

Randomized and Deterministic Simulations of PRAMs by Parallel Machines with Restricted Granularity of Parallel Memories^{*}

Kurt Mehlhorn¹ and Uzi Vishkin²

¹ Fachbereich 10, Universität des Saarlandes, D-6600 Saarbrücken (Fed. Rep.)

² Courant Institute, New York University, 251 Mercer Street, New York, NY 10012, USA

Summary. The present paper provides a comprehensive study of the following problem. Consider algorithms which are designed for shared memory models of parallel computation (PRAMs) in which processors are allowed to have fairly unrestricted access patterns to the shared memory. Consider also parallel machines in which the shared memory is organized in modules where only one cell of each module can be accessed at a time. *Problem.* Give general fast simulations of these algorithms by these parallel machines.

Each of our solutions answers two basic questions. (1) How to initially distribute the logical memory addresses of the PRAM, to be simulated, among the physical locations of the simulating machine? (2) How to compute the physical location of a logical address during the simulation?

We utilize two main ideas for the first question.

(a) Randomization. The logical addresses are randomly distributed among the memory modules. This is done using universal hashing.

(b) Copies. We keep copies of each logical address in several memory modules.

In a typical time cycle of the PRAM some number of memory requests has to be satisfied. As a primary objective, our simulations minimize the maximum number of memory requests which are assigned to the same module. Our solutions also optimize the following computational resources. They minimize the size of the physical memory, the time for computing the mapping from logical to physical addresses and the space for storing this mapping.

We discuss extensions of our solutions to various PRAMs and various shared memory parallel machines. Our solution is also applicable to synchronous distributed machines with no shared memory where the processors can communicate through a bounded degree network.

* A preliminary version of this paper was presented at the 9th Workshop on Graphtheoretic Concepts in Computer Science (WG-83), Fachbereich Mathematic, Universität Osnabrück, June 1983

I. Introduction

Consider algorithms designed for models of parallel computation in which processors have access to a shared memory. Consider also parallel machines in which this shared memory is organized in modules where only one cell of each module can be accessed at a time.

The problem of simulating these algorithms on these parallel machines is the problem of *granularity of parallel memories* (granularity, in short). Every intuitive idea for coping with the granularity problem has to be analyzed for different formal settings of assumptions for both the model of parallel computation and the parallel machine. In order to overcome this difficulty we present our main ideas on a setting of assumptions which enables us to simplify the presentation by limiting the discussion to the actual problem which is overcome by our ideas. Extensions of the ideas to other models of computation and machines are given later in the paper.

We study ways by which the second parallel machine model below can simulate the first. Both machine models employ p processing elements (PE's or processors) which operate synchronously and N common memory cells. In the first model each processor has access to each of the N cells in each time unit. We forbid only the case where two (or more) processors seek access to the same cell at the same time. This is the Exclusive-Read Exclusive-Write Parallel Random Access Machine (EREWPRAM). It is based on [11]. Our second model of computation is called Module Parallel Computer (MPC). The common memory of size N is partitioned into m memory modules. Say that at the beginning of a cycle of this model the processors issue R_j requests for addresses located in cells of module j , $0 \leq j \leq m-1$. Let $R_{\max} = \max \{R_j | 0 \leq j \leq m-1\}$. Then the requests for each module are queued in some order and satisfied one at a time. So a cycle takes R_{\max} time. We assume that immediately after a simulation of a cycle is finished every processor knows it. The problem of simulating efficiently one cycle of the EREWPRAM by the MPC is taken as the definition of the granularity problem in the next two chapters which include the main contribution of this paper. When $N=m$ the MPC can simulate a cycle of the EREWPRAM in one time unit while when $N \geq mp$ a naive simulation may result in R_{\max} as large as p .

The survey paper [10] emphasizes the importance of the granularity problem. It reports about the considerable attention this problem has received in the literature by mentioning fourteen papers that dealt with it. Most of these papers suggest strategies for partitioning the memory addresses among the modules for algorithms that either have access patterns which are known in advance or have access patterns in successive time units which satisfy some probabilistic assumptions. Our attitude is completely different. We present solutions and analyses for the general problem of simulation. They do not depend on the access behaviour of the algorithms being simulated. This is in sharp contrast to both classes of past research mentioned above. In this spirit [21] observed that in a few general cases the idea of dynamically changing location of addresses among modules throughout the performance of an algorithm enables efficient simulations utilizing only a moderate number of modules ($m=p$).

Our research is motivated by the Ultracomputer project. The NYU-Ultracomputer group [9] believes that a machine using 4096 processors and 4096 memory modules will be available by 1990. The MPC represents actually, an abstract Ultracomputer design which idealizes only one point: the interconnection of processors and memory modules. A significant part of this project involves heuristics for difficulties related to the granularity problem. The Ultracomputer is a general-purpose parallel computer that may be used for any parallel algorithm. Our general solutions are, therefore, of particular relevance to its design.

There is an even stronger relation between the present paper and the parallel-design distributed-implementation (PDDI) machine, proposed in [18]. The PDDI machine forms a counterpart to the Ultracomputer which differs from it mainly at the following point. Its interconnection network, between processors and memories, performs well in the worst case while the interconnection network of the Ultracomputer performs well in the average (as was indicated by many simulations). In order to explain how the PDDI paper and the present paper complement one another we need the following simulation problem. Consider algorithms designed for models of parallel computation in which processors have access to a shared memory. (We refer to these models as PRAMs in short. So far it is similar to the granularity problem). Consider also synchronous distributed machines with no shared memory where the processors can communicate through a bounded degree network which is fixed. The problem is to simulate these algorithms on such a distributed machine. Let us explain first why this problem is important. Many efficient algorithms for abstract PRAMs have been published in recent years. Because of the relative simplicity of the PRAMs it is widely accepted that these models provide a desirable setting for the design of parallel algorithms. On the other hand, this model of synchronous distributed machines represents technologically feasible machines. Therefore, the theory of parallel computation has to provide a solution for this problem. (This discussion on desirable and feasible models of parallelism follows [15].) Awerbuch, Israeli, Shiloach, Schwartz and Upfal solved this simulation problem directly.

However, we suggest a different strategy. It is suggested to combine our solutions for the granularity problem together with the PDDI computer into an efficient solution for the simulation problem. This combining can be done with a relatively small effort that does not require any new ideas.

The efficiency of our solution compares favorably with other solutions. Let us also mention a few advantages of our strategy over a direct solution. (1) The granularity problem and the interconnection problem of processors and memories (solved by the PDDI) seem to be each considerably simpler than the simulation problem. (2) Each of these problems has the flavor of a theoretically basic problem since they overcome a basic difficulty, and are likely to arise in other applications. (3) Our strategy demonstrates a methodological application of randomization in the following sense. We first "clean" the simulation problem from difficulties that have satisfactory deterministic solutions and only then apply randomization. Specifically, the PDDI gives an efficient deterministic solution for the interconnection and routing problem. We already indicated that a naive deterministic solution must fail to give a good worst case solution

for the granularity problem. So, we are remained with the granularity problem ready to demonstrate that randomization is provably helpful. This is unlike [17] who applied randomization to the whole simulation problem. For more on this application see Chapter IV.

Part of this research can also be motivated by some data base applications; where, for instance, there are p processes (transactions), m servers (resources, disks) and N files distributed among the servers, so that each server may serve at most one process at a time (a resource is locked by a transaction). The case where a process may require only one file at a time readily fits our framework. However, a few alternative assumptions regarding how many files can be required simultaneously by the same process may reflect different circumstances. It might be interesting to investigate which such assumptions fits our framework and possible extensions to others. We do not elaborate on this motivation any more in this paper.

Each of our solutions for the granularity problem deals with two basic questions. (1) How to initially distribute the logical memory addresses of the EREWPRAM among the physical locations of the MPC? (2) How to compute the physical location of a logical address during the simulation?

We utilize two main ideas for the first question.

(a) *Randomization.* The logical addresses are randomly distributed among the memory modules. It is explained later in this paragraph why we have to be very strict about efficiencies of algorithms for computation of physical locations of logical addresses (in response to Question 2 above). The key idea behind the proposed approach is to utilize universal hashing in the simulating machine. The MPC itself picks at random a hash function from an entire class of hash functions, instead of a specific hash function. This function is used in order to distribute the logical address among the memory modules.

(b) *Copies.* We keep several copies of each logical address in distinct memory modules.

The impact of these two ideas on R_{\max} is as follows. Randomization is shown to keep memory contention low in the average. The copies idea enables to decrease memory contention in the worst case.

The second question. Observe that the copies idea requires not only to find where there is a copy of a requested memory address. It is not less important to assign each memory request to the 'right' copy. Namely, the assignment of memory requests to copies should be done in a way that takes into account simultaneous memory requests by other processors or the memory contention will not be low. There is a typical difficulty that we had to cope with in algorithms for the second question. Every algorithm we suggest for finding the physical location of a logical address or for assignment to a physical copy of a logical address is beneficial only if it is an efficient parallel algorithm. By efficient we mean that it is very fast and does not use too much local or common memory. Note that since the worst R_{\max} that we want to improve is p , our algorithms have to be significantly faster than that. We would also like that the size of the physical memory, and the space for storing the mapping from logical to physical addresses will be minimized.

The paper provides a three stage study of the granularity problem. The ideas of each stage can be applied separately or in conjunction with the others. Chapter II studies the first stage which involves applications of the randomization idea. The first stage is designed to keep us 'out of trouble', in the first place, in the average case. Chapter III studies the second stage which uses the copies idea. This idea, in conjunction with fast algorithms for picking the 'right' copy of each address request, is shown to decrease memory contention in the worst case, for the less fortunate cases of the first stage. Our above definition of the granularity problem made the third stage somewhat indistinct. In simulations of other models than the EREWPRAM by the MPC or other machines, the problem of simulation is not completely solved by specifying for each address request the module that satisfies it. Problems like scheduling the requests for a module (in case queues are not available) or combining simultaneous requests for the same address in the same module may arise. They are solved in the third stage which is described in Chapters IV.

II.1. A Probabilistic Simulation

In this section we begin to study a simple probabilistic simulation of PRAMs on MPCs. Consider a PRAM with p processing elements and a shared memory of size N . Also consider an MPC with p processing elements, and a shared memory of size N which is divided into m modules. More precisely, let memory module MM_j , $0 \leq j < m$, contain all (physical) addresses a with $0 \leq a < N$ and $a \bmod m = j$.

Our probabilistic simulation is based on *universal hashing* as introduced by Carter and Wegman. Let H be a subset of S_N , the full set of permutations of $[0 \dots N-1]$. We use elements of H to make the connection between logical and physical addresses. More precisely, we proceed as follows:

Initialization. Choose $h \in H$ at random and store h in every processing element of the MPC. The i -th PE of the MPC will run the same program as the i -th PE of the PRAM to be simulated. We maintain the invariant that cell $h(a)$ of the MPC has the same content as cell a of the PRAM for $0 \leq a \leq N-1$.

Step by Step Simulation. Let a_i be the (logical) address generated by the i -th PE of the MPC. Apply h to a_i and obtain (physical) address $b_i = h(a_i)$. Issue a request for memory cell b_i . This describes the behavior of the i -th PE of the MPC, $1 \leq i \leq p$. Memory module MM_j , $0 \leq j < m$, collects all requests for cells in MM_j and serves them sequentially. When all requests are served the next cycle of the PRAM is simulated.

Of course, the quality of the simulation described above depends crucially on class H of permutation used in the simulation. Note that the simulation is probabilistic because h is chosen at random from class H . We want

- 1) H to be small; because every PE needs additional local memory of at least $O(\log |H|)$ bits (assuming a suitable encoding) to store an element $h \in H$.
- 2) random elements of H to be easy to generate; because this will hold the cost of the initialization phase small.

- 3) elements $h \in H$ to be easy to evaluate; because this determines the cost of translating from logical to physical addresses.
- 4) the length of the queues arising in the simulation to be short; because they essentially determine the quality of the simulation.

We will next study expected queue length in more detail. Let $S = \{a_1, a_2, \dots, a_p\} \subseteq [0 \dots N - 1]$ be a set of p addresses. Let $h \in H$ be a permutation. Define

$$R_j(h, S) = |\{a \in S; h(a) \bmod m = j\}|, \quad 0 \leq j < m,$$

$$R_{\max}(h, S) = \max_{0 \leq j < m} R_j(h, S), \quad \text{and}$$

$$R_{\max} = \max_{\substack{S \subseteq [0 \dots N - 1] \\ |S| = p}} \sum_{h \in H} R_{\max}(h, S) / |H|$$

$R_j(h, S)$ is the length of the queue in front of memory module MM_j when permutation $h \in H$ is used and set S of addresses is issued by the processors. $R_{\max}(h, S)$ is the length of the longest queue in front of any memory module under the same conditions. Next $\sum_{h \in H} R_{\max}(h, S) / |H|$ is the expected value of $R_{\max}(\cdot, S)$. Finally, R_{\max} is the worst case of that value taken with respect to all possible sets of p addresses. In other words, R_{\max} is the worst case (with respect to addresses) expected (with respect to random elements of H) length of the longest queue. Or even less formally, we (the good guys) have to make sure that no matter which subset S of logical addresses will be specified by the bad guys, the expected performance of a random element h of H is still not too bad.

We want R_{\max} to be small. What can we expect? In order to get a feeling for R_{\max} we will briefly study a limiting case first: $H = S_N$, the full set of permutations of N elements. Note that class S_N is much too large ($N \log N$ bits local memory would be required in every processor) to be practically useful.

Theorem 1. *Let $H = S_N$ and let $m \geq p$ (suppose $m \geq 16$). Then*

$$R_{\max} \leq \min(\log m / \log(m/p), \log m / \log \log m) + 1$$

Proof. Let $S = \{a_1, \dots, a_p\} \subseteq [0 \dots N - 1]$ be arbitrary. Let $p_{k,j}$ be the probability that at least k elements of S are mapped into memory module j by a random element $h \in S_N$. Then $p_{k,j} \leq \binom{p}{k} (1/m)^k$ since images of different addresses are independent and uniformly distributed. Let p_k be the probability that at least k elements of S are mapped into some memory module. Then

$$p_k \leq p_{k,0} + p_{k,1} + \dots + p_{k,m-1} \leq \binom{p}{k} (1/m)^{k-1}.$$

Hence

$$R_{\max} = \sum_{k \geq 1} p_k \leq \sum_{k \geq 1} k(p_k - p_{k+1}) \leq \sum_{k \geq 1} \min\left(1, \binom{p}{k} (1/m)^{k-1}\right)$$

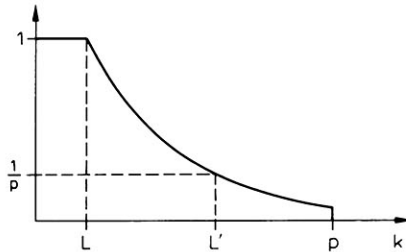
$$\leq \sum_{k \geq 1} \min(1, (m/k!)(p/m)^k).$$

Since for $m \geq 16$ we have $(p/m)^k(m/k!) \leq 1$ for $k \geq \min(\log m/\log(m/p), \log m/\log \log m)$ and decreases exponentially as a function of k for larger k , we have

$$R_{\max} \leq \min(\log m/\log(m/p), \log m/\log \log m) + 1 \quad \square$$

We infer from Theorem 1 that $R_{\max} \leq 1 + \log m/\log \log m$ if $m \geq p$ and $H = S_N$ and that $R_{\max} \leq 2 + 1/\varepsilon$ if $m = p^{1+\varepsilon}$ and $H = S_N$. Unfortunately, S_N is too large for our purposes. However, close inspection of the proof of Theorem 1 suggests methods for finding smaller classes of H which yield essentially the same value of R_{\max} as $H = S_N$. The diagram below shows a plot of a majoring function of p_k as a function of k . Our bound on p_k decreases exponentially for $k \geq L := \min(\log m/\log(m/p), \log m/\log \log m)$. In particular $p_k \leq 1/p$ for values of k exceeding L' where L' is only slightly larger than L . Quantity R_{\max} is the area under curve p_k . The following simple observation is crucial for the sequel: quantity R_{\max} increases by at most one if p_k were equal to $1/p$ for $k \geq L'$. Thus the bound on R_{\max} shown in Theorem 1 stays essentially true if we only know that $p_{k,j} \leq \binom{p}{k} (1/m)^k$ for $k \leq L'$ and that $p_{k,j}$ is non-increasing as a function of k .

We will explore this approach in the next section.



II.2. On the Expected Length of the Longest Chain in Universal Hashing

Universal hashing was introduced by Carter and Wegman. They showed that very small classes of hash function suffice to obtain an expected case behavior which is similar to the one of ordinary hashing. We show in this section that universal hashing is also competitive with respect to expected worst case behavior.

Definition [5]. Let $c \in \mathbb{R}$, $k \in \mathbb{N}$, $N \in \mathbb{N}$, $m \in \mathbb{N}$. A multiset $H \subseteq \{h; h: [0 \dots N - 1] \rightarrow [0 \dots m - 1]\}$ is c strongly k universal if for all $a_1, \dots, a_k \in [0 \dots N - 1]$, pairwise distinct, and all $b_1, \dots, b_k \in [0 \dots m - 1]$,

$$|\{h \in H; h(a_i) = b_i \text{ for } 1 \leq i \leq k\}| \leq c \cdot |H|/m^k. \quad \square$$

As in the preceding section, let $p \in \mathbb{N}$ and let $S \subseteq [0 \dots N - 1]$, $|S| = p$. For $h \in H$ let

$$R_{\max}(h, S) = \max_{0 \leq j < m} |\{a \in S; h(a) = j\}|,$$

and let

$$R_{\max}^p = \max_{\substack{S \subseteq [0 \dots N-1] \\ |S|=p}} \sum_{h \in H} R_{\max}(h, S)/|H|$$

We are now in a position to state the observation following Theorem 1 as

Theorem 2. *Let H be a c strongly k universal multiset of functions from $[0 \dots N-1]$ to $[0 \dots m-1]$. Then*

$$R_{\max}^p \leq k + cpm(p/m)^k/k!$$

for all $p \in \mathbb{N}$.

Proof. Let $S \subseteq [0 \dots N-1]$, $|S|=p$ be arbitrary. Let $p_i(S)$ be the probability that $R_{\max}(h, S) \geq i$, i.e. $p_i(S) = |\{h \in H; R_{\max}(h, S) \geq i\}|/|H|$. Then $1 \geq p_1 \geq p_2 \geq p_3 \geq \dots$ and

$$\begin{aligned} R_{\max}^p &= \max_{\substack{S \subseteq [0 \dots N-1] \\ |S|=p}} \sum_{i=1}^p p_i(S) \\ &\leq k + p \max_{\substack{S \subseteq [0 \dots N-1] \\ |S|=p}} p_k(S). \end{aligned}$$

Next observe that $p_k(S) \leq p_{k,0}(S) + p_{k,1}(S) + \dots + p_{k,m-1}(S)$ where $p_{k,j}(S)$ is the probability that at least k elements of S are mapped onto j . For fixed $a_1, a_2, \dots, a_k \in S$ (pairwise distinct) we have

$$|\{h \in H; h(a_l) = j \text{ for } 1 \leq l \leq k\}| \leq c|H|/m^k$$

since H is c strongly k universal. Hence $p_{k,j}(S) \leq c \binom{p}{k} / m^k \leq c(p/m)^k/k!$ and $p_k(S) \leq cm(p/m)^k/k!$ for all S . \square

Before we give some examples of strongly universal classes we recall the following lemma.

Lemma [5]. *If $H \subseteq \{h; h: [0 \dots N-1] \rightarrow [0 \dots N-1]\}$ is c strongly k universal and $r: [0 \dots N-1] \rightarrow [0 \dots m-1]$ is such that $|r^{-1}(j)| \leq \lceil N/m \rceil$ for all j , $0 \leq j < m$, then multiset*

$$\hat{H} = \{roh; h \in H\}$$

is \hat{c} strongly k universal where $\hat{c} = (m \lceil N/m \rceil / N)^k c$.³

Proof. Let $a_1, \dots, a_k \in [0 \dots N-1]$ be pairwise distinct and let $b_1, \dots, b_k \in [0 \dots m-1]$. Then there are at most $\lceil N/m \rceil^k$ tuples $c_1, \dots, c_k \in [0 \dots N-1]$ with $r(c_i) = b_i$ for $1 \leq i \leq k$. For every such tuple (c_1, \dots, c_k) there are at most $c|H|/N^k$ functions $h \in H$ with $h(a_i) = c_i$ for $1 \leq i \leq k$. Thus H is \hat{c} strongly k universal with $\hat{c} = (\lceil N/m \rceil / (N/m))^k c$. \square

We will next give some examples. Applications 1 and 2 are based on the fact that there are small doubly and triply transitive permutation groups. A set H of permutations of set X is transitive if for all $a, b \in X$ there is $h \in H$ such that $h(a) = b$. It is doubly (triply) transitive if it contains a permutation replacing

³ Let x be a real number. $\lceil x \rceil$ denotes the smallest integer $\geq x$.

any whatever given ordered pair (triple) of elements in X by any whatever ordered pair (triple) of elements in X . The reader may consult [4] for a detailed discussion. Application 3 uses the fact that a polynomial of degree k is fixed by its values at $k+1$ points, i.e. a random polynomial of degree k maps a set of $k+1$ points into a random set.

Application 1. Let N be a prime and let $H_1 = \{h; h(x) = (ax + b) \bmod N \text{ for some } a, b \in [0 \dots N-1], a \neq 0\}$, let $r(x) = x \bmod m$ and let $\hat{H}_1 = \{\text{roh}; h \in H_1\}$. Since for every $x_1, x_2, y_1, y_2 \in [0 \dots N-1]$, $x_1 \neq x_2, y_1 \neq y_2$ there is exactly one pair $a, b \in [0 \dots N-1]$ such that $y_1 = ax_1 + b \bmod N$ and $y_2 = ax_2 + b \bmod N$ class H_1 is 1 strongly 2 universal and hence \hat{H}_1 is $(m \lfloor N/m \rfloor / N)^2$ strongly 2 universal. For $m = p^3$ we obtain $R_{\max}^p \leq 2 + 4/2! = 4$; note that $(m \lfloor N/m \rfloor / N) \leq 2$. Finally, observe that H_1 is a set of permutations.

Application 2. Let N be a prime. For $a, b, c, d \in [0 \dots N-1]$ with $ad - bc \neq 0 \bmod N$ define $h_{a,b,c,d}: [0 \dots N-1] \cup \{\infty\} \rightarrow [0 \dots N-1] \cup \{\infty\}$ by

$$h_{a,b,c,d}(x) = \begin{cases} a/c & \text{if } x = \infty \\ \infty & \text{if } x = d/c \\ (ax - b)/(cx - d) \bmod N & \text{otherwise} \end{cases}$$

Note that division is well-defined since the integers mod N are a field. Also the first two clauses in the definition coincide if $c = 0$. It is known (a proof can be found in [4]) that $h_{a,b,c,d}$ is a permutation and that set $H_2 = \{h_{a,b,c,d}; a, b, c, d \in [0 \dots N-1], ad - bc \neq 0\}$ is a triply transitive group of permutations, i.e. for all x_1, x_2, x_3 and $y_1, y_2, y_3 \in [0 \dots N-1] \cup \{\infty\}$, x_1, x_2, x_3 pairwise distinct, there is exactly one $h \in H_2$ with $h(x_i) = y_i$ for $1 \leq i \leq 3$. Thus H_2 is 1 strongly 3 universal and hence $\hat{H}_2 = \{\text{roh}; h \in H_2\}$ and $r(x) = x \bmod m$ is 8 strongly 3 universal. Note that $(m \lfloor N/m \rfloor / N)^3 \leq 8$. Thus for $m = p^2$ we have $R_{\max}^p \leq 3 + 8/3! = 13/3$. Finally, we want to mention that $h_{a,b,c,d}(x)$ can be evaluated efficiently using Euklid's algorithm (cf. [1], p. 300–302). More precisely, $h_{a,b,c,d}(x)$ can be computed in $O(\log N)$ arithmetic steps.

Application 3. Let N be a prime, let k be an integer and let $H_3 = \{h; h(x) = \sum_{0 \leq i < k} a_i x^i \bmod N \text{ for some } a_i \in [0 \dots N-1] \text{ and } a_i \neq 0 \text{ for some } i \geq 1\}$ be the set of all polynomials of degree at most $k-1$. Since for every x_1, \dots, x_k and $y_1, \dots, y_k \in [0 \dots N-1]$, x_1, \dots, x_k pairwise distinct, there is at most one non-trivial polynomial h of degree at most $k-1$ with $h(x_i) = y_i$ for $1 \leq i \leq k$ we conclude that H_3 is 1 strongly k universal. Also $H_3 = \{\text{roh}; h \in H_3\}$ is \hat{c} strongly k universal with $\hat{c} = (m \lfloor N/m \rfloor / N)^k \leq (1 + m/N)^k$. Thus for $m = p^{1 + 2/(k-1)}$ we have

$$\begin{aligned} R_{\max}^p &\leq k + \hat{c} p^{2 + 2/(k-1)} p^{-\lfloor 2/(k-1) \rfloor k / k!} \\ &\leq k + \hat{c} / k! \end{aligned}$$

Another interesting choice is $k = 3 \ln p / \ln \ln p$ and $m = p$. Then

$$R_{\max}^p \leq k + \hat{c} p^2 / k! = O(\ln p / \ln \ln p).$$

Application 3 should be compared with Theorem 1. It states that if the class of hash functions is restricted to the set of all polynomials of degree at most $3 \ln p / \ln \ln p$ (note that there are only $N^{3 \ln p / \ln \ln p}$ such functions) then the expected worst case behavior is as good as if we use all hash functions (and there are N^N of those). Similarly, if $m = p^{1+\epsilon}$ and $k = 1 + 2/\epsilon$ for some $\epsilon > 0$ and if the class of hash functions is restricted to the set of polynomials of degree at most $k - 1$ (note that there are only $N^{k-1} = N^{2/\epsilon}$ such functions) then the expected worst case behavior is almost as good (except for a multiplicative factor of two) as if we use all hash functions (and here are N^N of those). Unfortunately, class H_3 is not a class of permutations and hence class H_3 cannot be directly used in simulating PRAMs by MPCs.

Example 1. Example 1 is slightly more difficult to treat; it is not a direct application of Theorem 2. Let $N = 2^n$ and $m = 2^b$. Then $[0 \dots N - 1]$ can be identified with the bit vectors of length n , i.e. $[0 \dots N - 1] \simeq \{0, 1\}^n$. The bit vectors of length n form a vector space of dimension n over the field of two elements. In a vector space we can use linear transformations to map any set of (linearly independent) vectors into any other set of vectors. This suggests to consider $H_4 = \{h; h: \{0, 1\}^n \rightarrow \{0, 1\}^n \text{ and } h(x) = Mx \text{ for some } n \text{ by } n (0, 1) - \text{invertible matrix } M\}$.

Here matrix multiplication is over the field of two elements. It is important for the application in Sect. II 3 that H_4 is a set of permutations. As before, let $r(x) = x \bmod m$ and $\hat{H}_4 = \{\text{roh}; h \in H_4\}$.

Before we analyze the behavior of class \hat{H}_4 we show that elements of \hat{H}_4 are easy to find by a probabilistic algorithm. More precisely, we show that a significant fraction of all $(0, 1)$ -matrices is invertible.

Lemma 1. $|H_4| = \prod_{0 \leq i \leq n-1} (2^n - 2^i) \geq 2^{(n^2)} / e^{7/5}$.

Proof. Lemma 1 is well known and can be found for example in *E. Artin: Geometric Algebra*. We include the very short proof for the sake of completeness. Choose M column by column. When columns 1 to i have been chosen (and are linearly independent) then column $i + 1$ must be different from all linear combinations of columns 1 through i . Hence there are $2^n - 2^i$ choices for column $i + 1$. Thus

$$\begin{aligned} |H_4| &= (2^n - 1)(2^n - 2) \dots (2^n - 2^{n-1}) \\ &= 2^{(n^2)} \prod_{i=1}^n (1 - 2^{-i}). \end{aligned}$$

Next observe that

$$\begin{aligned} \prod_{i=1}^n (1 - 2^{-i}) &\geq e^{\sum_{i=1}^n \ln(1 - 2^{-i})} \\ &\geq e^{-\sum_{i=1}^n (7/5) 2^{-i}} \quad \text{since } \ln(1 - x) \geq -7x/5 \text{ for } 0 \leq x \leq 1/2 \\ &\geq e^{-7/5}. \quad \square \end{aligned}$$

Lemma 1 shows that at least 25% of all (0, 1)-matrices are invertible. Hence invertible (0, 1)-matrices can be found by taking a few random (0, 1)-matrices and checking for singularity by Gaussian elimination. We will next show that H_4 is a good class of hash functions.

For x_1, x_2, \dots, x_k a set of vectors we write $\dim(x_1, x_2, \dots, x_k)$ to denote the dimension of the space spanned by x_1, \dots, x_k . Also, if $x \in \{0, 1\}^n$ we write x_b for the b -dimensional vector consisting of the last b components of x .

Lemma 2. *Let $a_1, \dots, a_k \in \{0, 1\}^n$, pairwise different. Let $d = \dim(a_1 - a_2, \dots, a_1 - a_k)$. Then*

$$|\{M; (Ma_1)_b = (Ma_2)_b = \dots = (Ma_k)_b\}| \leq 2^{(n^2)/m^d}.$$

Proof. If $(Ma_1)_b = \dots = (Ma_k)_b$ then $(M(a_1 - a_i))_b = 0$ for $2 \leq i \leq k$. Assume w.l.o.g. that $a_1 - a_2, \dots, a_1 - a_{d+1}$ are linearly independent. Let X be the n by d matrix whose columns are $a_1 - a_2, \dots, a_1 - a_{d+1}$. Let X_1 be the first d rows of X and let X_2 be the remaining $n - d$ rows. Assume w.l.o.g. that X_1 is non-singular. Let M' be the matrix consisting of the first b rows of M . Let M_1 consist of the first d columns of M' and let M_2 consist of the remaining $n - d$ columns. Then

$$M_1 X_1 + M_2 X_2 = 0$$

or

$$M_1 = -M_2 X_2 X_1^{-1}.$$

Thus M_1 is determined by the choice of M_2 and hence there are at most $2^{n^2 - bd} = 2^{n^2}/m^d$ matrices M such that $(Ma_1)_b = \dots = (Ma_k)_b$. Recall that $m = 2^b$. \square

Lemma 3. *Let $S = \{a_1, \dots, a_p\} \subseteq \{0, 1\}^n$. Let $t_{k,l}$ be the number of subsets of S of cardinality k and dimension l . Then*

$$t_{k,l} \leq \binom{p}{l} \binom{2^l}{k-l}.$$

In particular, $t_{k,l} = 0$ for $2^l + l < k$.

Proof. Any subset of cardinality k and dimension l can be written as a set of l linearly independent elements plus a set of $k - l$ vectors which are linear combinations of the first l elements. There are only $\binom{p}{l}$ choices for the first set and only $\binom{2^l}{k-l}$ choices for the second set. \square

We are now in a position to estimate the behavior of class \hat{H}_4 . Let $S = \{a_1, a_2, \dots, a_p\} \subseteq \{0, 1\}^n$ be arbitrary. As above, let

$$p_k = \text{prob}(R_{\max}(h, S) \geq k) = |\{h \in \hat{H}_4; R_{\max}(h, S) \geq k\}| / |\hat{H}_4|$$

Lemma 4. a) $p_1 \geq p_2 \geq p_3 \geq \dots$

b) $p_k \leq cm \sum_{l \geq 0} t_{k,l} m^{-l}$ where $c = e^{7/5}$

c) $p_k \leq 2cm(p2^k/m)^{\log k - 1}$ for $2 \leq k \leq \log m/p - 1$.

Proof. a) obvious.

b) We have

$$\begin{aligned}
 p_k &\leq \sum_{l \geq 0} \sum_{\substack{A \subseteq S \\ |A|=k \\ \dim A=l}} |\{h \in \hat{H}_4; h \text{ maps all points of } A \text{ into the same location}\}| / |\hat{H}_4| \\
 &\leq \sum_{l \geq 0} \sum_{\substack{A \subseteq S, |A|=k \\ \dim A=l}} c/m^{l-1}
 \end{aligned}$$

by Lemma 1, 2 and the observation that $\dim A = l$, $A = \{a_1, \dots, a_k\}$ implies $\dim(a_1 - a_2, \dots, a_1 - a_k) \geq l - 1$

$$\leq c \sum_{l \geq 0} t_{k,l} m^{1-l}$$

by the definition of $t_{k,l}$. This proves part b).

c) For $k \geq 2$ we have $t_{k,l} = 0$ for $l < \log k - 1$. Hence by part b) and Lemma 3.

$$\begin{aligned}
 p_k &\leq c m \sum_{\log k - 1 \leq l \leq k} \binom{p}{l} \binom{2^l}{k-l} m^{-l} \\
 &\leq c m \sum_{\log k - 1 \leq l \leq k} (p 2^k / m)^l \\
 &\leq c m (p 2^k / m)^{\log k - 1} \sum_{l \geq 0} (p 2^k / m)^l \\
 &\leq 2 c m (p 2^k / m)^{\log k - 1}
 \end{aligned}$$

since $p 2^k / m \leq 1/2$ for $k \leq \log(m/p) - 1$. \square

Lemmas 1 to 4 directly lead to

Theorem 3. Let $N = 2^n$, $m = 2^b$ and let \hat{H}_4 be defined as above. a) there is a function $f: \mathbb{N} \rightarrow \mathbb{R}$, $f(p) = O(p^\epsilon)$ for all $\epsilon > 0$ such that for $m \geq p f(p)$ we have:

$$R_{\max}^p \leq 20 \log p / (\log(10 \log p))^2 + O(1).$$

b) Let $\epsilon > 0$. For $m = p^{1+\epsilon}$ we have:

$$R_{\max}^p \leq O(1).$$

Proof. Note first that

$$R_{\max}^p = \sum_{1 \leq k \leq p} p_k \leq k_1 + p p_{k_1}$$

for every k_1 , $1 \leq k_1 \leq p$. From Lemma 4, c we conclude further

$$R_{\max}^p \leq k_1 + 2 c m p (p 2^{k_1} / m)^{\log k_1 - 1}$$

where $c = e^{7/5}$.

a) Let $f(p) = \max(2^{2^{10}}, 2^{10 \log p / \log(10 \log p)})$. Then $f(p) = O(p^\epsilon)$ for all $\epsilon > 0$. Let $k_1 = \log f(p) / \log \log f(p)$. Then $k_1 \leq (\log f(p)) / 10$ and hence

$$p 2^{k_1} / m \leq p f(p)^{1/10} / p f(p) = f(p)^{-9/10}.$$

Also $\log k_1 - 1 = \log \log f(p) - \log \log \log f(p) - 1 \geq (\log \log f(p))/2$ since $\log \log f(p) \geq 10$. Thus

$$\begin{aligned} pp_{k_1} &\leq 2c m p f(p)^{-9 \log \log f(p)/20} \\ &\leq 2c p^2 f(p) f(p)^{-9 \log \log f(p)/20} \\ &\leq 2c 2^{2 \log p + \log f(p) - 9 \log f(p) \log \log f(p)/20} \\ &\leq 2c 2^{2 \log p - 7 \log f(p) \log \log f(p)/20} \text{ since } \log \log f(p) \geq 10 \\ &= O(1) \text{ since } 7 \log f(p) \log \log f(p)/20 \geq 2 \log p \end{aligned}$$

for sufficiently large p . Thus

$$\begin{aligned} R_{\max}^p &\leq \log f(p)/\log \log f(p) + O(1) \\ &\leq 20 \log p / (\log(10 \log p))^2 + O(1). \end{aligned}$$

b) Let $m = p^{1+\epsilon}$. Then

$$R_{\max}^p \leq k_1 + 2c p^{2+\epsilon} (2^{k_1}/p^\epsilon)^{\log k_1 - 1}$$

for all k_1 , $1 \leq k_1 \leq p$. Let k_1 be such that $2 + \epsilon = \epsilon(\log k_1 - 1)$, i.e. $k_1 = 2^{2+2/\epsilon}$. Then

$$\begin{aligned} R_{\max}^p &\leq k_1 + 2c p^{2+\epsilon - \epsilon(\log k_1 - 1)} 2^{k_1(\log k_1 - 1)} \\ &= 2^{2+2/\epsilon} + 2c 2^{(1+2/\epsilon)2^{2+2/\epsilon}} \\ &= O(1). \quad \square \end{aligned}$$

Theorem 3 states that the performance of class \hat{H}_4 is “close” to the performance of the full set of permutations. More precisely, if $m = p^{1+\epsilon}$ then $R_{\max}^p = O(1)$ in both cases. Of course, the constants involved are dramatically different. Also $R_{\max}^p = O(\log p / \log \log p)$ can be achieved by class \hat{H}_4 for $m = pf(p)$ where $f(p)$ grows slower than any root of p . If the full set of permutations is used there $R_{\max}^p = O(\log p / \log \log p)$ can be achieved for $m = p$.

Why doesn't class \hat{H}_4 behave well in the case $m = p$? The answer comes in two parts. Firstly, the bounds derived in Lemmas 3, 4 and hence in Theorem 3 are not sharp. Secondly and more significantly, class \hat{H}_4 has a certain deficiency. Multiplying by a random, invertible matrix hashes a set of *independent* vectors very well, however, it does not do that well on a set of linearly dependent vectors. A variant of Lemma 1 can be used to show that the expected dimension of a *random* set a_1, \dots, a_n of vectors in $\{0, 1\}^n$ is very large, namely $n - O(1)$. Unfortunately, this observation is of no use since we are dealing with worst case behavior.

We will now discuss a possibility to overcome this problem. As above, let $N = 2^n$, $m = 2^b$. Assume also that $q = N - c$ for some small constant c is a prime. For $a \in [1 \dots q - 1]$ let

$$h_a: [0 \dots 2^n - 1] \rightarrow [0 \dots 2^n - 1]$$

be defined by

$$h_a(x) = \begin{cases} (ax) \bmod q & \text{if } x < q \\ x & \text{if } q \leq x \leq N - 1. \end{cases}$$

Note that h_a is a permutation. We conjecture

Conjecture. Let $a_1, \dots, a_p \in \{0, 1\}^n \equiv [0 \dots 2^n - 1]$ be arbitrary. For $a \in [1 \dots q - 1]$ let \dim_a be the dimension of set $h_a(a_1), \dots, h_a(a_p)$. Then

$$\sum_{1 \leq a < q} \dim_a / (q - 1) \geq \min(n, p) - f$$

for some small constant f . Moreover,

$$|\{a; 1 \leq a < q \text{ and } \dim_a = \min(n, p) - f - i\}| \leq qp^{-i}$$

for all $i \geq 1$.

Informally, the conjecture states that a random function h_a turns a set of p vector a_1, \dots, a_p into a (nearly) independent set of vectors. Moreover, it is very unlikely that the dimension of the resulting set of vectors is much less than the expected dimension.

Consider class

$$\hat{H}_4 = \{h: \{0, 1\}^n \rightarrow \{0, 1\}^b; h(x) = (M \cdot h_a(x)) \bmod 2^b$$

for some $a, 1 \leq a < q$, and some invertible n by n $(0, 1)$ -matrix. $\}$. We next show that class \hat{H}_4 is a very good class of hash functions provided that the conjecture is true.

Theorem 4. Let \hat{H}_4 be defined as above and let $m \geq 2p$. If the conjecture above is true then

$$R_{\max}^p \leq k + 3cmp^{f+1}/k!$$

for all $k, f < k \leq \min(n, p)$. Here $c = e^{7/5}$.

Proof. Let $S = \{a_1, \dots, a_p\} \subseteq [0 \dots N - 1]$. As above, let $p_k = \text{prob}(R_{\max}(h, S) \geq k)$. Then

$$\begin{aligned} p_k &= \sum_{a=1}^{q-1} |\{h \in \hat{H}_4; h \text{ maps at least } k \text{ points} \\ &\quad \text{of } h_a(S) \text{ into the same location}\}| / |\hat{H}_4| \\ &\leq \sum_{a=1}^{q-1} \sum_{l \geq 0} \sum_{\substack{A \subseteq S \\ |A|=k \\ \dim h_a(A)=l}} |\{h \in \hat{H}_4; h \text{ maps all points of } h_a(A) \\ &\quad \text{into the same location}\}| / |\hat{H}_4| \\ &\leq \sum_{a=1}^{q-1} \sum_{l \geq 0} \sum_{\substack{A \subseteq S \\ |A|=k \\ \dim h_a(A)=l}} cm^{1-l}/q \end{aligned}$$

by Lemma 2

$$\begin{aligned} &= \sum_{\substack{A \subseteq S \\ |A|=k}} \sum_{l \geq 0} |\{a; \dim h_a(A)=l\}| cm^{1-l}/q \\ &\leq \sum_{\substack{a \subseteq S \\ |A|=k}} \sum_{l=0}^{\min(n, k) - f - 1} cm^{1-l} p^{-(\min(n, k) - f - l)} \end{aligned}$$

$$\begin{aligned}
 &+ \sum_{\substack{A \subseteq S \\ |A|=k}} c m^{1 - (\min(n, k) - f)} \\
 &\leq 2c \binom{p}{k} m \cdot p^{f-k} + c \binom{p}{k} m^{1+f-k}
 \end{aligned}$$

since $k \leq n$ and $m \geq 2p$ by assumption

$$\leq 3c \binom{p}{k} m p^{f-k} = 3c m p^f / k!$$

since $m \geq 2p$ and $k > f$ by assumption. The proof is now completed since $R_{\max}^p \leq k + p p_k$ for all k . \square

Theorem 4 has an interesting consequence. Assume that $m = 2p$ and $n \geq \log p$. The latter assumption is certainly realistic. Let $k = (f+2) \log p / \log \log p$. Then $R_{\max}^p \leq (f+2) \log p / \log \log p + 6$. This is basically the same behavior as the behavior of the full class of permutations; cf. Theorem 1.

II.3. Probabilistic Simulations Revisited

We will now apply the results of Sect. II.2 to the probabilistic simulation described in Sect. II.1. We assume throughout this section that operations on addresses, like multiplying by an integer, take unit time.

Theorem 1. *Let $m = p^3$. Then a $T(n)$ - time bounded PRAM with p PE's and N memory cells, can be simulated by a randomized MPC with p PE's and m memory modules and total memory of size $N + 2p$ in time $O(T(n))$.*

Proof. We use the simulation as described in Sect. II.1 except that permutation π is chosen from class H_1 . Every processor needs two additional storage cells to store π . Also one step of the PRAM takes expected time $O(1)$ on the MPC. \square

Theorem 2. *Let $m = p^2$. Then a $T(n)$ time bounded PRAM with p PE's and N memory cells can be simulated by a randomized MPC with p PE's m memory modules and total memory size $N + 4p$ in time $O(T(n) \log N)$.*

Proof. Replace H_1 by H_2 in the proof of Theorem 1 and observe that the functions in H_2 can be evaluated in time $O(\log N)$. \square

Theorem 3. *a) Let $m = p$. Then a $T(n)$ time bounded PRAM with p PE's and N memory cells can be simulated by a randomized MPC with p PE's, m memory modules and total memory size $(N + p) \log p$ in time $O(T(n) \log p)$.*

b) Let $\epsilon > 0$ and let $m = p^{1+\epsilon}$. Then a $T(n)$ time bounded PRAM with p PE's and N memory cells can be simulated by a randomized MPC with p PE's, m memory modules and total memory size $(1 + 2/\epsilon)(N + p)$ in time $O(T(n))$.

Proof. a) We use the simulation as directed in Sect. II.1 except that π is chosen from class H_3 with $k = 3 \log p / \log \log p$. Then every processor requires $O(\log p)$

cells to store π . There is one additional problem now: class H_3 is not a class of permutations. Rather $\pi \in H_3$ might map up to $\log p / \log \log p$ distinct points into the same cell. It can certainly map no more distinct into the same cell since $\pi \in H_3$ is a nontrivial polynomial of degree at most $\log p / \log \log p$. Therefore the simulating MPC has $\log p / \log \log p$ copies of every memory cell. A memory access is made by sending the original PRAM address a and the modified address $\pi(a)$ to the appropriate memory module. We can then build up a balanced tree (say) for all addresses which are mapped to the same address by π . Thus access time within a memory module might be as large as $O(\log \log p)$. Also the expected number of concurrent accesses to the same module is $O(\log p / \log \log p)$ by Sect. II.2, application 3 and hence it takes an expected number of $O(\log p)$ MPC-steps to simulate our PRAM step. Also evaluation of a hash function in H_3 takes $O(\log p)$ steps.

b) The proof is completely analogous to the proof of part a). We choose $k = 1 + 2/\varepsilon$. Then at most $2/\varepsilon$ distinct points are mapped into the same cell. Thus the total memory size is $(1 + 2/\varepsilon)(N + p)$ and it takes an expected number of $O(2/\varepsilon \log 2/\varepsilon)$ MPC steps to simulate one PRAM step. \square

Theorem 4. *a) There is a function $f: \mathbf{N} \rightarrow \mathbf{R}$ with $f(p) = O(p^\varepsilon)$ for all $\varepsilon > 0$ such that: if $N = 2^n$, $m = 2^b \geq pf(p)$ then a $T(n)$ time bounded PRAM with p PE's and memory size N can be simulated on a randomized MPC with p PE's, m memory modules and total memory size $N + p \log N$ in time*

$$O((\log N)^3 + T(n)(\log p / (\log \log p)^2 + \log N)).$$

b) Let $\varepsilon > 0$ be fixed. If $m = p^{1+\varepsilon}$ then the time bound reduces to $O((\log N)^3 + T(n) \log N)$.

Proof. Replace H_1 by H_4 in the proof of Theorem 1 and observe that an element of H_4 can be stored in $\log N$ words of length $(\log N)$ each. Thus every processor needs $\log N$ additional memory cells. A random element of H_4 can be chosen as follows. Generate $\log N$ random bitstrings of length $\log N$ each. (In time $O(\log N)^2$) and check whether the $(0-1)$ matrix M generated in this way is invertible. This takes time $O((\log N)^3)$ on a sequential machine. Also $O(1)$ tries suffice on the average. Thus choosing a random element of H_4 takes time $O((\log N)^3)$. Next observe that it takes time $O(\log N)$ to multiply an $\log N$ by $\log N$ matrix by a vector. The time bounded is now an easy consequence of Theorem 3 of Sect. II.2. \square

We do not feel that Theorems 1 to 4 are best possible. They complement each other in that they optimize different parameters of the problem. Theorems 1, 2 and 4 present solutions with only a very moderate increase in total memory size ($O(1)$ additional cells per processor in Theorems 1 and 2 and $O(\log N)$ in Theorem 4) and only a moderate increase in running time (a multiplicative factor of $O(1)$ in Theorem 1 and $O(\log N)$ in Theorems 2 and 4). However, all three theorems require a non-linear number of memory modules thus increasing the size of the interconnection network. Theorem 4a is our best result in that respect.

On the other hand, both parts of Theorems 3 provide us with solutions with a small number of memory modules (p in part *a* and $p^{1+\epsilon}$ in part *b*) and modest increases of running time ($O(\log p)$) in part *a* and $O(1)$ in part *b*). However, both parts force us to increase total memory size considerably.

Theorem 4 of the previous section has the potential (if the conjecture were true) of combining most advantages of the other schemes. With only two memory modules per processor and only $\log N$ additional memory cells per processor it achieves a slowdown of $O(\log p / \log \log p)$.

III. Efficient Deterministic Simulations

Say that each of the N addresses is contained in exactly c of the m memory modules for some integer $c(c > 1)$. Each such distribution of addresses is called a c -partitioning.

We may break each cycle of the EREWPRAM into two halves: one includes all read instructions from the common memory, while the other includes all the write instructions of the cycle. This enables us to classify the cycles into reading cycles and writing cycles without multiplying the running time by more than a factor of two. We argue, in the paper, that the worst case time for simulating writing cycles worsens only a little while the worst case time for simulating reading cycles improves a lot. Many simulations of programs for the Ultracomputer have been run. The ratio between read instructions and write instructions that relate to the common memory was around 8:1. This fact is important since in case a store instruction into the common memory is executed we need access to all copies of the memory address, while in case a fetch instruction is executed one of the copies would suffice. See more on writing cycles on the last chapter. Our main concern is to study the advantages and limitations of the copying approach for simulating reading cycles. Our analysis applies also to input addresses (which are only read).

We start this chapter by studying some limitations of the copying approach. Then, a specific c -partitioning is proposed. For this c -partitioning we give upper bounds for the optimal R_{\max} achievable. We include also two sections that discuss how to compute efficiently good assignments of address requests to modules. The first of these sections deals with polynomial-time sequential algorithms for computing optimal assignments, namely with minimum R_{\max} . The second section suggests fast parallel algorithms that give low (but not necessarily minimum) R_{\max} . For the purpose of fast simulation we only need algorithms of the second kind. The proposed c -partitioning is very efficient when this stage is applied separately. However, an alternative c -partitioning which combines well with the first stage is considered subsequently. The last section of this chapter demonstrates that a number of copies c is useful when $\binom{m}{c}$ is much larger than N .

Note that the number of simultaneous requests for memory addresses is $\leq p$. In order to use one parameter less, and thereby simplify the presentation,

we consider the most difficult case only; i.e., where this number is p . It will be straightforward to extend our results to cases where this number is $< p$.

III.1. Lower Bounds

Assume that some c -partitioning is given. Problem: Find a lower bound for the optimal worst case time delay (R_{\max}) for any p reading requests for memory addresses.

Let $1 \leq i_1 < i_2 \dots < i_k \leq m$ be k modules and A_{i_1, i_2, \dots, i_k} be the set of all memory addresses such that all their c copies are contained in memory modules i_1, i_2, \dots, i_k . We are looking for a (very unfortunate) set of p memory addresses such that all their copies are contained in a minimum number of modules. More formally, we are looking for a minimum size subset of modules $\{i_1, i_2, \dots, i_k\}$ such that $\#A_{i_1, i_2, \dots, i_k} \geq p$.

Claim. The equation $\sum \#A_{i_1, i_2, \dots, i_k} = \binom{m-c}{k-c} N$ holds for any $k, c \leq k \leq m$.

The summation on the left hand side is on all $\binom{m}{k}$ subsets of k modules. The right hand side gives an alternative evaluation which is based on the fact that each memory address is contained in c modules; fixing $k-c$ additional modules (in any of the $\binom{m-c}{k-c}$ possibilities) gives a subset of k modules containing this address. This completes the proof of the claim.

Hence, there is at least one subset of k modules that contains

$$\binom{m-c}{k-c} N / \binom{m}{k} = k \dots (k-(c-1)) N / (m \dots (m-(c-1)))$$

elements. We are looking for the minimum k such that

$$(1) \quad p \leq k \dots (k-(c-1)) N / (m \dots (m-(c-1))).$$

Denote this k by K_p . This implies that

Observation 1. $R_{\max} \geq p/K_p$. If $N = \binom{m}{c} x$ for some integer x , then from inequality (1) we get $p \leq k \dots (k-(c-1)) x/c!$ or $p \leq \binom{k}{c} x$. This implies

Observation 2. For $N = \binom{m}{c} x$, $R_{\max} \geq p/K_p$ where K_p is the smallest k satisfying $p \leq \binom{k}{c} x$.

It is shown later that the last lower bound meets *exactly* an upper bound on R_{\max} for a specific c -partitioning that we propose. Therefore, we delay presentation of explicit evaluations until this later discussion.

III.2. The Proposed c -Partitioning

Our suggested c -partitioning is simple. For $N \leq \binom{m}{c}$ and an address i , $0 \leq i < N$, take the i -th subset of c modules (out of $\binom{m}{c}$ in the lexicographic order) and put a copy of address i in each of the c modules. If $(x - 1) \binom{m}{c} < N \leq x \binom{m}{c}$, for some integer x , then partition the addresses into x approximately equal subsets (layers) and fix the c -partitioning of each of the x layers separately as for $N \leq \binom{m}{c}$.

Remark. For an address i it takes $O(c^2)$ time to compute the subset of c modules containing its copies. Clearly, $i'_1 = i \pmod{\left(\left\lceil \frac{N}{\binom{m}{c}} \right\rceil\right)}$ is the serial number of i in its layer. The minimum $i_1 (\geq c)$ such that $\binom{i_1}{c} \geq i'_1$ implies that module number $i_1 - 1$ contains a copy of address i . Let us explain this. We have m modules numbered $0, 1, \dots, m - 1$. The lexicographic order implies that all c copies of the first $\binom{i_1 - 1}{c}$ addresses of the present layer are distributed among modules $0, 1, \dots, i_1 - 2$. Then, the next $\binom{i_1}{c}$ addresses (including the address being considered) have a copy at module $i_1 - 1$ and their other $c - 1$ copies are distributed among modules $0, 1, \dots, i_1 - 2$. The computation of i_1 , takes $O(c)$ time. This is because we get a constant difference approximation to i_1 by Stirling formula and explicit presentation of i_1 .

Denote $i'_2 = i_1 - \binom{i_1 - 1}{c}$. The minimum $i_2 (\geq c - 1)$ such that $\binom{i_2}{c - 1} \geq i'_2$ implies that module number $i_2 - 1$ contains a copy of address i . This takes $O(c - 1)$ time; and so on.

III.3. Upper Bounds

Theorem 1. Let t be an integer, $N = x \binom{m}{c}$, $S_f = \min \left\{ s; \binom{s}{c} x \geq (t - 1)s + 1 \right\}$ and $r = S_f(t - 1) + 1$. If $S_f \leq m$ then: (a) There exist r address requests which cause memory contention of at least $t (R_{\max} \geq t)$ for any assignment of requests to copies. (b) For any $r - 1$ requests, however, it is possible to get $R_{\max} \leq t - 1$.

Proof. We show (b) first. Say, in contradiction, that a set of $r - 1$ requests is given such that the best R_{\max} obtainable is t . (It will be easy to modify our argument if the best R_{\max} that can be obtained is greater than t). Among all possible ways to partition requests among our modules which imply $R_{\max} = t$ choose these that assign minimum number of modules with exactly t requests. Then, restrict further the choice to a partition where a module with the

smallest serial number possible is assigned with t requests. So we have t requests for addresses in some module. All these addresses have other copies in other modules. Let S_1 denote a lower bound on the number of modules containing all copies of these t addresses and T_1 denote a lower bound on the total number of address requests assigned to these modules. The T_1 addresses may have copies in other modules. Let S_2 denote a lower bound on the number of modules containing all copies of these T_1 addresses and T_2 denote a lower bound on the total number of address requests assigned to these modules. It should be clear how to define S_3, S_4, \dots and T_3, T_4, \dots .

Claim 1. *Each module which enters the scene is assigned with at least $t-1$ requests.*

Proof. Our sequence satisfies the following. There exists a directed path of the following form from the first module (the one we started with, which is assigned with t requests) to each module which is counted by the S_i numbers; the nodes along the path are modules; and there is a directed edge from module A to module B if module A is assigned with some address request and another copy of this address is located in module B . By propagating requests along this path we may decrease the number of requests assigned to the first module by one, increase the number of requests assigned to the last module on the path by one and not change the number of requests assigned to any other module. Thus, the existence of a module as in Claim 1, which is assigned with less than $t-1$ requests contradicts the choice of the partitioning of requests among modules above.

Corollary 1. $T_i = S_i(t-1) + 1$, for $i \geq 1$. And since $S_{i+1} = \min \left\{ s; \binom{s}{c} x \geq T_i \right\}$ we get,

$$S_1 = \min \left\{ s; \binom{s}{c} x \geq t \right\}, \quad S_{i+1} = \min \left\{ s; \binom{s}{c} x \geq S_i(t-1) + 1 \right\}, \quad \text{for } i > 1.$$

Claim 2. *The sequence (of integers) $\{S_i\}$ converges to S_f .*

Proof. Recall $S_f \leq m$. If $S_i < S_f$ then S_{i+1} satisfies $S_i < S_{i+1} \leq S_f$. If $S_i = S_f$ then $S_{i+1} = S_f$.

So we could not start with less than $S_f(t-1) + 1$ requests. A contradiction.

Proving (a) is easy. Take $r (= S_f(t-1) + 1)$ requests that all their copies are in S_f modules. This completes the proof of Theorem 1. \square

The Connection with the Lower Bound

We showed that for any r address requests the smallest number of modules to contain their copies is the smallest k satisfying $r \leq \binom{k}{c} x$, which is exactly our S_f . Then we concluded that R_{\max} is at least the smallest integer satisfying $R_{\max} \geq r/S_f$. Here this integer is t which is exactly the best R_{\max} achievable. So, the upper bound and the lower bound are *exactly equal*.

For any p requests let us establish an explicit upper bound on the best R_{\max} achievable. Denote this R_{\max} by t . Let x be the integer such that

$$(x-1) \binom{m}{c} < N \leq \binom{m}{c}.$$

Similar to the proof of Theorem 1 we restrict the choice of the assignments of requests to modules. Then, we start with a module which is assigned with t requests. Now, s is the minimum number of modules that have to be assigned with at least $t-1$ requests in order not to enable us decrease t in the first module. The following inequality holds from similar reasons to the proof of

Theorem 1, $x \binom{s}{c} \geq (t-1)s + 1$. It implies

$$x s^c / c! \geq (t-1)s \rightarrow s \geq (t-1)c! / x^{1/c-1}$$

Together with the fact $p \geq (t-1)s + 1$ we get

$$(t-1)((t-1)c! / x)^{1/c-1} \leq p$$

implying

$$(t-1)^{c/c-1} \leq p / (c! / x)^{1/c-1},$$

$$(t-1)^c \leq x p^{c-1} / c!$$

and

$$t \leq 1 + (x p^{c-1} / c!)^{1/c}.$$

Thus,

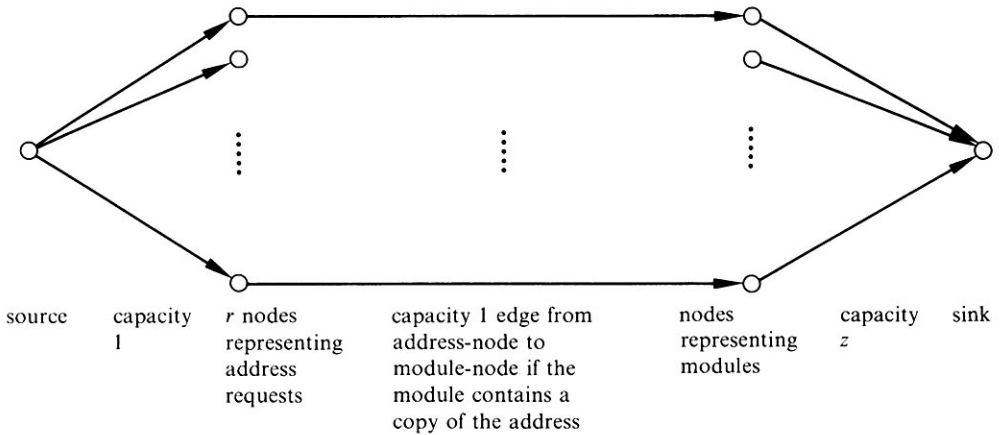
Theorem 2. $R_{\max} \leq 1 + \left(p^{c-1} / c! \left[N / \binom{m}{c} \right] \right)^{1/c}.$

III.4. Optimal Assignment

The problem of optimal assignment to the ‘right’ copy of each requested memory address, in order to minimize R_{\max} , is solvable in polynomial time. We actually solve the problem for any distribution of copies among modules (not only c -partitionings), in Fig. 1. For more on the Max-Flow problem, see [6].

An Alternative Solution

We have to assign p requests to modules. Assign one request at a time. Say that i out of the p address requests, for some $0 \leq i < p$, got already (temporary) assignment to modules. Assume that these i requests are assigned to modules in a way which minimizes R_{\max} with respect to them. Denote this R_{\max} by $R_{\max, i}$. We show how to extend this assignment to another temporary assignment for one more request, in a way which achieves minimum R_{\max} with respect to the $i+1$ requests. Our algorithm and proof are similar to the proof of Theorem 1. Consider the following auxiliary directed graph. Nodes represent modules. There is an edge from module A to module B if module A is assigned with a



The optimal assignment algorithm: find (by binary search) minimum z that enables flow of r .

Fig. 1. Finding optimal assignment to copies by utilizing instances of Max-flow

request for some address a while B contains another copy of a . The algorithm for extending the assignment of i requests to an assignment of $i + 1$ requests is as follows.

(1) Assign the $(i + 1)$ -st address request to a module M_1 containing one of the copies of the address. (Add to the auxiliary digraph $c - 1$ edges from this module to the other modules that contain copies of this address.)

(2) If M_1 is now assigned with $\leq R_{\max,i}$ requests then we are done. (It is impossible to add a request and get $R_{\max,i+1} < R_{\max,i}$.)

(3) If M_1 is now assigned with $R_{\max,i} + 1$ requests search the auxiliary digraph for a module which is assigned with less than $R_{\max,i}$ requests and is reachable through a directed path from M_1 . (Any efficient search can be utilized here. For instance Breadth First Search.) If such a module is found then propagate request assignments along a path in the digraph from M_1 to this module. (This results by the following change: the number of requests assigned to this module increases by one. Note, however, that it is still $\leq R_{\max,i}$. The number of request assigned to all other module is exactly the same as it was for the i requests.)

(4) If no such module is found in (3) do nothing (the request is assigned to M_1).

The only thing that has to be proved is that step (4) is correct. In other words: why is it impossible to assign these $i + 1$ requests with $R_{\max,i+1} = R_{\max,i}$ rather than $R_{\max,i+1} = R_{\max,i} + 1$? Assume, in contradiction, that there is an assignment of these $i + 1$ requests with $R_{\max,i+1} = R_{\max,i}$. In this contradictory assignment M_1 is assigned with $\leq R_{\max,i}$ requests. Therefore, there exists a request for some address a , which is assigned to M_1 in our assignment (the result of step (4)) but not in the contradictory assignment. Hence, the auxiliary graph contain an edge from M_1 to M_2 - the node representing the module that the request for a is assigned to it in the contradictory assignment. By step (3), M_2 has $R_{\max,i}$ request assigned to it in our assignment. Define S , an auxiliary set of modules, to include presently M_1 and M_2 . In our assignment both

modules of S have together $2R_{\max,i} + 1$ requests assigned to them while in the contradictory assignment there are at most $2R_{\max,i}$ such requests. Therefore, there exist a request for some address a_2 which is assigned to a module of S in our assignment but not in the contradictory assignment where this request is assigned to another module not in S , say M_3 . This implies the existence of an edge in the auxiliary graph from this module in S to M_3 . By step (3) M_3 has $R_{\max,i}$ requests assigned to it in our assignment. Add M_3 to S . Similarly we show that the set S can grow infinitely large. The number of modules is finite, so we got a contradiction. Therefore, the assignment achieved by our algorithm yields the minimum R_{\max} .

III.5. Fast Parallel Approximation Algorithms

The above mentioned optimal assignment algorithms are of general theoretical interest and might be relevant for data-base applications as mentioned in the introduction. However, the simulation requires fast parallel algorithms for the assignment problem even if the optimal result is not achieved (remember that we are after fast simulation of one EREWPRAM cycle). This poses a challenge of a new kind.

Lemma. *Let $N \leq \binom{m}{c}$ be the number of addresses. For p address requests we can achieve $R_{\max} \leq cp^{(c-1)/c}$ in parallel time $O(c^2 + c \log p)$. For $c=2$, we can do it in constant parallel time.*

Remark. The claim for $c=2$ needs a little bit stronger assumptions about the MPC; like, a memory request can be dropped by the requesting processor, or, alternatively, a memory module can discard memory requests sent to it if their location in the module's queue exceeds a certain point.

Proof of Lemma. By induction on c .

For $c=2$ start with any assignment of address requests to copies. For all modules where the number of requests is $\geq \sqrt{p}$ (there exist $\leq \sqrt{p}$ such modules) switch all requests above the \sqrt{p} line (by doing that for all requests that were not responded in the first \sqrt{p} time units of the cycle which is being simulated, we avoid computation of serial numbers for $c=2$) to the other copy. Second copies belong to pairwise distinct modules. Therefore, each module gets no more than \sqrt{p} requests for second copies. Assume by induction that the Lemma holds for $c-1$. We want to show that it holds for c . Assign all address requests to their first copy. Take a module which is assigned with $\geq p^{(c-1)/c}$ requests. There exist $\leq p^{(1/c)}$ such modules. The first $p^{(c-1)/c}$ among these requests will be satisfied by the module. The assignment of the other, say x , requests will be switched to other modules. The other $c-1$ copies of the x addresses relate like $c-1$ copies of a subset of $\binom{m-1}{c-1}$ addresses. Namely, it is impossible that there are two out of these x requests such that all their $c-1$ copies are in the same modules. We have to decide which are the $p^{(c-1)/c}$ first

requests assigned to and satisfied by the module. We do it in a way which utilizes sorting. We actually do more. We schedule the requests to the module by numbering them from one on. The scheduling is done as follows.

Say that processor i ($1 \leq i \leq p$) has a read request for some address at module j_i ($0 \leq j_i \leq m-1$). We represent it by a pair (j_i, i) .

Step 1. Sort these pairs according to the following lexicographic order:

$$(j_{i_1}, i_1) < (j_{i_2}, i_2) \quad \text{if } [j_{i_1} < j_{i_2}] \text{ or } [j_{i_1} = j_{i_2} \text{ and } i_1 < i_2].$$

The sorting can be done in $O(\log p)$ time on an EREWPRAM using the algorithms of [2] or [14].

Let $A[1], \dots, A[p]$ be the output sorted vector. Step 1 groups all requests for the same module into the intervals in A . Namely, if there are any requests for some module j they appear as $A[\alpha_j], A[\alpha_j+1], \dots, A[\beta_j]$ for some $1 \leq \alpha_j \leq \beta_j \leq p$ ($0 \leq j \leq m-1$). Let $A[\gamma]$ be a request for module j . The serial number of $A[\gamma]$ among the requests for module j is $\gamma - \alpha + 1$. We show how to compute this serial number into $R(\gamma)$.

Step 2. Initialization.

for $1 \leq \gamma \leq p$ **pardo**

$R(\gamma) := 1$

(*Comment.* The following instruction breaks the vector A into linked lists corresponding to modules. The linked list for module j will start at $A[\beta_j]$ and end at $A[\alpha_j]$.)

NEXT(1) := 'undefined'

for $2 \leq \gamma \leq p$ **pardo**

if the module of requests $A[\gamma]$ and $A[\gamma-1]$ are the same

then NEXT(γ) := $\gamma-1$

else NEXT(γ) := 'undefined'

Apply $\lceil \log p \rceil$ iterations

for $2 \leq \gamma \leq p$ **pardo**

if NEXT(γ) is not 'undefined'

then $R(\gamma) := R(\gamma) + R(\text{NEXT}(\gamma)); \text{NEXT}(\gamma) := \text{NEXT}(\text{NEXT}(\gamma))$

else abort.

This step takes $O(\log p)$ time on an EREWPRAM and completes the computation of the serial numbers into the vector R .

The scheduling of requests to second copies is done separately in each set of requests, that were switched from the same module, and so is the scheduling for later copies which is done separately for even more refined sets. The definition of such a refined set (at each transition from i to $i-1$ copies and application of the inductive step) includes all address requests that passed through the same sequence of modules but have not yet been satisfied. The computation of serial numbers for such set meets exactly the scheme of computation of serial numbers described above. (Here, we do not need the assumptions of the remark for $c=2$ since switched requests are not sent in the first place.) Therefore, we can apply the induction. Let A_i be sizes of switched

sets then

$$\sum (c-1)A_i^{((c-2)/(c-1))} \leq (c-1)p^{1/c} p^{((c-1)/c)((c-2)/(c-1))} = (c-1)p^{((c-1)/c)}$$

The inequality holds because the left hand side is maximized when all the A_i numbers are equal. The c^2 size is due to the need to compute the subset of modules that contain a copy of an address by each processor using its local memory. The $\log p$ is due to computations of serial numbers which is required at each of the c steps of the induction. \square

Theorem 3. *We can achieve $R_{\max} \leq c \left(p^{c-1} \left[N / \binom{m}{c} \right] \right)^{1/c}$ in parallel time $O(c^2 + c \log p)$.*

Proof. Let x be the integer such that $(x-1) \binom{m}{c} < N \leq x \binom{m}{c}$. Apply the lemma for each of the x layers, separately. Let p_i be the number of requests in the i -th layer, $1 \leq i \leq x$. Then we get

$$\begin{aligned} R_{\max} &\leq \sum_1^x R_{\max}(i) \leq \sum_1^x c p_i^{(c-1)/c} \leq x c (p/x)^{(c-1)/c} \\ &= c (p^{(c-1)} x)^{1/c} \leq c \left(p^{c-1} \left[N / \binom{m}{c} \right] \right)^{1/c}. \end{aligned}$$

$R_{\max}(i)$ is the R_{\max} obtained for each layer. The first inequality is obvious. The second is implied by the lemma. The third is because its left hand side is maximized when all the p_i are equal. The scheme of computation of serial numbers as above should be used here as well. \square

III.6. Connection with the First Stage

The proposed c -partitioning above poses some difficulties in combining it with the probabilistic simulation of the previous chapter. There, a copy of address a_i could have been found in module $h(a_i) \pmod m$. While in this chapter the set of modules that contain copies of address $h(a_i)$ was selected differently. Each address is assigned to a set of modules containing its copies (according to the definition of the proposed c -partitioning at the beginning of Sect. III.2) where no function like the remainder mod m seems to be involved in determining any member of this set. (For instance, module $h(a_i) \pmod m$ may not have a copy of address a_i). This is the reason why for the purpose of combining this stage with the previous one we propose an alternative c -partitioning. Its small disadvantage is that it gives a little bit inferior results than the first c -partitioning. Its big advantage is that it suits to be a second stage following Chapter II. Note that in the sequel we always omit the hashfunction h and refer to the address a_i only. When the solution of this section is set to follow Chapter II address a_i should be replaced by $h(a_i)$.

The Alternative c-Partitioning

For $N \leq m!/(m-c)!$ and address i , $0 \leq i < N$, the first copy of i is in module $i \pmod m$, the second copy is in module $i \pmod{(m-1)}$ of the remaining $m-1$ modules and so on. Namely the j -th copy, $1 \leq j \leq c$, is in module $i \pmod{(m-j+1)}$ of the $m-j+1$ modules not occupied by the first $j-1$ copies.

Example. Let $m=10$ $c=3$ and $i=15$. The first copy of i is in module 5 ($15 \pmod{10}=5$), the second is in module 7 ($15 \pmod{9}=6$ and since module 5 is occupied, module 7 corresponds to 6) and the third is in module 9 ($15 \pmod{8}=7$ and module 5 and 7 are occupied).

If $(x-1)(m!/(m-c)!) < N \leq x(m!/(m-c)!)$ for some integer x then partition the addresses into x approximately equal subsets (layer) and fix the c -partitioning of each of the x layers as for $N \leq m!/(m-c)!$.

Remark. For an address i it takes $O(c \log c)$ time to compute the subset of c modules containing its copies. Clearly, $i'_1 = i \pmod{[N(m-c)!/m!]}$ is the serial number of i in its layer. Find the modules one at a time. Create a 2-3 tree for the modules that were chosen so far. (For more on 2-3 trees, see [1].) By keeping in each internal node of the tree information about the number of unoccupied modules among its leaf-descendants we can identify the module of the next copy and update the tree in time $O(\log c)$. The simple additional details are omitted. In the cases where $N = xc! \binom{m}{c} (=xm!/(m-c)!)$ for some integer x , Observation 2 and Theorem 1 still hold since $\#A_{i_1, i_2, \dots, i_k} = xc!$ for any subset $\{i_1, i_2, \dots, i_k\}$ of k modules in both c -partitionings (using the notations of the lower bounds section). In order to shorten the paper we reconstructed here only the upper bound analogues to Theorem 2. This is done in an informal way since it follows the same lines as the proof of Theorem 2.

For any p requests we want to find an upper bound on the best R_{\max} achievable. Let x be the integer such that $(x-1)m!/(m-c)! < N \leq xm!/(m-c)!$. For some $R_{\max} = t$ and a module assigned with t requests denote by s the minimum number of modules (including the first module) that have to be assigned with at least $t-1$ request in order not to enable us to decrease t in the first module. The following two inequalities hold:

- (1) $xs!/(s-c)! \geq (t-1)s + 1.$
- (2) $p \geq (t-1)s + 1.$

(1) implies $xs^c \geq (t-1)s$ and $s \geq (t/x)^{c-1}$. This and (2) imply

$$(t-1)((t-1)/x)^{1/c-1} \leq p, \quad (t-1)^c \leq xp^{c-1} \quad \text{and} \quad t \leq 1 + (xp^{c-1})^{1/c}$$

Theorem 2'. *For the alternative c-partitioning*

$$R_{\max} \leq 1 + (p^{c-1} [N(m-c)!/m!])^{1/c}.$$

The analogue to Theorem 3 will be

Theorem 3'. *For the alternative c -partitioning we can achieve*

$$R_{\max} \leq c(p^{c-1} c! [N(m-c)!/m!])^{1/c}$$

in parallel time $O(c \log c + c \log p)$.

Proof. The only new idea in this proof is the following. We show that all addresses that fall in the same layer of our alternative c -partitioning can be efficiently further partitioned into $c!$ ‘sublayers’ as follows. Take a layer. It has at most $m!/(m-c)!$ addresses. Given any c modules there are at most $c!$ addresses (in the layer) such that all their copies are in these c modules. Recall that the definition of the alternative c -partitioning specifies for each address the module in which copy α , $1 \leq \alpha \leq c$, of this address is stored. Thus, we have to specify for each address i of layer l not only the c modules $0 \leq \beta_1 < \beta_2 < \dots < \beta_c \leq m-1$ in which this address has copies. But also a permutation on this list of modules that indicates in which module is the first, second, ..., c -th copy. It is simple to see that we can instead use a permutation on $[1, \dots, c]$ (together with the list of modules β_1, \dots, β_c). This permutation is the key for partitioning each layer into $c!$ sublayers. The discussion below fills in the details for computing these sublayers from the alternative c -partitioning. Let (i_1, i_2, \dots, i_c) be a c -tuple of modules that contain the respective first, second, ..., c -th copy of address i that belongs to layer l for some i and l . All addresses of layer l Denote by $N(j)$ the cardinality of $\{i_k | k < j \text{ and } i_k < i_j\}$. Note that these cardinalities can be observed upon searching the 2-3 tree mentioned above for i_j without changing the $O(\log c)$ time estimate. Obviously $0 \leq N(j) < j$.

Define $L(i_1, i_2, \dots, i_c) = \sum_1^c (j-1)! N(j)$ to be the sublayer of address i in layer l .

The only address in this sublayer which is contained in modules i_1, i_2, \dots, i_c is i . Therefore, the N addresses form altogether $c! [N(m-c)!/m!]$ sublayers. Applying the lemma similar to the proof of Theorem 3 to each of this sublayers completes the proof of the theorem. The $O(c \log p)$ size represents the computation of serial numbers as in the proof of the lemma.

III.7. Some Nonconstructive Upper Bounds

In this section we demonstrate a way for utilizing a number of copies c where $\binom{m}{c} \gg N$. A typical counting argument that provides for non-constructive c -partitionings having the desirable property of enabling small R_{\max} is presented. We leave the problem of constructing such efficient c -partitionings open.

We first need the following theorem:

Theorem 4. *Let p and m be as before. Let d be a positive integer, where $p \leq dm$. Assume that a c -partitioning is given such that for every subset of size p of the N addresses, the set of modules that contain all c copies of these p addresses is of size $\geq \lceil p/d \rceil$. Then it is possible to get $R_{\max} \leq d$ for every p memory requests.*

Proof. Assume, in contradiction, that there exists a set of p ($\leq dm$) requests that causes $R_{\max} > d$. By a similar technique to the proof of correctness of the alternative sequential algorithm, we choose an assignment such that some module is assigned with $\geq d + 1$ requests. They must have copies in one more module. This second module must be assigned with $\geq d$ request. So far we have $\geq 2d + 1$ requests. They must have copies in $\geq 3 = \left\lceil \frac{2d+1}{d} \right\rceil$ modules. The third module is also assigned with $\geq d$ requests, and so on. So we get that this set of requests was of size $\geq dm + 1$. A contradiction.

This is, actually an extension of Hall's Theorem (cf. [6]) which is given for the case $d=1$. An upper bound on the portion of c -partitionings that do not have the favorable property described in the theorem is

$$I_p = \sum_{1 \leq k \leq p} \binom{N}{k} \left(\left\lfloor \frac{k}{d} \right\rfloor^m - 1 \right) \frac{c^{\frac{N}{m}} \left(\left\lfloor \frac{k}{d} \right\rfloor - 1 \right)!}{\left(c^{\frac{N}{m}} \left(\left\lfloor \frac{k}{d} \right\rfloor - 1 \right) - ck \right)!} \frac{(cN - ck)!}{(cN)!}$$

assuming that N is divisible by m . (Note, that the definition of a c -partitioning is extended to the case where more than one copy of an address is contained in the same module.) This upper bound reminds to some extent the proof of the existence lemma of "expander bipartite graphs" given in p. 298 of [12].

If we prove for some values of N, m, p, d and c that I_p (the portion of "bad" c -partitionings) is smaller than one then we gave a (non-constructive) existence proof of a "good" c -partitioning.

Appendix I exemplifies a result that can be derived from this upper bound. We get that if $N = p^3, m = p^2, d = 1$ and p is "large enough" than there exists a "good" 10-partitioning.

In a similar way to the considerations above we may obtain an upper bound of $1 - b$ on the portion of bad c -partitionings, for some positive constant b . Since we do not have a way to construct a "good" c -partitioning, we can choose one at random and use it till the first time it fails. Namely, if we find for the first choice a set of requests for which it is impossible to get $R_{\max} = d$, then we choose another c -partitioning and so on. Obviously, with high probability, we find after a small number of trials a "good" c -partitioning.

IV. Efficient Low-Level Simulations

The purpose of this paper is to present ways for simulations of shared memory models of parallel computation which allow fairly unrestricted access patterns of processors to shared memory cells by machines in which the memory is organized in modules where only one cell of each module can be accessed at a time. As was mentioned in the introduction the paper envisions a three stage analysis and solution to the problem. The combination of the two earlier stages brings us to a point where each processor of the MPC simulating machine specifies in each cycle being simulated both an address request and

the module that was chosen to satisfy this request. The problem is how to complete the simulation. The choice of the EREWPRAM and the MPC for presentation of our ideas for the first two stages is due to the fact that the simulation of the former by the latter distinguished both the problems and solutions. However, this choice made the need for a third stage indistinct. This is the reason that here we switch to other models of computation. Instead of the EREWPRAM we take a more permissive model of computation, the concurrent-read concurrent-write (CRCW)PRAM. In this model several processors are allowed to read simultaneously from the same memory location. If several processors try to write simultaneously into the same memory location the lowest numbered processors succeeds. This model is based on [7] and [16]. We substitute the MPC machine by a weaker machine named non-queued (NQ)-MPC. Here only one address request may arrive at each module at each time unit. Besides that the NQ-MPC is similar to the MPC.

We wish to simulate one cycle of the CRCWPRAM. Assume that each processor of the NQ-MPC is assigned already both with an address request from the common memory and the memory module which has to satisfy this address request. Assume also that $m \geq p$, namely the number of modules is at least as big as the number of processors. Our proposed solution resembles the simulation of a CRCWPRAM by a EREWPRAM in [19]. Unfortunately, our present problem requires not only the circumvention of access conflicts to the same memory location but also to the same module. Our simulation resembles also the proof of the Lemma in Sect. III.5. Whenever the considerations are similar to this proof we shorten the presentation. The first step of the simulation is:

(1) Sort in parallel the p triples specified below in the lexicographic order.

Each processor enters the triple: (the serial number of the module assigned to its address request, the serial number of the address request itself, its own serial number).

The NQ-MPC sorts them using [2] sorting algorithm in time $O(\log p)$. It does not need more than size p shared memory. This is achieved by using one cell of each module. The result is that all requests for the same module (resp. the same address of the same module) appear in successive locations of the sorted vector. Let us call the set of such successive locations interval (resp. subinterval) denoted $I(M)$ (resp. $SI(M, a)$). M and a represent indeterminants corresponding to modules and addresses, respectively. Note that subintervals are sorted by the serial number of the processors.

(2) The serial number of each subinterval relative to the other subintervals in the sorted vector is computed. Denote it by $\#SI(M, a)$.

A processor is allocated to each triple. By comparison with the triple in the preceding place of the sorted vector the processor finds out whether its triple is the smallest in its subinterval. If yes it is chosen to 'represent' the subinterval as will be seen later; let us call the triple in this case, a subinterval-triple. Now, we apply the following computation. It is similar to Step 2 in the proof of the Lemma (Sect. III.5).

Initialization.

```

for each triple  $x$ ,  $1 \leq x \leq p$ , pardo
  if  $x$  is a subinterval-triple
  then  $R(x) := 1$ 
  else  $R(x) := 0$ 
NEXT(1) := 'Undefined'
for each triple  $x$ ,  $2 \leq x \leq p$ , pardo
  NEXT( $x$ ) :=  $x - 1$ 
Apply  $\lceil \log p \rceil$  iterations
for  $2 \leq x \leq p$  pardo
  if NEXT( $x$ ) is not 'undefined'
  then  $R(x) := R(x) + R(\text{NEXT}(x))$ ; NEXT( $x$ ) := NEXT(NEXT( $x$ ))
  else abort.

```

This computation takes $O(\log p)$ time. It results with the required subinterval serial number associated with the subinterval triple.

(3) The serial number of each subinterval which is smallest in its interval ($\min_b (\#SI(M, b))$) is 'broadcasted' to all other triples of the interval.

The number $\#SI(M, a) - \min_b (\#SI(M, b)) + 1$ is the serial number of the subinterval $SI(M, a)$ relative to other subintervals of interval $I(M)$. The broadcasting is done in $\lceil \log p \rceil$ pulses. In pulse i , $1 \leq i \leq \lceil \log p \rceil$, each processor that knows already, the serial number of the smallest subinterval of the interval of its triple writes it into the 2^{i-1} successor of the triple in the sorted vector if it belongs to the same interval. It is simple to see (see [19]) that all triple get the appropriate message and each module is accessed by one processor at a time.

(4) The processor of each subinterval-triple performs the requested access to the right module in time corresponding to the serial number computed in Step (3) above.

In case, we simulate a writing cycle (recall the classification of the previous chapter into reading and writing cycles) we are done. In case a reading cycle is simulated we finish by:

(5) The content read by the processor of each subinterval-triple is broadcasted to all other triple of this subinterval.

The broadcasting is done by a similar technique to Step(3). Broadcasting for the same purpose is used in [19]. It is appropriate to mention at this point that the $O(\log p)$ time sorting algorithms of [2] has a large constant in front of the $\log p$. It is possible to use instead the $O(\log p)$ time (with very high probability) algorithm of [14]. Another alternative is Batcher's $O(\log^2 p)$ time algorithm which is probably best for small values of p .

Simulations by a Synchronous Distributed Machine

Our next goal is to replace the NQ-MPC by the Parallel-Design Distributed-Implementation (PDDI) general-purpose computer of [18]. Every algorithm for the CRCW PRAM can be run on the PDDI. (Even algorithms for a model of computation which is stronger than the CRCW PRAM can be used but we

prefer not to discuss it here.) Each such algorithm is automatically simulated by the network which is described briefly below. The network has three kinds of processors. S_1, \dots, S_s are called *super-processors*. Each simulates the behaviour of (*several*) processors of the PRAM. M_1, \dots, M_m are called *memory-processors*. Each simulates a memory module. An interconnection network connects the super-processors with the memory-processors. It consists of a sorting network followed by a merging network. The ‘switches’ in this network are called *comparator-processors* since they correspond to comparator modules in the sorting and merging network. There are $f(s, m)$ comparator-processors, where the function f depends on the specific selection of sorting and merging network. Let $l(s, m)$ denote the longest directed path in the sorting and merging networks that were selected ($l(s, m)$ corresponds to the worst case execution time of a network). A typical time cycle of the CRCW PRAM includes instructions for reading from or writing into its shared memory. The PDDI simulates it by routing messages between the super-processors and the memory-processors. The main result concerning the PDDI is the following.

Theorem [18]. *Suppose we are given an algorithm for the CRCW PRAM whose running time is $O(t/p)$ for any number of $1 \leq p \leq x$ processors and N common memory locations, where t, x and N are some numbers. The interconnection network of the PDDI machine can simulate this algorithm using s super-processors, $m = N$ memory-processors and $f(s, m)$ comparator-processors in time $O(t/s)$ for any $s \leq x/l(s, m)$.*

Remark. For some selection of sorting and merging networks for the PDDI we get $f(s, m) = O(s \log s + m \log m)$ and $l(s, m) = O(\log s + \log m)$.

This theorem assumes $N = m$. We also refer the reader to a discussion in [18] which implies how to extend the proof of the theorem for cases where $m < N$, by queuing requests for different cells of the same module. This discussion implies essentially that the time $O(t/s)$ in the conclusion of the theorem should be replaced by $O(tR/s)$ where $R (= R_{\max})$ is the average maximum memory contention over all memory modules.

(Remark. We only explain briefly why such a result is possible. The communication between super-processors and memory-processors makes a heavy use of pipelining in the interconnection network. All computations required for the queuing are performed in the interconnection network itself as “side effects” of the communication between super-processors and memory-processors. It is done in a way which “absorbs” the time for these computations into the “big oh” of t/s in the theorem. We do not elaborate any more on how it is done since this would require a much longer description of the PDDI.)

Suppose now that we want to use the probabilistic solutions of the present paper in order to keep the memory contention R_{\max} low. Then we get the following.

Probabilistic Simulations Revisited for the Second Time

Corollary. *Suppose we are given an algorithm for the CRCW PRAM whose running time is $O(t/p)$ for any number of $1 \leq p \leq x$ processors and N common*

memroy locations, where t , x and N are some numbers. The interconnection network of the PDDI machine can simulate this algorithm using s super-processors, m memory-processors and $f(s, m)$ comparator-processors in time $O(tR/s)$ for any $s \leq x/l(s, m)$, where the following applies.

- (1) If $m = s^3$ then $R = 1$ and the simulation needs total memory size of $N + 2s$.
- (2) If $m = s^2$ then $R = \log N$ and the simulation needs total memory size of $N + 4s$.
- (3a) If $m = s$ then $R = \log s$ and the simulation needs total memory size of $(N + s) \log s$.
- (3b) If $\varepsilon > 0$ and $m = s^{1+\varepsilon}$ then $R = 1$ and the simulation needs total memory size of $(1 + 2/\varepsilon)(N + s)$.
- (4a) There is a function $f: \mathbb{N} \rightarrow \mathbb{R}$ with $f(s) = O(s^\varepsilon)$ for all $\varepsilon > 0$ such that: if $N = 2^n$ (n is the size of the input) and $m = 2^b \geq sf(s)$ the time of the simulation is $O((\log N)^3 + t(\log s/(\log \log s)^2 + \log N)/s)$ and the simulation needs total memory size of $N + s \log N$.
- (4b) Let $\varepsilon > 0$ be fixed. If $m = s^{1+\varepsilon}$ then the time of the simulation reduces to $O((\log N)^3 + t \log N/s)$.

[17] proved the following result. For any p , there exist synchronous distributed machines in which p processors can communicate with a bounded number of others such that: Given an algorithms for a PRAM which runs in $O(x)$ using p processors, it can be simulated in time $O(x \log^2 p)$ (almost surely) by the machine. Let us compare his result with ours. The literature contains many algorithms which make an efficient use of a large number of processors (for large enough input). Therefore, it is likely that the number of processors in a real machine will be far less than the number of processors employed by a PRAM algorithm. Let us demonstrate an on example how the efficiencies of the simulation results for the PDDI take advantage of that. Consider the problem of finding the k -th largest out of n elements. [20] gave an algorithm of time $O(n/p)$ using $p \leq n/(\log n \log \log n)$ processor and common memory of size $O(n)$. We wish to simulate this algorithm by a synchronous distributed machine. Observe that real processors should be pretty powerful in order to simulate processors of a PRAM including their full instruction set. Say that a machine can have \sqrt{n} such processors. So the PDDI will have \sqrt{n} super-processors. Say also that a machine can have m processors that simulate memory modules. Such processors need to have totally different properties than the super-processors. Indeed the PDDI does not link the number of super-processors to the number of memory-processors. The Corollary implies that while m ranges between \sqrt{n} and $n\sqrt{n}$ the running time of a PDDI ranges between $O(\sqrt{n})$ and $O(\sqrt{n} \log n)$ and the additional memory between $O(n)$ and $O(n \log n)$. We have not mentioned yet the $O(m \log n)$ comparator-processors. Observe, that this processors need not be as complex as the other processors. Actually they can be fairly degenerate: their instruction set and local memories are small. Therefore, the hardware cost of the comparator-processors is dominated by the hardware cost of the other processors.

On the other hand Upfal's solution implies a running time of $O(\sqrt{n} \log^2 n)$ of the simulating machine.

V. Summary

We gave above a detailed description of each component of our solution. Here we would like to give a high level description of a simulation of the CRCWPRAM by the NQ-MPC.

The CRCWPRAM (resp. NQ-MPC) is permissive (resp. restrictive) relative to the spectrum of other permissive (resp. restrictive) models of computation that appear in the literature. This enables us to go again through the main notions of our solution in order to summarize them in a uniform fashion.

Stage one assumes a class H of hashfunction. Pick at random one of them, say h . This hashfunction implies the location of the c copies of each address a_i . By the alternative c -partitioning of stage two the j -th copy is in module $h(a_i) \pmod{m - (j - 1)}$ of the modules that were not occupied by the first $j - 1$ copies ($1 \leq j \leq c$). The step by step simulation starts with each processor of the NQ-MPC machine specifying an address request for read or write like its corresponding CRCWPRAM processor. Now we have to split into reading and writing cycles. For reading cycles Theorem 3' gives an algorithm for allocation of address requests to modules. A scheduling of the requests to time units of the simulated cycle and a way to transmit the read request to the modules and the response back to the processors are described in the previous chapter. For writing cycles each address request is assigned to all its c copies. We do not do worse than c times (the time for one reading cycle). The reader is invited to conclude this by observing that we can access: first copies, then second copies, and so on.

In this paper we were able to apply a few powerful concepts developed by Theoretical Computer Science to a problem which arises in a natural way from practice. One concept is the use of randomization in the sense described in [13]. The second concept is the use of augmenting paths in the analysis and algorithms of Chapter III. The third concept is expander graphs in Sect. III.7.

Among the technical contributions of the paper we think that two are of particular interest: (1) The way universal hashing is used and analyzed. This is due to a different efficiency criterion than in [5] that has to be satisfied. (2) The fast parallel approximation algorithms. This is interesting since we used approximation for problems that were solved serially using augmenting paths. Augmenting path seems to be a procedure which is inherently serial. There are other problems whose algorithms use variations of augmenting paths. These include finding maximum matching in a graph and maximum flow in a network. So far, there are no fact parallel algorithms using a small number of processors for these problems. (By this we mean that the running time of the algorithm is poly-log in the length of the input for the problem; and, a simulation of such an algorithm by a single processor results in a serial algorithm which is as good, or "almost" as good, as the fastest serial algorithm for the problem.) It is possible that such algorithms do not exist. If this is true

it might still be possible to find fast parallel approximation algorithms, as we did here.

The formulation of the granularity problem together with the problem solved by the PDDI is an interesting exercise in a methodological solution to a “real problem”. That is the problem of simulating PRAMs by synchronous distributed machines. As was indicated in the introduction it was methodological to first identify these two problems and later attempt to solve them rather than to solve this simulation problem directly.

Appendix 1

The following example may illustrate what are the results that one may expect out of Sect. III.7. No special effort was made to get the best possible result in the subsequent computation.

Example. Let $N = p^3$, $m = p^2$ and $d = 1$. Then,

$$I_p \leq \sum_{1 \leq k \leq p} \binom{p^3}{k} \binom{p^2}{k} \frac{(cpk)!}{(cpk - ck)!} \frac{(cp^3 - ck)!}{(cp^3)!}$$

$$= \sum_{1 \leq k \leq p} \binom{p^3}{k} \binom{p^2}{k} \frac{(cpk)}{(ck)} \bigg/ \binom{cp^3}{ck}.$$

Since

$$\binom{cp^3}{ck} \geq \binom{p^3}{k} \binom{p^2}{k} \binom{(c-1)p^3 - p^2}{(c-2)k}$$

we get

$$\leq \sum_{1 \leq k \leq p} \frac{(cpk)}{(ck)} \bigg/ \binom{(c-1)p^3 - p^2}{(c-2)k}.$$

Since

$$\frac{(cpk)}{(ck)} \leq \frac{(cpk)^{ck}}{(ck)!}$$

and

$$\binom{(c-1)p^3 - p^2}{(c-2)k} \geq \frac{((c-1)p^3 - p^2 - (c-2)k)^{(c-2)k}}{((c-2)k)!}$$

we get

$$\leq \sum_{1 \leq k \leq p} \frac{(ck)! (cpk)^{ck}}{((c-2)k)! ((c-1)p^3 - p^2 - (c-2)k)^{(c-2)k}}.$$

Now $\frac{(ck)!}{((c-2)k)!} \leq (ck)^{2k}$; and for p “large enough” $p^3 > p^2 + (c-2)k$ where c is a constant and $k \leq p$. So

$$\leq \sum_{1 \leq k \leq p} \frac{(ck)^{2k} (cpk)^{ck}}{((c-2)p^3)^{(c-2)k}}.$$

Call each element L_k . Let $c=10$. For large enough p

$$L_1 = \frac{10^2(10p)^{10}}{(8p^3)^8} < 1/p.$$

Since it is easy to verify that the derivate of L_k with respect to k ($1 \leq k \leq p$) is negative for large enough values of p , we get $\sum_{1 \leq k \leq p} L_k < 1$. So, we proved the existence of a 10-partitioning such that for these parameters gives always a memory contention of one.

Acknowledgement. We are grateful to the referees for their thorough remarks. Discussions with A. Gottlieb, M. Snir, P. Spirakis and A. Wigderson are gratefully acknowledged. We thank Kevin McAuliffe for the information about the ratio between read and write instructions achieved in simulations of the Ultracomputer Project.

References

1. Aho, A.V., Hopcroft, J.E., Ullman, J.D.: The Design and Analysis of Computer Algorithms. Reading, MA: Addison-Wesley 1974
2. Ajtai, M., Komlos, J., Szemerédi, E.: An $O(n \log n)$ sorting network. Proc. Fifteenth ACM Symposium on Theory of Computing, pp. 1-9, 1983
3. Awerbuch, B., Israeli, A., Shiloach, Y.: Efficient simulations of PRAM by Ultracomputer. (Preprint). Dept. of Computer Science, Technion, Haifa, Israel, 1983
4. Carmichael, R.D.: Groups of finite orders. Dover: Dover Publications 1956
5. Carter, J.L., Wegman, M.N.: Universal classes of hash functions. Proc. Ninth ACM Symposium on Theory of Computing, pp. 106-112, 1977
6. Even, S.: Graph Algorithms. Potomac, MD: Computer Science Press 1979
7. Goldschlager, L.M.: A Unified Approach to Models of Synchronous Parallel Machines. Proc. Tenth ACM Symposium on Theory of Computing, pp. 89-94, 1978
8. Gonnet, G.H.: Expected length of the longest probe sequence in hash code searching. JACM **28**, 289-304 (1981)
9. Gottlieb, A., Grishman, R., Kruskal, C.P., McAuliffe, K.P., Rudolph, L., Snir, M.: The NYU Ultracomputer-Designing, a MIMD Shared Memory Parallel Machine. IEEE Trans. Comput. **c-32**, 175-189 (1983)
10. Kuck, D.J.: A survey of parallel machine organization and programming. Comput. Surveys **9**, 29-59 (1977)
11. Lev, G., Pippenger, N., Valiant, J.G.: A fast parallel algorithm for routing in permuting networks. IEEE Trans. Comput. **c-30**, 93-100 (1981)
12. Pippenger, N.: Superconcentrators. SIAM J. Comput. **6**, 298-304 (1977)
13. Rabin, M.O.: Probabilistic algorithms. In: Algorithms and Complexity, J.F. Traub (ed.). New York: Academic Press 1976
14. Reif, J., Valiant, L.J.: A logarithmic time sort for linear size networks. Proc. Fifteenth ACM Symp. Theory Comput. pp. 10-16, 1983
15. Schwartz, J.T.: Ultracomputers. ACM Trans. Progr. Lang. Syst. **2**, 484-521 (1980)
16. Shiloach, Y., Vishkin, U.: Finding the maximum, merging and sorting in a parallel computation model. J. Algorithms **2**, 88-102 (1981)
17. Upfal, E.: A probabilistic relation between desirable and feasible models of parallel computation. Proc. Sixteenth ACM Symp. Theory Comput. 1984 (To appear)
18. Vishkin, U.: Parallel-Design space Distributed - Implementation space (PDDI) general pur-

- pose computer. RC 9541, IBM T.J. Watson Research Center, Yorktown Heights, NY 10598, 1982. To appear in *Theoretical Computer Science*)
19. Vishkin, U.: Implementation of simultaneous memory address access in models that forbid it. *J. Algorithms* **4**, 45–50 (1983)
 20. Vishkin, U.: An optimal parallel algorithm for selection. (Preprint, 1983)
 21. Vishkin, U., Wigderson, A.: Dynamic parallel memories. *Information and Control* **56**, 174–182 (1983)

Received May 28, 1984 / July 2, 1984