

Elimination Trees and the Construction of Pools and Stacks¹

Nir Shavit*

Dan Touitou

MIT and
Tel-Aviv University

Tel-Aviv University

February 28, 1996

Abstract

Shared *pools* and *stacks* are two coordination structures with a history of applications ranging from simple producer/consumer buffers to job-schedulers and procedure stacks. This paper introduces *elimination trees*, a novel form of diffracting trees that offer pool and stack implementations with superior response (on average constant) under high loads, while guaranteeing logarithmic time “deterministic” termination under sparse request patterns.

¹A preliminary version of this paper appeared in the proceedings of the *7th Annual Symposium on Parallel Algorithms and Architectures (SPAA)*. Contact Author: E-mail: shanir@theory.lcs.mit.edu

1 Introduction

As multiprocessing breaks away from its traditional number crunching role, we are likely to see a growing need for highly distributed and parallel coordination structures. A real-time application such as a system of sensors and actuators will require fast response under both sparse and intense activity levels (typical examples could be a radar tracking system or a traffic flow controller). Shared *pools* offer a potential solution to such coordination problems, with a history of applications ranging from simple producer/consumer buffers to job-schedulers [7] and procedure stacks [26]. A *pool* [16] (also called a pile [22], global pool [7] or a producer/consumer buffer) is a concurrent data-type which supports the abstract operations: `enqueue(e)` – adds element e to the pool, and `dequeue` – deletes and returns some element e from the pool. A stack is a pool with a last-in-first-out (*LIFO*) ordering on enqueue and dequeue operations.

Since the formal introduction of the problem and its first solution by Manber [16], the literature has offered us a variety of possible pool implementations. On the one hand there are queue-lock based solutions such as of Anderson [2] and Mellor-Crummey and Scott [15], which offer good performance under sparse access patterns, but scale poorly since they offer little or no potential for parallelism in high load situations. On the other hand, there are a variety of that “load-balanced local pools” based algorithms like Manber’s *search tree* structure [16] and the simple and effective randomized *work-pile* and *job-stealing* techniques as designed by Kotz and Ellis [13], Rudolph, Slivkin-Allaluf, and Upfal [22], Lüling and B. Monien [21], and Blumofe and Leiserson [7]. These algorithms offer good *expected* response time under high loads, but very poor performance as access patterns become sparse (their expected response time becomes linear in n – the number of processors in the system – as opposed to that of a “deterministic” queue-lock based pool that is linear in the number of participating processors). This linear behaviour under sparse access patterns holds also for Manber’s tree based deterministic *job-stealing* method [16].

Shavit and Zemach’s *diffracting trees* [24] have recently been proposed as a reasonable middle-of-the-road solution to the problem. They guarantee termination within $O(\log w)$ time (where $w \ll n$) under sparse access patterns, and rather surprisingly manage to maintain similar average response time under heavy loads.

1.1 Elimination Trees

This paper introduces *elimination trees*, a novel form of diffracting trees that offers pool implementations with the same $O(\log w)$ termination guarantee under sparse patterns, but with far superior response (on average constant) under high loads. Our empirical results show that unlike diffracting trees, and in spite of the fact that elimination trees offer a “deterministic” guarantee

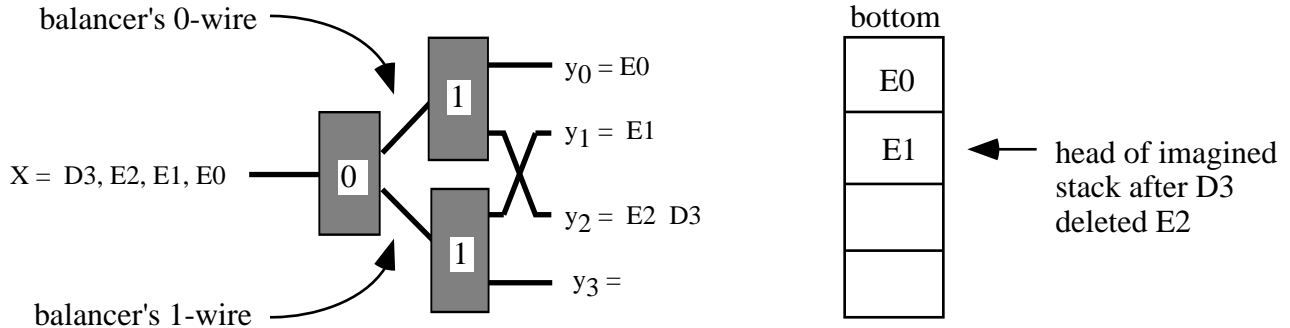


Figure 1: A sequential execution on a $\text{STACK}[4]$ elimination tree

of coordination,¹ they scale like the “randomized” methods [7, 13, 21, 22], providing improved response time as the load on them increases.

In a manner similar to diffracting trees, elimination trees are constructed from simple one-input two-output computing elements called *elimination balancers* that are connected to one another by wires to form a balanced binary tree with a single root input wire and multiple leaf output wires. While diffracting trees route *tokens*, elimination trees route both *tokens* and *anti-tokens*. These arrive on the balancer’s input wire at arbitrary times, and are output on its output wires. The balancer acts as a toggle mechanism, sending tokens and anti-tokens left and right in a balanced manner. For example, to create a pool implementation that has stack-like behavior, the balancer can consist of a single bit, with the rule that tokens toggle the bit and go to the 0 or 1 output wire according to its *old* value, while anti-tokens toggle the bit and go left or right according to its *new* value. Now, imagine that stack array entries are placed at the leaves of the tree, and think of tokens as enqueue (“push”) requests and anti-tokens as dequeue (“pop”) requests. Figure 1 shows a width four tree after 3 enqueues and a dequeue have completed. The reader is urged to try this sequence with toggles initially 0. The state of the balancers after the sequence is such that if next a token will enter it will see 0 and then 1 and end up on wire y_2 , while if the next to enter is an anti-token it will get a 1 and then a 0 and end up on wire y_1 , finding the value to be deleted. In fact, our tree construction is a novel form of a *counting network* [4] based counter, that allows decrement (anti-token) operations in addition to standard increment (token) operations.

However, this simple approach is bound to fail since the toggle bit at root of the tree will be a hot-spot [17, 18] and a sequential bottleneck that is no better than a centralized stack implementation. The problem is overcome by placing a *diffracting prism* [24] structure in front of the toggle bit inside every balancer. Pairs of tokens attempt to “collide” on independent locations in the prism,

¹They guarantee that a dequeue operation on a non-empty queue will always succeed.

diffracting in a coordinated manner one to the 0-wire and one to the 1-wire, thus leaving the balancer without ever having to toggle the shared bit. This is not a problem since in any case after both toggled it, the bit would return to its initial state. This bit will only be accessed by processors that did not succeed in colliding, and they will toggle it and be directed as before.

Our first observation is that the stack behavior will not be violated if pairs of anti-tokens, not only tokens, are diffracted. The second, more important fact, is that it will continue to work if collisions among a token and an anti-token result in the “elimination” of the pair, without requiring them to continue traversing the tree! In other words, a token and anti-token that meet on a prism location in a balancer can exchange enqueue/dequeue information and complete their operation without having to continue through $\log w$ balancers. In fact, our empirical tests show that under high loads, most tokens and anti-tokens are eliminated within two levels. Of course, the tree structure is needed since one could still have long sequences of enqueues only.

We compared the performance of elimination trees to other known methods using the Proteus Parallel Hardware Simulator [8] in a shared memory architecture similar to the Alewife machine of Agarwal et al. [3]. We first compared under high loads a variety of methods that can be used to implement a stack-like pool and are known to perform well under sparse access patterns. We found that elimination trees scale *substantially* better than all of these methods including queue-locks [15], Combining trees [10], and Diffracting Trees [24].

We then compared Elimination trees to the *load-balanced local pools* techniques [16, 13, 22, 21, 7] which cannot be used to implement a stack-like pool and theoretically provide only linear performance under sparse access patterns. We found that in many high load situations elimination trees are inferior to these methods (as explained in the sequel, we chose for the comparison a representative technique, the randomized technique of Rudolph, Slivkin, and Upfal [22]), especially for job distribution applications where a typical processor is the dequeuer of its latest enqueue (though in many cases not by much). However, our empirical evidence suggests that elimination trees provide up to a factor of 30 better response time than randomized methods under sparse loads. Finally, we present evidence that our new elimination balancer design offers a more scalable diffracting balancer construction even in cases where no collisions are possible.

2 Pools

We begin with our pool specification and implementations, later showing how to modify them to create stack-like pools.

A *pool* [16](also called a pile [22], centralized “pool” [7] or a producer/consumer buffer) is a concurrent data-type which maintains a multiset of values by supporting the abstract operations: `enqueue(e)` – adds element e to the multiset, and `dequeue` – deletes and returns some element e

from the multiset. For simplicity, assume that all enqueued elements e are unique, that is, multiset is simply a set. A pool is a relaxation of a first-in-first-out queue: apart from the queue’s basic safety properties, no causal order is imposed on the enqueued and dequeued values. However, it is required that:

P1 an enqueue operation always succeeds, and

P2 a dequeue operation succeeds if the pool is non-empty, that is, for every execution in which the number of enqueue operations is greater or equal to the number of dequeue operations, all the dequeue operations succeed.

A *successful* operation is one that is guaranteed to return an answer within finite (in our construction, *bounded*) time. Note that the randomized decentralized techniques of [7, 13, 21, 22] implement a weaker “probabilistic” pool definition, where condition *P2* is replaced by a *probabilistic* guarantee that dequeue operations succeed.

2.1 Elimination Trees

Our pool implementation is based on the abstract notion of an *elimination tree*, a special form of the diffracting tree data structures introduced by Shavit and Zemach in [24]. Our formal model follows that of Aspnes, Herlihy, and Shavit [4] I/O-automata of Lynch and Tuttle [20].

An *elimination balancer* is a routing element with one input wire x and two output wires y_0 and y_1 . *Tokens* and *anti-tokens* arrive on the balancer’s input wire at arbitrary times, and are output on its output wires. Every token carries a value. Whenever a token “meets” an anti-token in a balancer, it passes the value to the anti-token and both token and anti-token are eliminated and never output from the balancer. More formally, a pool balancer is a shared object that allows processors to execute `TokenTraverse(Token Type, v)` operations which have as input the token’s type, `TOKEN` or `ANTI-TOKEN`, and its input value v (which is non-empty in case of a `TOKEN` type traversal). Each such operation returns 0 or 1, depending on which of the output wires y_0 and y_1 the token should proceed, or the pair `(ELIMINATED, v)` meaning that the token (or anti-token) was eliminated and that the value v was exchanged. We slightly abuse our notation and denote by x and \bar{x} the number of tokens and anti-tokens ever received, and by y_i and \bar{y}_i , $i \in \{0, 1\}$, the number of tokens and anti-tokens ever output on the i th output wire. The pool balancer object must guarantee:

Quiescence Given a finite number of input tokens and anti-tokens, the balancer will reach a *quiescent* state, that is, a state in which all the tokens and anti-tokens traversal operation executions have completed.

Pairing In any quiescent state, there exists a *perfect matching* between eliminated tokens and eliminated anti-tokens, such that the value returned by an eliminated anti-token is matched with the value carried by its corresponding eliminated token.

Pool Balancing In any quiescent state, if $x \geq \bar{x}$ then for every output wire $i \in \{0, 1\}$, $y_i \geq \bar{y}_i$.

Let $\text{POOL}[w]$ be a binary tree of elimination balancers with a root input wire x and w designated output wires: y_0, y_1, \dots, y_{w-1} , constructed inductively by connecting the outputs of an elimination balancer to two $\text{POOL}[w/2]$ trees. From the quiescence property of the balancers, given a finite number of input tokens and anti tokens, $\text{POOL}[w]$ will reach a quiescence state in which all the tokens and anti tokens are either eliminated or have exited through one of $\text{POOL}[w]$ output. We extend pool balancing to trees in the natural way claiming that:

Lemma 2.1 *The outputs y_0, \dots, y_{w-1} of $\text{POOL}[w]$ satisfy the pool balancing property in any quiescent state.*

Proof: The proof is by induction on w . When $w = 2$ this follows directly from the balancer definition. Assume the claim for $\text{POOL}[w/2]$ and let us prove it for $\text{POOL}[w]$. If the number of tokens entering the root balancer of $\text{POOL}[w]$ is greater or equal to the number of anti-tokens, then, by definition this property is kept on the output wires of the root balancer, and by the induction hypothesis holds for the output wires of both $\text{POOL}[w/2]$ trees. ■

On a shared memory multiprocessor, one can implement an elimination tree as a shared data structure, where balancers are records, and wires are pointers from one record to another. Each of the machine’s asynchronous processors can run a program that repeatedly traverses the data structure from the root input pointer to some output pointer, each time shepherding a new “token” or “anti-token” through the network (see Figure 3). Constructing a *pool* object from a $\text{POOL}[w]$ tree is straightforward: each tree output wire is connected to a sequentially accessed “local” pool, a simple queue protected by a Mellor-Crummey and Scott MCS-queue-lock [15] will do. The MCS-queue-lock has the property of being “fair,” and so every access request to the queue will be granted within a bounded number of operations. A process performs an enqueue operation by shepherding a token “carrying” the value the down the tree. If the token reaches the output wire, the associated value is enqueued in the local pool connected to that wire. The dequeue operation is similarly implemented by carrying an anti-token through the network. If this anti-token collides with a token in a balancer, the dequeuing process returns the token’s value. Otherwise it exits on a wire and performs a dequeue operation on the anti-token’s local pool. Naturally if the local pool is empty the dequeuing process waits until the pool is filled and then access it. The elimination tree is thus a load-balanced *coordination* medium among a distributed collection of pools. It differs from elegant

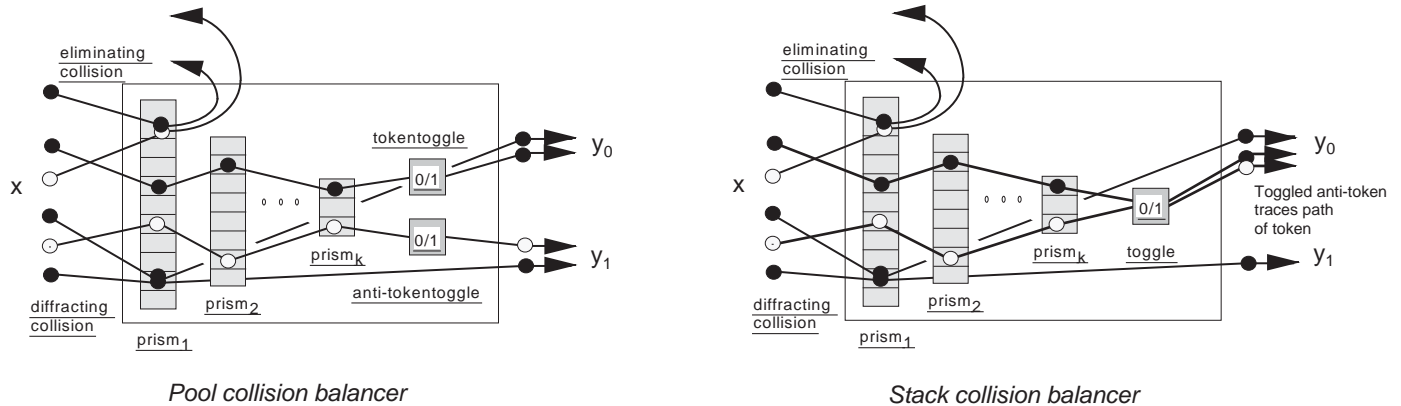


Figure 2: The structure of Pool and Stack elimination balancers

randomized constructions of [7, 13, 21, 22] in its deterministic dequeue termination guarantee and in performance. While work in an individual balancer is relatively high, each enqueue or dequeue request passes at most $\log w$ balancers both under high and under low loads.

Theorem 2.2 *The elimination tree based pool construction is a correct pool implementation.*

Proof: The basic safety properties of the pool are satisfied thanks to the perfect matching between eliminated tokens. By the quiescence property of the balancers all the tokens and anti-tokens will eventually reach the exits of the elimination tree. Since the MCS-queue-locks controlling access to the local pools are fair, all the enqueue operations will succeed in adding their value to the local pools within some bounded number of operations and property **P1** will be satisfied. Now, if the number of dequeue operations is greater than the number of enqueue operations, by Lemma 2.1 this will eventually be the case at each of the each of the local pools at the leaves. In that case no dequeue operation will never have to wait indefinitely at a leaf. This satisfies property **P2**. ■

2.2 Pool Elimination Balancers

The scalable performance of our pool constructions depends on providing an efficient implementation of an elimination balancer.

Diffracting balancers were introduced in [24]. Our shared memory construction of a diffracting elimination balancer, apart from providing a mechanism for token/anti-token elimination, also improves on the performance of the original diffracting balancer design. While a regular diffracting balancer [24] is constructed from a single prism array and a toggle bit, the elimination balancer we use in our pool construction (see lefthand side of Figure 2) has a sequence of prism arrays and two

toggle bits, one for tokens and one for anti-tokens². Each of the toggle bit locations is protected by an MCS-queue-lock [15]. A process shepherding a token or anti-token through the balancer decides on which wire to exit according to the value of the respective token or anti-token toggle bit, 0 to the left and 1 to the right, toggling the bit as it leaves. The toggle bits effectively balance the number of tokens (resp. anti-tokens) on the two output wires, so that there is in any quiescent state at most one token (resp. anti-token) more on the 0 output wire than on the 1 wire. The reader can easily convince herself that this suffices to guarantee the pool-balancing property. However, if many tokens were to attempt to access the same toggle bit concurrently, the bit would quickly become a hot spot. The solution presented in [24] is to add a *prism* array in front of each toggle bit. Before accessing the bit, the process shepherding the token selects a location l in the prism uniformly at random, hoping to “collide” with another token which selected l . If a collision occurs, then the tokens “agree” among themselves that one should be “diffracted” left and the other right (the exact mechanism is described in the sequel), without having to access the otherwise congested toggle bit. If such a *diffracting collision* does not occur, the process toggles the bit as above and leaves accordingly. As proved in [24], the combination of diffracted tokens and toggling tokens behaves exactly as if all tokens toggled the bit, because if any two diffracted tokens were to access the bit instead, after they both toggled it the bit state would anyhow return to its initial state. The same kind of prism could be constructed for anti-tokens.

The key to our new constructions is the observation that for data structures which have complementary operations (such as enqueues and dequeues), one can gain a substantial performance benefit from having a joined prism for both tokens and anti-tokens. In addition to toggling and diffracting of tokens and anti-tokens, if a collision between a token and anti-token occurs in the shared prism, they can be “eliminated” (exchanging the complementary information among themselves) without having to continue down the tree. We call this an *eliminating collision*. Unlike with diffracting collisions, if the eliminating collision had not occurred, each of the token and anti-token toggle bits would have changed. Nevertheless, the combination of toggling, diffracting and elimination preserves the pool-elimination balancer’s correctness properties, which by Lemma 2.1 guarantees pool-balancing.

The *size* of (number of locations in) the prism array has critical influence on the efficiency of the node. If it is too high, tokens will miss each other, lowering the number of successful eliminations, and causing contention on the toggle bits. If it is too low, too many processes will collide on the same prism entry, creating a hot-spot. We typically found the optimal performance was when the prism width at a balancer on a given level is the same as the width of the subtree below it (this conforms with recent projections based on steady-state analysis [25]). Moreover, unlike the single prism array of [24], we found it more effective to pass a token through a series of prisms of decreasing size, thus

²The two separate toggle locations are an artifact of the pool-balancing property. In our stack construction in Section 3 the elimination balancer uses a single toggle bit for both tokens and anti-tokens.


```

root : global ptr to root of elimination tree

procedure enqueue(v:value);
  b:= root
  while not leaf(b)
    r :=TokenTraverse(TOKEN,v) on balancer b;
    case r of
      ELIMINATED: return;
      0          : b := left child of b;
      1          : b := right child of b;
    endcase
  endwhile
  enqueue_local_pool(b,e)

function dequeue(): value;
  b:= root
  while not leaf(b)
    r:=TokenTraverse(ANTITOKEN,EMPTY) on balancer b;
    case r of
      <ELIMINATED,v> : return v;
      0              : b := left child of b;
      1              : b := right child of b;
    endcase
  endwhile
  return dequeue_local_pool(b);

```

Figure 3: Tree traversal code

increasing the chances of a collision. This way, at high contention levels most of the collisions will occur on the larger prisms while at low levels they happen on the smaller ones.

Figure 4 gives the code for traversing an elimination balancer. Note that for algorithmic simplicity we omitted input values and the code for their exchange, and have deferred a discussion of this issue to Section 2.4.

Apart from reading and writing memory, our implementation uses a hardware

- `register_to_memory_swap(addr, val)` operation, and a
- `compare_and_swap(addr, old, new)`, an operation which checks if the value at address `addr` is equal to `old`, and if so, replaces it with `new`, returning `TRUE` and otherwise `FALSE`.

```

Location: shared array[1..NUMPROCS];

Function TokenTraverse(b: ptr to bal, mytype: TokenType)
    returns (ptr to bal or ELIMINATED);

    Location[mypid] := <b,mytype>;
    /* Part 1 : attempt to collide with another token on k prism levels */
    for i:=1 to k do
        place := random(1,size_i);
        him := register_to_memory_swap(Prism_i[place],mypid);
        if not_empty(him) then
            <his_b,his_type> := Location[him];
            if his_b = b then
                if compare_and_swap(Location[mypid],<b,mytype>, <0,EMPTY>) then
                    if my_type = his_type then
                        if compare_and_swap(Location[him],<b,his_type>,<0,DIFFRACTED>) then
1.                 return b->OutputWire[1]
                        else Location[mypid] := <b,mytype>;
                        else if compare_and_swap(Location[him],<b,his_type>,<0,ELIMINATED>) then
2.                 return ELIMINATED;
                        else Location[mypid] := <b,mytype>;
                        else if Location[mypid]= <0,DIFFRACTED> return (b->OutputWire[0])
                        else return ELIMINATED
                repeat b->Spin times /* wait in hope of being collided with */
                    if Location[mypid] = <0,DIFFRACTED> then return b->OutputWire[0];
                    if Location[mypid] = <0,ELIMINATED> then return ELIMINATED;
    /* Part 2 access toggle the bits */
    AcquireLock(b->Locks[mytype]);
    if compare_and_swap(Location[mypid],<b,my_type>, <0,EMPTY>) then
        i:= b->Toggles[mytype];
        b->Toggles[mytype] := Not(i);
        ReleaseLock(b->Locks[mytype]);
3.     return b->OutputWire[i];
    else ReleaseLock(b->Locks[mytype]);
        if Location[mypid]= <0,DIFFRACTED> return (b->OutputWire[0])
        else return ELIMINATED

```

Figure 4: Traversing an eliminating balancer

Our implementation also uses standard `AcquireLock` and `ReleaseLock` procedures to enter and exit the MCS-queue-lock [15].

Initially, processor p announces the arrival of its token at node b , by writing b and its token type to `Location`[p]. It then chooses a location in the `Prism1` array uniformly at random (note that randomization here is used only to load-balance processors over the prism, and could be eliminated in many cases without a significant performance penalty) and swaps its own PID for the one written there. If it read a PID of an existing processor q (i.e. `not_empty(him)`), p attempts to collide with q . This collision is accomplished by first executing a `<his_b, his_type> := Location[him]` read operation to determine the type of token being collided with, and then performing two compare-and-swap operations on the `Location` array. The first clears p 's entry, assuring no other processor will collide with it during its collision attempt (this eliminates race conditions). The second attempts to mark q 's entry as "collided with p ," notifying q of the collision type: `DIFFRACTED` or `ELIMINATED`. If both compare-and-swap operations succeed, the collision is successful, and p decides based on collision type to either diffract through the right output wire or to be eliminated. If the first compare-and-swap fails, it follows that some other processor r has already managed to collide with p . In that case p diffracts through the left output wire or is eliminated, depending on the type of the processor that collided with it. If the first succeeds but the second fails, then the processor with whom p was trying to collide is no longer at balancer b , in which case p resets its `Location` entry to contain the balancer name and its token type, and having failed to "collide with" another processor, spins on `Location`[p] waiting for another processor to "collide with it." If after `spin` time units no collision occurs, p restarts the whole process at the next level `Prism2` and so on. If p has traversed all the prism levels without colliding, it acquires the lock on the toggle bit, clears its element, toggles the bit and releases the lock. If p 's element could not be erased, it follows that p has been collided with, in which case p releases the lock without changing the bit and diffracts or is eliminated accordingly.

2.3 Correctness Proof of Pool Balancer Implementation

Clearly if no diffractions and no eliminations occur during an execution, by the code all the tokens would access the toggle bits and the balancing property will easily be satisfied. Hence, in order to prove the correctness of our implementation we should focus on showing that eliminating and diffracting tokens are paired off correctly. For example, we must show that a scenario in which token T_1 diffracts with token T_2 and in which T_2 is not aware of it and still toggles the bit, will never happen. As a first step, let us assume that every token in a given execution has a unique virtual ID T_p , and let the subscript p denote the PID of the process shepherding the token. We use the "*" notation throughout the paper to denote an unspecified value. In the following lemma we show that if some process p reads `Location`[q]=`<b,*>`, then process q is currently shepherding a token through balancer b .

Lemma 2.3 *For every process p , if $\text{Location}[p]=\langle b,* \rangle$ then p is executing `TokenTraverse` on balancer b .*

Proof: Initially $\text{Location}[p]=0$. From the algorithm it is clear that only p can write a value different than 0 as a balancer name in $\text{Location}[p]$. Since p always writes 0 into $\text{Location}[p]$ (a successful `compare_and_swap`) before completing `TokenTraverse`, the claim follows. ■

We now define a token T_p traversing a balancer b as a *diffracting* token if p has executed Line 1 in the algorithm and thus “leaves on output wire 1.” Since for every diffracting token T_p , p executed a successful `compare_and_swap`($\text{Location}[\text{him}], \langle b,* \rangle, \langle 0, \text{DIFFRACTED} \rangle$), we know by Lemma 2.3 that at the same time process him was shepherding some token T_{him} through b . We designate T_{him} , which “leaves on output wire 0” as *diffracted by* T_p . We also define a token T_p as an *eliminating* token if p executed Line 2. In a similar way as for diffracting tokens we designate the token T_{him} as *eliminated by* T_p . Finally we define a token T_p as a *toggling* token if p has executed Line 3 in the algorithm. From the flow control of the algorithm it is clear than a token cannot be both toggling and eliminating, or toggling and diffracting, or eliminating and diffracting.

In the next two lemmas we show that tokens are paired off correctly during elimination and diffraction.

Lemma 2.4 *Every token traversing a balancer b can be diffracted or eliminated by at most one other token.*

Proof: By way of contradiction. Assume that a token T_p , while traversing b has been eliminated or diffracted by two other tokens T_q and T_r . In that case, both q and r have successfully executed `compare_and_swap`($\text{Location}[p], \langle b,* \rangle, \langle 0,* \rangle$). It follows that p must have written $\langle b,* \rangle$ in $\text{Location}[p]$ at least twice during the execution of the `TokenTraverse` carrying T_p through b . But in that case `compare_and_swap`($\text{Location}[p], \langle b,* \rangle, \langle 0, \text{EMPTY} \rangle$) was successfully executed by p before writing $\langle b,* \rangle$ on $\text{Location}[p]$ for the second time. A contradiction. ■

Lemma 2.5 *A toggling, eliminating, or diffracting token T_p cannot be eliminated or diffracted by some other token T_q .*

Proof: Follows since q executes Lines 1,2, or 3, or writes $\langle b,* \rangle$ on $\text{Location}[q]$, only after executing a successful `compare_and_swap`($\text{Location}[q], \langle b,* \rangle, \langle 0, \text{EMPTY} \rangle$), no other process will be able to execute a successful `compare_and_swap`($\text{Location}[q], \langle b,* \rangle, \langle 0, \text{EMPTY} \rangle$). ■

We now prove that:

Theorem 2.6 *The pool balancer implementation given in Figure 4 satisfies the pool balancing property.*

Proof: Given any execution of the pool implementation, let d_1 and \bar{d}_1 be the number of diffracting (leaving on wire 1) tokens and anti-tokens respectively and let d_0 and \bar{d}_0 be the number of diffracted (leaving on wire 0) tokens and anti-tokens. We designate by e the number of eliminated and eliminating tokens and by \bar{e} the number of eliminating and eliminated anti-tokens. Finally let t and \bar{t} be the number of toggling tokens and anti-tokens respectively.

By Lemma 2.5 $x = d_0 + d_1 + e + t$ and $\bar{x} = \bar{d}_0 + \bar{d}_1 + \bar{e} + \bar{t}$. By Lemma 2.4, $\bar{e} = e$, $d_0 = d_1$ and $\bar{d}_0 = \bar{d}_1$. Now, if $x \geq \bar{x}$ then $t + d_0 + d_1 = x - e \geq \bar{x} - \bar{e} = \bar{t} + \bar{d}_0 + \bar{d}_1$. Consequently

$$\lceil \frac{t + d_0 + d_1}{2} \rceil \geq \lceil \frac{\bar{t} + \bar{d}_0 + \bar{d}_1}{2} \rceil,$$

and since $d_0 = d_1$ and $\bar{d}_0 = \bar{d}_1$ then $\lceil \frac{t}{2} \rceil + d_0 \geq \lceil \frac{\bar{t}}{2} \rceil + \bar{d}_0$. Therefore $y_0 \geq \bar{y}_0$. Using the same arguments, one can show that $\lfloor \frac{t}{2} \rfloor + d_1 \geq \lfloor \frac{\bar{t}}{2} \rfloor + \bar{d}_1$ and therefore $y_1 \geq \bar{y}_1$. ■

2.4 Exchanging Values in Eliminating Collisions

The purpose of the eliminating collisions is to allow enqueueers and dequeuers to exchange values and to leave the pool. The algorithm in Figure 4 can be easily modified to handle value exchanges: every process writes and reads from `Location[mypid]` a triplet `<b,mytype,value>` instead of just the pair `<b,mytype>`. To eliminate an anti-token, a token writes `<0,ELIMINATED,value>` in the anti-token's `Location`. Note that it knows this is an anti-token following the preliminary `<his_b,his_type> := Location[him]` read operation. In this way the eliminated anti-token will find this value and return it. On the other hand, an eliminating anti-token returns the value it has read from the eliminated token's `Location` entry. Since, the triplets stored in `Location` are written and updated atomically, only minor modifications are needed in the correctness proof: we just have to show that an eliminating (or eliminated) anti-token returns the value carried by the token it has eliminated (or was eliminated by). The proof of this lemma is identical to the proof of Lemma 2.3.

Lemma 2.7 *For every process p , if `Location[p]=<b,TOKEN,v>` then p is shepherding a token carrying value v on balancer b .*

We have shown in Lemmas 2.4 and 2.5 that eliminated tokens and anti-tokens are paired off correctly. We prove now that eliminated or eliminating anti-tokens exchange values in a proper way.

Lemma 2.8 *Every eliminated anti-token returns the value carried by the token that has eliminated it. Every eliminating anti-token returns the value carried by the token it has eliminated.*

Proof: Assume that T_p is an eliminated anti-token. Let T_q be the token which eliminated T_p . By the modified algorithm `compare_and_swap(Location[p], <b, ANTI-TOKEN, NULL>, <0, ELIMINATED, v>)` was successfully executed by q , where v is the value carried by T_q . Since only p can change the content of `Location[p]`, and it could not, it must have returned v .

Assume that T_q is an eliminating anti-token which returned a value v and let T_p be the token it eliminated. Process q executed `compare_and_swap(Location[p], <b, TOKEN, v>, <0, ELIMINATED, NULL>)` successfully, and therefore by Lemma 2.7, v must be the value carried by T_p . ■

2.5 Performance of the Elimination Tree Based Pool

We evaluated the performance of our *elimination tree* based pool construction relative to other known methods by running a collection of benchmarks on a simulated 256 processor distributed-shared-memory machine similar to the MIT *Alewife* machine [3] of Agarwal et. al. The presented results hopefully exemplify the potential in using elimination trees, but in no way claim to be a comprehensive study of their performance.

Our simulations were performed using *Proteus* a multiprocessor simulator developed by Brewer, Dellarocas, Colbrook and Weihl [8]. *Proteus* simulates parallel code by multiplexing several parallel threads on a single CPU. Each thread runs on its own virtual CPU with accompanying local memory, cache and communications hardware, keeping track of how much time is spent using each component. In order to facilitate fast simulations, *Proteus* does not complete *cycle per cycle* hardware simulations. Instead, local operations (that do not interact with the parallel environment) are run uninterrupted on the simulating machine's CPU. The amount of time used for local calculations is added to the time spent performing simulated globally visible operations to derive each thread's notion of the current time. *Proteus* makes sure a thread can only see global events within the scope of its local time.

Our simulated *Alewife* like machine has 256 processors, each at a node of a Torus shaped communication grid. Each node also contains a cache memory, a router, and a portion of the globally-addressable memory. The cost of switching or wiring in the *Alewife* architecture is 1 cycle/packet. Each processor has a cache with 2048 lines of 8 bytes. The cache coherence is provided using a using a version of Chaiken's directory-based cache-coherence protocol [9].

2.5.1 The Produce-Consume Benchmark

We begin by comparing under various loads *deterministic* pool constructions which are known to guarantee good enqueue/dequeue time when the load is low (sparse access patterns). These methods are also the ones that can be modified to provide stack-like pool behaviour. In the produce-consume

```

Pool: array[1..N] of elements; - initially set to NULL -- N must be chosen optimally
headcounter, tailcounter:integer; - initially set to 0

Procedure Enqueue(el:elements);
  i:= fetch_and_increment(headcounter);
  repeat
    flag:= compare_and_swap(Pool[i],NULL,el);
  until flag= TRUE;

Function Dequeue() returns elements;
  i:= fetch_and_increment(tailcounter);
  repeat
    repeat el := Pool[i] until el <> NULL;
    flag := compare_and_swap(Pool[i],el,NULL);
  until flag= TRUE;
  return el;

```

Figure 5: A pool based on a cyclic array and shared counters.

benchmark each processor alternately enqueues a new element in the pool, dequeues a value from the pool and then waits a random number of cycles between 0 and `Workload` (see Figure 6).

```

repeat
  produce(val);
  val := consume;
  w := random(0..Workload);
  wait w cycles;
until 10^6 cycles elapsed

```

Figure 6: Produce-Consume Benchmark.

We ran this benchmark varying the number of processors participating in the simulation during 10^6 cycles, measuring: *latency*, the average amount of time spent per produce and consume operation, and *throughput*, the number of produce and consume operations executed during 10^6 cycles.

In preliminary tests we found that the most efficient pool implementations are attained when using shared counting to load balance and control access to a shared array (see Figure 5).

We thus realized the centralized pool in the style of [4], given in Figure 5, where the `headcounter` and `tailcounter` are implemented using two counters of the following type:

MCS The MCS-queue-lock of [15], whose response time is linear in the number of concurrent requests. Each processor locks the shared counter, increments it, and then unlocks it. The code was taken directly from the article, and implemented using atomic operations: `register_to_memory_swap` and `compare_and_swap` operations.

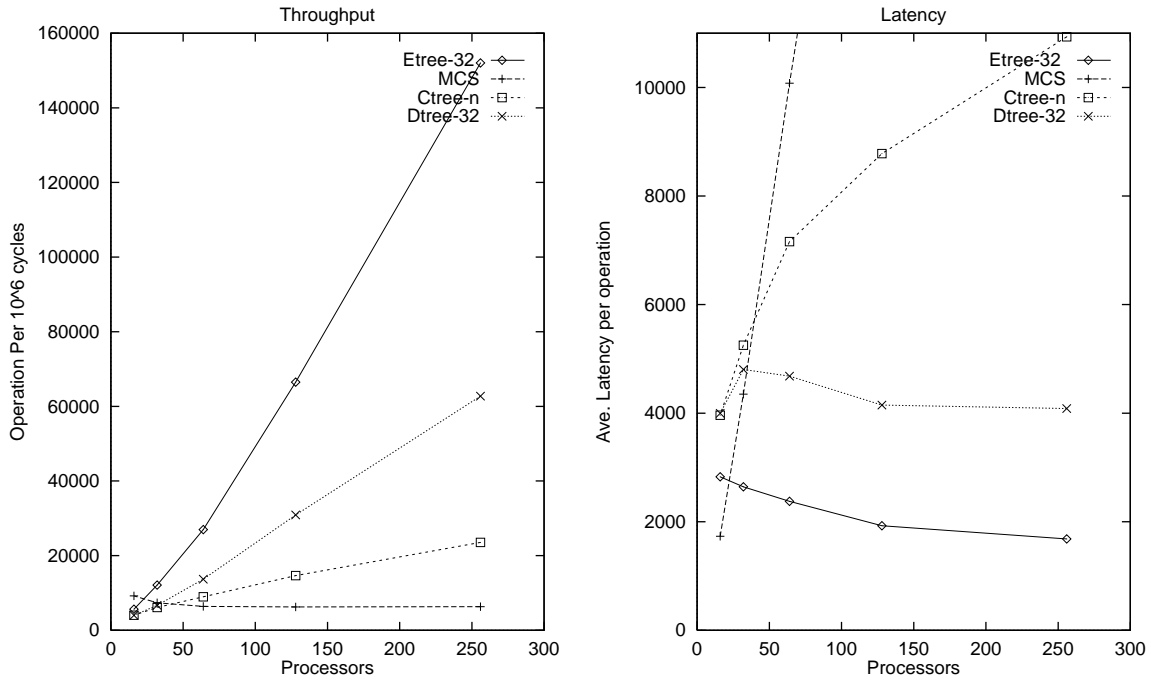


Figure 7: Produce-Consume: Throughput and Latency with `Workload= 0`

CTree A *Fetch&Inc* using an optimal width software combining tree following the protocol of Goodman et al. [10], modified according to [11]. The tree’s response time is logarithmic in the maximal number of processors. Optimal width means that when n processors participate in the simulation, a tree of width $n/2$ will be used [11].

DTree A Diffracting Tree of width 32, using the optimized parameters of [24], whose response time is logarithmic in $w = 32$ which is smaller than the maximal number of processors. The prism `sizes` were 8,4,2,2 and 1 for levels 1, . . . , 5 respectively. The `spin` is equal to 32,16,8,4 and 2 for balancers at depths 0,1,2,3,4 and 5 respectively.

and compared it to:

ETree A POOL[32] elimination tree based pool, whose response time is logarithmic in $w = 32$ which is smaller than the maximal number of processors. This size was chosen based on empirical testing. The root node and its children contain two prisms of size 32 and 8 for the root and 16 and 4 its children. The nodes at depths 3,4 and 5 have a single prism of size 2,1, and 1 respectively. The `spin` is equal to 32,16,8,4 and 2 for balancers at depths 0,1,2,3,4 and 5 respectively.

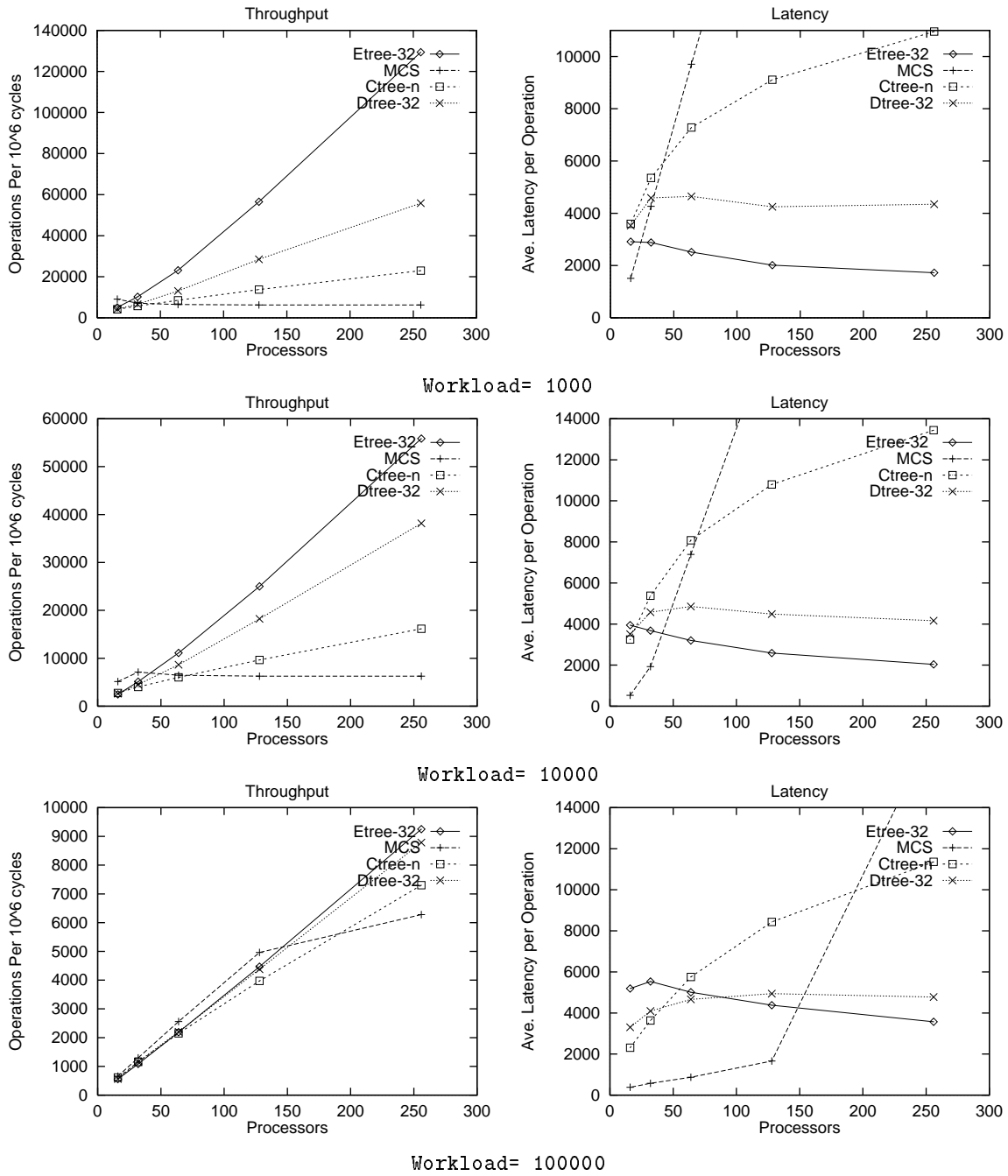


Figure 8: Produce-Consume: Throughput and Latency with Workload > 0

From Figure 7 we learn that under high loads diffracting and elimination trees provide the most scalable high load performance. However, as observed by Shavit and Zemach [24], as the level of

concurrency increases, while the diffracting tree manages only to keep the average latency constant, the average latency in the elimination tree continues to *decrease* due to the increased numbers of successful eliminating collisions taking place on the top levels of tree. The effect on the throughput is an up to 2.5 times increase in requests that are answered by the elimination tree! The fraction of eliminated tokens at the root varies between 44.7% when only 16 processors are participating and up to 49.7% for 256 processors. In fact, as can be seen from Table 1, most enqueue/dequeue requests never reach the lower level balancers, and the expected number of balancers traversed (including the pool at the leaf) for 16 processors is 3.14 nodes (38.9% of the request access the leaf pools) and for 256 processors 2.082 (only 8.95% of the request eventually access the pools at the leaves). As seen in Figure 7, at such high levels of concurrency the elimination tree is almost as fast as the MCS-queue-lock is when there are just a few processes.

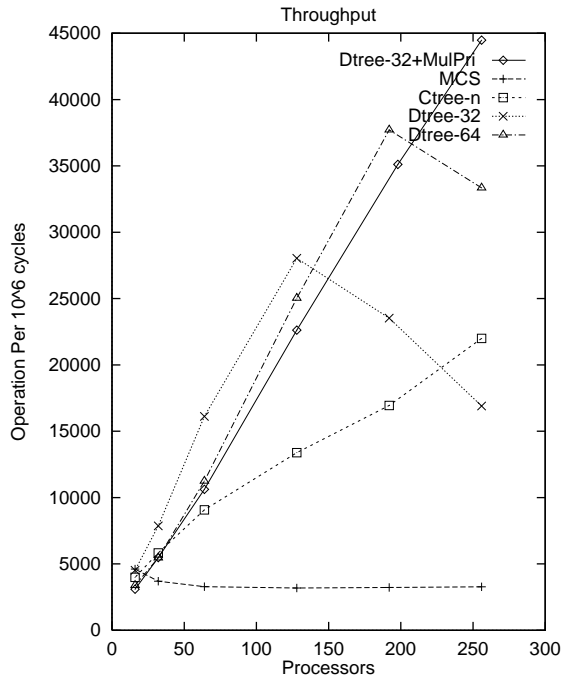
In Figure 8 we compared the various methods as access patterns become more sparse. The MCS lock outperforms all others when the number of processes is small, and unlike in the high load case of Figure 8, even with a high number of processes the elimination tree cannot match its low latencies because of the low levels of elimination on the root balancer. As the chances of combining, diffraction and elimination drop, the depth of the structures comes more into play. For 256 processes the optimal combining tree requires $2 \log n = 16$ node traversals (up and down the tree), while the optimal width 32 diffracting and elimination trees have depth 5 and thus require far fewer operations. It follows that the elimination and diffracting tree performance graphs converge, and at sufficiently high levels of concurrency remain far better than the combining tree.

2.5.2 Counting Benchmark

Our new multi-layered prism approach is slightly more costly but scales better than the original single prism construction of Shavit and Zemach [24], since it increases the likelihood of successful collisions. This conforms with the steady-state modeling of diffracting trees by Shavit, Upfal, and Zemach [25]. As can be seen from Figure 9, when running a benchmark of *fetch&increment* operations where no *eliminating collisions* can occur, the DTREE[32] and DTREE[64] with original

	16 procs	256 procs
level 0	44.7%	49.8%
level 1	24%	49.1%
level 2	5.8%	45.2%
level 3	1.9%	32.9%
level 4	0%	6.8%

Table 1: Fraction of Tokens Eliminated Per Tree Level



```

repeat
    fetch_and_inc();
until 10^6 cycles elapsed

```

Figure 9: Counting Benchmark

single `Prism` balancers outperform a `DTREE`[32] with our new multi-layered balancers in almost all the levels of concurrency which could be incurred in the 256 processor produce-consume benchmark (on average each `DTREE`[32] has 128 or so concurrent enqueues). However, unlike our the multi-layered balancer constructions, they do not continue to scale well at higher levels of concurrency.

2.5.3 Response Time Benchmark

We compared elimination trees to the randomized method of Rudolph, Silvkin-Allalouf, and Upfal (RSU) [22], which we chose as a representative of the class of *load-balanced local pools* methods, which also include the randomized methods of Kotz and Ellis [13] (RSU is a refinement of this method), of Lüling and B. Monien [21] (this method is a refinement of RSU), and the job-stealing method of Blumofe and Leiserson [7]. We also did not compare to Manber’s deterministic method [16] as Kotz and Ellis [13] have shown empirically that the randomized methods tend to give better overall performance. One should keep in mind that there are various situations in which any one

of these techniques outperforms all the others and vice versa.

The RSU scheme is surprisingly simple:

RSU A processor enqueues tasks in its private task queue. Before dequeuing a task, every processor flips a coin and executes a *load balancing* procedure with probability $1/l$ where l is the size of its private task queue. Load balancing is achieved by first choosing a random processor and then moving tasks from the longer task queue to the smaller so as to equalize their sizes.

We note that under high loads, and especially in applications such as job-distribution where each process performs both enqueues and dequeues, these methods are by far superior to elimination trees and all other presented methods. (The 10-queens benchmark in the lefthand side of Figures 11 and 10 is a lesser example of RSU's performance. Initially one processor, generates 10 tasks of depth 1 simultaneously. Each one of n processor repeatedly dequeues a task and if the task's depth is smaller than 3 it waits $work = 8000$ cycles and enqueue 10 new tasks of depth increased by one.) However, as we know from theoretical analysis, their drawback is the rather poor $\Theta(n)$ expected latency when there are sparse access patterns by producers and consumers that are trying to pass information from one to the other, as could happen say, in an application coordinating sensors and actuators.

The righthand side of Figures 11 and 10 show the results of an experiment attempting to evaluate (in a synthetic setting of course) how much this actually hampers performance, by measuring the average latency incurred by a dequeue operation trying to find an element to return. We do so by running a 256 processor machine in which $n/2$ processors are enqueueers and $n/2$ are dequeuers where n varies between 2 and 256. Each one of the enqueueing processors repeatedly enqueues an element in the pool and waits until the element has been dequeued by some dequeuing process. Each time we measured the time elapsed between the beginning of the benchmark until 2560 elements were dequeued, and normalized by the number of dequeue operations per process. Note that because of the way it is constructed, there is no real pipelining of enqueue operations, and this benchmark does not generate the high work-load of the produce-consume benchmark for large numbers of participants.

As can be seen, RSU does indeed have a drawback since it is almost 100 times slower than the queue-lock and 30 times slower than an elimination tree for sparse access patterns. This is mostly due to the fact that the elimination tree even without eliminating collisions will direct tokens and anti-tokens to the same local piles within $O(\log w)$ steps. RSU reaches a crossover point when about a quarter of all local piles are being enqueued into. In summary, elimination trees seem to offer a reasonable middle-of-the-way response time over all ranges of concurrency.

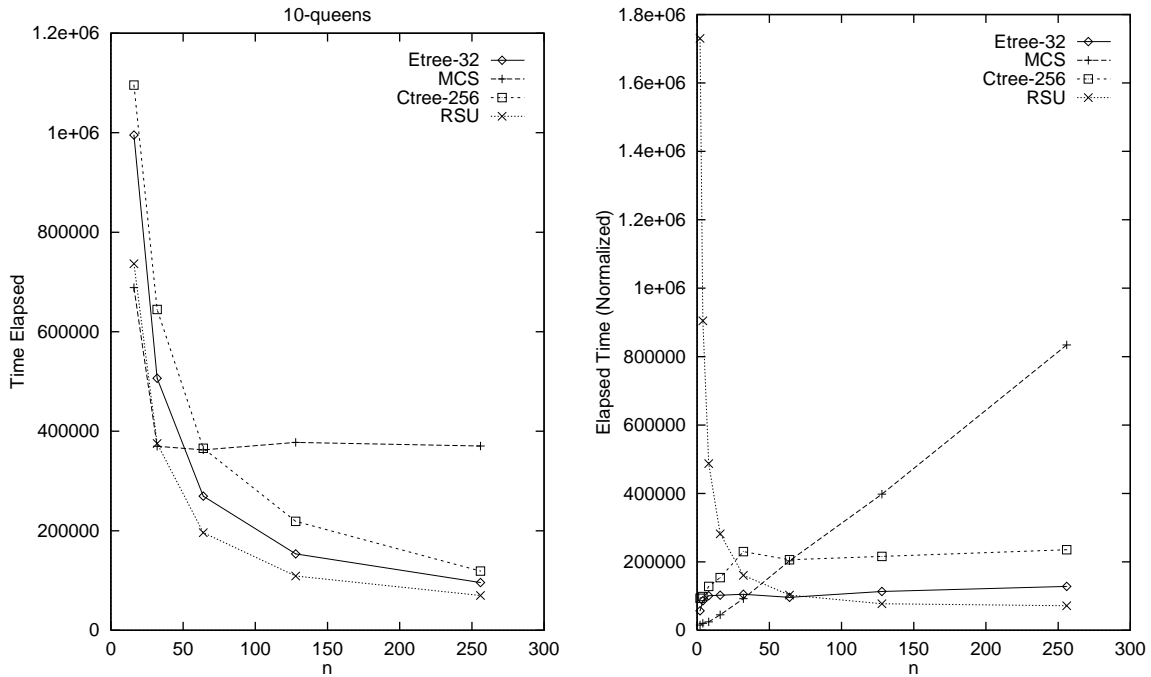


Figure 10: 10-Queens and Response Time Graphs

```

Initialization
  produce one instance with depth=0
repeat
  instance = consume();
  wait 8000 cycles;

  if instance's depth < 3 then
    produce 10 instances with depth greater by 1
until all instances have been consumed

producer:
  repeat
    produce(val);
    wait until the element is consumed;
    until a total of 2560 elements are consumed

consumer:
  repeat
    consume()
    until a total of 2560 elements are consumed

```

Figure 11: 10-Queens and Response Time Benchmarks

3 Stack-like Pools

Many applications in the literature that benefit by keeping elements in LIFO order would perform just as well if LIFO would be kept among all but a small fraction of operations. LIFO-based scheduling will not only eliminate in many cases excessive task creation, but it will also prevent

processors from attempting to dequeue and execute a task which depends on the results of other tasks [26]. Blumofe and Leiserson [7] provide a scheduler based on a randomized distributed pool having stack-like behavior on the level of local pools. We present here a construction of a pool that globally behaves like a stack. Our construction is based on the use of an elimination tree to create a single counter that can be both incremented and decremented concurrently, and can thus serve as high bandwidth pointer to the head of the stack.

3.1 Increment-Decrement Counting Trees

We define a new type of balancer, the *gap elimination balancer*, that allows both tokens and anti-tokens as inputs, and balances the “difference” between them (the surplus of tokens over anti-tokens) on its output wires. We use *gap elimination balancers* to construct counting trees that allow both increments and decrements. It has recently been shown by two independent teams, Busch and Mavronicolas [6] and Aiello, Herlihy, Shavit, and Tuitou [5] that the increment/decrement properties we describe hold for counting networks in general, not only for trees.

A *gap elimination balancer* is a *elimination balancer* that in addition to the Quiescence and Pairing property must satisfy the additional requirement that:

Gap Step Property In any quiescent state $0 \leq (y_0 - \bar{y}_0) - (y_1 - \bar{y}_1) \leq 1$.

In other words, any surplus of tokens over anti-tokens on the balancers output wires is distributed so that there is a gap of no more than one token on wire 0 relative to wire 1 in any quiescent state. Clearly, the gap step property implies the *pool balancing* property on the balancer’s output wires.

Claim 3.1 *Every gap elimination balancer satisfies the pool balancing property.*

We design $\text{INCDEC COUNTER}[w]$ as a *counting tree* [24] (a special case of the structure with regular token routing balancers replaced by token/anti-token routing *gap elimination balancers*). For w a power of two, $\text{INCDEC COUNTER}[2k]$ is just a root gap balancer connecting to two $\text{INCDEC COUNTER}[k]$ trees with the output wires y_0, y_1, \dots, y_{k-1} of the tree hanging from wire “0” re-designated as the even output wires $y_0, y_2, \dots, y_{2k-2}$ of $\text{INCDEC COUNTER}[2k]$, and the wires of the tree extending from the root’s “1” output wire re-designated as the odd output wires $y_1, y_3, \dots, y_{2k-1}$.

Lemma 3.2 *The $\text{INCDEC COUNTER}[w]$ tree constructed from gap elimination balancers has the gap step property on its output wires, that is, in any quiescent state:*

$$0 \leq (y_i - \bar{y}_i) - (y_j - \bar{y}_j) \leq 1$$

for any $i < j$.

Proof: We use that fact that the layout of the INCDECOUNTER is identical to that of a counting-tree [24], in order to show that if for some execution the INCDECOUNTER reaches a quiescent state which does not satisfies the gap step property, then there is an execution of the counting tree in which the step property is violated too. This is a contradiction to Theorem 5.5 of [24]. Let T^g be an INCDECOUNTER constructed from gap balancers g , and let T^b be the isomorphic counting tree which is the result of replacing every gap balancer g in INCDECOUNTER by a regular balancer b . Given an execution history h^g of T^g , for every gap balancer g , let h_x^g be the gap between tokens and anti-tokens on g 's input wire x , and let h_0^g and h_1^g be the gap at each of g 's output wires y_0 and y_1 . Define h_x^b , h_0^b , and h_1^b for h^b of T^b analogously.

Assume that for some execution history h^g of T^g , the gap step property is violated in a quiescent state. Assume first that the total difference between the number of tokens and anti-token accessing T^g is some non-negative number G . Let h^b be an execution of T^b in which G tokens access the tree T^b . By a simple inductive argument using on the depth of the trees, one can show that for every gap balancer g in T^g and its matching balancer b in T^b , the following holds: $h_x^g = h_x^b \wedge h_0^g = h_0^b \wedge h_1^g = h_1^b$. Consequently, it follows that:

Claim 3.3 *If for some execution history h^g of T^g , where G is non-negative, the gap step property is violated in a quiescent state, then it is violated also for the matching history h^b of T^b .*

Assume now that for h^g , the difference G between the total number of tokens and anti-tokens is negative. Let k be the smallest number such that $2^d * k + G \geq 0$ where d is the depth of the tree. Let $h1^g$ be an execution of T^g , in which after the completion of h^g , $2^d * k$ tokens were pushed through T^g . Using a simple inductive argument on the depth of the tree, one can show that for every node g of depth d' in T^g , $h_x^g + k * 2^{d-d'} = h1_x^g$. Therefore, since k tokens will have been equally added to all the exits of T^g , the gap step property will be violated in $h1^g$ too. Since in $h1^g$, the gap at the entrance of the tree is non-negative, the claim follows by applying Claim 3.3. ■

A *Stack-like Pool* is constructed, as with the pool data structure, by placing sequentially accessed “local stacks” at the leaves of a INCDECOUNTER[w] tree. The following theorem is a corollary of Theorem 2.2 and Claim 3.1:

Theorem 3.4 *The stack-like pool construction is a correct pool implementation.*

The next theorem, which explicates the the LIFOish behaviour of stack-like pool is a direct corollary from from step property of Lemma 3.2, and is left to the interested reader.

Theorem 3.5 *In any sequential execution the stack-like pool provides a last-in-first-out order on enqueues and dequeues.*

In Section 3.5 we present empirical evidence that suggests that even though the stack-like pool is not linearizable [12] to a sequential stack, it is linearizable in executions without severe timing anomalies, hence our use of the term “stack-like.”

3.2 Implementing the Gap Elimination Balancer

One can modify the pool elimination balancer construction from the former section so that it satisfies the gap step property. This is done by replacing Part 2 of the code in Figure 4 with the following:

```

AcquireLock(b->Lock);
if compare_and_swap(Location[mypid],<b,my_type>, <0,EMPTY>) then
    i:= b->INCDECToggle;
    b->INCDECToggle := Not(i);
    ReleaseLock(b->Tokens[mytype]);
    return b->OutputWire[i];
else
    ReleaseLock(b->Lock);
    if Location[mypid]= <0,DIFFRACTED> return (b->OutputWire[0])
    else return ELIMINATED

```

Instead of accessing two different toggle bits, both tokens and anti-tokens use the same toggle bit `INCDECToggle`. If a token does not collide in the prisms, it toggles `INCDECToggle` and chooses an output wire according to the old value of the bit. An anti-token similarly toggles `INCDECToggle`, but it chooses an output wire according to the *new* value of `INCDECToggle` (using machine language notation, tokens perform a `fetch&complement` and anti-tokens a `complement&fetch`). On an intuitive level, this combination causes an anti-token to “trace” the last inserted token.

3.3 Correctness Proof of Gap Balancer Implementation

In order to prove the correctness of our gap balancer implementation we first show that all the tokens that have accessed the toggle bit satisfy the gap step property. As before, let t_i and \bar{t}_i be the number of toggling tokens and anti-tokens exiting the balancer on wire i .

Lemma 3.6 *In any quiescent state $0 \leq (t_0 - \bar{t}_0) - (t_1 - \bar{t}_1) \leq 1$.*

Proof: The proof is by induction on the length of the history h of accesses to the toggle bit. If history h contains only token transitions or only anti-token transitions then the property holds

trivially. If h consists of transitions of both token types, there must be at least one token transition τ and one anti-token access $\bar{\tau}$ which followed one other in the history. Let us define h' to be the history h without τ and $\bar{\tau}$. Since following τ and $\bar{\tau}$ the `INCDECtoggle` bit returns to the same state it was before these transitions accessed it, h' is a possible history of the access to `INCDECtoggle` and by induction hypothesis satisfies the step property. Now, since both τ and $\bar{\tau}$ leave on the same output wire, h also satisfies the balancing property. ■

Since the elimination protocols are identical in both the pool and gap elimination balancer implementations, the proof of the following 3 lemmas are identical to the proofs of Lemmas 2.3,2.4 and 2.5 respectively, and are therefore omitted.

Lemma 3.7 *For every process p , if in a given state `Location`[p]= $\langle b, * \rangle$, then p is executing `TokenTraverse` on balancer b .*

Lemma 3.8 *Every token traversing a balancer b can be diffracted or eliminated by at most one other token.*

Lemma 3.9 *A toggling, eliminating, or diffracting token T_p cannot be eliminated or diffracted by some other token T_q .*

We can now conclude the correctness proof of our gap balancer implementation:

Theorem 3.10 *The gap eliminating balancer implementation satisfies the gap step property.*

Proof: Using the same notations as in the correctness proof of the pool balancer, we know from Lemmas 3.7,3.8 and 3.9 that $\bar{e} = e$, $d_0 = d_1$ and $\bar{d}_0 = \bar{d}_1$. Therefore $(t_0 - \bar{t}_0) - (t_1 - \bar{t}_1) = ((t_0 + d_0) - (\bar{t}_0 + \bar{d}_0)) - ((t_1 + d_1) - (\bar{t}_1 + \bar{d}_1))$. Since, $y_0 = t_0 + d_0, y_1 = t_1 + d_1, \bar{y}_0 = \bar{t}_0 + \bar{d}_0$ and $\bar{y}_1 = \bar{t}_1 + \bar{d}_1$ we may conclude that $0 \leq (y_0 - \bar{y}_0) - (y_1 - \bar{y}_1) \leq 1$. ■

3.4 Performance of the Stack-like Pool

We tested the performance of the stack-like pool for the produce-consume benchmark from Section 2. We implemented a `INCDEC COUNTER`[32] with prism sizes and spin times as in the `POOL`[32]. In Figure 12 we present the result of a comparison between an `INCDEC COUNTER`[32] based stack-like pool and a `POOL`[32] in the producer-consumer benchmark under high load `Workload = 0`. As can be seen, though tokens are accessing a shared toggle bit instead of two separate ones, high elimination rates on the prisms allow the efficiency of the stack-like pool to fall from that of the `POOL`[32] only slightly.

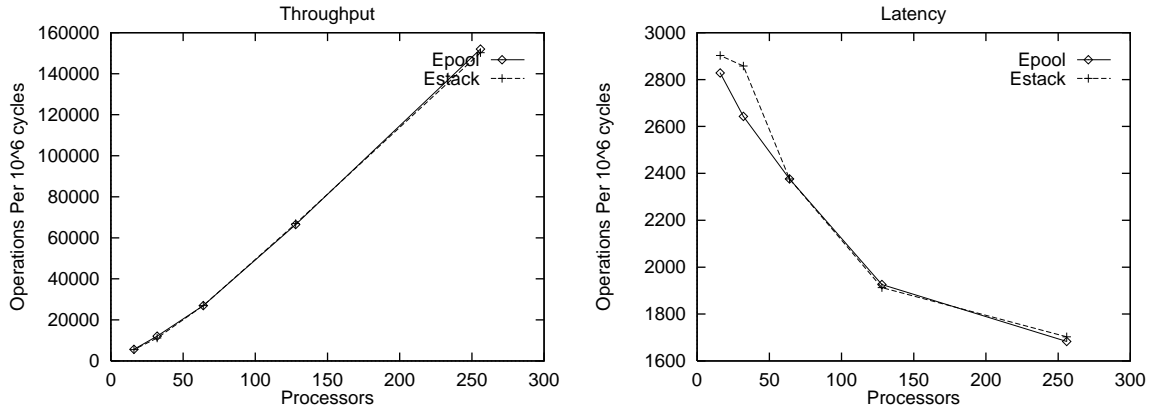


Figure 12: Comparison between a Pool and a Stack-like Pool

3.5 Almost Linearizability

Herlihy and Wing’s Linearizability [12] is a consistency condition that specifies the allowable concurrent behaviours of an object by way of a mapping to a sequentially specified object whose behaviours are easy to state. A linearization mapping exists if one can pick a point within the execution interval of every concurrent operation so that the collection of operations executed sequentially according to the order among these points, meets the sequential object specification. We present some empirical evidence that suggests that even though the stack-like pool is not always linearizable to a sequential stack, it behaves very much like one.

Given a stack-like pool implementation, let $E(e)$ and $D(e)$ respectively denote an enqueue operation of e and a dequeue operation returning e . Let \rightarrow be the real time order between the operations ($OP_1 \rightarrow OP_2$ iff OP_1 has terminated before OP_2 has started). We say that the operation $D(x)$ in an execution e is *not linearizable* if there are $E(y), E(x)$ such that $E(x) \rightarrow E(y) \rightarrow D(x)$ and either $D(y)$ does not exist in e or $D(y)$ exists in e and $E(x) \rightarrow E(y) \rightarrow D(x) \rightarrow D(y)$. A stack-like pool implementation is *linearizable* [12] if it ensures that every execution does not contain a dequeue operation that is not linearizable .

Our elimination tree based `INCDECOUNTER`[w] is easily shown not to be linearizable to a sequential counter with increments and decrements. However, we present in Figure 13 empirical evidence suggesting that scenarios in which the linearizability of our stack-like pool is violated require extreme timing anomalies that one might argue are not likely to occur frequently. We ran the producer-consumer benchmark where each processor, after traversing a balancer node, waits a random number of cycles between 0 and $W = 0, 1000, 10000, 100000$ until 2000 dequeue operations are executed. The graph presented plots the fraction (%) of dequeue operations that are not linearizable. Note that for tightly synchronized executions ($w = 0$), our stack-like implementation

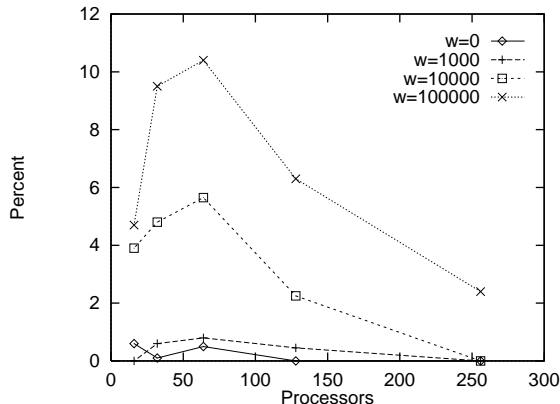


Figure 13: Produce-Consume: Percentage of Dequeue operations that are not linearizable.

is linearizable to a stack at almost all levels of concurrency.

4 Conclusions and Further Research

Our paper introduces the notion of “anti-tokens” to allow decrement operations on a counting-tree [24]. Two independent research teams, Busch and Mavronicolas [6] and Aiello, Herlihy, Shavit, and Touitou [5], have recently extended our proofs to show that counting networks [4] in general, not only trees, work with anti-tokens (Busch and Mavronicolas [6] show this also for multi-balancers [1, 14], that is, balancers with multiple inputs and output wires).

In summary, *elimination trees* represent a new class of concurrent algorithms that we hope will prove an effective alternative to existing solutions for produce/consume coordination problems. This paper presents shared memory implementations of elimination trees, and uses them for constructing pools and and stack-like pools.

There is clearly room for experimentation on real machines and networks. Given the hardware *fetch-and-complement* operation to be added to the Alewife machine’s Sparcle chip’s set of colored load/store operations [19], one will be able to implement a shared memory elimination-tree in a wait-free manner, that is, without any locks. Our plan is to test such “hardware supported” elimination-tree performance. We also plan to develop better measures and methods for setting the tree parameters such as prism `size` and balancer `spin`, and are currently developing message passing versions of our algorithms.

5 Acknowledgements

We would like to thank Yehuda Afek, Bill Aiello, Maurice Herlihy, and Asaph Zemach for their many helpful comments.

References

- [1] E. Aharonson and H. Attiya. Counting networks with arbitrary fan out. In *Proceedings of the 3rd Symposium on Discrete Algorithms*, Orlando, Florida, January 1992. Also: Technical Report 679, The Technion, June 1991.
- [2] T.E. Anderson. The Performance of Spin Lock Alternatives for Shared-Memory Multiprocessors. *IEEE Transactions on Parallel and Distributed Systems*, 1(1):6–16, January 1990.
- [3] A. Agarwal et al. The MIT Alewife Machine: A Large-Scale Distributed-Memory Multiprocessor. In *Proceedings of Workshop on Scalable Shared Memory Multiprocessors*. Kluwer Academic Publishers, 1991. An extended version of this paper has been submitted for publication, and appears as MIT/LCS Memo TM-454, 1991.
- [4] J. Aspnes, M.P. Herlihy, and N. Shavit. Counting Networks. *Journal of the ACM*, Vol. 41, No. 5 (September 1994), pp. 1020-1048.
- [5] W. Aiello, M. Herlihy, N. Shavit and D. Touitou. Inc/Dec Counting Networks. Manuscript, December 1995.
- [6] C. Busch and M. Mavronicolas. The Strength of Counting Networks. Proceedings of the 15th Annual ACM Symposium on Principles of Distributed Computing, to appear, May 1996.
- [7] R.D. Blumofe, and C.E. Leiserson. Sheduling Multithreaded Computations by Work Stealing. In *Proceeding of the 35th Symposium on Foundations of Computer Science*, pages 365-368, November 1994.
- [8] E.A. Brewer, C.N. Dellarocas, A. Colbrook and W.E. Weihl. PROTEUS: A High-Performance Parallel-Architecture Simulator. MIT Technical Report /MIT/LCS/TR-561, September 1991.
- [9] D. Chaiken. Cache Coherence Protocols for Large-Scale Multiprocessors. S.M. thesis, Massachusetts Institute of Technology, Laboratory for Computer Science Technical Report MIT/LCS/TR-489, September 1990.

- [10] J.R. Goodman, M.K. Vernon, and P.J. Woest. Efficient Synchronization Primitives for Large-Scale Cache-Coherent multiprocessors. In *Proceedings of the 3rd ASPLOS*, pages 64–75. ACM, April 1989.
- [11] M. Herlihy, B.H. Lim and N. Shavit. Low Contention Load Balancing on Large Scale Multiprocessors. *Proceedings of the 3rd Annual ASM Symposium on Parallel Algorithms and Architectures*, July 1992, San Diego, CA. Full version available as a DEC TR.
- [12] M. Herlihy and J.M. Wing. Linearizability: A correctness condition for concurrent objects. In *ACM Transaction on Programming Languages and Systems*, 12(3), pages 463-492, July 1991.
- [13] D. Kotz and C. S. Ellis. Evaluation of Concurrent Pools. In *Proceedings of the International Conference on Distributed Computing Systems*, pages 378-385, June 1989.
- [14] E.W. Felten, A. LaMarca, R. Ladner. Building Counting Networks from Larger Balancers. University of Washington T.R. #93-04-09.
- [15] J.M. Mellor-Crummey and M.L. Scott Synchronization without Contention. In *Proceedings of the 4th International Conference on Architecture Support for Programming Languages and Operating Systems*, April 1991.
- [16] Udi Manber. On maintaining dynamic information in a concurrent environment *SIAM J. Computing* 15(4), pages 1130–1142, November 1986.
- [17] G.H. Pfister and A. Norton. ‘Hot Spot’ contention and combining in multistage interconnection networks. *IEEE Transactions on Computers*, C-34(11):933–938, November 1985.
- [18] D. Gawlick. Processing ‘hot spots’ in high performance systems. In *Proceedings COMPCON’85*, 1985.
- [19] J. Kubiawicz. Personal communication (February 1995).
- [20] N.A. Lynch and M.R. Tuttle. Hierarchical Correctness Proofs for Distributed Algorithms. In *Sixth ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing*, August 1987, pp. 137–151. Full version available as MIT Technical Report MIT/LCS/TR–387.
- [21] R. Lüling, and B. Monien. A Dynamic Distributed Load Balancing Algorithm with Provable Good Performance. In *Proceedings of the 5rd ACM Symposium on Parallel Algorithms and Architectures*, pages 164-173, June 1993.
- [22] L. Rudolph, M. Slivkin, and E. Upfal. A Simple Load Balancing Scheme for Task Allocation in Parallel Machines. In *Proceedings of the 3rd ACM Symposium on Parallel Algorithms and Architectures*, pages 237–245, July 1991.

- [23] N. Shavit, and D. Touitou. Elimination Trees and the Construction of Pools and Stack. In *Proceedings of the 7th Annual Symposium on Parallel Algorithms and Architectures (SPAA)*, pages 54-63, July 1995.
- [24] N. Shavit and A. Zemach. Diffracting Trees. In *Proceedings of the Annual Symposium on Parallel Algorithms and Architectures (SPAA)*, June 1994.
- [25] N. Shavit, E. Upfal, and A. Zemach. A Steady-State Analysis of Diffracting Trees. Unpublished manuscript. Tel-Aviv University. October 1995.
- [26] K. Taura, S. Matsuoka, and A. Yonezawa. An Efficient Implementation Scheme of Concurrent Object-Oriented Languages on Stock Multicomputers. In *Proceedings of the 4th Symposium on Principles and Practice of Parallel Programming*, pages 218–228, May 1993.