# Proposal and Evaluation of Realization Approach for a Shared Memory System in λ Computing Environment

Hirohisa NAKAMOTO[†], *Nonmember*, Ken-ichi BABA[††], *and* Masayuki MURATA[†], *Members*

**SUMMARY**   In this paper, we propose a new computing environment (which we will refer to λ computing environment) that provides a base for parallel computing among nodes distributed in the wide area. In our concept, virtual channels are provided utilizing optical networks connecting computing nodes. It can offer high-speed and reliable connection pipe among nodes, so that it is efficiently applicable to SAN (Storage Area Network) and/or Grid computing. In the environment, shared memory is constituted on a virtual ring of the photonic network. Consequently, it is not necessary to distinguish shared memory in a wide-area distributed system from a communication channel; thus high-speed data exchange between computing nodes on the ring can be achieved. The key to realizing such a computing environment is how to construct a shared memory system on the photonic ring. In this paper, we propose and evaluate a shared memory system suitable to the virtual optical ring network, which takes into account contention resolution, cache coherency, and synchronization methods because the propagation delays among nodes are much larger than the conventional shared memory system. Through simulation experiments of using three benchmark programs as representative parallel computing applications, we show the applicability of our shared memory system on the wide-area virtual photonic ring.
*key words:   λ Computing Environment, Shared Memory, Cache Coherency, Memory Access Contention*

## 1.   Introduction

As users of networks such as the Internet have increased, so has the amount of traffic steadily increased. Various applications that utilize images have come to be used, and the demands made on the technology that enables the high speed and large scale transmission in a network have increased. To satisfy these demands, research into optical transmission technology has been actively pursued. Research into WDM technologies that use multiplexed light wavelengths have been the main target for development and technology from new WDM research that can use 1000 wavelengths has also been advanced [1]. In recent years, IP over a WDM network has been studied and developed to provide high-speed transmission on the Internet based on WDM technology. Moreover, standardization of the routing technology of the Internet, called GMPLS, which is the com-

Manuscript received January 1, 2003.
Manuscript revised January 1, 2003.
Final manuscript received January 1, 2003.
    [†]The author is with Department of Information Networking, Graduate School of Information Science and Technology, Osaka University, Suita, Osaka 565-0871, Japan
    [††]The author is with Cybermedia Center, Osaka University, Ibaraki, Osaka 567-0047, Japan

munication technology that uses various optical technologies for a lower layer than the WDM technology, has also been advanced in IETF [2]. Further, aiming to realize the true IP communication of a photonic network, research into optical packet switches based on optical technology has also begun [3]–[9].

However, many such technologies presuppose the existing Internet technology. That is, an IP packet is treated as a degree of granule treating information, and it is made into the target for research and development of how to carry it at high speed on a network. Therefore, as long as architecture based on packet switching technology is focused on, realization of high quality communication to each connection will be very difficult. New technologies such as SAN and Grid computing need to be applied to provide end users with a high speed and reliable communication pipe; for that, a mass wavelength path needs to be set up between end users and provided for users. That is, it is possible to provide an end user with a ultra high-speed and high quality communication pipe by building a photonic network that uses established fibers, or newly laid fiber if needed, and by utilizing wavelengths multiplexed in the fiber as the minimum particle size for information exchanges.

As middleware aimed at the realization of a high-speed distributed computation environment using an optical network, OptIPuter is proposed [10]. It has been studied and developed to build a Grid environment established on optical networks. It also provides virtual communication paths. However, this is based on present Internet technology and treats a packet as the informational particle size; so that the problem mentioned previously regarding packet processing arises again.

Thus we propose a new architecture, the λ computing environment that has virtual channels utilizing optical fibers connecting computing nodes. In the conventional Grid environment, data is exchanged with the message passing using TCP/IP. In the λ computing environment, by realizing communication between nodes on the Grid not by conventional TCP/IP but by established wavelength paths, we can achieve highly reliable high-speed communication. Then, by making virtual channels on a mesh upon the photonic network that is connected to the network nodes and the computer nodes with optical fibers, distributed computa-

tion on a high speed channel is enabled. Moreover, it is possible to utilize wavelengths as a shared memory by constituting a virtual ring in the $\lambda$ computing environment [11]. As a result, it is not necessary to distinguish shared memory from communication channels in a wide-area distributed system; and we expect that the high-speed data exchange between computing nodes can be achieved (see Fig. 1). In our research group, we also use wavelengths as a high-speed transmission channel and implemented high-speed access method to the shared memory that exists on each computing node [12].

In this paper, we propose and evaluate an approach to realize a shared memory system using a virtual optical ring network. Specifically our shared memory system uses a level-1 cache in a CPU of each computer group as the cache of such a shared memory. When using the virtual ring as a shared memory, it is necessary, unlike a bus between a CPU and the shared memory in a computer, to consider restrictions in timing and the frequency of access, since the shared memory is spread out on a long-distance optical fiber; so to take into consideration coherency between the shared memory in the virtual ring and the cache of each computer group more strictly than the conventional shared memory system. Next we have to solve the contention of shared memory access like a conventional shared memory system. This problem arises in cases where a processor has not finished writing access to the data on the shared memory, and another processor tries to read or write to the same data. When we perform parallel computation on shared memory using a virtual optical ring network, synchronization is needed to collaborate between computing nodes. However, it is not necessary to distinguish the shared memory in a wide-area distributed system from a communication channel so that it appears that the high-speed data exchange between computing nodes is achieved. As noted above, after considering such features, we propose a shared memory access method for the $\lambda$ computing environment, and evaluate the method through simulations.

The rest of the paper is organized as follows. In Section 2, we describe the cache coherency problem, memory access contention, and synchronization between computing nodes of the conventional shared memory system. In Section 3, we propose a realization approach for a shared memory system in the $\lambda$ computing environment. In Section 4, we evaluate our approach using a benchmark program for parallel computing. Finally, we conclude the paper and describe future work in Section 5.

## 2. Conventional Shared Memory System

To realize a shared memory system, it is necessary to avoid memory access contention, maintain cache coherency [13], [14], and realize synchronization between
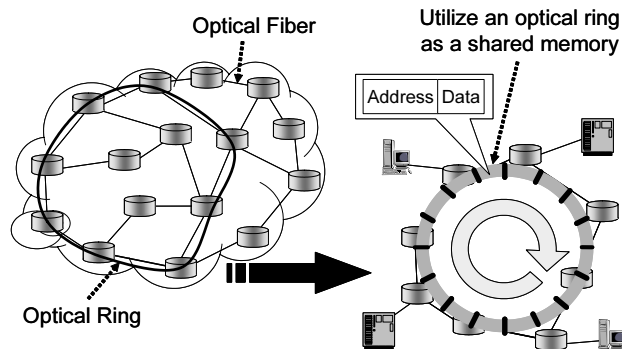


**Fig. 1** $\lambda$ Computing Environment.

computing nodes to enable collaboration. In this section, we briefly explain these techniques.

### 2.1 Contention Avoidance in a Shared Memory System

A problem exists when a processor has not finished writing access to data on the shared memory, and another processor tries to read or write to the same data. This kind of contention is solved by a lock mechanism. Each processor has to lock the shared memory before it tries to write the data on the shared memory. By using a lock mechanism, write access to the shared memory is protected. When we realize a shared memory system in the $\lambda$ computing environment, we have to resolve this problem.

### 2.2 Cache Coherency

When each processor has a cache, it is necessary to fully take into consideration the consistency between the data on the cache and the data on the shared memory. Two ways, a directory method and a snoop cache method, are techniques for generally maintaining cache coherency. In realizing shared memory on the $\lambda$ computing environment, access to a directory table may become a bottleneck when a directory method is adopted. So a snoop cache method is adopted in this paper.

The snoop cache method offers cache consistency between data on caches and on the shared memory. It snoops memory access on a share bus, and performs consistency control to a local cache block by the distributed technique if needed. When a processor tries to read data, firstly it searches in the local cache, and when data does not exist there, it accesses to the shared memory. When a processor tries to read or write to a shared memory, consistency control is not needed if another processor does not the same data in its local cache. However, if another processor has the same data in its local cache, some methods can be considered using a method to keep cache consistency. Moreover, when a processor writes or updates the data on its cache, keeping cache consistency becomes still more complicated
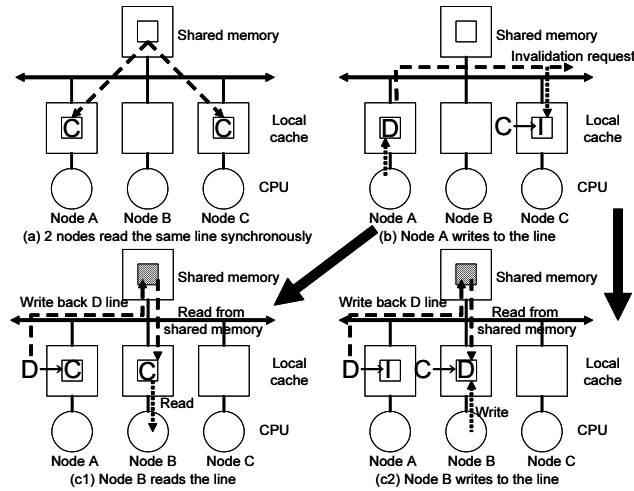
**Fig. 2** Behavior of Write-back Invalidation Protocol.

and there are some ways to keep consistency.

Such cache consistency protocols are classified into four types according to the timing (write-through, write-back) and the method (invalidation, updating). Among those, a write-back invalidation protocol has the least shared memory access. When the access delay time to a shared memory is large, this protocol is very effective. So we adopt this protocol when we realize a shared memory system in the $\lambda$ computing environment. The write-back invalidation protocol is simply explained below.

In the write-back invalidation protocol, data in the local cache has three states; Invalid (I), Clean (C), and Dirty (D). The I state means that data is invalidated and cannot be used, the C state means that the data on the cache is the same value compared to the data on the shared memory, and the D state means that the data on cache is not the same value compared to the data on shared memory.

When one or some computers refer to an address, the data is copied to A cache from the shared memory and that cache's state becomes a clean state as shown in Fig. 2 (a). Since the value of the data on the shared memory and the data in the C state is the same, read access to this cache block does not need bus operation. If a processor writes to data in a C state, the state will become a D state. At this time, the control message requesting invalidation of the relevant data is sent on a bus. Since the cache controllers of other processors snoop a bus, they receive the control message and invalidate the relevant data in their local cache (Fig. 2 (b)). Henceforth, read and write accesses to the data in a D state do not need bus operation. When other processors read to the data in a D state, the data in a D state is written back to the shared memory, and cache consistency is completed. Next, the data is sent by bus to the processor that requested the read access and states of cache blocks on both processor's cache

will become a C state (Fig. 2 (c1)). On the other hand, when other processors write to the address of the data in a D state, like reading, writing back of the data in a D state to the shared memory takes place, and the data is sent to the processor which sent the demand message. Finally, the demanding processor writes the data on a local cache, and the state of the data will be D. The cache data on the processor that has the original data are invalidated (Fig. 2 (c2)).

## 2.3 Synchronization between Computing Nodes

When we perform application programs of parallel computation on a shared memory system, synchronization between computing nodes is needed to collaborate. Kinds of synchronization include atomic operation, caching of synchronous variable, blocking control on a shared memory, a memory lock, and the barrier synchronization method. Barrier synchronization is used when each processor needs to wait until all processors reach to the same break point.

Application programs used for evaluation in this paper first calculate locally and then after local calculation perform synchronization processing. In barrier synchronization, each computing node reads and updates the data of synchronous variable by turns. So a synchronous memory on an optical ring is suited to barrier synchronization because data is circling on the optical ring and each computing node can easily read and update the data on a synchronous memory by turn. So in this paper we adopt the barrier synchronization method. The method of realizing barrier synchronization using shared memory on an optical ring is shown in Sec. 3.4.

## 3. Realization Approach for a Shared Memory System in the $\lambda$ Computing Environment

In this section, we explain an approach for realizing a shared memory system in the $\lambda$ computing environment. Firstly, we describe our network model. Like a conventional shared memory system, we have to realize cache coherency, avoidance of memory access contention, and synchronization between computing nodes to enable collaboration. We also describe how to resolve these points needing consideration.

## 3.1 Network Model on Consideration

We show a network model in Fig. 3. Computing nodes that compose the $\lambda$ computing environment are connected with optical fibers that make a virtual ring network. In this paper, I presuppose that each computing node has one CPU, a level-1 cache, and a local memory. A local memory is used for storage of programming codes and local data, a shared memory is used for storage of the shared data that all computing nodes use
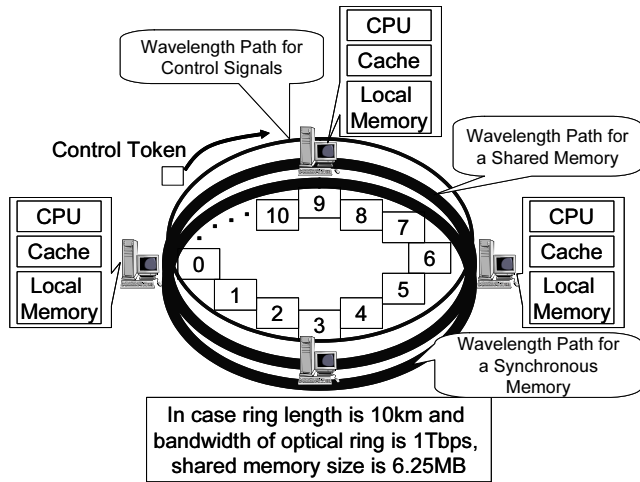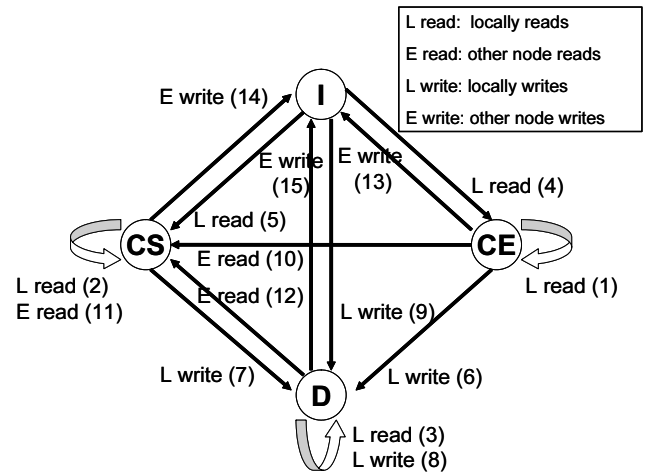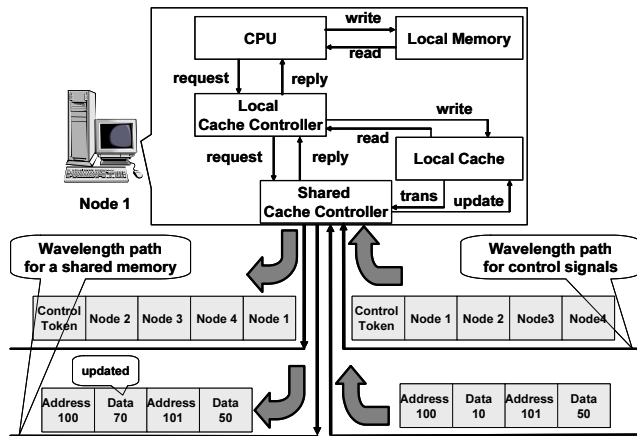
**Fig. 3** Network Model.



**Fig. 4** Cache Control Scheme and Data Flows.



**Fig. 5** State Transition Diagram.

in computing and a synchronous memory is used for synchronization between computing nodes. An optical ring network has a wavelength path for shared memory, a wavelength path for control signals and a wavelength path for synchronous memory. The bandwidth of an optical ring for shared memory is set to 1Tbps. Since propagation delay time is 5 ns/m, we can use an optical ring network as a shared memory, of which the capacity is equivalent to 6.25MBytes. We usually use these parameters except otherwise explicitly stated. For comparative evaluation, we also use the model with 1km ring length and 10 number of rings in parallel, and with 100m ring length and 100 number of rings while these models might be unrealistic even in the near future in Sec. 4.3. The processing delay time in the interface of each computing node and middle nodes, such as network devices which constitute an optical ring network, is not explicitly taken into account here. Indeed, we assume that it is included in the propagation delay time.

Next, we show a configuration of each computing node and data flows between a computing node and

wavelengths of an optical ring in Fig. 4. The local cache controller searches the data on the local cache when it receives read or write request to the data on the shared memory and the synchronous memory from the CPU. If the data requested by the CPU is not found in the local cache, the local cache controller commits the request from the CPU to the shared cache controller. The shared cache controller monitors wavelengths for control signals, a shared memory, and a synchronous memory. When it receives the request from the local cache controller, it reads or updates the data on the shared memory and the synchronous memory. And when it receives the control messages from other nodes, it changes the state of the cache block on a local cache to keep cache coherency.

## 3.2 Contention Avoidance

Like a conventional shared memory system, we have to solve contention to the shared and synchronous memories to realize a shared memory system in the $\lambda$ computing environment. To resolve this problem, we adopt a lock mechanism. That is, each computing node has to send a lock message using a control token before it tries to write or update data in state C on a local cache, and data on the shared memory or the synchronous memory. By using this mechanism, write access to the shared memory is protected. The approach to realize contention avoidance is explained in Sec. 3.3 along with cache coherency.

## 3.3 Cache Coherency

As mentioned in Sec. 2.2, the write-back invalidation protocol has the least shared memory access. It is very effective because in our shared memory system, the access delay time to a shared memory is large. Thus, there are some write-back invalidation protocols; such as the Illinois [15] and Symmetry protocols [16]. We

adopt the Illinois protocol because its consistency control is the simplest. In the Illinois protocol, data of a local cache has four states; Invalid (I), Clean Exclusive (CE), Clean Shared (CS) and Dirty (D). I and D states are the same states like in the conventional protocol. The CE state means that the data on the cache is the same value compared to the data on the shared memory and another processor does not have the same data, and the CS state means that another processor has the same data in its local cache. However, the Illinois protocol presupposes the shared memory system using a shared bus. So we have to adapt Illinois protocol to the shared memory system in the $\lambda$ computing environment.

So we propose a cache coherency protocol that solves cache coherency and contention to the shared and synchronous memories in the $\lambda$ computing environment on the basis of the Illinois protocol. We show the state transition diagram of our proposed cache coherency protocol in Fig. 5. Hereafter, we first explain each arrowhead of the transition in Fig. 5. Second, we explain the behavior of the computing nodes that received a control message. In our cache coherency protocol, five control messages are used. A line copy request message is sent by a computing node that tries to read the data when it does not have the data in its local cache. A line move request message is sent by a computing node that tries to write the data when it does not have the data in its local cache. Lock and invalidation request messages are sent by a computing node that tries to update the data in the CE or CS states or write the data in a I state. A lock request message is used to avoid contention to shared and synchronous memories. And an invalidation request message is used to keep that the cache block in the D state is only one among all computing nodes for maintaining cache coherency. An unlock request message is sent by a computing node that has finished updating or writing the data.

At first, we explain the case a computing node that reads the data from A local cache. In this case, the state of the cache block (Fig. 5 (1), (2), (3)) does not change.

And, we explain how a computing node reads data from another computing node's cache or shared memory. A computing node attaches a line copy request message to a control token after catching the control token. It waits until the control token returns. If it receives a line copy ack (Fig. 5 (5)), it has to wait until the copy of the cache block is sent by another computing node and reads the cache block from the control token. Then it changes the state of the cache block $(I \rightarrow CS)$. If it does not receive a line copy ack (Fig. 5 (4)), it directly reads the cache block from the shared memory. Then it changes the state of the cache block in its local cache $(I \rightarrow CE)$.

Next, we explain how a computing node updates the data on a local cache. If the state of the cache block is D (Fig. 5 (8)), it only updates the data on the local cache. If the state of the cache block is CE (Fig. 5 (6)), it updates the data on the local cache and changes the state of the cache block $(CE \rightarrow D)$. If the state of cache block is CS (Fig. 5 (7)), contention avoidance processing is needed. After catching the control token, it has to check whether a lock request message of another computing node is attached to the control token. If a lock message is attached, it has to wait until an unlock message is sent and then restart the write processing. If a lock message is not attached, it attaches a lock request message and an invalidation request message to the control token. Then it updates the data on the local cache and changes the state of the cache block $(CS \rightarrow D)$. Finally, it attaches an unlock request message to the control token after catching the control token.

Next, we explain how a computing node tries to write when it does not have the copy cache block on the local cache (Fig. 5 (9)). In this case, contention avoidance processing is also needed. A computing node searches the lock table. If the address of the cache block that it tries to write is registered in the lock table, it has to wait until an unlock request message is sent. If the address is not registered, it has to check whether a lock request message of other computing node is attached to the control token after catching the control token. If a lock message is attached, it has to wait until an unlock message is sent and then restart the write processing. If a lock message is not attached, it attaches a lock request message and a line move request message to the control token. It waits until the control token returns. If it receives a line copy ack, it has to wait until a cache block is sent by another computing node and read the cache block from the control token. Then it updates the cache block and changes the state of the cache block $(I \rightarrow D)$. If it does not receive a line copy ack, it directly reads the cache block from the shared memory. Then it updates the cache block and changes the state of the cache block $(I \rightarrow D)$. After updating the cache block, it attaches an unlock request message to the control token after catching the control token.

Next, we explain the behavior of computing nodes that received a control message. At first, we explain the behavior of a computing node that receives a line copy request message. If a computing node has a cache block of the relevant address in the CE state (Fig. 5 (10)), it attaches a line copy back and a copy of cache block to the control token after catching the control token. Then, it changes the state of the cache block $(CE \rightarrow CS)$. If a computing node has a cache block of the relevant address in the CS state (Fig. 5 (11)), it attaches a line copy back and a copy of the cache block to the control token after catching the control token. If a computing node has a cache block of the relevant address in the D state (Fig. 5 (12)), it attaches aline copy back to the control token after catching the control

token. Then, it writes back the cache block in the D state to the shared memory and changes the state of the cache $(D \rightarrow CS)$. It attaches a copy of the cache block to the control token after catching the control token.

Next, we explain the behavior of a computing node that receives a line move request message. If a computing node has a cache block of the relevant address in the CE state (Fig. 5 (13)), it attaches a line move ack and copy of the cache block to the control token after catching the control token. Then, it changes the state of the cache block $(CE \rightarrow I)$. If a computing node has a cache block of the relevant address in the CS state (Fig. 5 (14)), it attaches a line move ack and a copy of the cache block to the control token after catching the control token. Then, it changes the state of the cache block $(CS \rightarrow I)$. If a computing node has a cache block of the relevant address in D the state (Fig. 5 (15)), it attaches a line move ack to the control token after catching the control token. Then, it writes back the cache block in the D state to the shared memory and changes the state of the cache block $(D \rightarrow I)$. It attaches a copy of the cache block to the control token after catching the control token.

Next, we explain the behavior of a computing node that receives lock, unlock and invalidation request messages. If a computing node receives a lock request message, it registers the relevant address to its lock table. If a computing node receives an unlock request message, it removes the relevant address from its lock table. If a computing node receives an invalidation request message, it searches the relevant address in the local cache. If a computing node has a relevant cache block in its local cache, it changes the state of the cache block to the I state.

## 3.4 Synchronization between Computing Nodes

We explain the method for realizing barrier synchronization in the shared memory system using an optical ring network. First, part of the wavelength paths of an optical ring are allocated to the synchronous memory area. When a synchronous memory is accessed, a Fetch & Decrement operation like in the conventional method is indivisibly performed. That is, it ensures that access to a synchronous memory indivisibly causes a subtraction processing of the relevant data. Since only one computing node can simultaneously access the synchronous memory, when using an optical ring network for a synchronous memory, execution of an atomic operation is easy. With an application program, in bringing about synchronization among some computing nodes, each node accesses to the synchronous memory. The number of processors is set in the synchronous memory as the initial value. The value of the synchronous memory will be set to zero if all nodes access. If the value of a synchronous memory is set to zero, all nodes will finish the synchronous process and begin the next
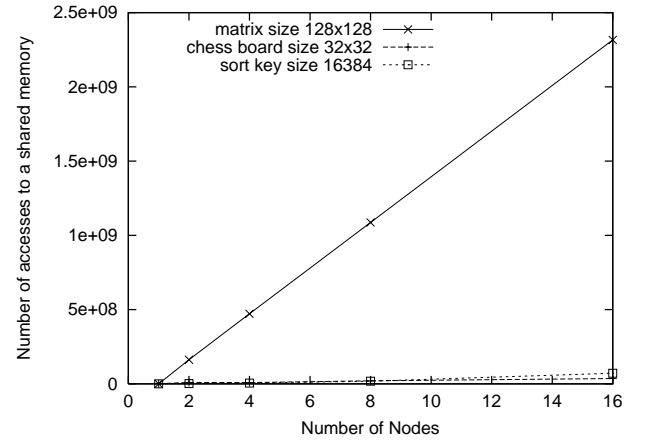


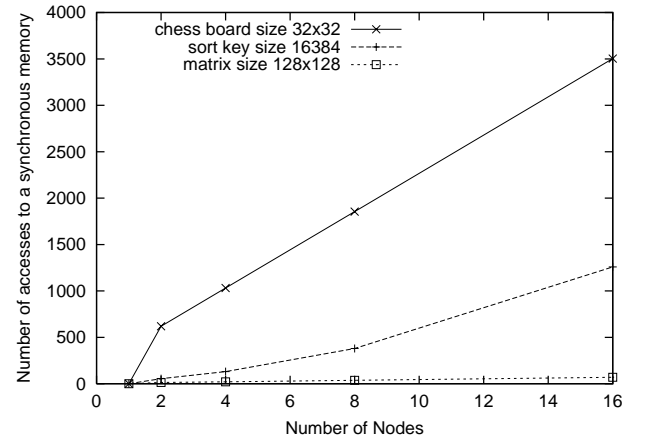**Fig. 6** Number of accesses to a shared memory (ring length 10km).



**Fig. 7** Number of accesses to a synchronous memory (ring length 10km).

processing.

## 4. Performance Evaluation

In this section, we evaluate through simulation the performance of the shared memory access method proposed in the previous section. We utilized the ISIS library [17], [18] currently developed at the Amano Laboratory of Keio University in coding the simulation program.

## 4.1 Simulation Model

We used the following network model. Each computing node in the $\lambda$ computing environment is interconnected with optical fibers and nodes are configured to virtually form the ring topology. Each computing node has one CPU, a level-1 cache, and a local memory. Clock frequency of a CPU is 1GHz, the capacity of a level-1 cache is 512KB, and the capacity of a local memory is 2GByte. It assumes that computing nodes are put

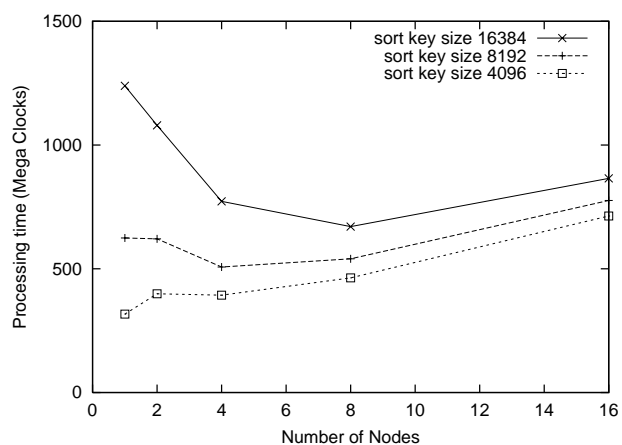**Table 1**  Characteristics of Memory Accesses (ring length 10km).

| | nodes | clocks | shared accesses | synchronous accesses | L1 hit |
|---|---|---|---|---|---|
| "radix sort" | 1 | 1239504563 | 116224 | 0 | 59.9622% |
| (16384 keys) | 2 | 1079750949 | 1296305 | 55 | 17.2424% |
| | 4 | 772450953 | 4822926 | 131 | 16.6878% |
| | 8 | 670200953 | 17821287 | 379 | 16.7660% |
| | 16 | 865300953 | 71520173 | 1259 | 16.8651% |
| "product of a matrix" | 1 | 2048357893 | 81920 | 0 | 60.2431% |
| (128×128) | 2 | 2665750970 | 162124605 | 13 | 79.5476% |
| | 4 | 2615251057 | 471528629 | 21 | 90.0370% |
| | 8 | 2539751057 | 1086505926 | 37 | 95.2501% |
| | 16 | 2506851057 | 2316505719 | 69 | 97.5408% |
| "queen problem" | 1 | 110045031 | 499250 | 0 | 30.4880% |
| (32×32) | 2 | 2494650912 | 8964784 | 619 | 15.6490% |
| | 4 | 2325700916 | 8869410 | 1031 | 17.0453% |
| | 8 | 2074550916 | 20263999 | 1855 | 15.9305% |
| | 16 | 2107500916 | 35978138 | 3503 | 17.6546% |

on the optical ring network with equal distances. An optical ring network has wavelength paths for shared memory, control signals, and synchronous memory.

The bandwidth of an optical ring for shared memory is set to 1Tbps. Since propagation delay time is 5 ns/m, we can use an optical ring network as a shared memory, of which the capacity is equivalent to 6.25MBytes. We usually use these parameters except otherwise explicitly stated. For comparative evaluation, we also use the model with 1km ring length and 10 number of rings in parallel, and with 100m ring length and 100 number of rings while these models might be unrealistic even in the near future. The processing delay time in the interface of each computing node and middle nodes, such as network devices which constitute an optical ring network, is not explicitly taken into account here. Indeed, we assume that it is included in the propagation delay time.

To evaluate the performance, we use a Splash2 benchmark program [19], such as the "radix sort" program that sorts the sequence of an integer value using a radix sort algorithm. We also use the "product of a matrix" program that calculates the product of $n \times n$ matrix and the "queen problem" program that solves the n-queen problem. See Table 1 for comparing the characteristic of memory accesses for these programs. The numbers show the order of frequencies in memory accesses of programs. See also Figs. 6 and 7 for actual data.

We note that sample programs that we are using in this paper are just intended to see how levels of parallelism affects the performance. In actual, the program could be clearly divided into independent tasks when applied to parallel computation, and we can enjoy a parallelism of computation by an increasing number of nodes, but some part of the entire program needs synchronization to an extent, and to effect on the performance depends on the problem. The three programs that we have chosen here are typical examples and have different characteristics as indicated in Table 1. Our



**Fig. 8**  Processing time of "radix sort" program (ring length 10km).

intention here is to test whether the typical parallel algorithm can work well in a sense that total parallel execution time is not unacceptably increased and an introduction of parallelism does not only result in the increasing execution time.

### 4.2  Comparisons: Basic Results

We show results of execution time for each application program by setting the ring length to 10km. The number of execution clocks in CPU for the "radix sort" program is first shown in Fig. 8. The numbers of keys for sorting are set at 4096, 8192, and 16384. When the number of key is 4096, the advantage of parallel computation cannot be observed even if the number of computing nodes is increased. This is because the ratio of synchronous operation to total operation is large. However, as problem size becomes larger such as 8192, it turns out that the advantage of parallel computation appears. As the number of nodes exceeds some number (4 in the case that the key size is 8192), the execution time is gradually increased because the number of synchronization becomes large by an increasing number of
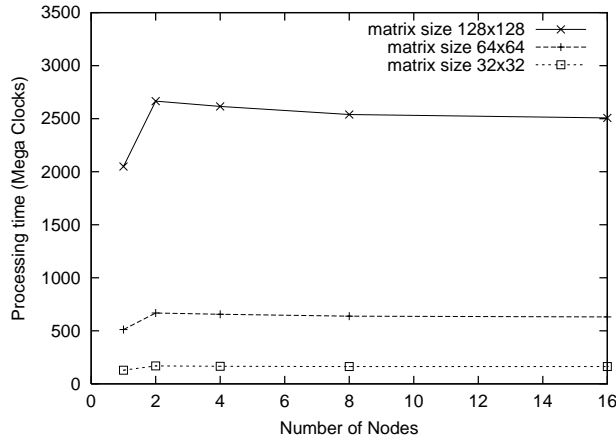
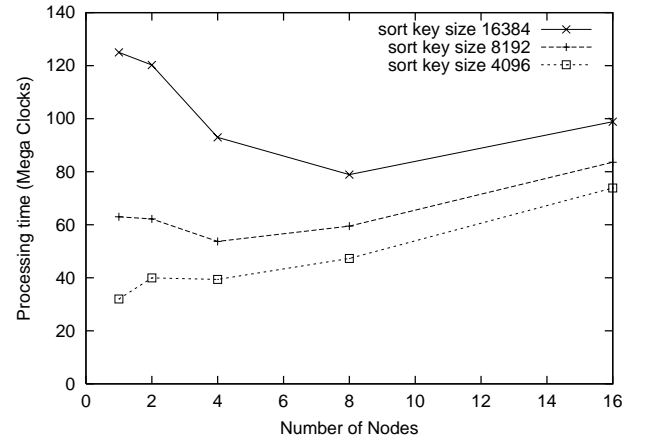**Fig. 9**  Processing time of "product of a matrix" program (ring length 10km).



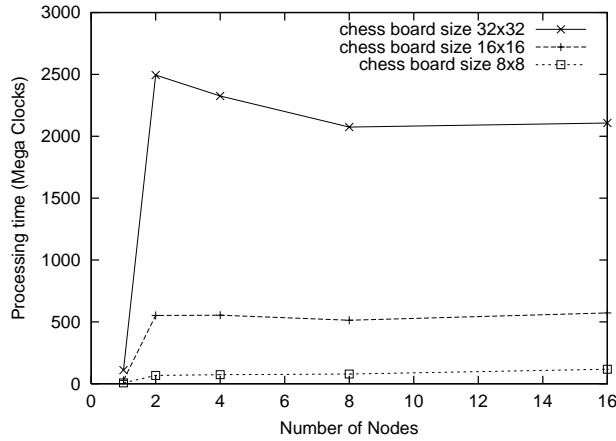**Fig. 11**  Processing time of "radix sort" program (ring length 1km).



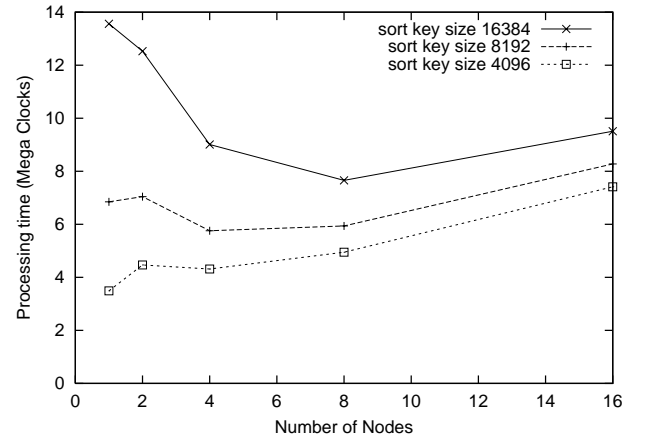**Fig. 10**  Processing time of "queen problem" program (ring length 10km).



**Fig. 12**  Processing time of "radix sort" program (ring length 100m).

nodes. When the sort key size is 16384, such a tendency becomes clearer. From these results, we found that the shared memory and access method for $\lambda$ computing environment are effective in parallel computation for the "radix sort" program when the number of parallel nodes is not so large.

The case of "product of a matrix" program is next shown in Fig. 9. The matrix sizes are changed from $32 \times 32$ to $128 \times 128$. The advantage of parallel computation becomes smaller in this case. It is because the number of accesses to the shared memory is large compared to other application programs as shown in Fig. 6, where the numbers of accesses to the shared memory are compared in three programs, and it compensates for introduction of parallel computation. However, it is still true that it does not introduce the additional delay if the "product of a matrix program" does not occupy the large portion of the entire program.

The same tendency can be observed in Fig. 10 where the "queen problem" program is considered. Its

problem size is changed from $8 \times 8$ to $32 \times 32$. In this case the number of synchronous accesses is much larger than other programs as shown in Fig. 7, where the numbers of accesses to the synchronous memory are compared in three programs. However, we can again see that the execution time is at least not increased even if the number of nodes becomes large.

### 4.3 Effect of Increasing the Optical Ring Length for Parallel Computation

When ring length is 10km long, the access delay time to the shared memory and the synchronous memory are large and may be the main factor that compensates for introduction of parallel computation. And when ring length becomes shorter, the advantage of parallel computation may become clearer. Accordingly, we next investigate the effect of decreasing the ring length on parallel computation time. We use three values of ring lengths: 100m, 1km and 10km.

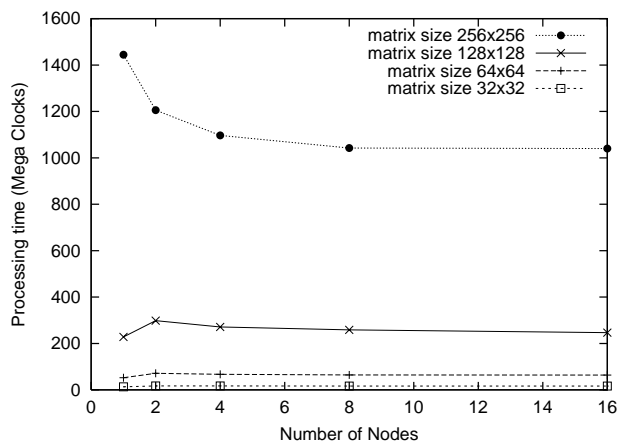In Figs. 11 and 12, the execution times of the

**Fig. 13** Processing time of "product of a matrix" program (ring length 1km).



**Fig. 15** Speed up ratio of "queen problem" program by changing the caching policy of the synchronous memory (ring length 10km).
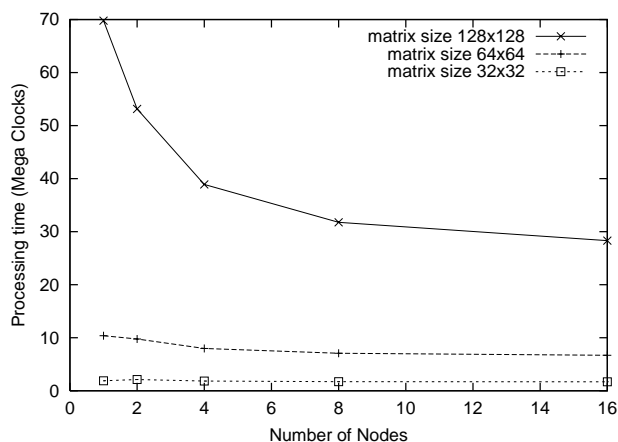


**Fig. 14** Processing time of "product of a matrix" program (ring length 100m).

"radix sort" program are shown against 100m and 1km of ring length. Compare with Fig. 8 where the case of 10km ring length was shown. We can here see that the results are almost same in three cases.

A different behavior is observed in using the "product of a matrix" program. See Figs. 9, 13, and 14. When the ring length was 10km (Fig. 9), there is no effect of parallel computation due to the access delay time is large much larger than the shared memory time. However, when the ring length is 1km (Fig. 13) and the matrix size is large enough ($256 \times 256$), the effect clearly appears when the number of nodes is less than eight. It is because the hit ratio of the level-1 cache becomes high as the matrix size becomes large. Also, when the ring length is 100m (Fig. 14), the effect of parallelism is attained even if the matrix size is small ($64 \times 64$ or $128 \times 128$). As we have already mentioned, the "product of a matrix program" has a characteristic that the number of accesses to the shared memory is the largest among all the application programs. See Fig. 6. Then
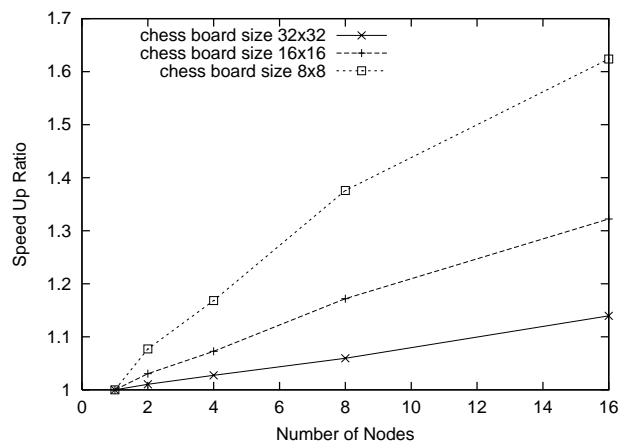
when ring length is short, the effect of parallel computation is obtained. Moreover, the number of accesses to the synchronous memory is the smallest, and therefore, parallel computation can be easily improved as the number of nodes increases.

In the "queen problem" program, on the other hand, no advantage of parallel computation is obtained even if the ring length is changed. It is due to the largest number of synchronous accesses in three programs even though the accesses to the shared memory do not frequently occur.

## 4.4 Synchronization Improvement

As mentioned in Sec. 4.3, synchronization has a great influence on the performance of parallel computation. We thus propose a method to improve synchronization time. As mentioned in Sec. 3.4, we use barrier synchronization to enable collaboration between parallel nodes. The data used for barrier synchronization are stored on synchronous memory and the copy of data are stored on the local cache of the node. Since data used for barrier synchronization are referred by each node only once, the performance must be improved by not storing data on synchronous memory to the local cache. This is because cache coherency processing is not needed in this case, and it is sufficient for each node to read or write to synchronous memory only when data on synchronous memory are not stored in the local cache.

Lastly, we compare the performance of the above–mentioned caching policy to that of the original one. For this purpose, we present the speedup ratio of the new policy compared to the original one. Noting that the speedup ratio of the "radix sort" program and the "product of a matrix" program are small, we only show the speedup ratio for the "queen problem" program in Fig. 15. The ring length is assumed to be 10km. The larger speeding ratio can be detained by the increas-

ing number of nodes. It is due to the fact that the "queen problem" program requires the largest number of accesses to the synchronous memory among the programs that we have tested. Then, the effect of the synchronization mechanism that we have introduced in this subsection becomes clear.
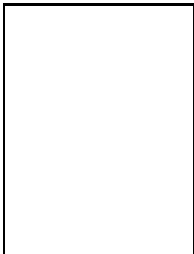
## 5. Conclusion

In this paper, we have proposed the shared memory access method in realizing the shared memory on photonic network. Moreover, we have evaluated the performance of the proposed method using the benchmark program for parallel computing. As a result, we show that the effectiveness of using optical ring as a shared memory and of parallel processing by the increase in the number of nodes when number of synchronous processing is small. We can see the future possibility of all-optical parallel computing environment in the wide-area. An efficient shared memory access method and a practical use of a local memory is due to be considered in the future.
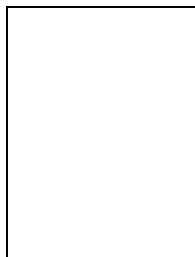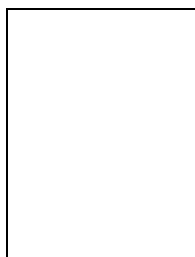
## Acknowledgements

## References

[1] M. Murata and K. Kitayama, "A 1,000-channel WDM network can resolve network bottleneck," in *Proceedings of the 7th Asia-Pacific Conference on Communications (APCC 2001) (Tokyo)*, pp. 113–116, Sept. 2001.

[2] E. L. Berger, "Generalized multi-protocol label switching (GMPLS) signaling functional description," *IETF RFC3471*, Jan. 2003.

[3] T. Yamaguchi, K. Baba, M. Murata, and K. Kitayama, "Packet scheduling for WDM fiber delay line buffers in photonic packet switches," in *Proceedings of OptiComm 2002*, vol. 4874, pp. 262–273, July 2002.

[4] K. Baba, R. Takemori, M. Murata, and K. Kitayama, "A packet scheduling algorithm for the 2x2 photonic packet switch with FDL buffers," in *Proceedings of 28th Europian Conference on Optical Communication 2002 (ECOC2002)*, Sept. 2002.

[5] S. L. Danielsen, B. Mikkelsen, C. Joergesen, T. Durhuus, and K. E. Stubkjaer, "WDM packet switch architectures and analysis of the influence of tuneable wavelength converters on the performance," *IEEE Jounal of Lightwave Technology*, vol. 15, pp. 219–227, Feb. 1997.

[6] S. L. Danielsen, C. Joergesen, B. Mikkelsen, and K. E. Stubkjaer,, "Analysis of a WDM packet switch with improved performance under bursty traffic conditions due to tuneable wavelength converters," *IEEE Journal of Lightwave Technology*, vol. 16, pp. 729–735, May 1998.

[7] D. Hunter, M. C. Chia, and I. Andonovic, "Buffering in optical packet switches," *IEEE Journal of Lightwave Technology*, vol. 16, pp. 2081–2094, Dec. 1998.

[8] K. L. Hall and K. A. Rauschenbach, "All-optical buffering of 40-gb/s data packets," in *IEEE Photonic Technology Letters*, vol. 10, pp. 442–444, 1998.

[9] T. Yamaguchi, K. Baba, M. Murata, and K. Kitayama, "Scheduling algorithm with consideration to void space reduction in photonic packet switch," *IEICE Transactions on Communications*, vol. E86-B, pp. 2310–2318, Aug. 2003.

[10] T. DeFanti, M. Brown, J. Leigh, O. Yu, E. He, J. Mambretti, D. Lillethun, and J. Weinberger, "Optical Switching Middleware for the OptIPuter," *IEICE Transaction on Communication*, vol. E86-B, Aug. 2003.

[11] H. Nakamoto, K. Baba, and M. Murata, "Shared memory access method for a    computing environment," in *Proceedings of IFIP Optical Networks and Technologies Conference (OpNeTec)*, pp. 210–217, Oct. 2004.

[12] E. Taniguchi, K. Baba, and M. Murata, "Implementation and evaluation of shared memory system for establishing lambda computing environment," Tech. Rep. 255, IEICE, Aug. 2004.

[13] H. Amano, *Parallel Computer*. Shoukoudou, June 1996.

[14] N. Suzuki, S. Shimizu, and N. Yamanouchi, *An Implemantation of a Shared Memory Multiprocessor*. Koronasha, Mar. 1993.

[15] M. S. Papamarcos and J. H. Patel, "A low-overhead coherence solution for multiprocessors with private cache memories," in *Proceedings of IEEE 11th Annual International Symposium on Computer Architecture*, (New York), pp. 348–354, 1984.

[16] S. S. Thakkar and M. Sweiger, "Performace of an OLTP application on symmetry multiprocessor system," *17th Annual International Symposium on Computer Architecture*, pp. 228–238, May 1990.

[17] M. Wakabayashi, K. Inoue, and H. Amano, "ISIS: Multiprocessor simulator library," *Applied Informatics AI'99*, pp. 198–200, Feb. 1999.

[18] M. Wakabayashi and H. Amano, "Environment for multiprocessor simulator development," *I-SPAN 2000*, pp. 64–71, Dec. 2000.

[19] S. C. Woo, M. Ohara, E. Torrie, J. P. Singh, and A. Gupta, "The SPLASH-2 programs: Characterization and methodological considerations," in *Proceedings of the 22nd International Symposium on Computer Architecture*, pp. 24–36, June 1995.

**Hirohisa Nakamoto**    received the B.E. degree in Informatics and Mathematical Science from Osaka University, Osaka, Japan, in 2003. He is currently the master student at the Graduate School of Information Science and Technology, Osaka University, Japan. His research work is in the area of Grid computing and photonic netowork systems.

**Ken-ichi Baba**   received the B.E. and M.E. degrees in Information and Computer Sciences from Osaka University, Osaka, Japan, in 1990 and 1992, respectively. After quitting the doctoral course in 1992, he was a Research Associate with Education Center for Information Processing, Osaka University until 1997. He was an Assistant Professor with Electronic and Photonic Systems Engineering, Kochi University of Technology at Kochi in Japan from 1997. He has been an Associate Professor with Cybermedia Center (then Computation Center) at Osaka University since December 1998. His research interests include broadband communication network, computer communication network, and photonic network systems. He received the D.E. degree from Osaka University in 1995.

**Masayuki Murata**   received the M.E. and D.E. degrees in Information and Computer Sciences from Osaka University, Japan, in 1984 and 1988, respectively. In April, 1984, he joined Tokyo Research Laboratory, IBM Japan, as a Researcher. From 14 September 1987 to January 1989, he was an Assistant Professor with Computation Center, Osaka University. In February 1989, he moved to the Department of Information and Computer Sciences, Faculty of Engineering Science, Osaka University. From 1992 to 1999, he was an Associate Professor in the Graduate School of Engineering Science, Osaka University, and from April 1999, he has been a Professor of Osaka University. He moved to Advanced Networked Environment Division, Cyber-media Center, Osaka University in April 2000. In March 2004, he moved to Graduate School of Information Science and Technology, Osaka University. He has more than four hundred papers of international and domestic journals and conferences. His research interests include computer communication networks, performance modeling and evaluation. He is a member of IEEE, ACM, The Internet Society, IEICE and IPSJ.