

# A Key-based Adaptive Transactional Memory Executor\*

Tongxin Bai    Xipeng Shen<sup>‡</sup>    Chengliang Zhang  
William N. Scherer III<sup>†</sup>    Chen Ding    Michael L. Scott

Computer Science Department, University of Rochester

<sup>‡</sup>Computer Science Department, The College of William and Mary

<sup>†</sup>Computer Science Department, Rice University

{bai,zhangchl,cding,scott}@cs.rochester.edu

xshen@cs.wm.edu

bill.scherer@cs.rice.edu

URCS TR 909

December 2006

## Abstract

Software transactional memory systems enable a programmer to easily write concurrent data structures such as lists, trees, hash tables, and graphs, where non-conflicting operations proceed in parallel. Many of these structures take the abstract form of a *dictionary*, in which each transaction is associated with a search key. By regrouping transactions based on their keys, one may improve locality and reduce conflicts among parallel transactions.

In this paper, we present an executor that partitions transactions among available processors. Our key-based adaptive partitioning monitors incoming transactions, estimates the probability distribution of their keys, and adaptively determines the (usually nonuniform) partitions. By comparing the adaptive partitioning with uniform partitioning and round-robin keyless partitioning on a 16-processor SunFire 6800 machine, we demonstrate that key-based adaptive partitioning significantly improves the throughput of fine-grained parallel operations on concurrent data structures.

---

\*This research was supported in part by NSF grants CCR-0204344, CNS-0411127, CNS-0509270, and CNS-0615139; financial and equipment support from Sun Microsystems Laboratories; and financial support from Intel Corporation.

# 1 Introduction

Task-level parallelism is becoming increasingly important for general-purpose programs as modern desktop and server machines are equipped with multicore and multithreaded processors. Though essential to increased throughput on such machines, concurrent programs are notoriously difficult to write, from both performance and correctness perspectives. One performance issue with highly concurrent multithreaded programs is the overhead of thread creation and destruction. This overhead becomes particularly apparent when a server program is written in a thread-per-request model with a large number of very short requests. The overhead problem may be addressed by arranging for a pool of threads to share a “bag of tasks,” as implemented, for example, by the Java 5 `java.util.concurrent Executor` framework.

Although a thread pool helps to efficiently execute short concurrent tasks, it shares with the thread-per-request model the usual pitfalls of lock-based synchronization, including deadlock, priority inversion, and “wrong lock” bugs. *Transactional memory* [9] addresses these pitfalls by allowing programmers simply to label blocks of code *atomic*, and to count on the underlying system to execute them in parallel whenever possible.

For applications consisting of many small tasks, it seems natural to combine executors and transactional memory. The former minimizes the overhead of threads; the latter ensures atomicity, consistency, and mutual isolation. In this paper, we present a *key-based* transactional memory executor that takes software transactions as inputs and schedules them across a pool of worker threads. The main contribution lies in the scheduling policy, which uses search keys or other transaction parameters to increase memory locality, reduce transaction conflicts, and adaptively balance load. We focus in particular on *dictionary* operations, in which the likelihood that two transactions will access similar data correlates with the numerical proximity of keys.

To improve temporal and spatial locality, we can execute transactions with similar keys on the same processor, so successive transactions find data in the cache, reducing access latency and load on the processor-memory interconnect. On a machine with multicore processors and shared chip-level caches, locality among transactions that run on the same chip (if not on the same core) may still improve performance by reducing the number of loads that go to memory.

Transaction conflicts happen when concurrent transactions violate Bernstein’s condition [1, 2]: they access the same data, and at least one of them is a writer. To reduce such conflicts, we can avoid concurrent execution of transactions with similar keys by scheduling them to the same work thread. Fewer conflicts mean less time wasted on aborted transactions, conflict detection and contention management.

To balance load among worker threads, we first dynamically sample keys of transactions. Then we estimate the distribution of keys and adaptively partition keys across the available worker threads.

The remainder of this paper is organized as follows: Section 2 presents background information on transactional memory and the Java `Executor` framework. Section 3 describes our executor model. Section 4 describes the adaptive partitioning of keys. Section 5.1 discusses our implementation. We empirically validate the feasibility and benefits of key-based adaptation in Section 5.4; we discuss related work in Section 6 and conclude in Section 7.

## 2 Background

### 2.1 Transactional Memory

Transactional memory (TM), originally proposed by Herlihy and Moss [9], borrows the basic execution model of database systems. With transactional memory, a programmer specifies atomic regions and the underlying system guarantees that these are executed atomically (in their entirety or not at all), consistently (in a way that preserves system invariants), and in isolation (with no intermediate results visible to other transactions). The overall effect is indistinguishable from an execution in which the transactions occur one at a time, in some order. If two concurrently executing transactions cannot be so serialized, the system typically forces one to wait or to abort and restart automatically. The optimistic expectation is that if no conflicts occur during the course of the two transactions, we can achieve full parallelism.

Transactional memory can be implemented in hardware, in software, or in some combination of the two. Our executor mechanism can be used with any of these. The experiments reported in Section 5.4 employ the DSTM (dynamic software TM) system of Herlihy et al. [8].

A good TM implementation will execute concurrent, nonconflicting transactions in parallel whenever possible. It is fundamentally constrained, however, by the characteristics of the workload. Consider, for example, a two-processor system and an application whose transactions modify data object *A* or *B* with roughly equal probability. If both processors attempt to run an *A* transaction at the same time, one will almost surely have to abort and retry. On the other hand, if we are able to write the application in such a way that processor 1 performs *A* transactions most of the time and processor 2 performs *B* transactions most of the time, then they will usually be able to run in parallel. The motivating observation behind our work is that for important classes of applications we can often tell, based on the dynamic inputs of a transaction, which data it is likely to access. We can then take steps to dispatch transactions with similar access patterns to the same application thread. In this paper we look in particular at applications whose transactions access dictionary data structures, and we use the dictionary key to drive a prediction of data access pattern.

In addition to reducing conflict, a key-based transaction executor can be expected to lead to better temporal locality: an object is more likely to be in cache if the previous transaction to access it was executed by the same thread. If dictionary entries with similar keys are near each other in memory (as they are in many implementations) then a key-based executor may improve spatial locality as well.

### 2.2 Java Executor

For servers that execute many short tasks, a bag of tasks and thread pool serve to avoid unnecessary thread creation and destruction overhead. Each *worker* in the thread pool takes tasks (in our case, transactions) from the bag and executes them one at a time. If the bag is empty, a worker waits for an available task. Conversely, if the thread pool is empty, a task waits for an available worker.

In the `java.util.concurrent` package of JDK 1.5, the `Executor` framework manages task bags and thread pools. The `Executor` interface declares a method, `execute`, with the following form:

```
public interface Executor {
    public void execute (Runnable task);
}
```

In response to a call to `execute`, the `Executor` arranges for the specified task to be executed by a worker thread at some point in the future. Most programs use executors provided by the standard



transactions. More importantly for this work, the executor can reorder transactions to improve throughput.

**Parallel executors.** A single executor may be a scalability bottleneck. Parallel executors, as shown in Figure 1 (c), take tasks generated by producers and dispatch them to workers in parallel. If the number of executors is the same as the number of producers, the executor can be a part of the producer thread, so the data passing becomes thread local, as shown by the dotted lines in the Figure.

**Load balancing.** For maximum speedup, load on the worker threads should be balanced. Many factors affect the speed of a worker thread, including the size of the transactions, their temporal and spatial locality with respect to other transactions in the same worker, and the chance of conflicts with transactions in other workers. Many of these factors have been systematically studied for parallel tasks other than STM transactions. A common solution is *work stealing* [3], where a worker thread with an empty queue grabs work from other queues. An alternative method is for the executor to dynamically change the distribution of transactions.

For concurrent data structure updates, the size of transactions may not vary much. In such a case, an executor can control load balance by assigning the same number of transactions to each worker thread. When work is plentiful, the number transactions in each queue can be expected to be roughly equal. For transactions that have more complex sizes, some type of work stealing may be needed, although the overhead of the executor scheme may increase because the work queue is more complex to manage.

## 4 Key-based Scheduling

In preceding Sections we have implied that transactions on a dictionary data structure should be reordered based on dictionary keys. In actuality, it is valuable to distinguish between dictionary keys and “transaction keys”, which the executor uses for scheduling. Transaction keys are a more general concept, applicable to a wider variety of workloads. By creating a separate notion of transaction key we arrive at a two-phase scheduling scheme. The first phase maps transaction inputs, whatever they may be, into a linear key space. The second phase dispatches transactions based on their location in this space.

### 4.1 Mapping Transactions to Key Space

The central requirement for keys is that numerical proximity should correlate strongly (though not necessarily precisely) with data locality (and thus likelihood of conflict) among the corresponding transactions. One possible way to find a good mapping from inputs to keys is to profile an application on a representative set of inputs. Alternatively, users might provide appropriate mapping functions for each application. Given a priori knowledge, such manual specification may be both easier and more accurate than profiling. We use manual specification in our experiments.

In generating keys, we have two conflicting goals. On the one hand, we want our keys to be *accurate*: transactions with similar data access patterns should have similar keys. On the other hand, key generation should be *efficient*. Perfect accuracy could be obtained by fully executing a transaction and generating a Gödel number that encodes its data accesses, but this would obviously defeat the purpose of scheduling altogether! Instead, we seek the best estimate of a transaction’s accesses that we can obtain without incurring unreasonable overhead.

Some mappings may be trivial. In a stack, for example, we are more interested in the fact that every push and pop will begin by accessing the top-of-stack element than in the exact address of this element. So for this simple case, the hint we provide to the scheduler (the key) is constant for every transactional access to the same stack. This allows the executor to recognize that transactions will race for the same data and schedule them appropriately. This example may seem uninteresting, but it demonstrates how summaries of transactional data access patterns can be meaningfully used by a scheduler that requires no other information about the contents of transactions.

For the benchmark results shown in this paper, we have manually coded functions that generate keys for transactions. Although one would obviously prefer to generate them automatically, we focus for now on demonstrating the potential benefits to system throughput that can be obtained by managing locality. We believe that automatic key generation will prove to be an interesting challenge for future work.

## 4.2 Scheduling Given Keys

We have experimented with three schemes to schedule transactions, using the executor model of Figure 1 (c). The baseline scheme is a *round robin* scheduler that dispatches new transactions to the next task queue in cyclic order. The second scheme is a *key-based fixed scheduler* that addresses locality by dividing the key space into  $w$  equal-sized ranges, one for each of  $w$  workers. There is no guarantee, however, that a mapping function that captures data locality will lead to a uniform distribution of keys. The third scheme is a *key-based adaptive scheduler* that samples the input distribution and partitions the key space adaptively in order to balance load.

To simplify the explanation of the adaptive scheduler, we assume that a transaction is some operation on an integer, which we may treat as the transaction key. During the early part of program execution, the scheduler assigns transactions into worker queues according to a fixed partition. At the same time, it collects the distribution of key values. Once the number of transactions exceeds a predetermined *confidence threshold*, the scheduler switches to an adaptive partition in which the key ranges assigned to each worker are no longer of equal width, but contain roughly equal numbers of transactions.

More specifically, our executor employs a Probability Distribution-based partition (*PD-partition*) [19], which is similar to the extended *distributive partitioned sorting* of Janus and Lamagna [10]. It requires time linear in the number of samples. Figure 4.2 explains the *PD-partition* algorithm. Given an unknown distribution (a), the algorithm constructs a histogram (b) that counts the number of samples in ranges of equal width. It then adds the numbers in successive ranges to obtain cumulative probabilities (c) and uses these to obtain a piece-wise linear approximation of the cumulative distribution function (CDF) (d). Finally, by dividing the probability range into equal-size buckets and projecting down onto the  $x$  axis, it obtains the desired adaptive ranges (e).

The number of samples determines the accuracy of CDF estimation, which in turn determines the *PD-partition* balance. Shen and Ding show how to determine the sample size required for a guaranteed bound on accuracy by modeling the estimation as a multinomial proportion estimation problem [19]. In this work, we use a threshold of 10,000 samples, which guarantees with 95% confidence that the CDF is 99% accurate.

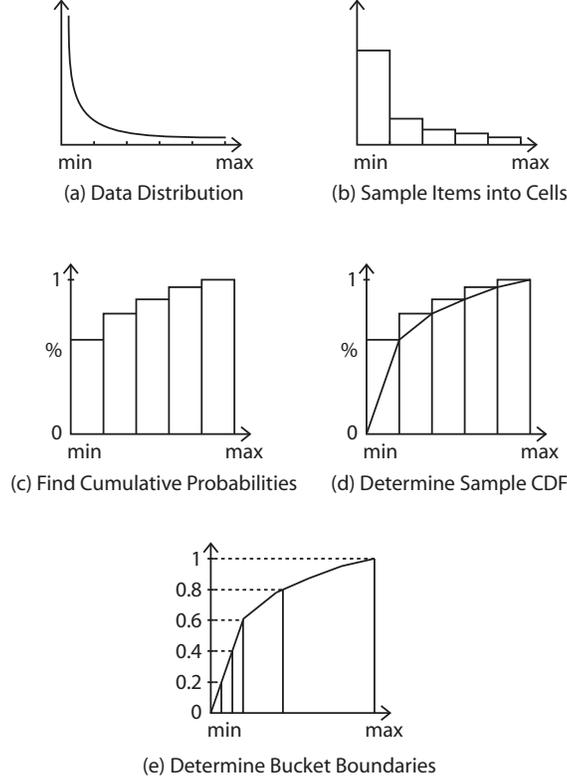


Figure 2: Adaptive calculation of key ranges.

## 5 Evaluation

### 5.1 Implementation

Our implementation is based on the Java dynamic software transactional memory (DSTM) system of Herlihy et al. [8]. We have implemented the parallel executor scheme of Figure 1(c) with two separate classes of threads: producer threads, which generate transactions (the producer part in Figure 1(c)) and then dispatch them (the executor part in Figure 1(c)); and worker threads, which execute transactions within the DSTM. These are connected by task queues, one for each worker. The queues are instances of `java.util.concurrent.ConcurrentLinkedQueue`, which implements the Michael & Scott concurrent queue [14].

The entire system is orchestrated by a test driver thread, which selects the designated benchmark, starts the producer threads, records the starting time, starts the worker threads, and stops the producer and worker threads after the test period (10 seconds in our experiments). After the test is stopped, the driver thread collects local statistics from the worker threads and reports the cumulative throughput.

A producer thread operates in a loop. It first checks with the test driver ensuring that it should continue. Then it generates the next transaction. In our experiments, we use four parallel producers (eight for the hash table benchmark to prevent worker threads being hungry). For efficiency we insert the parameters of a transaction rather than the transaction itself into the task queue.

A worker thread follows a simple regimen: check with the test driver for permission to continue, obtain the next transaction, execute it, retry in the case of a failure until successful, and then increment the local counter of complete transactions. We assume that there is no need to report

the result of a transaction back to the producer. The effect has already taken place in the shared data structure.

## 5.2 Microbenchmarks

We use the following four benchmarks, each of which performs insertions and deletions on concurrent data structures that implement an abstract dictionary with 16-bit search keys. We name the benchmarks by the type of data structure. To emphasize the impact of conflicts we do not include lookup operations. Although lookup operations would exhibit locality, they do not conflict with one another.

**Hash table:** Uses external chaining from an array of 30031 buckets (a prime number close to half the value range). The hash function is the hash key modulo the number of buckets. The benchmark uses the same number of inserts and deletes, so the load factor at stable state is around 1. A conflict occurs when the same bucket is modified by two transactions. We use the output of the hash function (not the dictionary key) as the value of the transaction key. With this particular hash function the dictionary key would also ensure spatial locality among transactions with similar keys, but this would not be the case for other, more “random” hash functions. Because hash table transactions are extremely fast, we use twice the usual number of producers to keep the workers busy.

**Red-black tree:** A balanced binary search tree from the original DSTM paper [8], modified to use 16-bit keys instead of 8-bit keys. The program uses the same number of inserts and deletes and has tens of thousands of tree nodes by the end of a ten-second run. The typical transaction therefore accesses about 16 tree nodes. Most nodes are at the bottom of the tree, as are most insertions, deletions, and rotations. Two operations conflict if they change the same part of the tree. We use the unmodified search key as the transaction key. Clearly two transactions with the same key will conflict. Other transactions may also conflict, if they attempt to perform rotations in the same part of the tree. The chance of conflict correlates with the numerical closeness of the keys, because close key values typically have a close common ancestor. The numerical closeness of the keys is also a good indicator of spatial locality, because the closer two leaves are, the more ancestors they are likely to share.

**Sorted linked list:** A singly-linked list with numerically increasing key values, also from the original DSTM paper [8]. A transaction, when inserting or deleting a number  $i$ , traverses the list until finding the first number no smaller than  $i$ . The benchmark starts with an empty set and interleaves the same number of random inserts and deletes. Because transactions on the linked list are slower than transactions on the tree (and conflicts are more common), the benchmark would need more than 10 seconds to reach stability. In our tests the typical run ends with a list of several hundred nodes.

To avoid unnecessary conflicts, the list code uses *early release* to “forget” its use of nodes after it has passed them by. This allows concurrent inserts and deletes in different parts of the list. A transaction modifies at most three nodes, and two transactions conflict if their write sets overlap or if they pass through the same part of the list at exactly the same time. As in the red-black tree, we use the unmodified search key as the transaction key. In this case, however, numerical closeness is a comparatively poor indicator of spatial locality: two transactions have overlapping write sets only if their keys are within two nodes of each other in the list.

The hash table and the red-black tree should benefit from the key-based executor because of better locality and lower conflicts, though it is difficult to predict the degree of improvement a priori,

given that the programs run in parallel over the DSTM and the Java virtual machine. Given the weaker connection between keys and locality, the benefit of key-based scheduling is harder to predict for the sorted list. Adaptive partitioning should prove valuable whenever the key distribution is highly nonuniform.

### 5.3 Data Collection

All results were obtained on a SunFire 6800, a cache-coherent symmetric multiprocessor with 16 1.2Ghz UltraSPARC III processors and 8MB of L2 cache for each processor. We tested in Sun’s Java 1.5.0.06, server version, using the HotSpot JVM augmented with a JSR 166 update `jar` file from Doug Lea’s web site [11]. The dynamic software transactional memory (DSTM) implementation came from Herlihy et al. [8]. We used the “Polka” contention manager to arbitrate among conflicting transactions. It combines randomized exponential back-off with a priority accumulation mechanism that favors transactions in which the implementation has already invested significant resources [18].

Memory management has a non-trivial cost in Java. We use the default generational mark-sweep garbage collector. The test driver thread explicitly invokes a full-heap garbage collection after starting the producers and the DSTM system and before recording the starting time. A typical 10-second test period invokes stop-the-world garbage collection two to five times at a total cost of 0.1–0.5 seconds.

For each executor on each benchmark, we take the mean throughput of ten runs. Though we have endeavored to keep the environment as uniform as possible, random factors cannot be fully eliminated. Randomizing factors include memory management, the execution of kernel daemons, and the shutdown of threads at the end of the test. With the exception of a few outliers, all experiments consumed between 10.0 and 10.1 seconds.

### 5.4 Performance Comparison

To test each microbenchmark, we generate transactions of three distributions in a 17-bit integer space. The first 16 bits are for the transaction content (i.e., the dictionary key) and the last is the transaction type (insert or delete). The first distribution is *uniform*. The second is *Gaussian*, with a mean of 65,536 and a variance of 12,000, which means that 99% of the generated values lie among the 72,000 (55%) possibilities in the center of the range. The third is *exponential*. It first generates a random double-precision floating-point number  $r$  in range  $[0, 1)$  and then takes the last 17 bits of  $-\log(1 - r)/0.001$ . The parameter 0.001 determines the narrowness of the distribution. In this case, 99% of the generated values lie between 0 and 6907, which is a little over 5% of the range.

Figures 3–5 present throughput, in transactions per second, for each of our three microbenchmarks and each of the three distributions. Each individual graph includes three curves, one each for the round-robin, fixed key-based and adaptive key-based executors. As supporting data, Figures 6–8 show the total number of “contention instances” encountered in each experiment. A “contention instance” occurs when a transaction first identifies a conflict, and again each time the conflict still exists after backoff. The data should be taken as suggestive only, since it is sensitive to backoff strategy.

**Hash table.** For the uniform distribution, both key-based executors outperform round robin, clearly showing the benefit of partitioning. The main effect comes from locality since the number of conflicts is negligible in a table of this size and the load balance of round robin is perfect.

The throughput of the fixed executor is 25% higher than round robin for two workers, and the gap increases with additional workers. Note, however, that the fixed executor does not completely

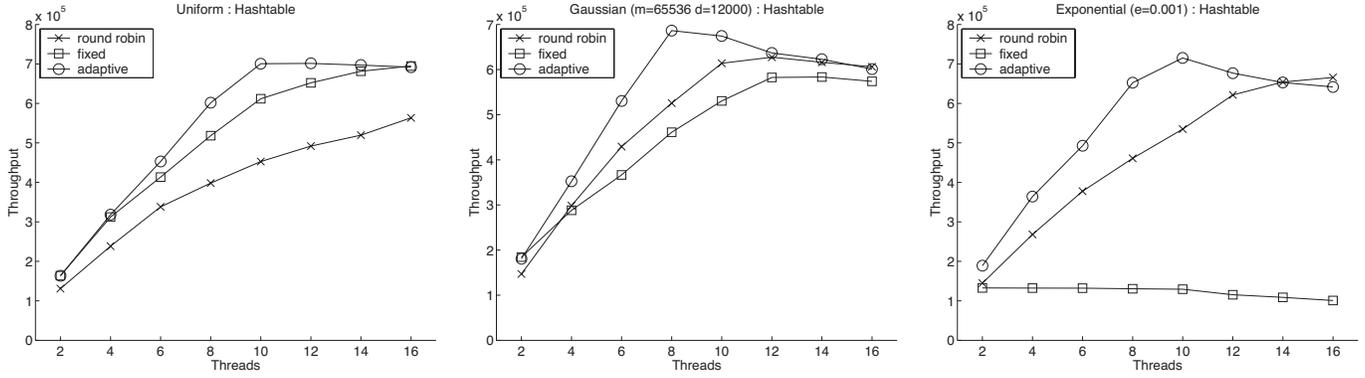


Figure 3: Throughput (txn/s) of the hash table microbenchmark with a uniform (left), Gaussian (middle), or exponential (right) distribution of transaction keys.

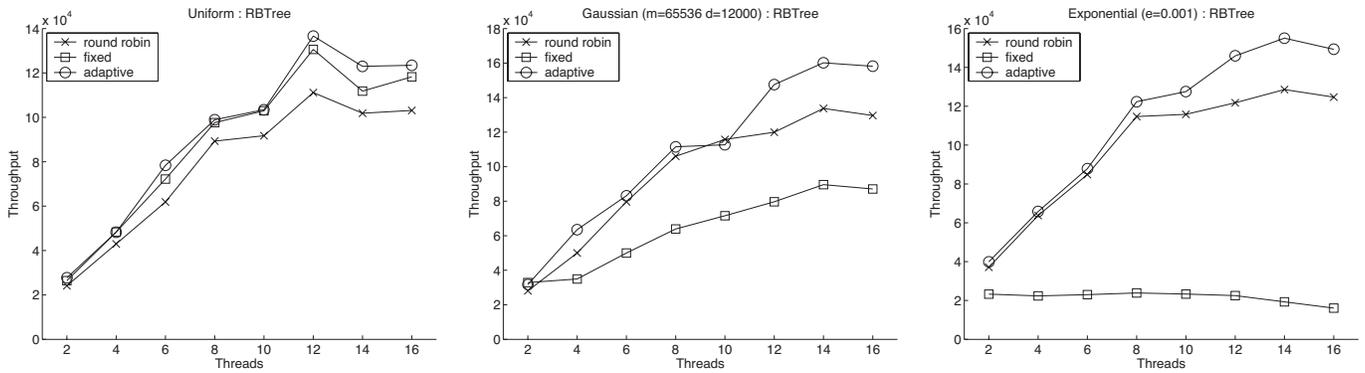


Figure 4: Throughput (txn/s) of the red-black tree microbenchmark with a uniform (left), Gaussian (middle), or exponential (right) distribution of transaction keys.

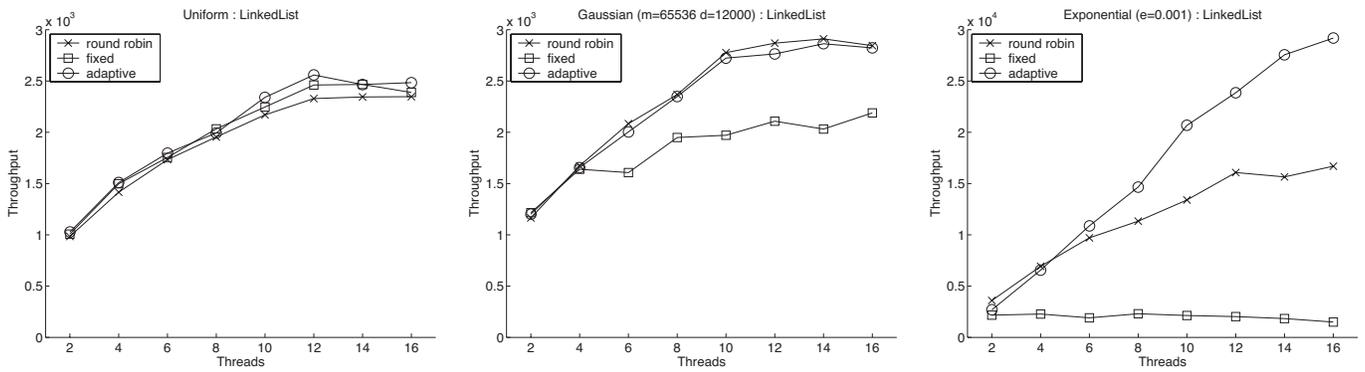


Figure 5: Throughput (txn/s) of the sorted list microbenchmark with a uniform (left), Gaussian (middle), or exponential (right) distribution of transaction keys.

balance load, because the modulo function produces 50% “too many” values at the low end of the range. The adaptive executor balances the load via uneven partitioning of the index range. As a result, the adaptive executor outperforms the fixed executor up through about ten threads; at

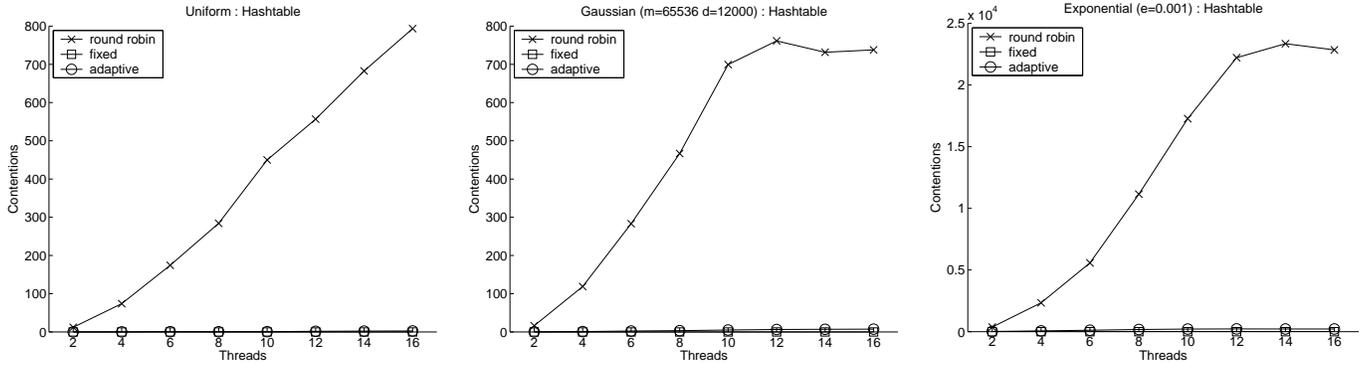


Figure 6: Hash table contention instances (10-second run).

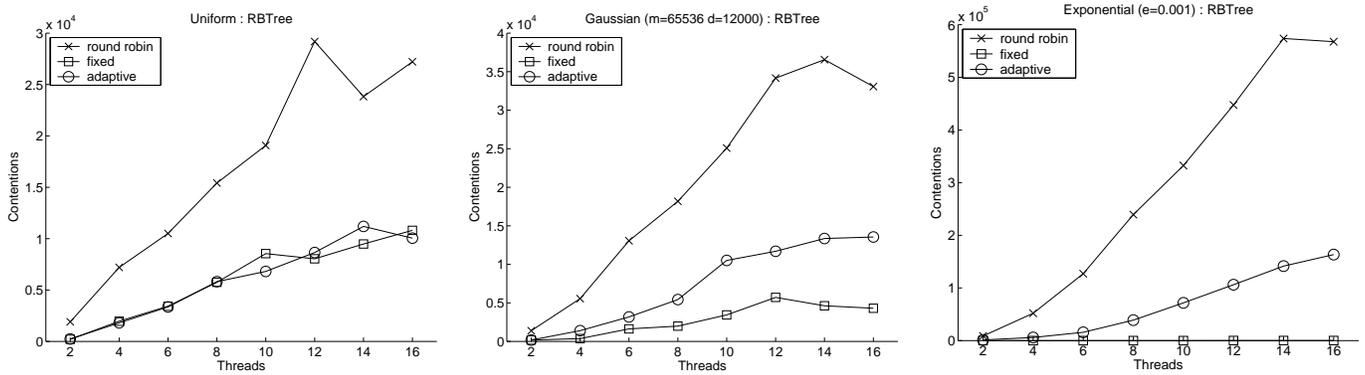


Figure 7: Red-black tree contention instances (10-second run).

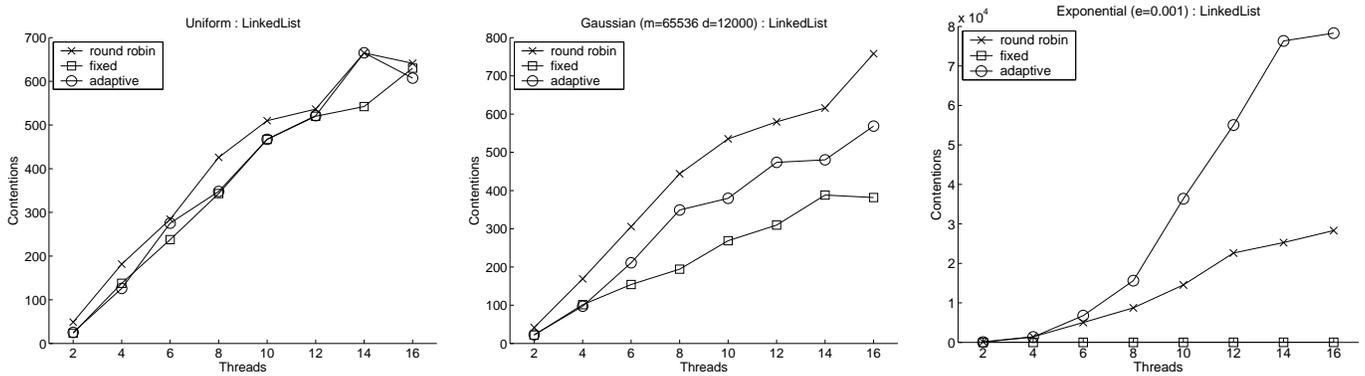


Figure 8: Sorted list contention instances (10-second run).

this point the fixed number of producers are unable to satisfy the processing capacity of additional workers.

For the nonuniform distributions, the throughput of the fixed executor suffers as transactions are concentrated in a narrower range. In the extreme (exponential) case, the fixed executor shows no speedup beyond two workers. The adaptive executor, by contrast, shows mostly scalable performance as long as producers are able to keep up with new workers. Clearly the adaptive executor's advantage over round robin also comes from better locality.

**Red-black tree.** Results for the balanced search tree are qualitatively similar to those of the hash table. The adaptive key-based executor provides the best throughput for all three distributions—up to 20% better than round robin. The performance gain is due to good locality. Since nearly half of all the nodes are leaves, partitioning helps threads repeatedly access the same groups of data while at the same time reducing the conflicts between them.

**Sorted list.** Here locality and contention are more dynamic and not as well captured by the transaction key. All three executors provide roughly the same throughput with the uniform key distribution. As in the previous two benchmarks, the throughput of the fixed executor drops below the others in the nonuniform distributions. Round robin delivers similar performance to adaptive for the Gaussian key distribution. For the exponential distribution, however, the adaptive executor outperforms round robin by 20–75% when more than six workers are active. The combination of the exponential distribution and the adaptive executor ensures that workers at the low end of the key range work on a very concentrated region of the list. A worker that successfully completes a small transaction has an excellent chance of completing additional transactions before losing a needed cache line to another worker. With round-robin scheduling, by contrast, every worker will see occasional requests in the tail of the distribution, and even requests near the low end will be scattered around the beginning of the list. After completing one transaction, a worker is likely to suffer a cache miss in its next transaction, during the servicing of which it is likely to lose the lines it used in the previous transaction, without being able to leverage their presence for additional requests. This explanation is supported by experiments that count the number of transactions completed in various key ranges with the exponential and round-robin schedulers. With the latter, workers complete roughly equal numbers of transactions in each key range; with the former, transactions with small keys see dramatically higher throughput.

#### 5.4.1 Performance Factors

**Locality and Load Balance.** Since load is equally balanced in the round robin and the adaptive executors, the main factor separating their performance would appear to be locality. The hash table here is a clear evidence. Adaptive partitioning also reduces conflicts, as shown in Figure 6, but the total amount of contention is so low even in round robin as to be essentially irrelevant. It seems reasonable to expect that conflict reduction will pay off in high-contention applications.

The nonuniform input distribution results show the advantage of adaptive partitioning. By sampling key values, the executor can balance load among workers.

**Contention.** In the hash table and the uniform and Gaussian distributions of the sorted list, the total number of contention instances is small enough (less than 1/100th the number of completed transactions) to have a negligible effect on performance. Even in the red black tree and the exponential distribution of the sorted list, fewer than one in four transactions encounters contention. In less scalable applications, contention may be more of an issue for executor throughput.

It is perhaps worth noting that the exponential distribution case of the sorted list provides the only example in which a better performing executor suffers *higher* contention. As discussed above, we conjecture that the adaptive executor is allowing workers whose transactions are clustered near the beginning of the list to complete multiple transactions before losing a needed cache line to another worker. In the 16-worker case, though the adaptive executor sees more than 50,000 additional instances of contention, relative to round robin, it completes 150,000 additional transactions. Cache misses, not contention, are the principal factor in throughput.

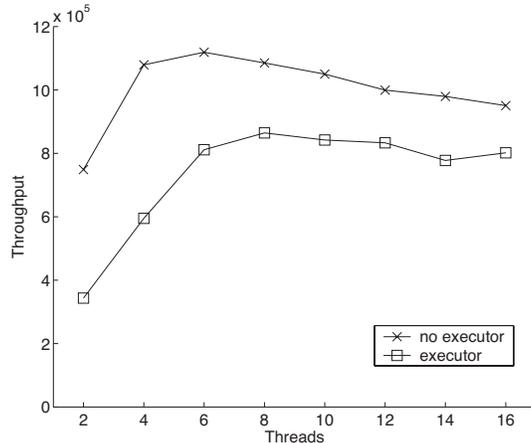


Figure 9: Throughput of empty threads and executor tasks.

**Executor Overhead** Executor overhead is a potential factor in overall performance. It includes the cost of queuing synchronization, both blocking and non-blocking; cache misses due to queue operations in different threads; and the garbage collection required to recover queue nodes and other metadata. To measure this overhead in isolation we devised an extreme case test that compares the throughput of a simple transactional executor with  $k$  workers to that of a collection of  $k$  threads executing transactions in a loop. For executor mode, we constantly use six producers. Figure 9 shows the result. With only two worker threads, the executor doubles the overhead of transactions. At higher thread counts the ratio is less pronounced. Naturally it is much less pronounced when transactions are of nontrivial length. In our experiments, executor overhead has a noticeable impact on the hash table benchmark, but almost none on the sorted list.

#### 5.4.2 Summary of Performance Comparison

Overall, our results show a clear performance benefit from key-based adaptive scheduling in the hash table and red-black tree benchmarks, where the transaction key accurately predicts a transaction’s data access pattern. The advantage is smaller in the sorted list benchmark, where the prediction is significantly weaker. We still see a significant benefit in the exponential distribution, however, and in all cases the throughput of the adaptive executor is either the best or essentially tied for the best. Fixed partitioning is clearly unscalable for nonuniform key distributions.

## 6 Related Work

The basic framework of the executor bears a close resemblance to various flavors of locality-based thread scheduling. In memory-conscious scheduling, Markatos and LeBlanc used application knowledge to give priority to threads that referenced the same data and ran them on the same processor [12]. In CacheMiner [23], Yan et al. defined a general programming interface and relied on the programmer to partition the computation into fine-grain tasks and declare their access pattern as the starting address of data access in each array. The performance was significantly better for parallel sparse matrix multiply, for which static compiler optimization such as tiling was not applicable. Philbin et al. used a similar approach to improve uniprocessor cache performance for N-body simulation [16]. While both the latter two methods used equal-size partitions of array data

and assigned the tasks using a hash table, the size of the partition was determined by the number of processors in CacheMiner and by the cache size in Philbin et al.

In comparison, we use adaptive partitioning based on a run-time estimate of the probability distribution. The context is very different because transactions are executed on a stack of software layers, they use dynamically allocated data, and they are not completely parallel.

We were originally motivated by ideas from dynamic computation regrouping. For N-body simulation, irregular mesh, and sparse matrix computations, a number of studies have showed that run-time reordering can improve sequential running time despite its cost. Ding and Kennedy used lexicographic grouping [6]. Han and Tseng used lexicographic sorting [7]. Mitchell et al. divided the data space into fixed-size tiles and grouped computations in cache-sized data sub-regions [15]. Mellor-Crummey et al. organized computation into space-filling curves for computations whose data had spatial coordinates [13]. Strout et al. compared lexicographic grouping, sorting, and sparse tiling and gave a run-time framework for automating them [20]. Strout and Hovland gave a cost model based on hyper graphs [21]. Pingali et al. manually regrouped work in integer computations, including tree traversals [17]. Dynamic regrouping is also part of the inspector-executor model, originally studied by Saltz and others for partitioning irregular computations on distributed-memory parallel machines [5]. Most of these studies considered data transformation in conjunction with computation reordering.

Since transactions are generated concurrently, a separate inspection step is impractical. Key-based executors interleave inspection with transaction processing. The adaptive partitioning has low run-time overhead. We characterize locality with a key, which is readily determined for dictionary operations and should, we conjecture, be readily available for other important classes of applications. The partitioning is fully automatic because of the transaction semantics. In this work we do not address the issue of data layout. The adaptive scheme uses data sampling to estimate the distribution of the workload. Other groups have used code-based sampling to identify data streams and strides [4, 22]. While most past techniques exploit dynamic locality for sequential execution, we focus on parallel execution.

## 7 Conclusions

In this paper, we have augmented the dynamic software transactional memory system of Herlihy et al. with locality-aware partitioning. We presented an adaptive partitioning method based on run-time estimation of the probability distribution of the input. The adaptive executor shows clear improvement when the key accurately predicts the data access by a transaction. Otherwise, the adaptive executor performs comparable to round robin transaction scheduling. Traditional, fixed partitioning, on the other hand, scales poorly when the input distribution is not uniform. Locality-aware partitioning also has the potential to significantly reduce the rate at which transactions conflict with one another, though this effect was not a principal determinant of performance in our experiments.

Although our results are preliminary, we believe that we have demonstrated significant potential for further exploration in this area. The notion of transactional executors, in particular, seems worth of additional study. Data conflicts among tasks are potentially problematic with traditional (nontransactional) executors, given the potential of lock-based code to trigger deadlock. The automatic atomicity of transactions would seem to broaden the class of applications to which executors are applicable. At the same time, while our results were obtained in the context of transactional memory, we believe that adaptive scheduling has significant promise for nontransactional executors as well.

Further work is clearly needed in the generation of transaction keys, both manual and automatic (possibly profile-driven). It may also be possible to use transaction keys for additional purposes, e.g., contention management [18]. Finally, the relative importance of locality and conflict bears reconsideration in high-contention workloads.

## Acknowledgements

Our thanks to Virendra Marathe and Mike Spear for their participation in early phases of this work.

## References

- [1] R. Allen and K. Kennedy. *Optimizing Compilers for Modern Architectures: A Dependence-based Approach*. Morgan Kaufmann Publishers, October 2001.
- [2] A. J. Bernstein. Analysis of programs for parallel processing. *IEEE Transactions on Electronic Computers*, 15(5):757–763, Oct. 1966.
- [3] R. D. Blumofe, C. F. Joerg, B. C. Kuszmaul, C. E. Leiserson, K. H. Randall, and Y. Zhou. Cilk: An efficient multithreaded runtime system. In *Proceedings of the ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, Santa Barbara, CA, 1995.
- [4] T. M. Chilimbi and M. Hirzel. Dynamic hot data stream prefetching for general-purpose programs. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation*, Berlin, Germany, June 2002.
- [5] R. Das, M. Uysal, J. Saltz, and Y.-S. Hwang. Communication optimizations for irregular scientific computations on distributed memory architectures. *Journal of Parallel and Distributed Computing*, 22(3):462–479, Sept. 1994.
- [6] C. Ding and K. Kennedy. Improving cache performance in dynamic applications through data and computation reorganization at run time. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation*, Atlanta, GA, May 1999.
- [7] H. Han and C. W. Tseng. A comparison of locality transformations for irregular codes. In *Lecture notes in computer science 1915 (Proceedings of LCR'2000)*, pages 70–84, Rochester, NY, 2000. Springer-Verlag.
- [8] M. Herlihy, V. Luchangco, M. Moir, and W. N. Scherer III. Software transactional memory for dynamic-sized data structures. In *Proceedings of the 22nd ACM Symposium on Principles of Distributed Computing*, pages 92–101, Boston, MA, July 2003.
- [9] M. Herlihy and J. E. Moss. Transactional memory: Architectural support for lock-free data structures. In *Proceedings of the 20th International Symposium on Computer Architecture*, pages 289–300, San Diego, CA, May 1993.
- [10] P. J. Janus and E. A. Lamagna. An adaptive method for unknown distributions in distributive partitioned sorting. *IEEE Transactions on Computers*, c-34(4):367–372, April 1985.
- [11] D. Lea. Concurrency JSR-166 interest site. <http://gee.cs.oswego.edu/dl/concurrency-interest/>.
- [12] E. P. Markatos and T. J. LeBlanc. Using processor affinity in loop scheduling on shared-memory multiprocessors. *IEEE Transactions on Parallel Distributed Systems*, 5(4):379–400, 1994.
- [13] J. Mellor-Crummey, D. Whalley, and K. Kennedy. Improving memory hierarchy performance for irregular applications. *International Journal of Parallel Programming*, 29(3), June 2001.

- [14] M. M. Michael and M. L. Scott. An efficient algorithm for concurrent priority queue heaps. In *Proceedings of the 15th ACM Symposium on Principles of Distributed Computing*, pages 267–275, 1996.
- [15] N. Mitchell, L. Carter, and J. Ferrante. Localizing non-affine array references. In *Proceedings of International Conference on Parallel Architectures and Compilation Techniques*, Newport Beach, CA, October 1999.
- [16] J. Philbin, J. Edler, O. J. Anshus, C. C. Douglas, and K. Li. Thread scheduling for cache locality. In *Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems*, 1996.
- [17] V. S. Pingali, S. A. McKee, W. C. Hsieh, and J. B. Carter. Computation Regrouping: Restructuring Programs for Temporal Data Cache Locality. In *International Conference on Supercomputing*, New York, NY, June 2002.
- [18] W. N. Scherer III and M. L. Scott. Advanced contention management for dynamic software transactional memory. In *Proceedings of the 24th ACM Symposium on Principles of Distributed Computing*, Las Vegas, NV, July 2005.
- [19] X. Shen and C. Ding. Adaptive data partition for sorting using probability distribution. In *Proceedings of the International Conference on Parallel Processing*, Montréal, Canada, August 2004.
- [20] M. M. Strout, L. Carter, and J. Ferrante. Compile-time composition of run-time data and iteration reorderings. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation*, San Diego, CA, June 2003.
- [21] M. M. Strout and P. Hovland. Metrics and models for reordering transformations. In *Proceedings of the 2nd ACM SIGPLAN Workshop on Memory System Performance*, Washington, DC, June 2004.
- [22] Y. Wu. Efficient discovery of regular stride patterns in irregular programs. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation*, Berlin, Germany, June 2002.
- [23] Y. Yan, X. Zhang, and Z. Zhang. Cacheminer: A runtime approach to exploit cache locality on SMP. *IEEE Transactions on Parallel Distributed Systems*, 11(4):357–374, 2000.