#### 2017 Practical Parallel Computing (実践的並列コンピューティング) No. 11

Distributed Memory Parallel Programming with MPI (4)

#### Toshio Endo

School of Computing & GSIC

endo@is.titech.ac.jp

## **Considering Performance of MPI Programs**

(Simplified) Execution time of an MPI program =

**Computation time** 

- + Communication time
- + Others



← including congestion

← load imbalance, I/O...





## Computation Time & Communication Time (1)



How are they determined? (very simplified discussion) <u>1. Aspect of software</u>

**Computation time** 

- Longer if computation costs are larger
  - O(mnk/p) in matmul,
  - O(NX NY NT/p) in diffusion per process

Communication time

- Longer if communication costs are larger
  - O(mk) in memory reduced matmul
  - O(NX NT) in diffusion per process

## Computation Time & Communication Time (2)

#### 2. Aspect of hardware

**Computation time** 

- Shorter if processor speed is faster
  - 140GFlops per node on TSUBAME2
  - Actually, memory access costs are important

#### **Communication time**

- Shorter if network speed is faster
  - 80Gbps per node on TSUBAME2





4

## **Parameters for Network Speed**

What parameters describes network speed?

- Bandwidth : Data amounts that network can transport per unit time → Larger is better
  - bps: X bits per second
  - B/s: X Bytes per second
  - On TSUBAME2, 80Gbps = 10GB/s per node
- Network latency: Time to transport minimum data (1bit, for example) → Smaller is better
  - On TSUBAME2, <10us

X Additionally, communication time may suffer from effects of network topology: how nodes/switches are connected to each other

### **Bandwidth and Latency**

Is "latency" reciprocal of "bandwidth"?  $\rightarrow$  No, because data are transported successfully



- T: Communication time
- M: Data size
- B: Bandwidth
  - L: Network latency

Be aware of difference between"Byte" and "bit": 1Byte=8bit

X In some contexts, T, not L, may be called "latency"

# Why L (Latency) > 0?

1. Overhead when data passes network switches



- 2. Software overhead
  - Cf) Socket library, MPI library performs data copy
- 3. Transfer speed of data cannot exceed speed of light (3x10<sup>8</sup> m/s)

Considering T = M / B + L,

batching communication may improve communication time

cf) Sending <u>1Gbytes at once</u> is much faster than sending <u>1Kbytes for 1,000,000</u> times



# How to Improve Performance of MPI Programs?

- Reduce computation time
  - Reduce computation amount
  - Using cache memory efficiently
- Reduce communication time
  - Reduce communication amount
  - Batch communication
  - Using collective communication is also good
- Reduce other time
  - Improve load balancing
  - Reconsider I/O

... And overlap computation and communication



## Idea of Overlapping

If "some computations" do not require contents of message, we may start them beforehand



#### **Overlapping in Stencil Computation** (related to [M1], but not requied)

When we consider data dependency in detail, we can find <u>computations that do not need data from other processes</u>



Rows C, D, E do not need data from other processes → They can be computed without waiting for finishing communication

On the other hand, rows B, F need received data

For such purposes, <u>non-blocking communications</u> (MPI\_Isend, MPI\_Irecv...) are helpful



## Implementation <u>without</u> Overlapping



## Implementation <u>with</u> Overlapping



for (t = 0; t < nt; t++) { Start Send B to rank-1, Start Send F to rank+1 (MPI\_Isend) Start Recv A from rank-1, Start Recv G from rank-1 (MPI Irecv) Compute rows C--E Waits for finishing all communications (MPI\_Wait) Compute rows B, F computations are Switch old and new arrays divided

$$T=max(T_N,T_{P1})+T_{P2}$$

#### Another Improvement: Reducing Communication Amounts



Multi-dimensional division may reduce communication





Each process communicate with upper/lower/right/left processes

- Comp: O(mn/p)
- Comm: O(n)

per 1 process, 1 iteration

- Comp: O(mn/p)
- Comm: O((m+n)/p<sup>1/2</sup>)
- per 1 process, 1 iteration
- → Comm is reduced

## Multi-dimensional division and Non-contiguous data (1)

 MD division may need communication of noncontiguous data
Comm



In Row-major format, we need send/recv of noncontiguous data for left/right borders

But "fragmented communication" degrades performance! (since Latency > 0) How do we do?

## Multi-dimensional division and Non-contiguous data (2)



Solution (1):

- Before sending, copy non-contiguous data into another contiguous buffer
- After receiving, copy contiguous buffer to noncontiguous area
- Solution (2):
- Use MPI\_Datatype
  - Skipped in the class; you may use Google :-p

#### It is Better to Use Collective Communications if Appropriate

• Comparing MPI\_Bcast and MPI\_Send&Recv In the graph, rank 0 called MPI\_Send for p-1 times to other processes



In most cases, MPI\_Bcast is faster

#### Why are Collective Communications Fast?



 Since Scalable communication algorithms are used inside MPI library



#### Binomial tree algorithm



## One of Scalable "Bcast" Algorithms

- Scatter&Allgather algorithm
  - Message is divided into p parts
  - Better than "binomial tree" if M is larger



# pL + M/B + (log p)L + M/B

R. Thakur and W. Gropp. Improving the performance of collective operations in mpich. EuroPVM/MPI conference, 2003.





- We have finished
  - Part 1: OpenMP for shared memory parallel programming
  - Part 2: MPI for distributed memory parallel programming
- Why are "parallel programs" slower than expectation?
  - Ideal: "p times speed-up with p processor cores"

#### **Too Many Factors that Limit Performance of Programs**

- Factors in algorithm
  - Load imbalance between threads, processes
  - Bottlenecks due to mutual exclusions
  - Communication costs
- Factors related to OpenMP/MPI system
  - Too many parallel region
  - Too many message
- Factors related to hardware
  - Memory access costs
  - Congestion in network

and many, many factors



## How Should We Tackle Performance Limiting Factors?



- It is important to know "why it is slow now"
- Consider what should be measured in order to specify current problem
  - Measuring time part by part may be helpful
  - Comparing computation time and communication time separately
  - Comparing 1-node performance and multi-node performance may be helpful
- It is good to use knowledge of computer hardware

## **Assignments in this Course**

- There is homework for each part. Submissions of reports for 2 parts are required
- Also attendances will be considered



## Assignments in MPI Part (Abstract)



Choose <u>one of</u> [M1]—[M3], and submit a report Due date: May 29 (Monday)

[M1] Parallelize "diffusion" sample program by MPI.[M2] Improve mm-mpi sample in order to reduce memory consumption.

[M3] (Freestyle) Parallelize *any* program by MPI.

For more detail, please see <u>Apr 27 slides</u> or <u>OCW-i</u>.

#### **Next Class**

Part 3 starts

• GPU parallel programming using CUDA

