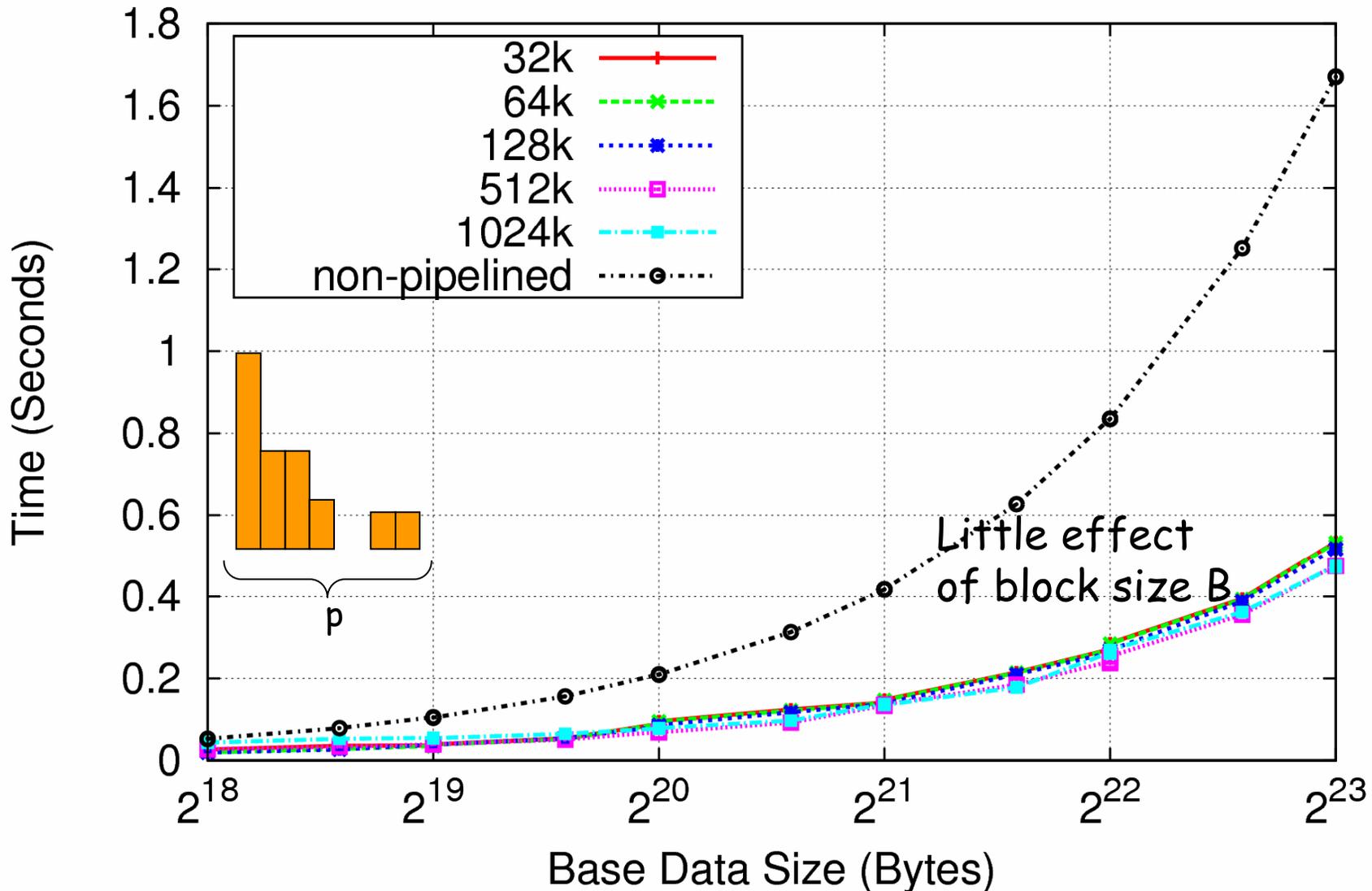
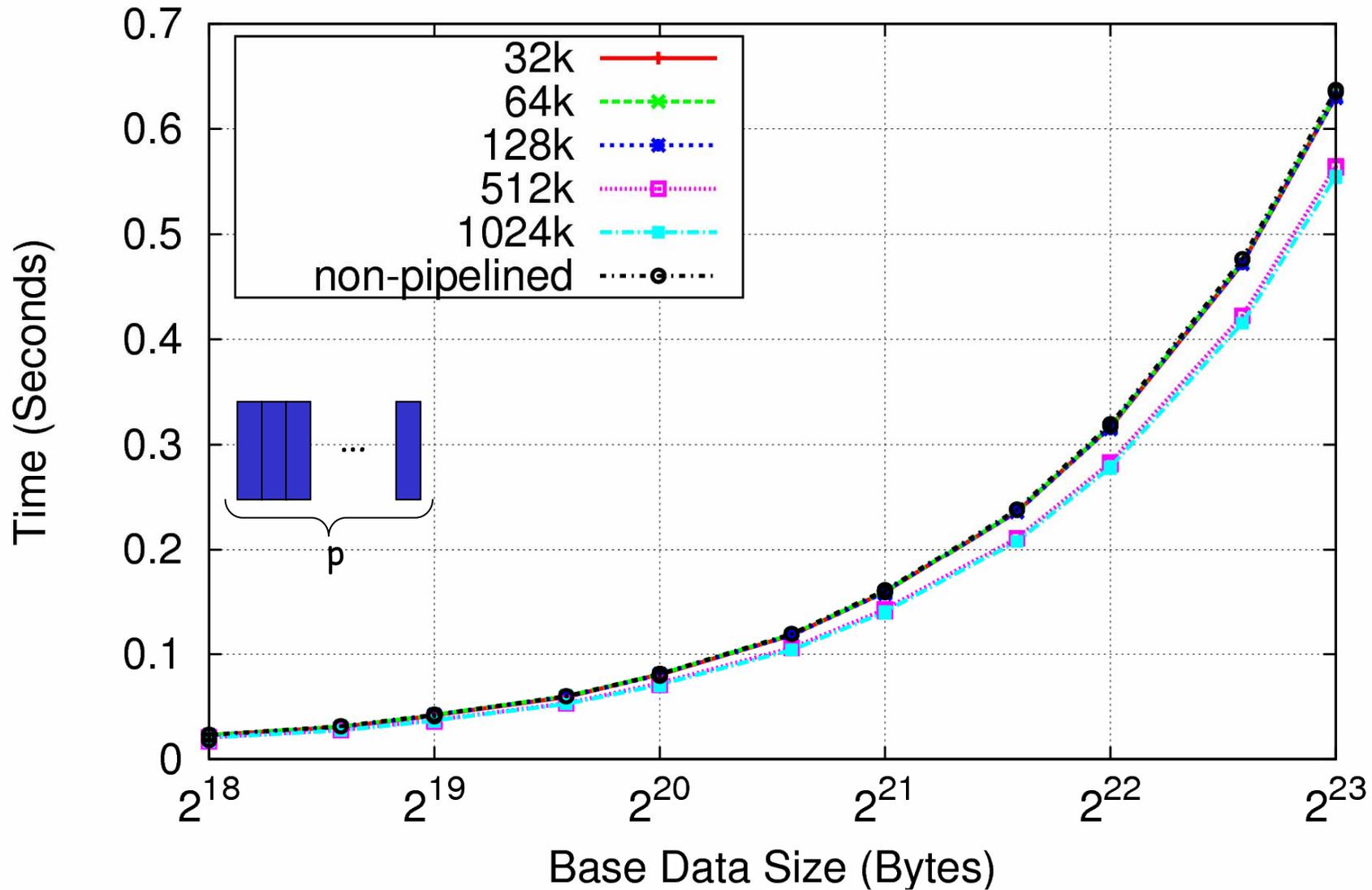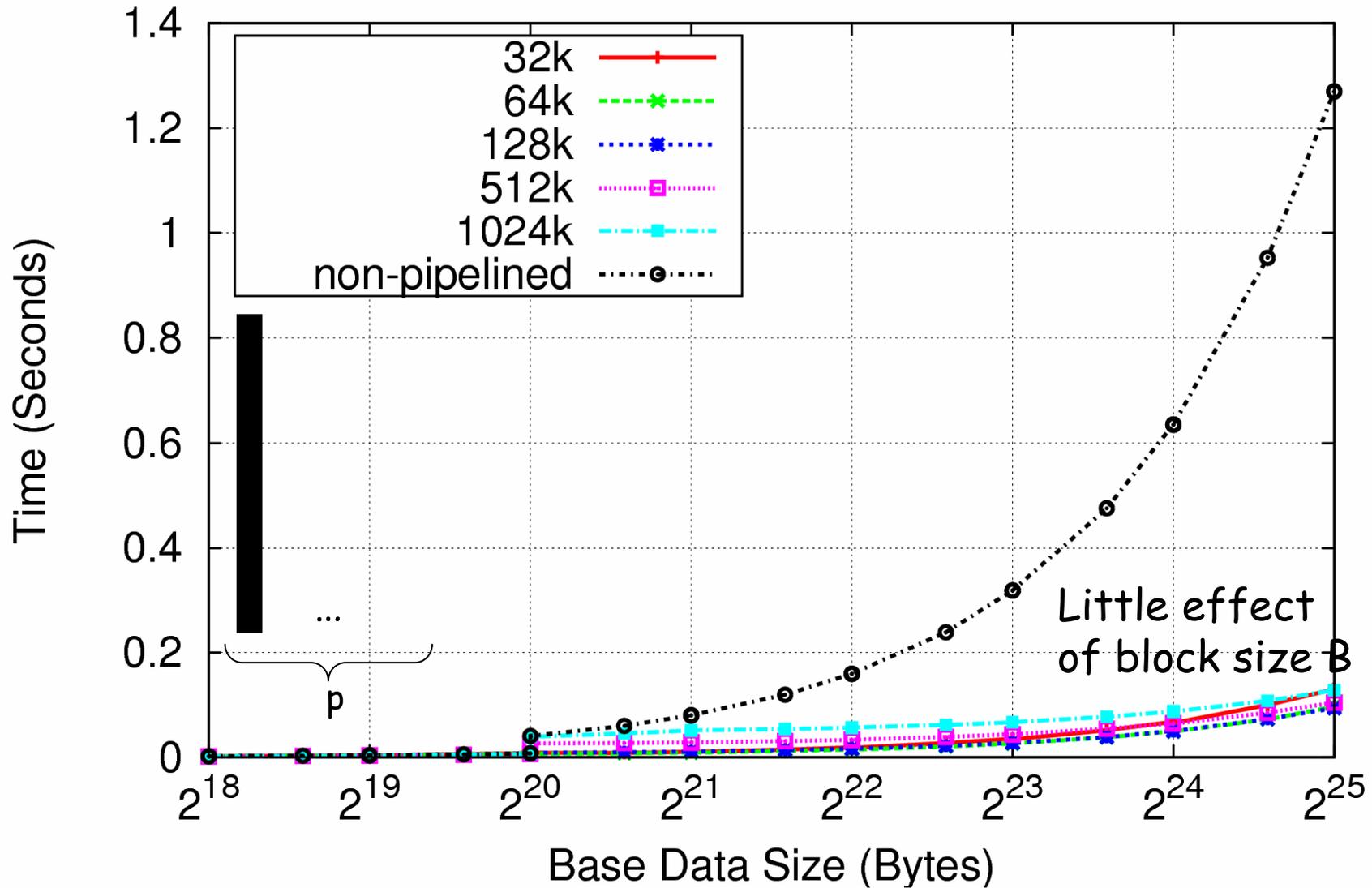Linux cluster, 16x2 procs

MPI_Allgatherv (Geometric curve)

Little effect of block size B

Linux cluster, 16x2 procs

MPI_Allgatherv (Regular)

Linux cluster, 16x2 procs

MPI_Allgatherv (Bcast)

Little effect of block size $B$

September 8-10, 2008    EuroPVM/MPI 2008, Dublin
© NEC Laboratories Europe

NEC

## Summary

- Simple, blocked linear ring algorithm for MPI_Allgatherv

  - NEW? Observation not found in literature

- Large performance gains for large problems on different systems

- Good limit behavior: identical to linear ring for regular problems, similar to pipelined broadcast for extreme distributions

- Tuning of block size: dependent on data distribution, linear model inadequate, experimental work needed

- There are relationships between regular and irregular collectives (on processes and nodes) that can (sometimes) be exploited for design of new algorithms

NEC