

# Parallel Programming

Prof. **Paolo Bientinesi**

`pauldj@aices.rwth-aachen.de`

WS 16/17



# Collective Communication: Lower Bounds

---

Cost of communication:  $\alpha + n\beta$

Cost of computation:  $\gamma \#ops$

$\alpha$  = “latency”, “startup”

$\beta$  = 1/“bandwidth”

$n$  = size of the message

$\gamma$  = cost of 1 flop

$p$  = # of processes

# Collective Communication: Lower Bounds

Cost of communication:  $\alpha + n\beta$

Cost of computation:  $\gamma \#ops$

$\alpha$  = “latency”, “startup”

$\beta$  = 1/“bandwidth”

$n$  = size of the message

$\gamma$  = cost of 1 flop

$p$  = # of processes

Primitive	Latency	Bandwidth	Computation
Broadcast	$\lceil \log_2(p) \rceil \alpha$	$n\beta$	-
Reduce	$\lceil \log_2(p) \rceil \alpha$	$n\beta$	$\frac{p-1}{p} n\gamma$

# Collective Communication: Lower Bounds

Cost of communication:  $\alpha + n\beta$

Cost of computation:  $\gamma \#ops$

$\alpha$  = “latency”, “startup”

$\beta$  = 1/“bandwidth”

$n$  = size of the message

$\gamma$  = cost of 1 flop

$p$  = # of processes

Primitive	Latency	Bandwidth	Computation
Broadcast	$\lceil \log_2(p) \rceil \alpha$	$n\beta$	-
Reduce	$\lceil \log_2(p) \rceil \alpha$	$n\beta$	$\frac{p-1}{p}n\gamma$
Scatter	$\lceil \log_2(p) \rceil \alpha$	$\frac{p-1}{p}n\beta$	-
Gather	$\lceil \log_2(p) \rceil \alpha$	$\frac{p-1}{p}n\beta$	-

# Collective Communication: Lower Bounds

Cost of communication:  $\alpha + n\beta$

Cost of computation:  $\gamma \#ops$

$\alpha$  = “latency”, “startup”

$\beta$  = 1/“bandwidth”

$n$  = size of the message

$\gamma$  = cost of 1 flop

$p$  = # of processes

Primitive	Latency	Bandwidth	Computation
Broadcast	$\lceil \log_2(p) \rceil \alpha$	$n\beta$	-
Reduce	$\lceil \log_2(p) \rceil \alpha$	$n\beta$	$\frac{p-1}{p}n\gamma$
Scatter	$\lceil \log_2(p) \rceil \alpha$	$\frac{p-1}{p}n\beta$	-
Gather	$\lceil \log_2(p) \rceil \alpha$	$\frac{p-1}{p}n\beta$	-
Allgather	$\lceil \log_2(p) \rceil \alpha$	$\frac{p-1}{p}n\beta$	-
Reduce-Scatter	$\lceil \log_2(p) \rceil \alpha$	$\frac{p-1}{p}n\beta$	$\frac{p-1}{p}n\gamma$

# Implementation of Bcast and Reduce

---

- IDEA: recursive doubling / “Minimum Spanning Tree” (MST)  
At each step, double the number of active processes.
- How to map the idea to the specific topology?

# Implementation of Bcast and Reduce

---

- IDEA: recursive doubling / “Minimum Spanning Tree” (MST)  
At each step, double the number of active processes.
- How to map the idea to the specific topology?
  - ring: linear doubling
  - (2d) mesh: 1 dimension first, then another, then another ...
  - hypercube: obvious, same as mesh
- Cost?

# Implementation of Bcast and Reduce

---

- IDEA: recursive doubling / “Minimum Spanning Tree” (MST)  
At each step, double the number of active processes.
- How to map the idea to the specific topology?
  - ring: linear doubling
  - (2d) mesh: 1 dimension first, then another, then another ...
  - hypercube: obvious, same as mesh
- Cost?
  - # steps:  $\log_2 p$
  - cost(step):  $\alpha + n\beta$
  - total time:  $\log_2(p)\alpha + \log_2(p)n\beta$  lower bound:  $\log_2(p)\alpha + n\beta$
  - note:  $\text{cost}(p^2) = 2 \text{cost}(p)!$



# Implementation of Bcast and Reduce

- IDEA: recursive doubling / “Minimum Spanning Tree” (MST)  
At each step, double the number of active processes.
- How to map the idea to the specific topology?
  - ring: linear doubling
  - (2d) mesh: 1 dimension first, then another, then another ...
  - hypercube: obvious, same as mesh
- Cost?
  - # steps:  $\log_2 p$
  - cost(step):  $\alpha + n\beta$
  - total time:  $\log_2(p)\alpha + \log_2(p)n\beta$  lower bound:  $\log_2(p)\alpha + n\beta$
  - note:  $\text{cost}(p^2) = 2 \text{cost}(p)!$
- Reduce  
BCast in reverse; cost(computation) ?

# Implementation of Scatter (and Gather)

---

- IDEA: MST again  
At step  $i$ , only  $\frac{1}{2^i}$ -th of the message is sent
- # steps:  $\log_2 p$
- cost(step $_i$ ):  $\alpha + \frac{n}{2^i}\beta$
- total time: 
$$\sum_{i=1}^{\log_2(p)} \alpha + \frac{n}{2^i}\beta = \log_2(p)\alpha + \frac{p-1}{p}n\beta$$
- lower bound:  $\log_2(p)\alpha + \frac{p-1}{p}n\beta$  optimal!

# A different implementation of Bcast

---

- IDEA: Scatter + cyclic algorithm (e.g., pass to the right)
- Cost?

# Implementation of Allgather (and Reduce-scatter)

- IDEA: “Recursive-doubling” (bidirectional exchange)  
Recursive allgather of half data + exchange data between disjoint nodes.

Node <sub>0</sub>	Node <sub>1</sub>	Node <sub>2</sub>	Node <sub>3</sub>
v[0]	v[1]	v[2]	v[3]



Node <sub>0</sub>	Node <sub>1</sub>	Node <sub>2</sub>	Node <sub>3</sub>
v[0] v[1]	v[0] v[1]	v[2] v[3]	v[2] v[3]



Node <sub>0</sub>	Node <sub>1</sub>	Node <sub>2</sub>	Node <sub>3</sub>
v[0]	v[0]	v[0]	v[0]
v[1]	v[1]	v[1]	v[1]
v[2]	v[2]	v[2]	v[2]
v[3]	v[3]	v[3]	v[3]

- # steps:  $\log_2 p$

- total time:

$$\sum_{i=1}^{\log_2(p)} \alpha + \frac{n}{2^i} \beta = \log_2(p) \alpha + \frac{p-1}{p} n \beta$$

# Another implementation of Allgather

- IDEA: Cyclic algorithm

Node <sub>0</sub>	Node <sub>1</sub>	Node <sub>2</sub>	Node <sub>3</sub>
v[0]			
	v[1]		
		v[2]	
			v[3]



Node <sub>0</sub>	Node <sub>1</sub>	Node <sub>2</sub>	Node <sub>3</sub>
v[0]	v[0]		
	v[1]	v[1]	
		v[2]	v[2]
v[3]			v[3]



Node <sub>0</sub>	Node <sub>1</sub>	Node <sub>2</sub>	Node <sub>3</sub>
v[0]	v[0]	v[0]	
	v[1]	v[1]	v[1]
v[2]		v[2]	v[2]
v[3]	v[3]		v[3]

- # steps:  $p - 1$

- total time:

$$\sum_{i=1}^{p-1} \alpha + \frac{n}{p} \beta = (p-1)\alpha + \frac{p-1}{p} n \beta$$