

하이퍼큐브 분산 시스템에서 향상된 부하분산에 관한 연구

A Study on Advanced Load Balancing for Hypercube distributed System

劉宰旭*, 朴仁甲*, 金重民*

Jae-Wook Yu*, In-Kap Park*, Joong-Min Kim*

요 약

본 논문에서는 n 차원의 하이퍼큐브(hypercube) 구조의 분산 시스템에서 새로운 부하분산 알고리즘을 제안하였다. 제안된 알고리즘은 인접노드에게만 과부하 노드의 부하를 전송하여 발생되었던 blocking phenomenon을 억제하였으며 processor thrashing을 방지하기 위하여 한 프로세서에 모든 프로세서의 정보를 모아 실행하는 중앙집중형 알고리즘의 장점을 응용하였다. 제안된 알고리즘의 성능 평가를 위하여 컴퓨터 모의실험을 수행하고 다른 알고리즘과 비교하였다. 2차원 하이퍼큐브 분산 시스템의 경우에는 각 알고리즘의 부하분산 결과가 큰 차이를 보이지 않았으나 3차원과 4차원의 경우에는 제안된 알고리즘이 보다 나은 결과를 보여주었다.

Abstract

In this paper, an advanced load balancing algorithm in n th order Hypercube distributed system has been proposed. The new algorithm uses centralized load-balancing to avoid blocking phenomenon and processor thrashing, and shows the results which makes loads to approach average value of loads. The new algorithm is compared with several other algorithm and it shows a merit in cost function value.

1. 서 론

점차적으로 고성능을 필요로 하는 응용분야의 증대로 고성능 시스템이 필요하게 되었다. 필요한 고성능을 얻기 위해 하나의 프로세서의 성능을 높이는 방법과 다수개의 프로세서를 연결하여 구현하는 분산 멀티프로세서 시스템(distributed multiprocessor system)을 사용할 수 있다. 그러나 하나의 프로세서만을 이용하여 고성능을 얻고자 할 경우 물리적인 한계로 인해 프로세서의 성능 향상에 한계가 있다. 그러므로 프로

세서들을 병렬로 구성한 분산 시스템을 사용하게 되었다. 이상적인 분산 시스템은 각각 개별적인 부하들을 처리하므로 시스템의 성능은 프로세서의 개수에 따라 선형적으로 증가한다. 그러나 실제로는 실행 중 발생하는 부하 불균형 현상 등의 이유로 시스템의 규모가 커지는 만큼의 성능향상이 이루어지지 않는다.

본 논문은 다른 구조에 비하여 구조적인 적응성이 매우 높고, 노드와 노드사이의 통신로의 개수가 많으며, 이동거리가 비교적 짧은 n 차원의 하이퍼큐브(hypercube) 구조의 분산 시스템에서 새로운 부하분산 알고리즘을 제안하고 제안된 알고리즘의 성능 평가를 위하여 컴퓨터 모의 실험을 수행하였다. 제안된

* 建國大學校 電子情報通信工學部

(Dept. of Electronic Engr., Konkuk Univ)

接受日: 2002年 5月 1日, 修正完了日: 2002年 7月 28日

알고리즘은 인접노드에게만 과부하 노드의 부하를 전송하여 발생되었던 blocking phenomenon을 억제하였으며 processor thrashing을 방지하기 위하여 한 프로세서에 모든 프로세서의 정보를 모아 실행하는 중앙 집중형 알고리즘의 장점을 응용하였다. 또한 부하량 sorting후에 부하의 최대값을 찾아서 최대값 노드의 부하만을 전송해 나감으로써 전체 노드간의 부하량의 차이를 줄여감으로 전체 노드의 부하량이 평균값에 가깝도록 하였다.

II. 부하분산 정의 및 방법

2.1 부하분산의 정의

부하분산은 통신부담을 예측하기 어려운 환경이나 동적인 실행 환경하에서 프로세서들에 할당된 부하의 양을 균일하게 만들어 전체적인 실행시간을 단축하는 방법이다. 그림 1은 부하분산 시스템을 설명한 것으로 한 개 이상의 프로세서로 구성되어 있는 분산 시스템에서 각 프로세서로 부하가 할당되어 실행대기 상태에 있을 때, 부하가 특정의 프로세서에 과다하게 주어졌을 경우 이러한 프로세서의 부하 중 일부를 부하가 적은 프로세서로 이동하여 실행함으로써 시스템 전체의 처리효율을 높이고 부하의 실행시간을 단축하는데 그 목적이 있다.

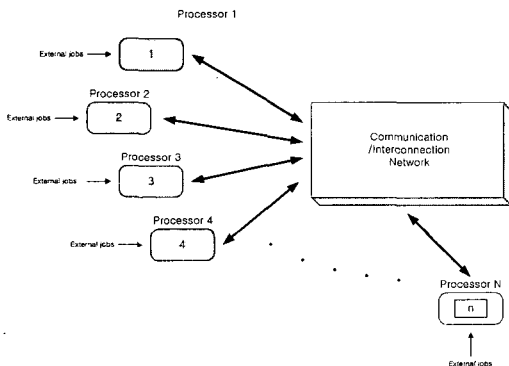


그림 1. 부하분산 시스템
Fig. 1 Load balancing system

2.2 부하분산의 방법

(1) 정적 부하분산

정적인 부하 이동은 시간의 경과에 따른 각 프로세서의 부하의 상태에 관계없이 영구적으로 특정한 프로세서를 그 부하분산 영역상의 서비스 프로세서로 설정하여 놓는 방법이다.

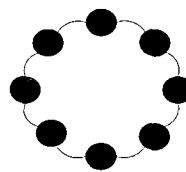
(2) 동적 부하분산

동적인 부하 이동은 시간의 경과에 따라서 부하분산 영역내의 프로세서의 역할이 부하의 상태 변화에 따라서 과부하 상태에 있을 때는 클라이언트로 작용하고 저부하 상태에 있을 때는 서버로 작용한다. 일반적으로 분산 시스템의 부하분산은 동적 부하분산을 말한다.

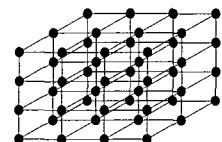
2.3 프로세서간의 연결 방법

각 노드간 기하학적 배치나 노드간의 연결을 topology라 한다. 각각의 노드는 하나 또는 여러 개의 프로세서로 구성되고 배치에는 그림 2와 같은 고리(ring), 격자(mesh), 토러스(torus), 팻 트리(fat-tree), 하이퍼큐브(hypercube) 등의 방법들이 있다.

이 중 하이퍼큐브(hypercube)는 $N(=2^n)$ 개의 노드(Node : 분산 시스템에서는 프로세서에 해당)로써 n 차원의 이진 큐브를 형성하므로 binary n-cube 라고도 말한다. 하이퍼큐브에는 여러 가지 특성이 있다. 각 프로세서는 자신의 로컬 메모리를 가지고 있으며, n 차원의 하이퍼큐브일 경우 n 개의 다른 프로세서와 직접 연결되는 통신로를 갖는다. 하이퍼큐브의 차원이 3 이상일 경우에는 tree, grid, ring, torus 등 임의의 연결구조를 갖는 시스템으로 노드 및 링크의 대응에 의한 전환이 가능하다. 따라서 하이퍼큐브는 다른 구조에 비하여 구조적인 적응성이 매우 높다.



Ring



3D Mesh

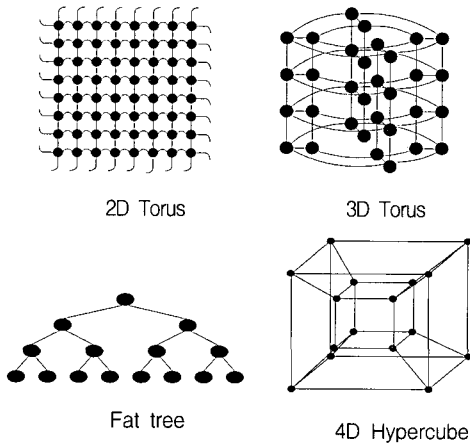


그림 2. 연결 방법
Fig. 2 Topology

III. 기존의 부하분산 알고리즘

3.1 집중형 알고리즘

시스템이 커질 경우에 많은 문제점이 있기 때문에, 순수하게 집중형 알고리즘만을 사용하는 연구보다는 집중형 알고리즘과 분산형 알고리즘을 혼합한 알고리즘이 주종을 이룬다.

모든 프로세서들의 부하정보를 관리하는 information collector를 두어서 부하정보를 관리하는 방법이다. 부하분산이 이루어지는 범위에 따라 local range, global range로 구분하고, 부하분산이 이루어지는 방법에 따른 LRM(Least Recently Migrated), MLM(Minimum Load Maintained)로 구분한다.

3.2 분산형 알고리즘

분산형 알고리즘은 모든 프로세서가 되도록 넓은 범위의 프로세서들의 정보를 참조하고, 부하상태를 상, 중, 하로 나누고 상에서 하로만 부하 전송이 일어나도록 하는 방식이다.

IV. 개선된 부하분산 알고리즘

4.1 기존 알고리즘의 문제점 및 해결방안

(1) 기존의 알고리즘은 과부하된 노드의 부하량을 저부하 노드에게 모두 전송하거나 전혀 전송하지 않는 방법을 택하였다.

→ 과부하된 부하를 모두 전송하지 못하더라도 프로세서 trashing 현상이 발생하지 않는 범위 내에서 최대의 과부하 양을 전송한다.

(2) 최대 부하량을 가진 노드의 부하의 양이 전체 부하의 실행에 필요한 비용함수에 가장 큰 영향을 미친다. 부하분산이 모두 행하여지지 못할 경우(부하를 전송할 노드를 찾지 못하는 경우) 비용함수가 크게 증가하게 된다.

→ 부하량 sorting후에 부하의 최대값을 찾아서 최대값 노드의 부하만을 전송해 나간다. 그러므로 전체 노드간의 부하 량의 차이를 줄여감으로 전체 노드의 부하의 값이 평균값에 가깝도록 한다.

(3) 한 과부하 노드의 부하들은 정해진 경로의 저부하 노드에게만 과부하 노드의 부하를 전송한다.

→ 매 iteration(반복)마다 아직 저부하 상태인 노드 가운데 최단 거리인 곳으로 부하 이동하므로 동일 과부하 노드의 부하는 다수개의 저부하 노드로 분배되어 전송될 수 있다.

4.2 비용 함수

분산 시스템의 비용 함수 H는 식(4-2)의 부하의 실행에 필요한 비용과 식(4-3)의 데이터의 통신에 드는 비용의 합인 식(4-1)로 표현할 수 있다. 표현될 수 있는데, 시스템의 성능을 향상시키기 위해서는 이 에너지, 즉 비용 함수를 최소화하도록 부하분산을 시켜야 한다.

$$H_{total} = H_{execute} + \tau H_{comm} \quad (4-1)$$

$$H_{execute} = \text{Max}[W_1, W_2, \dots, W_i] \quad (4-2)$$

$$H_{comm} = \sum_i W_{i,comm} \quad (4-3)$$

여기서 W_i 은 i 번 노드의 부하량을, $W_{i,comm}$ 은 i 번 노드의 이동되어야 할 부하량을 의미한다. $H_{execute}$ 는 각각의 프로세서들이 같은 부하를 가질 때 최소화되며 H_{comm} 은 통신시간이 작을 때 최대화된다. 하이퍼큐브의 차원이 n , 노드의 개수가 $N(=2n)$ 이라면, 부하분산 이후의 비용 함수를 구하면 다음과 같다.

① $H_{execute}$ 는 노드의 부하량 W_i 에서 저부하 노드로의 이동되는 부하량을 제외하거나 다른 노드로부터의 부하이동을 합한 것 중 최대 값이 되어 식(4-4)로 표현될 수 있다.

$$H_{execute} = \text{Max}[W_i - \sum_j M_{ij} \text{ for } i=0,1,\dots,n] \quad (4-4)$$

② 노드마다의 거리가 다르므로 그 거리를 i 노드로부터 j 노드까지의 거리를 D_{ij} 라고 하면 H_{comm} 은 다음과 같다.

$$H_{comm} = \sum_i \sum_j M_{ij} \times D_{ij} \quad (4-5)$$

여기서 M_{ij} 는 i node로부터 j 노드로의 이동된 부하량을 의미한다.

4.3 제안하는 개선된 알고리즘

초기의 부하량을 W_i 이라고 하고 반복되는 iteration에 따라 계산되는 nW_i 라고 하자. 부하분산을 취하기 전에는 $nW_i = W_i$ 이며 $N=2^n$ 이다.

$$L_i = f(T_L - W_i) \quad (4-6)$$

$$\text{Order}_i = \begin{cases} 1, & W_i = \text{Max}[W_1, W_2, W_3, \dots, W_n] \\ 0, & \text{otherwise} \end{cases} \quad (4-7)$$

$$\text{if } \sum_i g(W_i - T_H) > 0 \text{ and } \sum_i \text{Order}_i \leq \sum_i L_i \cdot g(T_H - W_i) \quad (4-8)$$

then $W_i = W_i - \text{Order}_i$,
 $W_j = W_j + 1$
 {if $\text{Min}(D_{ij} \text{ in } f(T_H - W_j) \cdot L_j = 1)$ },
 ($i, j = 1, 2, \dots, N$)

여기서

$$L_i = g(T_H - W_i) \cdot f(T_L - W_i) \quad (4-9)$$

$$f(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases} \quad (4-10)$$

$$g(x) = \begin{cases} x, & x \geq 0 \\ 0, & x < 0 \end{cases} \quad (4-11)$$

식(4-6)은 초기의 부하량이 저부하인지 아닌지를 결정하는 것으로 L_i 는 저부하 상태이면 1, 그렇지 않으면 0의 값을 가진다. 식(4-7)은 부하량의 Sorting 후에 최대의 부하량을 가진 노드의 Order_i 는 1, 그렇지 않으면 0을 만드는 식이다. 식(4-8)은 부하의 이동을 나타내는 식으로 아직 이동될 과부하량이 남아 있고, 최고 부하량을 가진 노드의 개수보다 저부하 노드의 부하이동 허용량이 크거나 같은 경우에 최고의 부하량을 가진 노드에서 최단 거리의 노드로 1개씩 부하를 이동하는 식이다.

알고리즘의 종료조건은 식(4-8)을 만족하지 않을 때이며, 종료 전까지는 식(4-7)과 식(4-8)을 계속하여 반복한다.

4.4 개선된 알고리즘의 처리순서

- (1) 모든 노드의 부하량을 계산한다.
- (2) 최대 부하량을 가진 노드를 찾는다.
- (3) 최대 과부하 노드의 부하중 1개를 최단 거리에 위치한 저부하 노드로 전송한다. 상태는 T_H , T_L 경계로 과부하, 중부하, 저부하의 3가지 상태로 나누고 부하의 불규칙한 이동을 방지하기 위하여 과부하 노드에서 저부하 노드로의 이동만 허용한다.
- (4) 만일 전체의 과부하량이 저부하량 보다 크다면 부하분산을 끝마치고 작다면 (1)부터 다시 반복한다.

4.5 개선된 알고리즘의 부하분산 실험 결과

제안된 부하분산 알고리즘은 하이퍼큐브의 차원 n 이 2, 3, 4 일 때의 부하분산의 예제들을 가지고 모의실험되었다. 노드의 숫자는 각 프로세서의 로컬 큐에 실행 대기중인 부하의 수이며 각 노드의 주소는 "[]"안에 표시되었다. 각 노드의 부하상태는 T_H 와 T_L 에 따라 과부하(11이상), 중부하(6~10), 저부하(5이하) 3단계로 색으로 구별하여 나타내었다. 다음 그림 4, 그림 5 그리고 표 1의 결과는 차원 $n=3$ 일 때의 결과이다.

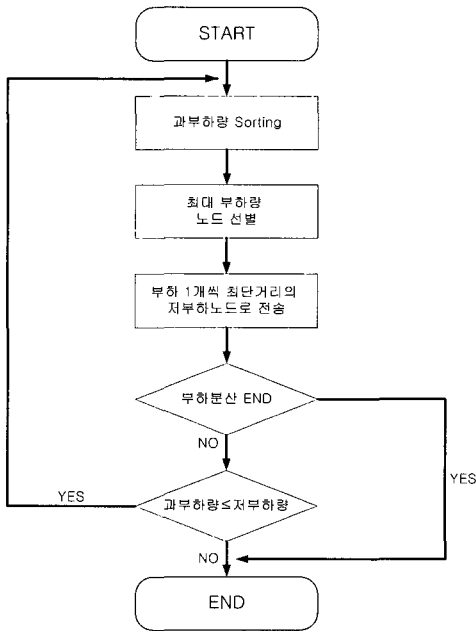


그림 3. 개선된 알고리즘 flow chart
Fig. 3 New algorithm flow chart

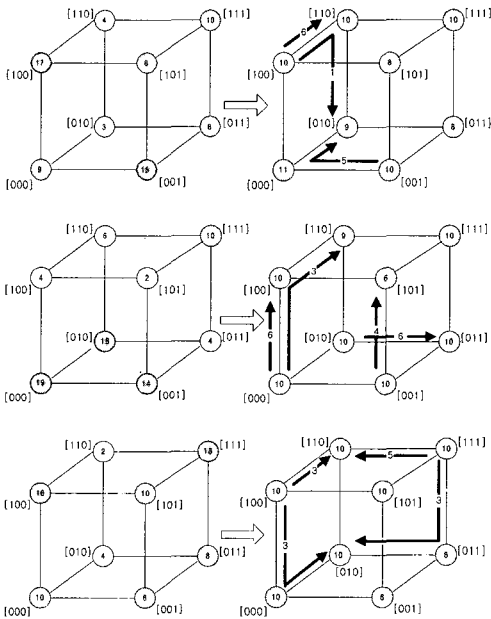


그림 4. 3차원 부하분산 예제(a),(b),(c)
Fig. 4 Example of 3rd order load balancing(a),(b),(c)

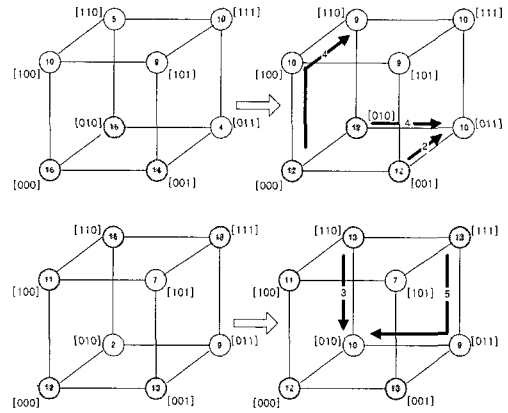


그림 5. 3차원 부하분산 예제(d),(e)
Fig. 5 Example of 3rd order load balancing(d),(e)

표 1은 다른 알고리즘들과 제안된 알고리즘에 의한 부하 분산 결과의 비교이다.

표 1. 3차원 부하분산 알고리즘 비교
Table 1. Comparison of 3rd order load balancing algorithms

그림	Cost Function	Shortest Method	Gradient Method	제안 알고리즘
4(a)	$H_{execute}$	10	11	10
	H_{comm}	18	17	18
	H_{total}	11.8	12.7	11.8
4(b)	$H_{execute}$	10	13	10
	H_{comm}	22	19	22
	H_{total}	12.8	14.9	12.8
4(c)	$H_{execute}$	10	16	10
	H_{comm}	20	14	20
	H_{total}	12	17.4	12
5(d)	$H_{execute}$	14	14	12
	H_{comm}	16	22	14
	H_{total}	15.6	16.2	13.4
5(e)	$H_{execute}$	16	18	13
	H_{comm}	16	6	13
	H_{total}	17.6	18.6	14.3

다음 그림 6, 그림 7 그리고 표 2의 결과는 차원 n=4일 때의 결과이다.

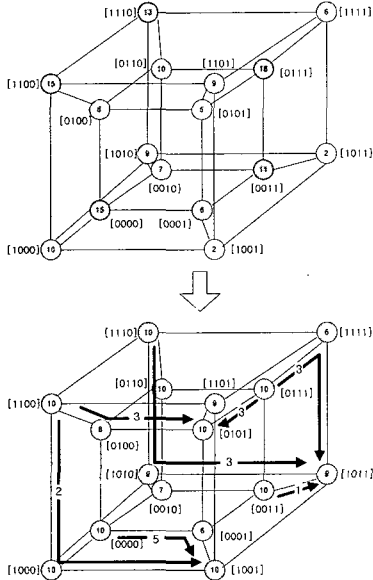


그림 6. 4차원 부하분산 예제(a)
Fig. 6 Example of 4th order load balancing(a)

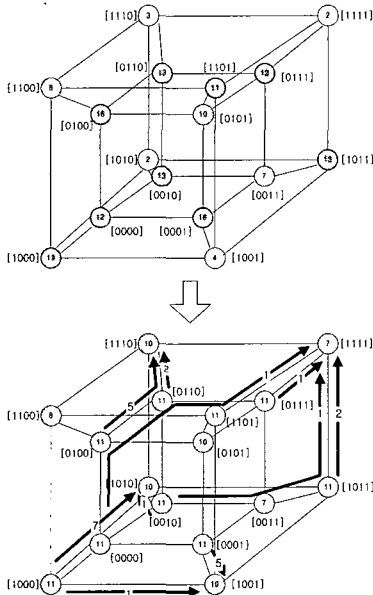


그림 7. 4차원 부하분산 예제(b)
Fig. 7 Example of 4th order load balancing(b)

표 2는 다른 알고리즘들과 제안된 알고리즘에 의한 부하 분산 결과의 비교이다.

표 2. 4차원 부하분산 알고리즘 비교
Table 2. Comparison of 4th order load balancing algorithms

그림	Cost Function	Gradient Method	Preferred List Method	제안 알고리즘
6	$H_{execute}$	10	15	10
	H_{comm}	39	34	36
	H_{total}	13.9	18.4	13.6
7	$H_{execute}$	19	19	11
	H_{comm}	41	41	28
	H_{total}	23.1	23.1	13.8

표 1, 표 2는 다른 알고리즘들과 제안된 알고리즘에 의한 부하분산 결과의 비교이다. 제안된 부하분산 알고리즘을 이용한 부하분산 결과와 Random Method, Shortest Method, Gradient Method, 그리고 우선 순위 부하분산 알고리즘에 의한 부하분산 결과를 비교하였다. 식(3-2)의 τ 는 0.1로 하여 비용함수를 계산하였다.

Random Method는 비용함수가 낮게 나오더라도 불필요한 부하 이동 현상이나 상태 불규칙 현상을 예측할 수가 없다. Shortest Method와 Gradient Method는 낮은 차원에서는 적당한 부하분산을 하지만 프로세서 쏠림 현상을 방지할 수 없다. 우선 순위 부하분산 알고리즘 또한 많은 노드에서 과부하 현상이 나타날 때는 프로세서 쏠림 현상이 발생한다. 제안된 알고리즘은 프로세서 쏠림 현상이 발생하지 않으면서 가장 가까운 거리의 저부하 노드로 부하 이동을 한다. 총 비용함수를 비교해보면, 2차원 하이퍼큐브 분산 시스템의 경우 각 알고리즘의 부하분산 결과가 큰 차이를 보이지 않았지만 3차원과 4차원의 경우에는 제안된 알고리즘이 보다 나은 결과를 보여준다.

V. 결 론

본 논문은 하이퍼큐브 구조의 병렬 프로세서 시스템에서 적합한 부하분산 알고리즘을 제안하였다. 제안된 알고리즘은 과부하 노드에서 인접노드에게만 부하를 전송하여 인접하지 않은 노드는 저부하이지만 부하를 받지 못하여 발생했던 blocking phenomenon을 억제하였으며 저부하 노드가 다수의 과부하 노드에서 너무 많은 부하를 받아 오히려 processor trashing 현상을 유발하는 문제점을 방지하기 위하여 한 프로세서에 모든 프로세서의 정보를 모아 실행하는 중앙집중형 알고리즘의 장점을 응용하였다. 또한 기존의 집중형 알고리즘에서는 최대 과부하노드가 가장 가까이 있는 저부하 노드로 부하를 전송시켜 나머지 과부하 노드들은 전송하지 못하는 문제가 발생하였다. 하지만 제안된 알고리즘은 부하량 sorting 후에 부하의 최대 값을 찾아서 최대 값 노드의 부하만을 저부하 노드로 전송하는 단계를 반복함으로써 전체 노드의 부하의 값이 평균값에 가깝도록 하였다.

제안된 알고리즘은 컴퓨터 모의실험을 통하여 기존의 부하분산 알고리즘과 제안한 알고리즘을 비교함으로써 부하분산의 효율성 증가를 검증하였으며 기존 알고리즘에서 발생했던 여러 문제점들을 개선하였다.

참 고 문 헌

- [1] Ming-Syan Chen and Kang G. Shin, "Processor Allocation in an N-Cube Multiprocessor Using Gray Codes," IEEE Trans. Computers, vol. C-36, no. 12, December 1987.
- [2] Sanjay Ranka, Youngju Won, and Sartaj Sahni, "Programming a Hypercube Multicomputer," IEEE Trans. Software Eng., pp. 69-77, September 1988.
- [3] D. E. Rumelhart, G. E. Hinton and R. J. Williams, "Learning Internal Representation by Error Propagation : in Parallel Distributed Processing (Eds. David E. Rumelhart and James

L. McClelland)", The MIT Press, Vol. 1, 1986.

- [4] Phillip Krueger and Niranjan G. Shivaratri, "Adaptive location policies for global scheduling", IEEE Transactions on Software Engineering, Vol. 20, No. 6, pp. 432-444, June 1994

저 자 소 개

劉宰旭 (學生會員)



2002년 2월 : 건국대학교 전자정보통신공학과 (공학사)
2002년 3월 ~ 현재 : 건국대학교 전자정보통신공학과 대학원 석사과정

朴仁甲 (正會員)



1973년 2월 : 고려대학교 전자공학과 (공학사)
1976년 2월 : 고려대학교 전자공학과 대학원 (공학석사)
1986년 2월 : 고려대학교 전자공학과 대학원 (공학박사)
1980년 8월 ~ 현재 : 건국대학교 전자정보통신공학과 교수

金重民 (正會員)



2000년 2월 : 경기대학교 전자공학과 (공학사)
2002년 2월 : 건국대학교 전자정보통신공학과 대학원 (공학석사)
2002년 3월 ~ 현재 : 건국대학교 전자정보통신공학과 대학원 박사과정