

Simulation-based Comparison of Hash Functions for Emulated Shared Memory*

Curd Engelmann¹ and Jörg Keller²

¹ Universität des Saarlandes, Computer Science Department
Im Stadtwald, 6600 Saarbrücken, Germany

² Centrum voor Wiskunde en Informatica
Postbus 4079, 1009 AB Amsterdam, The Netherlands

Abstract. The influence of several hash functions on the distribution of a shared address space onto p distributed memory modules is compared by simulations. Both synthetic workloads and address traces of applications are investigated. It turns out that on all workloads linear hash functions, although proven to be asymptotically worse, perform better than theoretically optimal polynomials of degree $O(\log p)$. The latter are also worse than hash functions that use boolean matrices. The performance measurements are done by an expected worst case analysis. Thus linear hash functions provide an efficient and easy to implement way to emulate shared memory.

1 Introduction

Users of parallel machines more and more tend to program with the view of a global shared memory. Commercial machines (with more than 16 processors) however usually have distributed memory modules. Therefore the address space has to be mapped onto memory modules, memory access is simulated by packet routing on a network connecting processors and memory modules. This has to be done in a way that for (almost) all access patterns the requests are distributed almost evenly among the memory modules. The reason to demand this is obvious: if cases happen where the number of requests per module (the so called *module congestion*) is too high, then performance gets very poor.

Several kinds of hash functions have been proposed. But their theoretically provable properties are asymptotical results. As currently available machines are quite small (the number p of processors and memory modules usually is less than 1000) the actual behaviour of the chosen hash function can differ quite a lot from these theoretical properties. The lack of experimental data makes the selection of a particular hashing scheme difficult in practice. We are not aware of comparisons of hash functions based on simulated behaviour.

* This work was supported by the German Science Foundation (DFG) in SFB 124, TP D4, and by the Dutch Science Foundation (NWO) through NFI Project ALADDIN under Contract number NF 62-376. Part of this work was done while the second author was working at Universität des Saarlandes, Computer Science Department, Saarbrücken, Germany.

The goal of this investigation is to provide these data by comparing four kinds of hash functions by simulations. In Sect. 2 the most common kinds of hash functions are introduced. Section 3 describes the types of synthetic and real access patterns that were used as workloads. Section 4 sketches the experiments made and Sect. 5 presents and discusses the results.

2 Hash Functions

As already mentioned, a hash function serves to map a global address space onto distributed memory modules. More formally, for an address space M of size $m = 2^v$ and a set N of $p = 2^u$ memory modules, the mapping is a function $h : M \rightarrow M$ that maps addresses to memory cells. The function $mod : M \rightarrow N, mod(x) = x \text{ div } m/p$ specifies the module of a memory cell x , the function $loc : M \rightarrow M', loc(x) = x \text{ mod } m/p$ specifies the local address of cell x .

An optimal mapping function h should guarantee low module congestion for almost all possible access patterns (if all addresses of one pattern are distinct). This is achieved by using classes of functions in which each function has low module congestion for almost all patterns. A particular function is randomly chosen. This guarantees with very high probability that the current application does not exhibit the patterns on which the chosen function produces hot spots.

An additional problem consists in patterns with several processors concurrently accessing one cell. This problem cannot be solved by hashing. However there exist routing algorithms that perform *combining*. Requests that access the same cell are merged during routing, answers are duplicated. Ranade's emulation algorithm [10] is a good example. Therefore, concurrent access does not increase module congestion.

A class that restricts module congestion to $O(\log p)$ is

$$\mathcal{H} = \left\{ p(x) = \left(\sum_{i=0}^{\xi} a_i \cdot x^i \right) \text{ mod } P \text{ mod } m : 0 \leq a_i < P \right\} .$$

P is a prime larger than m , $\xi = O(\log p)$. A function of \mathcal{H} is obtained by randomly choosing the values for a_i . This class was used in several theoretical investigations [6, 8, 10] to emulate shared memory on a processor network. The module congestion of $O(\log p)$ is sufficient because access from processors to memory modules across a constant-degree interconnection network needs time $\Omega(\log p)$ anyway.

However the functions in \mathcal{H} are not bijective. This means that several addresses of the shared memory could be mapped onto the same cell. This requires secondary hashing on each memory module. Ranade [10] describes a method that performs secondary hashing in constant time and increases the size of the memory module only by a constant factor.

In practice however one should avoid secondary hashing and waste of memory because a constant factor of performance loss can destroy an asymptotically good result. Furthermore, the time to evaluate the hash function should be short. The

functions in \mathcal{H} require $\xi = O(\log p)$ multiplications and additions and a modulo division by a prime which needs a lengthy computation.

Therefore some alternatives were proposed:

1. For $\xi = 1$ one obtains a linear function. This reduces evaluation time to one multiplication, one addition and one modulo division. The function is still not bijective.
2. Furthermore if the modulo division by a prime is skipped and the coefficient a_0 is set to zero, the evaluation time is reduced to one multiplication. The operation modulo m is not counted because m is a power of two. If only odd values are chosen for a_1 the function also is bijective.
3. If the binary representation of an address is seen as a boolean vector, the hash function consists of multiplying this vector with an invertible boolean matrix. The time to evaluate this function is shorter than one multiplication.

Dietzfelbinger et. al. prove that the first alternative is asymptotically equivalent to the second [5]. Furthermore he proves that linear functions can result in a module congestion of $\Theta(\sqrt{p})$ for patterns with addresses of the form $b + s \cdot i$ where $i = 0, \dots, n - 1$ [4]. The constants b and s are called *base* and *stride*. This means that linear functions modulo a power of two are asymptotically worse than polynomials.

The third alternative was used in the design of the IBM RP3. Norton and Melton [9] introduce a class of boolean matrices where all matrices are invertible (which means bijectivity). Optimal distribution can be guaranteed for patterns with strides where s is a power of two and where in the binary representation of base b bits s to $s + \log n - 1$ are zero. For other bases the module congestion is at most 2. No theoretical results are given for other patterns, but their simulations hint that distribution is acceptable for other patterns, too. One particular matrix is obtained by randomly choosing several bits of the matrix and then computing all the other bits with respect to the above properties.

3 Workloads

The workloads are chosen to compare the hash functions with respect to known differences, especially behaviour on access patterns with strides, and with respect to patterns taken from applications. Therefore both synthetically generated patterns and application traces were taken.

The synthetic traces consist of randomly chosen patterns as a reference and strides with $s = 1, 13, 32$. The strides were chosen to compare matrix hashing and the other hash functions and to check whether linear functions get worse on these patterns. For $s = 32$ and $s = 1$ matrix hashing is optimal [9]. Theoretical results about the performance on the others are not known.

The traces were taken from three application programs: list ranking, matrix multiplication and connected components. The reasons for taking traces from applications are the variety of produced patterns and the structure of single patterns that often is more complex and less regular than in synthetical traces.

The three applications are chosen to represent a large variety of algorithms. Matrix multiplication is an example of a class of algorithms where the access patterns are regular and do not depend on the particular input values. Many other numerical algorithms behave that way, especially as many of them are originally designed to work on a processor network with a fixed interconnection structure (see e.g. [3]).

List ranking represents combinatorial algorithms where access patterns depend on the actual data. An example technique is pointer doubling. Processor i loads or stores $F[F[i]]$, where F is an array in the shared memory. Part of the accesses to shared memory still are regular. If processor i loads or stores $F[i]$, the access pattern is a stride with $s = 1$. Many PRAM algorithms working on lists and graphs are of this type (see e.g. [7]).

The connected components algorithm represents algorithms where access patterns depend on the actual data, but not all processors may participate in the access. This together with concurrent accesses to some cells, which get combined, makes module congestion smaller. Thus, connected components and similar algorithms are remarkable exceptions compared to list ranking type algorithms.

The list ranking algorithm is taken from a survey [7]. For a given linked list of n elements, the distance (or *rank*) to the end of the list is computed for each element. The algorithm needs n processors and $O(\log n)$ time. The list is represented as an array F , where $F[i]$ means successor of i in the list. For the last element of the list, $F[i] = i$. The rank is contained in array R . The PARDO code is shown in Fig. 1(a). The access patterns of this algorithm partly depend on the structure of the list and partly are strides with $s = 1$.

In the matrix multiplication algorithm $C = A \cdot B$, each processor computes one element of the destination matrix C . In order to avoid concurrent accesses, all processors start at different rows and columns of the matrices A and B . The PARDO code is shown in Fig. 1(b). Matrices A and C consist of $n = w2^{2z}$ elements and have dimension $2^z \times w2^z$, matrix B has dimension $w2^z \times w2^z$. The algorithm needs n processors and takes time $O(n^{1/2})$. The access patterns of this algorithm only depend on the dimensions of the matrices.

The connected components algorithm was adapted from Shiloach and Vishkin [11]. For a given undirected graph $G = (V, E)$, the connected components are computed. The algorithm needs $n = \max(|V|, 2|E|)$ processors and takes time $O(\log n)$. The graph is represented by two arrays HEAD and TAIL. For a given edge e , HEAD[e] and TAIL[e] contain the nodes to which e is adjacent. The components are represented by an array F . Two nodes u, v are in the same component if and only if $F[u] = F[v]$ after running the program. The PARDO code is shown in Fig. 1(c). The access patterns partly depend on the structure of the input graph and partly are strides with $s = 1$. Not all processors participate in every access.

```

(* Init rank R *)
for i := 1 to n pardo
  if F[i] = i then R[i] := 0 else R[i] := 1
od ;
(* Compute rank R *)
for t := 1 to ⌈log n⌉ do
  for i := 1 to n pardo
    R[i] := R[i] + R[F[i]] ;
    F[i] := F[F[i]] (* Pointer doubling *)
  od
od ;

(a) list ranking

(* n = w22z *)
k := 2z; m := w2z; l := w2z;
for (i, j) := (1, 1) to (k, m) pardo
  C[i, j] := 0 (* Init C *)
od ;
for r := 1 to l do
  for (i, j) := (1, 1) to (k, m) pardo
    t := (i + j + r) mod l ;
    C[i, j] := C[i, j] + A[i, t] · B[t, j]
  od
od ;

(b) matrix multiplication

for u ∈ V pardo F[u] := u od;
for t := 1 to 2 log |V| do
  for u ∈ V pardo change[u] := 0 od;
  starcheck;
  for all (u, w) with {u, w} ∈ E pardo
    if star[u] and F[w] < F[u] then
      F[F[u]] := F[w];
      change[F[u]] := 1;
      change[F[w]] := 1
    fi
  od;
  starcheck;
  for all (u, w) with {u, w} ∈ E pardo
    if star[u] and not change[F[u]]
      and F[w] ≠ F[u] then
      F[F[u]] := F[w]
    fi;
    F[u] := F[F[u]]
  od
od.

proc starcheck ;
begin
  for i ∈ V pardo
    star[i] := true;
    if F[F[i]] ≠ F[i] then
      star[F[F[i]]] := false
    fi;
    star[i] := star[F[F[i]]]
  od
end;

(c) connected components

```

Fig. 1. Code of applications

4 Experiments

To obtain the input data for the experiments, all applications are simulated by sequential programs, only the address traces are extracted. This frees us from considering a particular microprocessor instruction set and compiler. The address traces of the synthetic workloads are generated by a program, that simulates 4 steps of the machine. In the workloads with strides, the base b is increased each step by ns .

We are only interested in the resulting module congestion and not in the time to route the requesting packets from processors to memory modules. Therefore we can neglect the structure of the interconnection network. We only model it by a latency term because the processors perform latency hiding (see below).

All experiments are carried out for $m = 2^{22}$, the prime P is chosen closest to m . We simulate machines with $p = 2^u$, $u = 5, \dots, 10$ processors. We run multiple processes per processor to hide the network latency from processors. The processes are executed in a round-robin manner, one instruction per turn. The exact number c of necessary processes per processor is depending on p , e.g. $O(\log p)$ in a butterfly network. We choose a fixed c to obtain comparable results and take $c = 5$ as an average from a machine size of $p = 128$ [2]. Therefore in each step $5p$ requests are made. Step in this context means synchronous execution of one instruction on each of the $5p$ processes.

As polynomials we used functions of degree $\xi = 2, 10, 20$. Each of the experiments was done 5 times with randomly chosen hash functions. More exactly, for each class five functions were randomly chosen and then used for all workloads and machine sizes.

As input for list ranking a list of length $n = 10p$ was randomly chosen. As input for connected components, a graph with $n = 10p$ nodes and $5p$ edges was randomly chosen. The problem size n is twice as large as the number of processes in these applications. Each process simulates two program processors step by step. A problem size larger than $5p$ is needed to obtain access patterns depending on the list or graph.

In matrix multiplication, the dimensions of the matrices are as follows: if $p = 2^{2z}$ then $w = c = 5$, if $p = 2^{2z+1}$ then $w = 2c = 10$.

In each experiment we measured for each step of the trace the maximum module congestion c_{max} and then computed the expected value of all c_{max} averaged over all steps. The analysis is a kind of (expected) worst case analysis. Each expected value was checked for significance by looking at the variance. The five values obtained by using five functions of one class for each experiment were checked against significant differences. In case there were none, the average was taken. In case there were some, ten additional hash functions were chosen and the average was taken from these 15 values. Significant differences appeared only for stride $s = 13$, $p = 2^7, \dots, 2^9$ in both linear functions and for stride $s = 32$, $p = 2^9, 2^{10}$ in the linear function modulo power of two.

Because of mapping $5p$ requests per step onto p memory modules, $E(c_{max}) \geq 5$. The only exception is connected components, because not necessarily all processors make accesses in IF statements (see Sect. 3).

5 Results

The results of the experiments are presented in two ways. First we show the performance of the hash functions sorted by benchmarks. In Fig. 2 the performance on random patterns is given as a reference. The legend of the hash functions is shown in Fig. 3, which shows all other benchmarks. Second we show the performance sorted by hash functions in Fig. 4.

All figures are built as follows: the x -axis shows $\log p$ in range $5 \dots 10$, the y -axis shows the expected value of the maximum module congestions in range $4 \dots 14$.

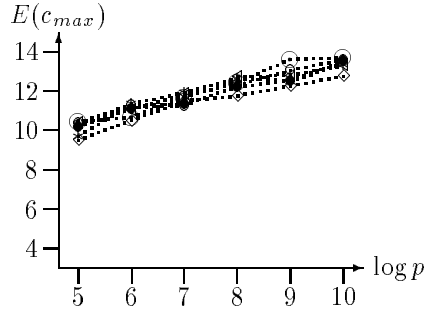


Fig. 2. Performance on random patterns

The performance on random patterns (see Fig. 2) is similar for all hash functions. Thus none of the hash functions is bad in an obvious way. The maximum module congestion rises from 10 for $p = 32$ to 12 for $p = 1024$. This will serve as a reference to analyse the performance on the other benchmarks.

5.1 Analysis of Benchmarks

The curves of Fig. 3 show similar shapes for all benchmarks: the polynomials of different degrees behave in a similar way and so do the three other hash functions. The behaviour of the polynomials furthermore is on all workloads worse than the behaviour of the simpler hash functions. Among the linear functions, the one modulo a prime always behaves a little bit worse than the linear function modulo a power of two. Thus the most interesting part is the comparison of our simple linear function with the boolean matrix hashing.

For strides that are a power of two, the boolean matrix hashes values optimally (see (a) and (c)) and reaches a module congestion of 6. The module congestion reached by the linear function lies between 6.5 and 7.5, so it is not far away.

A similar behaviour of linear function and boolean matrix can be seen in (d) and (f). This results from the fact that part of the accesses in these workloads are strides 1, when processors load or store values in arrays in the manner that processor i reads or writes $F[i]$.

However, as soon as we obtain other patterns, the boolean matrix hashing gets worse than the linear function (see (b) and (e)). Even for the matrix multiplication workload, where accesses always consist of $5 \cdot p^{1/2}$ strides with $s = 1$ and $p^{1/2}$ processors involved in each stride, the linear function is better.

5.2 Analysis of Hash Functions

Figure 4 shows the performance of the different hash functions. Because the connected components benchmark is not comparable to the others as explained

in Sect. 3, it is not shown here. The first observation is, that all hash functions behave on all workloads not worse than on random patterns. The second observation is that the polynomials show roughly the same behaviour on all workloads as they do on random patterns (see (d) to (f)). We conclude that their performance is independent of the application. That is what we expected. But this performance is bad in comparison to what is reached by the other functions that behave better than on random patterns on all workloads.

The linear function (see (a)) shows almost uniform behaviour on all workloads, too, but it varies between 6.5 and 8, which is significantly better than the behaviour on random patterns.

The behaviour of the linear function modulo a prime is not uniform and varies between 6.3 and 10.

The behaviour of the boolean matrix hashing function can be divided in an expected optimal behaviour for strides with s a power of two and a significantly higher module congestion for other patterns, which is however still below the one produced by random patterns.

6 Conclusions

The above experiments show surprisingly that linear functions modulo a power of two and boolean matrix functions show best performance for practical use. Both have the additional properties of bijectivity and short evaluation time. The choice between these two depends on the expected user profile (if such exists) and the surrounding machine architecture. For machines that already contain a hardware multiplier this could be used to perform hashing in the case of the linear functions. Moreover, the use of matrix hashing is restricted by the fact that an implementation needs $(\log m)^2$ bit register hardware to store the boolean matrix. Therefore, if no user profile is known and chip area is restricted (or a multiplier already available), the use of the linear function is preferable.

The observations presented here lead to the use of linear hash functions in the prototype design of the SB-PRAM [1, 2] which emulates a synchronous shared memory machine with $p = 128$ physical processors and provides hardware support for hashing and packet routing including combining.

Unfortunately, some open questions remain. First, there is no theoretical framework to explain why simple hash functions work better than complex ones. Also, the exact relationship between linear functions with and without “modulo prime” is still unknown.

Acknowledgements

The authors would like to thank Martin Dietzfelbinger for helpful discussions.

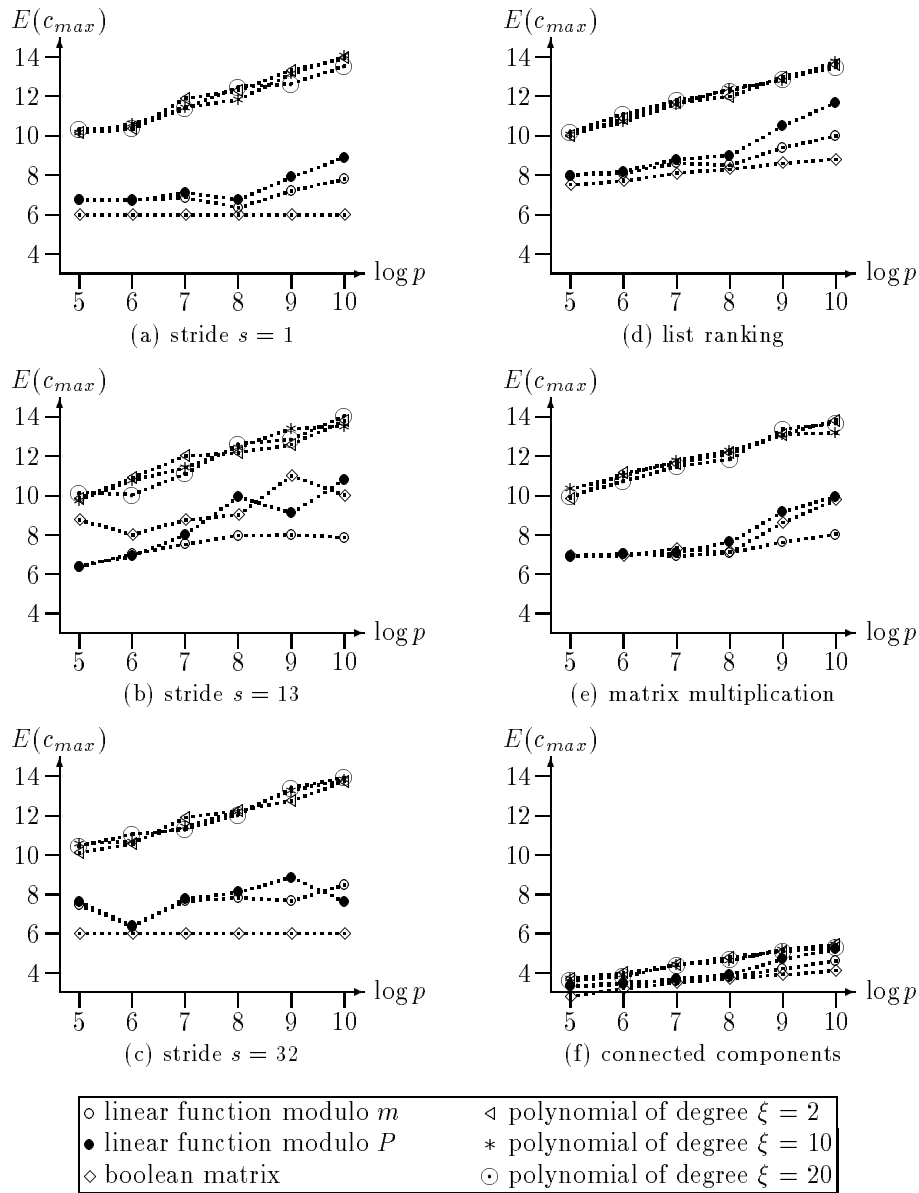


Fig. 3. Performance on benchmarks

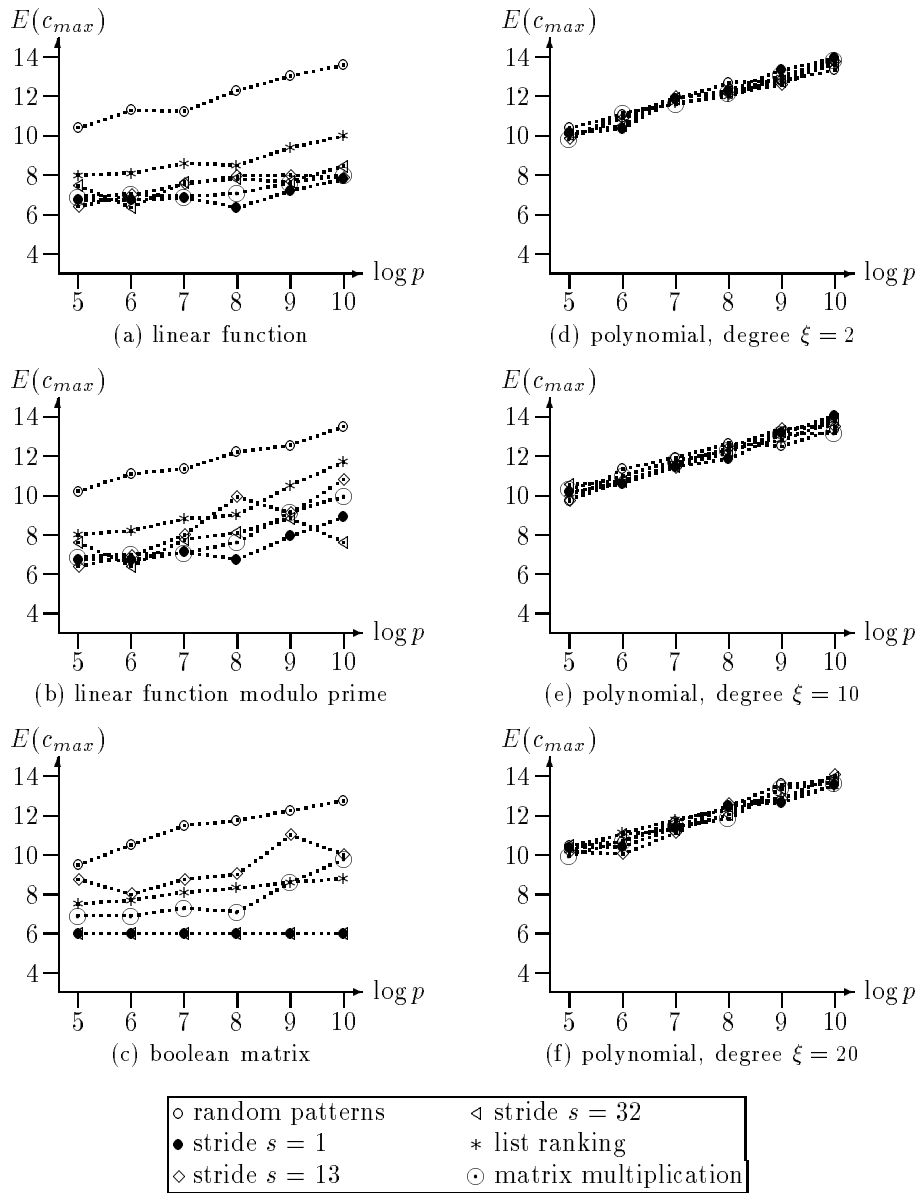


Fig. 4. Performance of hash functions

References

1. Abolhassan, F., Drefenstedt, R., Keller, J., Paul, W. J., Scheerer, D.: On the Physical Design of PRAMs. Informatik — Festschrift zum 60. Geburtstag von Günter Hotz (Teubner, 1992) 1–19
2. Abolhassan, F., Keller, J., Paul, W. J.: On the cost–effectiveness of PRAMs. Proc. 3rd IEEE Symp. on Parallel and Distributed Processing (1991) 2–9
3. Akl, S. G.: The Design and Analysis of Parallel Algorithms. (Prentice-Hall, 1989)
4. Dietzfelbinger, M.: On limitations of the performance of universal hashing with linear functions. Reihe Informatik Bericht Nr. 84 Universität–GH Paderborn (1991)
5. Dietzfelbinger, M., Hagerup, T., Katajainen, J., Penttonen, M.: A reliable randomized algorithm for the closest-pair problem. Manuscript (1992)
6. Karlin, A. R., Upfal, E.: Parallel hashing: An efficient implementation of shared memory. J. Assoc. Comput. Mach. **35** (1988) 876–892
7. Karp, R. M., Ramachandran, V. L.: A survey of parallel algorithms for shared–memory machines. Handbook of Theoretical Computer Science Vol. A (Elsevier, 1990) 869–941
8. Mehlhorn, K., Vishkin, U.: Randomized and deterministic simulations of PRAMs by parallel machines with restricted granularity of parallel memories. Acta Inform. **21** (1984) 339–374
9. Norton, A., Melton, E.: A class of boolean linear transformations for conflict–free power–of–two stride access. Proc. Internat. Conf. on Parallel Processing (1987) 247–254
10. Ranade, A. G.: How to emulate shared memory. Proc. 28th IEEE Symp. on Foundations of Computer Science (1987) 185–194
11. Shiloach, Y., Vishkin, U.: An $O(\log n)$ parallel connectivity algorithm. J. Algorithms **3** (1982) 57–67