



Implementation and Evaluation of MPI Library with Globus Toolkit for Establishing λ Computing Environment

Mai IMOTO
Osaka University, Japan

Outline

- 4 Research background
 - λ computing environment
- 4 Objective of our research
- 4 Implementation of MPI library utilizing shared memory in λ computing environment
 - Experimental system with AWG-STAR system
- 4 Evaluation
- 4 Conclusion and future work

Research Background

- 4 Grid computing
 - Connect distributed computing nodes
 - Share resources and storages
 - Large scale computing distributed in wide area
 - Transmitting huge data
- 4 Overhead of transmission over TCP/IP
 - Overhead from packet processing
 - Re-transmission due to packet loss

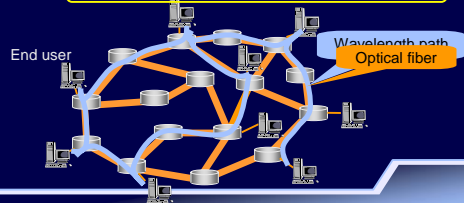


New technology which provides highly reliable, high speed connections for end users is needed

λ Computing Environment (1/2)

- 4 Connect computing nodes and optical routers by optical fibers
- 4 Establish wavelength paths on the fibers
- 4 Do not use TCT/IP
- 4 Data is exchanged on wavelength paths which are treated as a granular units

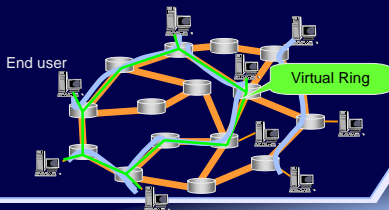
Release overhead of packet processing



λ Computing Environment (2/2)

- 4 Establish a virtual ring by connecting wavelength paths
- 4 By building a virtual ring, the optical ring is used as exclusive high speed transmission channel

Achieve highly reliable, high speed communication between end users



Establishing our Grid Environment with Globus Toolkit

- 4 Globus Toolkit
 - **Middleware** of the grid environment for communication, authentication and job control
 - Provides users with interface which is independent of implementation
 - De-facto standard of grid middleware
- 4 Adopt the Globus Toolkit as a upper layer of the λ computing environment

Users can perform high speed distributed computation without changing their original program

Objective of our Research

- 4 Adopt the Globus Toolkit to the λ computing environment
- 4 Implement and evaluate MPI library in the λ computing environment
 - Use the AWG-STAR system developed by NTT Photonics Laboratory

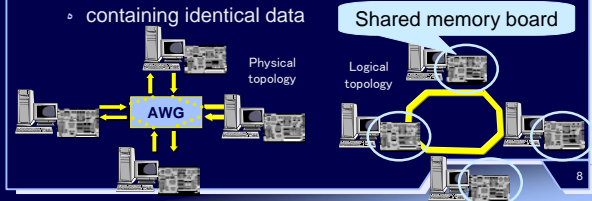
2005/11/11

APSITT 2005

7

AWG-STAR System

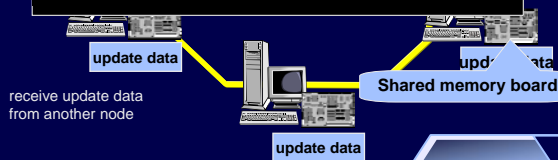
- 4 Information sharing network platform
- 4 Computing nodes connected to the AWG router in ring topology
- 4 Each node is equipped with a shared memory board
 - containing identical data



8

Data Sharing with the AWG-STAR System

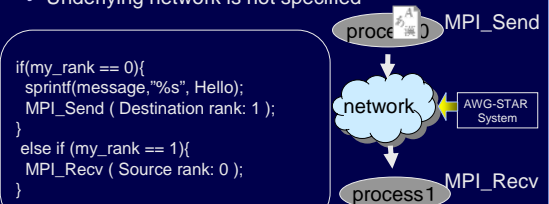
- 4 Data is shared by writing on the memory board
- 4 Reading from the shared memory does not need transmission
- 4 Writing delay time to access the memory board and delay time to go around the optical ring



9

MPI (Message Passing Interface)

- 4 Interface specification for data transmission in parallel computing
 - Underlying network is not specified



2005/11/11

APSITT 2005

10

Implementation of MPI library

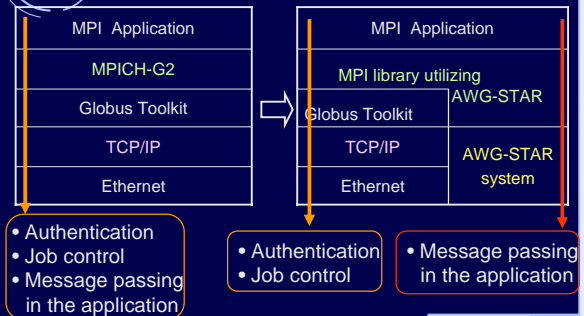
- 4 Implementation strategy
 - Create MPI library from scratch
 - Create MPI library based on another MPI library
- 4 Our MPI library bases on MPICH-G2 which works on the Globus Toolkit
 - Authentication and job control are same as MPICH-G2
 - Data exchange during the application is switched from TCP/IP to the AWG-STAR system

2005/11/11

APSITT 2005

11

Parallel Computing in the λ Computing Environment



2005/11/11

APSITT 2005

12

Implementation of MPI Library utilizing Shared Memory

- Dynamic memory allocation is not yet supported
 - Number of process: n
 - Shared memory is divided into $n \times n$ areas
 - Each area is used as a queue for one pair of sending and receiving processes

Area for one pair of sending and receiving process to transmit data

shared memory

2005/11/11 APSITT 2005 13

Message Passing

- Whole transmit data is written to the shared memory at a time

Packet division Header addition

TCP / IP

Packet reconstruction Header deletion

Sending node's local memory

Receiving node's local memory

Shared memory

2005/11/11 14

Sending Process

- When a sending function is called, process enqueues transmit data in the shared memory
- After writing, the process sends a signal to the receiving process
 - The signal is provided as a function of the AWG-STAR system

Data A process

Shared Memory

Receiving process

2005/11/11 APSITT 2005 15

Receiving Process

- Timings when the process receives the data and when the receiving function is called are different
 - Implementing two buffers in the local memory
 - Data buffer
 - Request buffer

Data A

Reading from shared memory

Shared memory

Data C Data B

Request buffer

Request A

Receive function

2005/11/11 16

Receiving Process

- Timings when the process receives the data and when the receiving function are different
 - Implementing two buffers on the local memory
 - Data buffer
 - Request buffer

Shared memory

Data C Data B

Request D

Request buffer

Request D

Receive function

2005/11/11 APSITT 2005 17

Evaluation - Application

- Evaluation by Himeno Benchmark
 - It measures the computational time in MFLOPS by solving Jacobi method
 - Three dimensional array is partitioned into smaller arrays which are assigned to all computing nodes
 - Every process calculates their partitioned array and transmits the boundary region to the next process by utilizing MPI.
 - Problem size is variable
 - Transmit data size are proportion to problem size
 - Number of message-passing is inverse proportion to problem size

2005/11/11 APSITT 2005 18

Evaluation - Simulation model

- Maximum number of computing nodes are 4
 - One node executes a single process
- Specification of the computing nodes
 - Distance between computing nodes: 20m
 - Both optical fibers for AWG-STAR system and Ethernet cables for TCP/IP
 - CPU: Xeon 3.0 GHz
 - OS: Redhat 7.3 Linux
- Specification of the shared memory board
 - Network interface speed: 2Gbps
 - Processing delay of a token on each node: 500ns
 - Access speed of shared memory from local memory: 60MB/s

2005/11/11 APSITT 2005 19

Evaluation by Comparing to the Ethernet

Executed by 2 processes

Executed by 4 processes

Lower performance increases and compared to Ethernet performance become better board is the bottleneck

MFLOPS

Problem size

2005/11/11 APSITT 2005 20

Communication and Calculation Time

AWG-STAR system

Ethernet

As number of processes increase, communication time become longer. Message passing are increased

100% Communication

0% Calculation

2005/11/11 APSITT 2005 21

Evaluation by the Number of Processes

MFLOPS

Number of processes

Compared to 2 processes and 4 processes. Because of overhead of message passing, 2 processes execute faster. For small data size, 2 processes execute faster

2005/11/11 APSITT 2005 22

Evaluation by the Number of Processes

MFLOPS

Number of processes

Large data size problem are solved when the executing processes increase

2005/11/11 APSITT 2005 23

Domain Decomposition

- 3-dimension array is parallelized by domain decomposition
 - Transmit data sizes are almost proportional to the size of the boundary region
 - In the case of decomposition 1x4x1, one process transmit twice the size of the case of decomposition 1x1x4

1x4x1

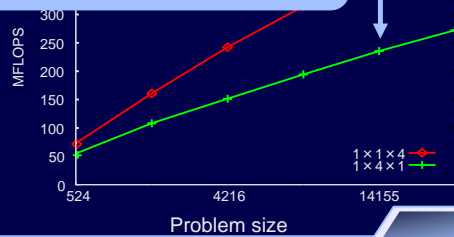
1x1x4

24

Evaluation by the Decomposition

Performance is affected by the transmit data size

- In the case when the problem size is 14155, performance is about 60%



25

Conclusion and Future Work

4. Conclusion

- Adopted the Globus Toolkit into the λ computing environment
- Implemented and evaluating MPI library in the λ computing environment
- Delay of access to the shared memory board is the bottleneck

4. Future work

- New architecture of the shared memory board is now investigated

2005/11/11

APSITT 2005

26