

TASK MIGRATION IN N -DIMENSIONAL WORMHOLE-ROUTED MESH MULTICOMPUTERS

*Gwo-Jong Yu**, *Chih-Yung Chang*** and *Tzung-Shi Chen****

* Department of Computer Science and Information Engineering
National Central University, Chung-Li, Taiwan

** Department of Computer and Information Science
Aletheia University, Taipei, Taiwan

*** Department of Information Management
Chang Jung University, Tainan, Taiwan

Email: {yu, johny, shi}@axp2.csie.ncu.edu.tw

ABSTRACT

In a mesh multicomputer, performing jobs needs to schedule submeshes according to some processor allocation scheme. In order to assign the incoming jobs to a free submesh, a task compaction scheme is needed to generate a larger contiguous free region. The overhead of compaction depends on the efficiency of the task migration scheme. In this paper, two simple task migration schemes are first proposed in n -dimensional mesh multicomputers with supporting dimension-ordered wormhole routing in one port communication model. Then, a hybrid scheme which combines advantages of the two schemes is discussed. Finally, we evaluate the performance of all of these proposed approaches.

Keywords: Dimension-ordered routing, mesh multicomputers, gather-and-scatter, task migration, wormhole routing.

1. INTRODUCTION

The simplicity and regularity of mesh topology make it popular in multiprocessor systems. Because it is suitable for VLSI implementation, several mesh-based commercial machines have been built in recent years. Examples includes Intel Touchstone Delta[18], Intel Paragon [19], and the Tera computer system[1]. In a mesh multicomputer, to perform a sequence of jobs needs to schedule submeshes according to some processor allocation strategy, that allocates each job to a submesh with appropriate size. Number of researches aimed at processor allocation and job scheduling schemes in hypercubes [8] [9] and mesh multicomputers[2] [4] [5] [7] [10] [11] [14] [16] [21].

Given a job, the mesh system first allocates required processors, then execute job, and finally free processors. Although current submesh allocation schemes try to find an efficient way to allocate processor and to reduce the condition of fragment, however after a lot of processor allocations and deallocations, the mesh system still have a lot of fragments due to the unpredictable of the size and execution time of new incoming jobs. In such condition, a new job can not be scheduled to execute on this mesh system due to the lack of large enough contiguous

area of free processors, even if the number of free processors is sufficient. So, a task migration which moves running jobs from source processors to target processors is needed for mesh subsystem to increase the system utilization. This work assumes that the locations of source and destination submesh are determined by some processor management scheme of the mesh system. Efficiency of task compaction is measured by the transmission latency of task migration. A lot of researchers design fast task migration schemes in hypercube by exploring disjoint paths between two subcubes. There are some difference between hypercube and mesh computer. The number of parallel paths is propotional to the degree of hypercube. However, the degree of each node is fixed to be $2n$ in an n -dimensional mesh system. In most case, it is impossible to find out the parallel paths between two submeshes. The objective of this paper is to minimize the transmission latency of task migration by means of transferring the job in a way of several phases.

In recent parallel computer machines, wormhole routing [12] [13] is the most important switching technique. In this paper, the target machine we discuss is an n -dimensional mesh multicomputers supporting one port communication with dimension-ordered wormhole routing. We first present two task migration schemes in n -dimensional mesh multicomputers. Then a hybrid task migration scheme is presented. Finally, we use performance analysis to compare all of the proposed task migration schemes and inspect factors which affect the performance of those proposed schemes.

In the next section, we first introduce the system model of mesh multicomputer and describe the problem to be solved in this paper. In Section 3, two task migration schemes are presented in n -dimensional mesh multicomputers. We also proposed a hybrid approach which integrates these two schemes. Performance analysis of various cases are demonstrated in Section 4. Section 5 is the conclusion of this paper.

2. PRELIMINARIES

In this section, we introduce our system model for designing the task migration schemes on an n -dimensional mesh multicomputer. The mesh multicomputer system

consists of nodes, each node is a computer with its own processor, local memory, and communication links. Each directed link connects to two neighboring nodes through network [12] [13]. A common component of nodes in a new generation multicomputer is a router. It can handle and control the message communication entering, leaving, and passing through the node. The architecture of the n -dimensional mesh network system used in this paper is to provide dimension-ordered wormhole routing with one-port communication, in which one node can only send and, simultaneously, receive a worm from the other respective node at the same time. The dimension-ordered routing used in this paper is assumed to route messages to destination nodes first along X_1 -direction, then X_2 -direction, and finally X_n -direction on the mesh.

An n -dimensional mesh system $M(\mathbf{s})$ consists of $s_1 \times s_2 \times \dots \times s_n$ number of processors arranged in an n -dimensional grid, where $\mathbf{s} = (s_1, s_2, \dots, s_n)$ and s_i denotes the number of processors in dimension i . A processor in the grid is denoted by the coordinate $\mathbf{x} = (x_1, x_2, \dots, x_n)$, where $0 \leq x_i \leq s_i - 1$. Let $M(\mathbf{s})$ denote the set of processors $\{(0, 0, \dots, 0), \dots, (s_1 - 1, s_2 - 1, \dots, s_n - 1)\}$. We define the submesh of $M(\mathbf{s})$ as $SM(\mathbf{x}, \mathbf{m})$, where $\mathbf{x} = (x_1, x_2, \dots, x_n)$ denotes starting position (the smallest address of processors in submesh in lexicographical order) and $\mathbf{m} = (m_1, m_2, \dots, m_n)$ denotes the size, $m_1 \times m_2 \times \dots \times m_n$, of submesh. In other words, $SM(\mathbf{x}, \mathbf{m})$ denote the set of processors with coordinates in $\{(x_1, x_2, \dots, x_n), \dots, (x_1 + m_1, x_2 + m_2, \dots, x_n + m_n)\}$. Here we denote the source submesh by $SM(\mathbf{x}^s, \mathbf{m}) = \{(x_1^s, x_2^s, \dots, x_n^s), \dots, (x_1^s + m_1, x_2^s + m_2, \dots, x_n^s + m_n)\}$ and the destination submesh by $SM'(\mathbf{x}^d, \mathbf{m}) = \{(x_1^d, x_2^d, \dots, x_n^d), \dots, (x_1^d + m_1, x_2^d + m_2, \dots, x_n^d + m_n)\}$. These two submeshes with both of the same shape and size, $m_1 \times m_2 \times \dots \times m_n$, are located in different starting locations, $\mathbf{x}^s = (x_1^s, x_2^s, \dots, x_n^s)$ and $\mathbf{x}^d = (x_1^d, x_2^d, \dots, x_n^d)$ with allowing them to be partially overlapped. One node $\mathbf{x} = (x_1, x_2, \dots, x_n)$ in $SM(\mathbf{x}^s, \mathbf{m})$ is needed to route its assigned subtask to the corresponding destination node $\mathbf{x}' = (x_1', x_2', \dots, x_n')$ in $SM'(\mathbf{x}^d, \mathbf{m})$, where

$$x_i' = x_i^d + (x_i - x_i^s), i = 1, \dots, n. \quad (1)$$

3. TASK MIGRATION

In this section, two task migration schemes are first proposed in n -dimensional mesh multicomputers based on dimension-ordered wormhole routing. Then, we propose a hybrid scheme to minimize the total routing latency.

3.1. Diagonal Scheme

In this subsection, a task migration approach which explores disjoint paths in mesh multicomputer with dimension-ordered wormhole routing is proposed. First, we use the following example of a 3-dimensional mesh to illustrate the main idea of the developed task migration scheme. A task with several subtasks allocated to a $5 \times 4 \times 3$ submesh is needed to be migrated to another $5 \times 4 \times 3$ submesh on a mesh $\{(0, 0, 0), \dots, (11, 11, 11)\}$ as depicted in Figure 1. The task executed in the source sub-

mesh $\{(1, 1, 1), \dots, (5, 4, 3)\}$ is needed to be migrated to the destination submesh $\{(6, 7, 8), \dots, (10, 10, 10)\}$ in a $12 \times 12 \times 12$ mesh systems. The proposed scheme use five phases to migrate the tasks distributed in nodes located in various X_1, X_2 , and X_3 axis to the corresponding destination nodes. As shown in Figure 2, 12 nodes migrate their subtasks to corresponding destinations respectively in parallel. Note that each node migrates its subtasks first along X_1 -direction, then X_2 -direction, and finally X_3 -direction. As we can check, no common links are used or shared in each phases. For example, in Phase 1 if we project those selected nodes (the solid nodes) to the X_2 - X_3 plane, then all nodes will fall in different position. We can conclude that no two nodes will share or use the same link along X_1 -direction. Similarly, it is easy to verify that no two nodes will share or use the same link along X_2 -direction and X_3 -direction. Therefore, in each phase, the routing is congestion-free based on dimension-ordered routing in one-port communication, in which one node can simultaneously send out and receive from one worm at a time. There are totally five phases needed to complete the task migration in this example.

Next, we generalize the above descriptions of the migration routing scheme as follows. In order to avoid congestion during performing dimension-ordered routing in some phase, the nodes in the source submesh should have different positions in $(n-1)$ -dimensional subspace when we project a processor from n -dimensional space to $n-1$ subspace along direction X_i for all $1 \leq i \leq n$. In what follows, we propose the *diagonal scheme* to migrate these subtasks. Here we only illustrate how to schedule the source nodes in some transmission phases. It is easy for node determine their routing path since the dimension ordered routing protocol is applied. The source submesh $SM(\mathbf{x}^s, \mathbf{m})$ is assumed to be $\{(x_1^s, x_2^s, \dots, x_n^s), \dots, (x_1^s + m_1, x_2^s + m_2, \dots, x_n^s + m_n)\}$. We state how to arrange the routing phases P_k for each k . The routing during migration is one-by-one from the first phase to the last phase. Let B be set to

$$B = \max_i \{m_i\}.$$

We will show that there are exactly B phases needed to migrate task to destination submesh. For each phase, P_k , all of nodes $\mathbf{x} = (x_1, x_2, \dots, x_n)$, $x_i^s \leq x_i \leq x_i^s + m_i$, that are lying on hyperplanes

$$\sum_{i=1}^n m_i = c \cdot B + k, \quad (2)$$

are scheduled to simultaneously migrate their subtasks to their corresponding destinations, where c is integers and $1 \leq k \leq B$. Clearly, we have a lot of parallel hyperplanes

$$\sum_{i=1}^n x_i = H',$$

where $\sum_{i=1}^n x_i^s \leq H' \leq \sum_{j=1}^n (x_j^s + m_j)$. From the above description, the entire task distributed to all nodes in the source submesh is migrated to destination submesh

in B phases. We prove that the routing is congestion-free in each phase below. In addition, we prove that the number of routing phases is minimum with congestion-free routing in one-port communication model.

Theorem 1 *By applying the proposed diagonal scheme, the routing in each phase is congestion free.*

Proof: We will show that in each phase, there is exactly one node occupying the X_1 -direction communication channels. By projecting nodes that are scheduled in the same phase to hyperplane $X_1 = 0$, all nodes will have different positions. So we ensure that all these nodes are congestion free along X_1 -direction. We prove these arguments in what follows. We have hyperplanes

$$\sum_{i=1}^n x_i = H$$

for $x_i^s \leq x_i \leq x_i^s + m_i, 1 \leq i \leq n$, where

$$\sum_{i=1}^n x_i^s \leq H \leq \sum_{j=1}^n (x_j^s + m_j).$$

In phase P_k , nodes in hyperplane

$$\sum_{i=1}^n x_i = c \cdot B + k, \quad (3)$$

are scheduled to migrate subtasks to destination, where $B = \max_i \{m_i\}$. Assume $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)$ is a node on the hyperplane in Equation (3). Projecting $(x_1^*, x_2^*, \dots, x_n^*)$ onto the hyperplane $x_1 = 0$ generates a corresponding position $(0, x_2^*, \dots, x_n^*)$. If there are two or more nodes projected to the same location $(0, x_2^*, \dots, x_n^*)$, these nodes will cause congestion during using X_1 -direction communication channels. So, the goal we have to prove is that there are exactly one node in phase P_k projected onto the specific position $(0, x_2^*, \dots, x_n^*)$. Assume nodes $(x_1, x_2^*, \dots, x_n^*)$ are projected onto $(0, x_2^*, \dots, x_n^*)$, then

$$x_1 + \sum_{i=2}^n x_i^* = c \cdot B + k$$

That is,

$$x_1 = c \cdot B + (k - \sum_{i=2}^n x_i^*).$$

Because $x_1^s \leq x_1 \leq x_1^s + m_1$, the following inequations holds:

$$x_1^s \leq c \cdot B + (k - \sum_{i=2}^n x_i^*) < x_1^s + m_1 \quad (4)$$

$$\Rightarrow 0 \leq c \cdot B + (k - \sum_{i=2}^n x_i^* - x_1) < m_1. \quad (5)$$

$$\Rightarrow 0 \leq c \cdot B + (k - \sum_{i=2}^n x_i^* - x_1) < B. \quad (6)$$

Because $(k - \sum_{i=2}^n x_i^* - x_1)$ is a constant, there will be only one integer c satisfies (6). So, for fixed $x_i^*, 2 \leq i \leq n$, only one x_1^* fulfills equation (3). If the projected positions of all of these nodes are all different, we can conclude that all of these nodes are congestion free along X_1 -direction. The same arguments can be applied to X_i -direction for $2 \leq i \leq n$. If all of n directions are congestion free, we conclude that the migration process is congestion free.

⊥

Theorem 2 *The number of phases, B , with congestion-free routing is minimum based on dimension-ordered routing in one port communication model.*

Proof: Based on dimension-ordered routing, there are at most $m_1 \times m_2 \times \dots \times m_{i-1} \times m_{i+1} \times \dots \times m_n$ nodes delivering data along x_i -direction simultaneously, for $i = 1, \dots, n$. For a submesh $SM(\mathbf{x}^s, \mathbf{m})$, therefore, there are at most

$$\min_i \left\{ \prod_{j=1, j \neq i}^n m_j \right\}$$

nodes being simultaneously able to route their data in a way of congestion-free transmission. This makes there is a limitation on the size in each dimension of submesh $SM(\mathbf{x}^s, \mathbf{m})$. Therefore, the minimum number of routing phases are

$$\frac{\prod_{i=1}^n m_i}{\min_j \left\{ \prod_{k=1, k \neq j}^n m_k \right\}},$$

where $\prod_{i=1}^n m_i$ is the total number of nodes on the submesh $SM(\mathbf{x}^s, \mathbf{m})$. According to Theorem 1, each phase is congestion-free. We need exactly $B = \max_i \{m_i\}$ phases. Thus, the minimum number of routing phases, $\frac{\prod_{i=1}^n m_i}{\min_j \left\{ \prod_{k=1, k \neq j}^n m_k \right\}} = \max_l \{m_l\}$, whose value is minimum.

⊥

3.2. Gathering-Routing-Scattering Scheme

We next describe our second task migration scheme based on two collective communication schemes, gathering and scattering operations[12]. An illustration of gathering and scattering operation in eight processors is shown in Figures 3 and 4 respectively, where the solid nodes represent target processor in gathering operation and source processor in scattering operation. Without loss of generality, we assume the source submesh has maximum length in X_1 -direction. We apply the gathering operation on nodes in the same X_1 -direction to collect subtasks into one node. After the node migrates the combined subtasks to the corresponding destination node, we use the scattering operation on nodes in a line along X_1 -direction to disperse a couple of subtasks assigned to the destination node, to their respective nodes

in the destination submesh. The gathering and scattering schemes can be referred to [12] and [16]. In Figure 3, after the node in the end of line along X_1 -direction receives all the subtasks, it sends all these subtasks to the scheduled solid node in the source submesh. In the scattering operation, an additional routing step is needed to route the data from solid node to the node in the end of a line along X_1 -direction in the destination submesh. Figure 4 shows the operation of scattering.

Without loss of generality, we assume the submesh has maximum length in X_1 -direction. All of nodes (x_1, x_2, \dots, x_n) located on hyperplanes $\sum_{i=1}^n x_i = c \cdot B$ are gathering subtasks distributed within all of nodes in X_1 -direction of (x_1, x_2, \dots, x_n) . The node (x_1, x_2, \dots, x_n) next migrates the combined subtasks to the corresponding node $(x'_1, x'_2, \dots, x'_n)$. The node $(x'_1, x'_2, \dots, x'_n)$ finally scatters the combined subtask to the respective destination nodes on lines along X_1 -direction to complete the task migration. An illustration of gathering-routing-scattering operation is shown in Figure 5

Theorem 3 *The gathering-routing-scattering scheme is congestion-free in each phase.*

Proof: In this approach, all of the gathering and scattering steps are congestion-free [12][16]. The middle step, the solid nodes send the combined subtasks to the corresponding destination nodes, is also congestion-free, which has been shown in Theorem 1. Thus, the gathering-routing-scattering scheme is congestion free in each phase. \perp

3.3. Hybrid Scheme

In this subsection, we present a hybrid task migration scheme. We combine two schemes stated in the previous subsections to a hybrid one. We first partition the source submesh $SM(\mathbf{x}, \mathbf{m})$ into several n -dimensional subpartitions, which is also of submesh form with the size of $p_1 \times p_2 \times \dots \times p_n$. That is, the number of subpartitions is $\left\lceil \frac{m_1}{p_1} \right\rceil \times \left\lceil \frac{m_2}{p_2} \right\rceil \times \dots \times \left\lceil \frac{m_n}{p_n} \right\rceil$, where $\left\lceil \frac{m_i}{p_i} \right\rceil$ is the number of partitions in dimension i of the partitioned submesh. In the following, we show how the hybrid scheme works. We first partition the source submesh into subpartitions. We then use the gathering scheme to collect subtasks into nodes, located at the designated positions, depending on the size of $m_i, i = 1, \dots, n$. This step is similar to the gathering-routing-scattering scheme we proposed. After that, we use the diagonal scheme to schedule $\max\left(\left\lceil \frac{m_1}{p_1} \right\rceil, \left\lceil \frac{m_2}{p_2} \right\rceil, \dots, \left\lceil \frac{m_n}{p_n} \right\rceil\right)$ phases, to route the aggregated subtasks to their corresponding nodes. In this step, each subpartitions is represented as a supernode to transfer their subtasks on each node to its corresponding destination subpartition. Finally, we use the scattering scheme to distribute subtasks on the direction, which is with the maximum length in each partition, to complete the task migration. Figure 6 shows the hybrid approach in a 2D mesh system.

Theorem 4 *The hybrid scheme is congestion-free in each phase.*

Proof: The gathering scheme is first used to collect subtasks in each subpartition on a specific X_i -direction; hence, it is congestion-free. It is also congestion-free that all nodes scheduled in the same phase in each subpartition migrates its combined subtasks to its corresponding node. By Theorem 1 and 3, we know that this step is also congestion-free. The next step is to route the combined subtasks, scheduled with $\max\left(\left\lceil \frac{m_1}{p_1} \right\rceil, \left\lceil \frac{m_2}{p_2} \right\rceil, \dots, \left\lceil \frac{m_n}{p_n} \right\rceil\right)$ phases, to their corresponding nodes. Finally, the scattering scheme is applied to distribute each subtask in each subpartition in X_i direction; hence, it is congestion-free. Therefore, the proposed hybrid scheme is congestion-free in each phase. \perp

4. PERFORMANCE ANALYSIS

In this section, we will evaluate the developed schemes by performance analysis. Some parameters used in performance analysis are described in below. We assume that the startup latency is t_s and the transmission time of a flit (or byte) is t_x on a link between two neighboring nodes. Due to the delay in the intermediate routers is small, commercial systems [13] demonstrate that wormhole routing is distance insensitivity. We also assume that the size of subtask distributed in each node is the same, l flits, for analyzed convenience. In general, in the wormhole routing model, the latency for a node sending one message with l flits to another node is $t_s + t_x \times l$ [12] [13].

Because the hybrid-scheme is a combination of the diagonal scheme and the gathering-routing-scattering scheme, we first discuss the total transmission latency related to the hybrid task migration scheme in below. The hybrid task migration schemes needs gathering, diagonal routing, and scattering steps. In the gathering, it takes $\max_i \{\lceil \log p_i \rceil + 1\}$ time steps to collect these subtasks into one node. The combined message is with a double size compared with the previous one in each step and we need one extra step to transmit the final combined subtasks to the node located at the position in the designated phase. Thus, the size of the collected message is $\max_i \{(l \times 2^{\lceil \log_2 p_i \rceil})\}$ in final. This concludes the gathering step takes the time in below which can be obtained by the above derivations.

We observe that the quantities $\max_i \{p_i\}$ and $\max_j \left\{ \left\lceil \frac{m_j}{p_j} \right\rceil \right\}$ are two important factors of the proposed task migrations schemes. Let $P = \max_i \{p_i\}$ and $Q = \max_j \left\{ \left\lceil \frac{m_j}{p_j} \right\rceil \right\}$. That is, symbol P stands for the maximum number of processors in each dimension of each subpartition (the number of processors collected by gathering and scattering operation). The transmission latency in gathering operation depends on the maximum dimension P as follows:

$$T_{gather} = (\lceil \log_2 P \rceil + 1) \times t_s + (l \times 2^{\lceil \log_2 P \rceil}) \times t_x.$$

The cost $T_{diagonal}$ depends on both P and Q as follows,

$$T_{diagonal} = Q \times (t_s + l \times P \times t_x).$$

The cost of $T_{scatter}$ is the same as T_{gather} , i.e.

$$T_{scatter} = (\lceil \log_2 P \rceil + 1) \times t_s + (l \times 2^{\lceil \log_2 P \rceil}) \times t_x,$$

Thus, we have the total transmission time

$$T_{hybrid} = T_{gather} + T_{diagonal} + T_{scatter}.$$

We know that the first two schemes proposed in the previous section are the special cases of the hybrid scheme. The first task migration scheme, diagonal scheme, takes the transmission time in below not containing the gather and scatter steps for $p_i = 1, i = 1, \dots, n$.

$$T_1 = \max_i \{m_i\} \times (t_s + l \times t_x).$$

The second task migration scheme, gathering-routing-scattering scheme, takes the transmission time in below, for $p_1 = m_1, p_2 = m_2, \dots$, and $p_n = m_n$.

$$\begin{aligned} T_2 &= 2 \times (\max_i \{\lceil \log_2 m_i \rceil + 1\}) \times t_s \\ &+ \max_j \{l \times 2^{\lceil \log_2 m_j \rceil}\} \times t_x \\ &+ (t_s + \max_k \{m_k\}) \times t_x. \end{aligned}$$

From the above analysis of time complexity, we know that the amount of startup latencies, $\max_i \{m_i\}$, in the diagonal scheme is larger than that of $2 \times (\max_i \{\lceil \log_2 m_i \rceil + 1\}) + 1$ in the gathering-routing-scattering scheme in general. That is the reason why we use the collective communication to reduce the total amount of startup latency in migrating a task. However, the transmission size of the message between two submeshes in the diagonal scheme is generally smaller than that in the gathering-routing-scattering scheme. The hybrid scheme, therefore, is proposed for optimizing the time of migrating one task. To minimize the total amount of transmission delay T_{hybrid} , it is necessary to derive the optimum partitioning with respect to the values of $p_i, i = 1, \dots, n$ for a specific type of mesh architecture.

Here we give some assumptions to the parameters of system architecture for the use of analyzing the routing performance. This analysis is made on 2D (64×64), 3D ($64 \times 64 \times 64$), or 4D ($64 \times 64 \times 64 \times 64$) submesh accommodating a job needed to be migrated to another location. We have the message startup latency t_s is 1.0 microsecond for the small startup latency and is 10.0 microseconds for the large startup latency. The transmission time of a flit t_x on a link is 20.0 nanoseconds. The amount of a task in one node is assumed to be l flits; here different values of 100, 300, and 600 flits are discussed.

According to formula of transmission latency in each step of hybrid approach, to minimize the total transmission latency, T_{hybrid} , both of P and Q must be as small as possible. In the first case, we perform the task migration using various configurations such that $P = \max_i \{p_i\}$ is set to 64 and $Q = \max_j \{\lceil \frac{m_j}{p_j} \rceil\}$ has various size. Because Q corresponds to the number of phases needed in diagonal routing step. When P is fixed, small Q is corresponding to small total latency. The simulation results is shown in Figure 7. From Figure 7, the configuration

$8 \times 8 \times 64$ corresponding to $P = 64$ and $Q = 8$ has the smallest transmission latency.

In the second case, we perform task migration in 3D submesh using various configuration such that all configuration corresponds to the same value of $Q = 64$ and various value of P . The value of P corresponds to the message size and the step needed in the gathering and scattering operation. To reduce transmission latency, the value of P is as small as better. It can be seen from Figure 8 that small P corresponding to small transmission latency.

In the third case, we discuss the impact of uniform partitioning size using large and small startup latency on the transmission latency in 3D mesh system, i.e. we have different values of $p_1 \times p_2 \times p_3$. In order to balance the number of routing phases, P , and the maximum dimension of subpartition, Q , the subpartitions are with the shapes of the size of $1 \times 1 \times 1, 2 \times 2 \times 2, 4 \times 4 \times 4, 8 \times 8 \times 8, 16 \times 16 \times 16, 32 \times 32 \times 32$ and $64 \times 64 \times 64$ used in this case of performance analysis. The configurations, $1 \times 1 \times 1$, corresponds to the diagonal scheme while the configuration, $64 \times 64 \times 64$, corresponds to the gathering-routing-scattering scheme. Figures 10 and 9 are the incurred transmission latency of small and large startup latency respectively. In small startup latency, as shown in Figure 10, we prefer to use the diagonal scheme for task migration in most cases. This is because in small startup latency, the impact of message size is larger than the startup latency. In small startup latency, and message size 100, 300, and 600, diagonal scheme requires smaller transmission latency than gathering-routing-scattering scheme. On the other hand, in large startup latency, when the message size is large ($l = 600$), diagonal scheme performs better than gather-and-scatter scheme. However, when the message size is small ($l = 100$), gather-and-scatter scheme performs better. In Figure 9, the configuration $4 \times 4 \times 4$ perform best when message size is 300 or 600. The configuration $8 \times 8 \times 8$ perform best in the case that the message size is 100.

To examine the system performance on various dimensions, we compute the transmission latency on 4D and 2D mesh system. In 4D mesh system, the size of source submesh is $64 \times 64 \times 64 \times 64$ and the size of subpartition ranges from $1 \times 1 \times 1 \times 1$ to $64 \times 64 \times 64 \times 64$. Figure 11 shows the performance of various configurations. Compare Figure 9 with Figure 11, we found that all configurations have the same transmission latency. This is because the wormhole routing is distance insensitive and the formula of T_{hybrid} is affected by P and Q , and is irrelevant to dimension. Figure 12 shows the transmission latency of the hybrid approach on source submesh with size 64×64 . We obtain the same results as in Figures 9 and 11.

Briefly, from the above analysis, we have the following comments and suggestions as performing task migration. Generally speaking, we use the diagonal scheme to perform task migration to gain the minimized transmission latency when the startup latency is smaller. The time by applying the second scheme is usually larger than that by applying other schemes as shown in Figure 6. Thus, the second one is more unsuitable for executing task migration. The reason is that the collective communication

used in each X_i -direction takes to take a lot of time so as to collect a large amount of messages. However, we have to take into account the factors which affect the transmission latency, including the message size assigned in each node, the startup latency, as well as the partitioning size and shape together as the hybrid scheme is used. By evaluating these factors together, we are able to get the optimum solution with having the minimum time of T_{hybrid} .

5. CONCLUSION

In this paper, two simple task migration schemes are proposed in n -dimensional mesh multicomputers with supporting dimension-ordered wormhole routing in one-port communication model. Furthermore, a hybrid task migration scheme was proposed with the attempt to minimizing the routing latency. Finally, we compare all of our proposed task migrations schemes via performance analysis. In addition, we discuss how to easily apply our proposed task allocation schemes to some different processor allocation schemes with contiguous methods. Our future work includes investigating the job scheduling approaches, integrating task migration and processor allocation together, and evaluating the entire system performance including system utilization, job response time and more.

6. REFERENCES

- [1] R. A. et. al., "The tera computer system," in *Proc. Int'l Conf. Supercomputing*, pp. 1–6, 1990.
- [2] S. Bhattacharya and W.-T. Tsai, "Lookahead processor allocation in mesh-connected massively parallel multicomputer," in *Proceeding of International Parallel Processing Symp.*, pp. 868–875, 1994.
- [3] G.-I. Chen and T.-H. Lai, "Constructing parallel paths between two subcubes," *IEEE Transactions on Computers*, vol. 41, pp. 118–123, January 1992.
- [4] G.-M. Chiu and S.-K. Chen, "An efficient submesh allocation scheme for two-dimensional meshes with little overhead," *IEEE Transaction on Parallel and Distributed Systems*, vol. 10, pp. 471–486, May 1999.
- [5] P.-J. Chuang and N.-F. Tzeng, "An efficient submesh allocation strategy for mesh computer systems," *Proceedings of 11th International Conference on Distributed Computing Systems*, pp. 256–262, May 1991.
- [6] P.-J. Chuang and N.-F. Tzeng, "Allocating precise submeshes in mesh connected systems," *IEEE Transaction on Computers*, vol. 5, pp. 211–217, February 1994.
- [7] J. Ding and L. Bhuyan, "An adaptive submesh allocation strategy for two-dimensional mesh connected systems," in *Proceeding of International Conference of Parallel Processing*, vol. II, pp. 193–200, 1993.
- [8] P.-J. Chuang and N.-F. Tzeng, "A fast recognition-complete processor allocation strategy for hypercube computers," *IEEE Transactions on Computers*, vol. 41, no. 4, pp. 467–479, 1992.
- [9] J. Kim, C. R. Das, and W. Lin, "A top-down processor allocation scheme for hypercube computers," *IEEE Transaction on Parallel and Distributed Systems*, vol. 2, pp. 20–30, January 1991.
- [10] K. Li and K. H. Cheng, "A two-dimensional buddy system for dynamic resource allocation in a partitionable mesh connected system," *Journal of Parallel and Distributed Computing*, vol. 31, pp. 79–83, December 1991.
- [11] V. Lo, K. J. Windisch, W. Liu, and B. Nitzberg, "Noncontiguous processor allocation algorithms for mesh-connected multicomputers," *IEEE Transaction on Parallel and Distributed Systems*, vol. 8, pp. 712–726, July 1997.
- [12] P. K. McKinley, Y.-J. Tsai, and D. F. Robinson, "Collective communication in wormhole-routed massively parallel computers," *IEEE Computers*, vol. 28, pp. 39–50, December 1995.
- [13] L. M. Ni and P. K. McKinley, "A survey of wormhole routing techniques in direct networks," *IEEE Computer*, vol. 26, pp. 62–76, February 1993.
- [14] D. D. Sharma and D. K. Pradhan, "Submesh allocation in mesh multicomputers using busy-list: A best-fit approach with complete recognition capability," *Journal of Parallel and Distributed Systems*, vol. 36, pp. 106–118, January 1998.
- [15] D. D. Sharma and D. K. Pradhan, "Job scheduling in mesh multicomputers," *IEEE Transactions on Parallel and Distributed Systems*, vol. 9, pp. 57–70, January 1998.
- [16] Y.-C. Tseng, S.-Y. Ni, and J.-P. Sheu, "Toward optimal complete exchange on wormhole-routed tori," in *Proceedings of the 1997 International Conference on Parallel and Distributed Systems*, (Seoul, Korea), pp. 96–103, December 1997.
- [17] N.-F. Tzeng and H.-L. Chen, "Fast compaction in hypercubes," *IEEE Transactions on Parallel and Distributed Systems*, vol. 9, pp. 50–56, January 1998.
- [18] I. Corp., "Paragon xp/s product overview," 1994.
- [19] I. Corp., "A touchstone delta system description," 1991.
- [20] D. D. Sharma and D. K. Pradhan, "A fast and efficient strategy for submesh allocation in mesh-connected parallel computers," in *Proc. IEEE Symp. Parallel and Distributed Processing*, pp. 682–689, 1993.
- [21] Y. Zhu, "Efficient processor allocation strategies for mesh-connected parallel computers," *Journal of Parallel and Distributed Computing*, vol. 16, pp. 328–337, December 1992.

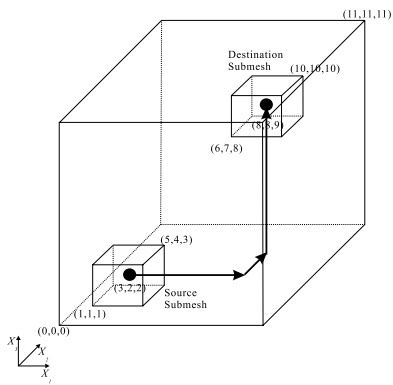


Fig. 1. Task migration between two $5 \times 4 \times 3$ submeshes in a 3-dimensional mesh with dimension-ordered worm-hole routing.

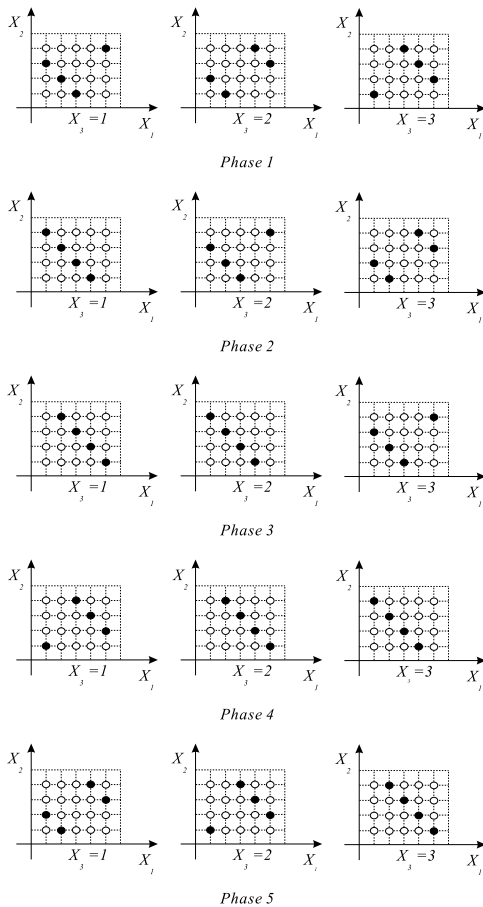


Fig. 2. Schedule of each phase in task migration process using the proposed diagonal scheme. Those solid nodes indicate processors that can be migrated to destination position in a congestion-free manner in each phase.

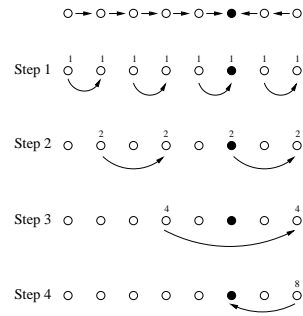


Fig. 3. Gather operation

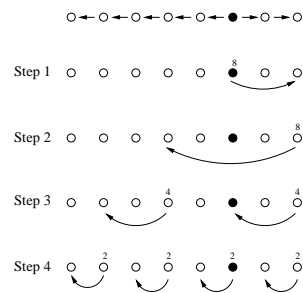


Fig. 4. Scatter operation

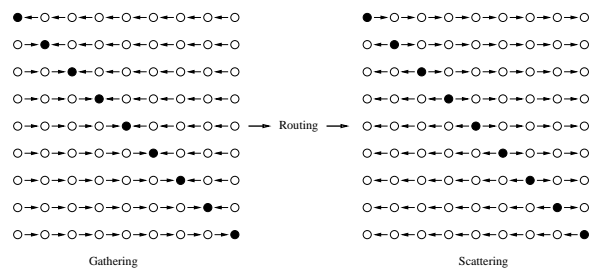


Fig. 5. The Gathering-Routing-Scattering operation

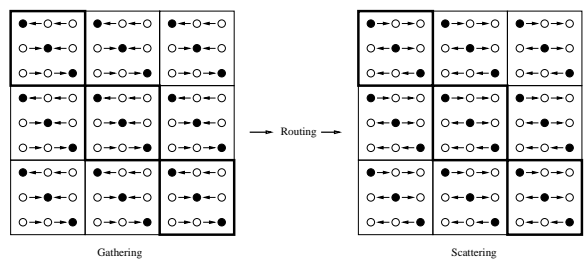


Fig. 6. The Hybrid Approach

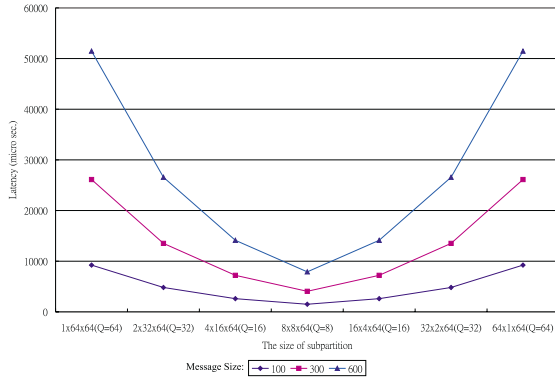


Fig. 7. Transmission Latency of 3D mesh with large startup latency. Configurations have the same value of $P = 64$ and various Q .

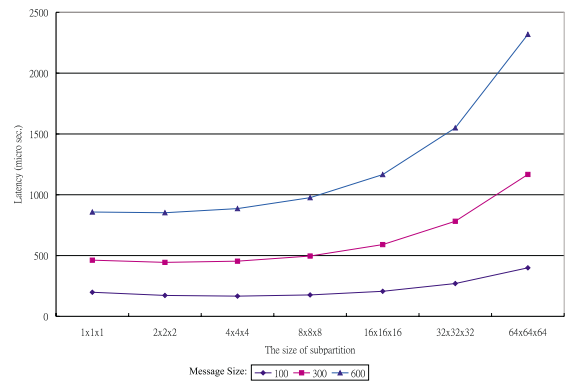


Fig. 10. Transmission Latency for small startup latency in 3D mesh with various subpartition size.

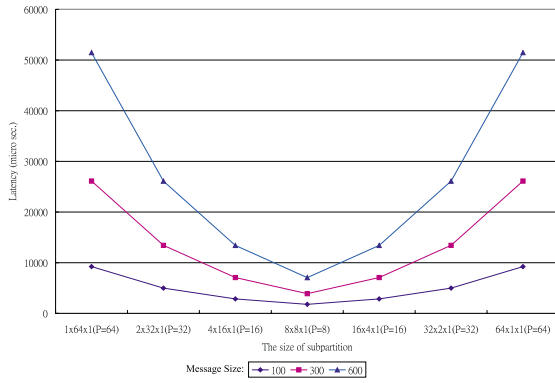


Fig. 8. Transmission Latency of 3D mesh with large startup latency. Configurations have the same value of Q but with different value of P .

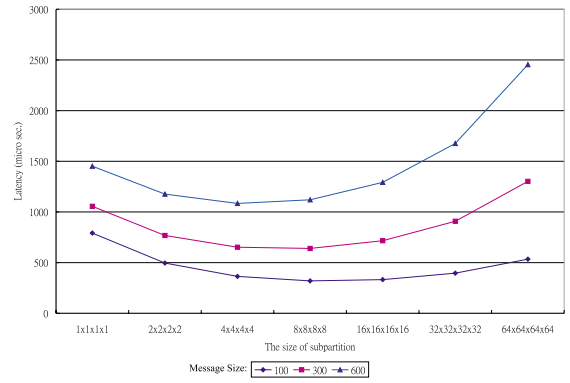


Fig. 11. Transmission Latency for large startup latency with in 4D mesh with various partition size.

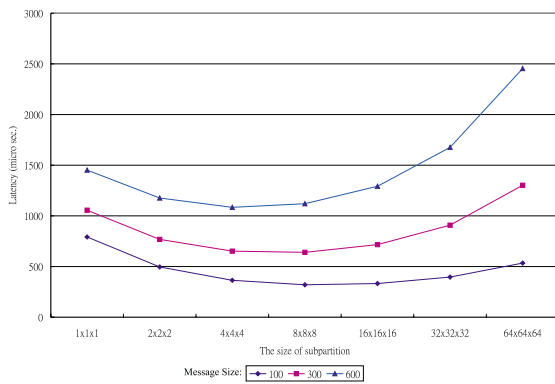


Fig. 9. Transmission Latency for large startup latency in 3D mesh with various subpartition size.

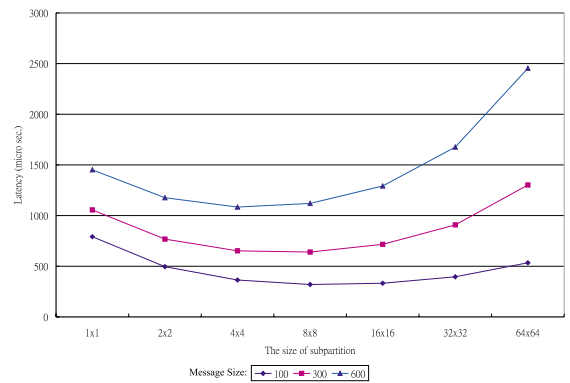


Fig. 12. Transmission Latency for large startup latency in 2D mesh with various partition size using hybrid approach.