

Data race: Tame the Beast

K. Leung, Z. Huang,* Q. Huang, P. Werstein

Department of Computer Science

University of Otago, Dunedin, New Zealand

Email: {kcleung; hzy; tim; werstein}@cs.otago.ac.nz

Abstract

Data races hamper parallel programming and threaten the reliability of future software. This paper proposes a data race prevention scheme, which can prevent data races in the View-Oriented Parallel Programming (VOPP) model. VOPP is a novel shared-memory data-centric parallel programming model, which uses views to bundle mutual exclusion with data access. We have implemented the data race prevention scheme with a memory protection mechanism. Experimental results show that the extra overhead of memory protection is trivial in our applications. The performance is evaluated and compared with modern programming models such as OpenMP and Cilk.

1 Introduction

Parallel programming has become inevitable with the advent of multicore and chip-multithreading (CMT) technologies [30]. These technologies allow multiple processors to be packed into a chip in a single computer, which often provides shared memory and cache. However, parallel programming with shared memory can be prone to errors such as data race, which is difficult to debug due to its non-determinism and thus can severely affect programmability and software reliability.

In a parallel multithreaded computation, a data race occurs if concurrent threads access the same memory location without mutual exclusion primitives such as locks, and at least one of the threads writes to the location. There have been many studies on debugging data races. Some perform a post-mortem analysis based on program execution traces [8, 11, 14, 23, 24], while others perform on-the-fly analysis during program execution [2, 10, 22, 29]. Among modern shared-memory parallel programming models [9, 25, 26, 28], only Cilk++ [9] provides a data race detector called Cilkscreen [2, 9, 19].

Even though race detectors can help debug some data races, they often have the following problems.

- Race detectors are often expensive to run, both in terms of computation and memory space. For example, Cilkscreen can take up to 30 times the nor-

mal execution time of the debugged program to run and the memory footprint can be “several times” the memory footprint of the original application [9].

- Race detectors can only detect data races for one given input of a program. If data races do not occur when the program is run with a given input, this does not imply the program is data race free. The reason is that a different input may result in threads being executed in different order, and the resultant interaction may cause data races.
- To a novice programmer, race detectors can be difficult to use. For example, Cilkscreen gives a detailed trace of memory addresses and their associated function names and line numbers, which can be very scary and confusing to inexperienced programmers. In addition, this trace is of little help to programmers about the dynamic nature of the data races, e.g. when and how the data races happen.

In this paper, instead of data race detection, we propose a data race prevention scheme, which can prevent data races from occurring in the first place. This scheme is implemented in our View-Oriented Parallel Programming (VOPP) model [15, 38]. In VOPP, shared data is partitioned into views. A view is a set of memory units (bytes or pages) in shared memory. Each view, with a unique identifier, can be created, merged, or destroyed at any time in a program. Before a view is accessed (read or written), it must be acquired (e.g. with *Vpp_acquire_view*); after the access of a view, it must be released (e.g. with *Vpp_release_view*). The most important property for views is that they do not intersect with each other (refer to [15, 38] for details).

VOPP is a data-centric programming model [3, 7, 35], which bundles mutual exclusion and data access together. With data-centric programming, the programmer only considers which data to access atomically, instead of the issues like mutual exclusion and data race. Accordingly, VOPP has the following advantages.

First, programmers are relieved from data race issues. VOPP requires that shared data of a program be partitioned into non-overlapping views according to the shared pattern of data. Since the partitioning of data is decided by the programmer as part of the parallel programming, the parallel algorithm can be finely tuned through careful

*Zhiyi Huang is currently a visiting scientist at the Cilk group of MIT CSAIL.

view allocation. Once views are allocated, the programmer is only concerned about which view (data) will be accessed, instead of worrying about the data races and mutual exclusion in lock-based programming. When a view is acquired, mutual exclusion is automatically achieved, so it is not possible for other processes to access the view at the same time. If a view is accessed without being acquired, either the programmer can be notified of the problem by the compiler with some VOPP related support, or the run-time system can report the problem with the support of the underlying virtual memory system, as we will describe in Section 2.

Second, VOPP encourages reduction of unnecessary sharing and enables high performance. We argue that shared memory is and will be a critical resource that needs to be used with care. As we know, shared memory/cache is and will still be a bottleneck in parallel computers [12], though memory space will keep increasing. The philosophy behind VOPP is to discourage data sharing through view allocation, rather than automatic sharing as in other shared memory programming models. Every time a view is created, the programmer is given a chance to justify if this sharing of data is necessary. It is worthwhile to spend time on careful view partitioning in order to reduce unnecessary data sharing.

Third, VOPP is portable over a range of parallel computers and can hide the differences of memory architectures. A view can be implemented efficiently on both distributed memory and shared memory, and can even be prefetched into the cache [38]. We have efficiently implemented VOPP on both cluster computers and multi-core computers [16, 17, 38], which is not easily attainable for programming models such as OpenMP and Cilk.

Fourth, VOPP enables optimizations such as prefetching. Since the base address and size of a view are known at view allocation, when a view is acquired, VOPP can pass the information to the system which can prefetch the view into the cache. Experiments show that prefetching with helper threads can significantly improve the performance of memory access [18, 21, 38].

Finally, VOPP is language independent. It proposes a general idea regarding how to use shared memory for inter-process communications. Any language can apply this idea of view partitioning. For example, for object-oriented languages, view classes can be defined and their instances can be created for inter-process communications. In this paper, we use the C language as an example to demonstrate the idea and performance of VOPP.

The rest of this paper is organized as follows. Section 2 describes a data race prevention scheme that can eliminate data races in VOPP. In Section 3, we briefly introduce the advanced features of Maotai 2.0 for improving programmability and performance. Maotai 2.0 is our up-to-date implementation of VOPP on CMT computers. Section 4 presents the performance evaluation of Maotai 2.0. Finally, our future work is suggested in Section 5.

2 Data Race Prevention

In VOPP, shared data is defined through views. Unlike most shared memory parallel programming models, variables are private to a process by default in VOPP. Shared objects must be *explicitly* defined as “views”.

Views can be created, destroyed, merged, or resized, but a process must acquire a view (read-only or read-write) before accessing it and must release it after finishing with the view. VOPP adopts the Single-Writer Multiple-Reader (SWMR) model. At any given time, a view can either be read/written by one process *or* allow read-only access to multiple processes. In our current implementation, a view uses a contiguous memory space to store shared variables. Below is a simple example of VOPP in C.

```

1 typedef struct {int a[ARRAY_SIZE];
2                 int result;} Foo;
3
4 Foo *ptr;
5 if (0 == Vpp_proc_id) {
6     /* master allocates view 0 with
7      * type SWV, which is a shared object
8      * with "Foo" type */
9     Vpp_alloc_view(0, sizeof(Foo), SWV);
10 }
11 Vpp_barrier();
12 ...
13 ptr = Vpp_acquire_view(0);
14 ptr->result += do_work(ptr->a);
15 Vpp_release_view(0);
16 ...

```

As illustrated in the above example, if a data structure should be shared by multiple processes, a view has to be created for it with *Vpp_alloc_view*. For exclusive access to the view, the view type is SWV, which means “Single Writer View”. However, we also provide other advanced views to enhance the programmability and flexibility of VOPP (refer to Section 3).

If a process wants access to a view, the view must be acquired with *Vpp_acquire_view* (or *Vpp_acquire_Rview* for read-only access). The view must be released with *Vpp_release_view* after accessing it.

2.1 Implementation

In our data race prevention scheme, data races are prevented by a memory protection mechanism available in most UNIX systems. All views are initially protected from access using system calls such as *mprotect()*. *mprotect()* can deny access to a page, or allows read-only access to a page, or allows read/write access to a page. We use this mechanism to prevent a view from illegal accesses. Only after a view is acquired is a process allowed to access the memory pages of the view via *mprotect()*. When a view is released, the process is again denied access to the view.

If a process accesses a view before *Vpp_acquire_view* or after *Vpp_release_view*, the pages of the view would not have the necessary access permission and thus a segmentation fault will occur. Our system will handle the fault,

send a warning message to the programmer that a view is accessed without acquisition, and quit the program execution.

In this way, a view can either be written to by one process or read by multiple processes at a time. Programmers do not need to worry about the data race bugs. If a view is accessed by calling `Vpp_acquire_view`, mutual exclusion of the view access is automatically done by the system. If a view is accessed without view acquisition, a segmentation fault will occur, and the system will alert the programmer about which view is accessed without acquisition. The programmer can easily fix the bug by inserting `Vpp_acquire_view` and `Vpp_release_view` into the faulted code section.

The extra cost of this data race prevention scheme is the overhead of the memory protection. In our VOPP implementation Maotai 2.0, this cost is very low. On a Sun T2000 Server equipped with a 1GHz UltraSPARC T1 processor [31], micro-benchmarking results demonstrate that the overhead of memory protection added to the view primitives is generally very low (around 2-3 μ s). The exception is `Vpp_acquire_view`, requiring up to 35 μ s extra, which covers the essential overhead of the memory protection mechanism (see Table 1). Note that `Vpp_acquire_Rview` and `Vpp_release_Rview` means acquiring and releasing views as read-only.

Table 1: Breakdown of view primitive costs (in μ s)

Primitive	no prot	prot	cost
<code>Vpp_acquire_view()</code>	3.14	39.08	35.94
<code>Vpp_acquire_Rview()</code>	3.60	6.32	2.72
<code>Vpp_release_view()</code>	1.91	4.54	2.63
<code>Vpp_release_Rview()</code>	1.99	4.64	2.65

However, in our application benchmarks, this overhead does not cause noticeable difference in application speedup. Table 2 shows the speedups (at 32 processes) of our applications with and without memory protection in Maotai 2.0. We have six benchmark applications: Successive Over-Relaxation (SOR), Gaussian Elimination (GE), Integer Sort (IS), Neural Network (NN), Mandelbrot, and Mergesort, which typically represent a wide variety of parallel applications. For details of these applications, refer to Section 4. As we can see from Table 2, in all 32-process benchmark cases, the difference is around 0.5%. Therefore, the overhead introduced by data race prevention is trivial.

Table 2: Effects of memory protection on benchmark application speedups with 32 processes / threads

Application	no prot	prot
SOR	16.82	16.77
GE	22.41	22.36
IS	16.51	16.47
NN	16.98	16.92
Mandelbrot	17.80	17.79
Mergesort	12.52	12.50

One issue about the implementation is that memory protection such as `mprotect` is page-based. Therefore, in order to protect view data properly, memory space allocated to a view is aligned by pages. This can result in memory space wastage. Table 3 shows the requested and actual sizes of the memory space allocated by VOPP in our benchmark applications. The page size is 8kB and 32 processes are used when the data of the table are collected. From this table, it can be seen that some applications like GE and Mandelbrot, which have many views that do not exactly fit a page, have a higher proportion of memory wastage (up to 51%), though other applications have less than 7% wastage. However, this memory wastage is much smaller than the memory footprint of race detectors, which can be "several times" the memory footprint of the original applications.

Table 3: Requested vs actual VOPP shared size (in Kbytes) in different applications

Algorithm	Requested	Actual	Wasted	Percent wasted
SOR	4,097,024	4,194,304	97,280	2.32
GE	64,016,004	98,328,576	34,312,572	34.9
IS	4,194,304	4,194,304	0	0
NN	271,612	294,912	23,300	7.90
Mandelbrot	2,000,000	4,096,000	2,096,000	51.2
Mergesort	1,600,001,280	1,600,274,432	273,152	0.0171

Fortunately, with architectural support of variable-size pages [6, 37], this memory wastage can be greatly reduced.

2.2 Related Work

Shared memory systems have different approaches to the data race issue. In most systems (such as OpenMP [26], Cilk [32], PThread [25] and UPC [33]), locks are not associated with shared objects and programmers are responsible for arranging locks properly to prevent data races, therefore these systems are prone to data races and deadlocks caused by programming errors.

Transactional memory systems are very convenient for parallel programming. However, its major goal is to guarantee atomicity of memory accesses without locking, instead of addressing the data race issue. It rolls back one or more conflicting transactions when atomicity may be violated. Therefore, it never removes data races. Also live-lock can be a new issue with transactional memory (all competing processes repeatedly roll back and make no progress).

Deterministic Parallel Java (DPJ) [4, 5] is an object-oriented shared-memory concurrent model based on a Java language extension. In DPJ, the compiler uses the "type and effect" system on classes and methods to statically check whether two concurrent code blocks can be executed concurrently, if not, then the tasks will be run serially instead in the order they are listed to ensure determinacy. The concept of "region" in DPJ is similar to "view" in VOPP. Both models bundle access management

into shared objects and relieve programmers from the responsibility of manually setting locks to prevent data race. However, the difference is that DPJ avoids the data race problem through falling back to serial execution, while VOPP detects the data races at runtime and helps the programmer fix the bugs.

3 Advanced Features in Maotai 2.0

In addition to data race prevention, Maotai 2.0 also offers primitives for acquiring multiple views in order to avoid deadlocks, producer/consumer views, and system queues to enhance programmability and performance. These features are discussed below.

3.1 Deadlock Avoidance

Similar to data race, deadlock is another pain that can happen easily but is difficult to debug in shared-memory parallel programming. In VOPP, deadlock can happen if views are acquired in a nested way and different processes acquire them in different orders.

To avoid deadlocks due to acquiring multiple views in different orders, Maotai 2.0 offers primitives for acquiring multiple views. Programmers can list all views to be acquired with these primitives which will acquire the views in a specific, same order. In this way, there is no chance for deadlocks to happen.

Below is an example illustrating the use of the primitives for acquiring multiple views:

```

1  /* acquire access to both view 0 and 1 */
2  Vpp_acquire_multiviews(0, &ptr0, 1, &ptr1);
3  ptr0->result += compute0(ptr0->a, ptr1->a);
4  ptr1->result += compute1(ptr1->a, ptr0->a);
5  Vpp_release_view(); /* release all views */

```

In the above example, the process acquires both view 0 and 1 with *Vpp_acquire_multiviews* which puts the view base addresses into *ptr0* and *ptr1*. Finally the process releases both views with *Vpp_release_view*.

Note that the above solution cannot eliminate deadlocks from VOPP programs as the data race prevention scheme does data races. There are two reasons: first, the programmer may choose not to use *Vpp_acquire_multiviews* for nested view acquisition; second, even if the programmer would like to use the primitive, it is difficult to know which views to acquire in advance in some programs where inner views can only be decided after the outer views are processed.

Nevertheless, the above primitives provide an avenue for novice programmers to avoid unnecessary deadlocks.

3.2 Producer/Consumer View

A Producer/Consumer View (PCV) is provided to allow direct expression of producer/consumer relationships in

parallel algorithms. Traditionally barriers are used to synchronize the producers and the consumers in shared memory parallel programming. With the introduction of PCV, programming with producer/consumer problem is more straightforward and thus increases programmability. In addition, PCV can avoid expensive barriers, which makes all processes wait and whose cost would increase with increasing number of processes.

PCV is implemented as a queue. The producer enqueues a new version of the view by acquiring the view, producing the data, and finally releasing the view. The consumer dequeues a version of the view by acquiring read-only access to the view. After it finishes with the view, it releases the view whose buffer may be recycled by the producer.

In our experiments, the SOR and GE benchmark applications demonstrate that PCVs give a better speedup than barrier based implementations. Figure 1 and 2 shows the speedup difference between applications using barriers and those using PCVs.

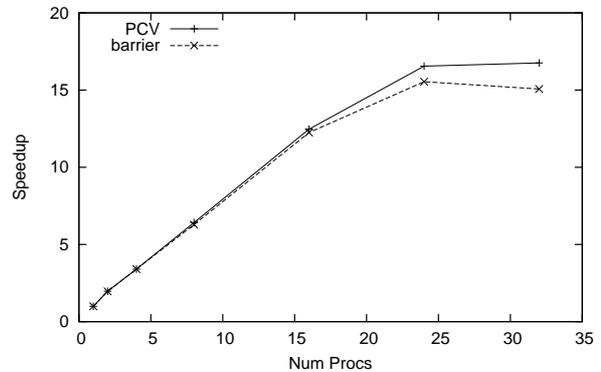


Figure 1: Speedup of SOR in VOPP

Figure 1 shows the speedup of SOR which uses PCV to improve its performance. Compared with its barrier implementation, the improvement of speedup is 11.2% at 32 processes.

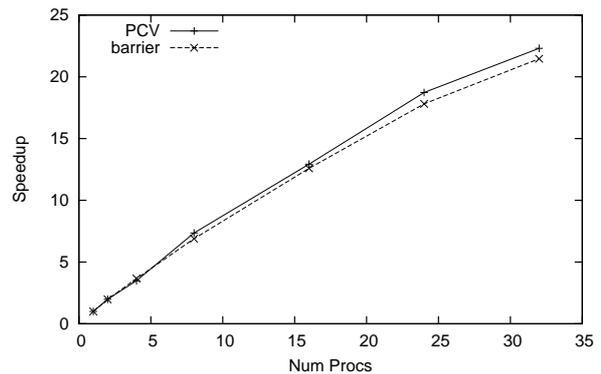


Figure 2: Speedup of GE in VOPP

Figure 2 shows the speedup of GE which uses PCV to improve its performance. Compared with its barrier

implementation, the improvement of speedup is 4.2% at 32 processes.

3.3 System Queues

System queues are provided in Maotai 2.0 to store view IDs. This facility allows easy implementations for task queues. Task queues are good for load balancing parallel applications (e.g. Mandelbrot and tree search algorithms), where the data for each job or node can be put in a view and its ID is simply enqueued in a system queue for other processes to work on.

In Maotai 2.0, the enqueue and dequeue calls are efficient. In a microbenchmark test on a Sun T2000 server, an enqueue call only takes $2.65\mu s$ and a dequeue call takes $2.56\mu s$.

4 Performance Evaluation with Other Models

In this section, we compare the performance of Maotai 2.0 with other modern shared memory parallel programming models OpenMP and Cilk. Our benchmark applications include Successive Over-Relaxation (SOR), Integer Sort (IS), Gaussian Elimination (GE), Neural Network (NN), Mandelbrot and Mergesort. The experiments are carried out on a Sun T2000 server with an UltraSPARC T1 processor and 16GB memory. The UltraSPARC T1 has eight cores, each of which is clocked at 1GHz and supports four hardware threads. In total, the UltraSPARC T1 processor supports up to 32 hardware threads [31]. Linux kernel 2.6.24-sparc64-smp and the compiler gcc-4.4 are used during benchmarking. The benchmark applications are implemented on Maotai 2.0, Cilk-5.4.6 [32], and OpenMP 3.0 [26], respectively. All programs are compiled with the optimization flag “-O2”. In each case, speedup is measured against the serial implementation of the benchmark algorithm. The elapsed time calculated in each case excludes initialization and finalization costs, because they are one-off and are difficult to measure within the program in models that involve source-translation, such as Cilk and OpenMP. Instead, startup and finalization times for each model are measured separately. Runtime of functions that are irrelevant to the original application, such as generation of random sequences and result-verification, are also excluded.

Successive Over-relaxation (SOR) is a multiple-iteration algorithm where each element is updated by the values of the neighbouring elements from the last iteration. In this experiment, the implementation is adapted from [38]. Matrix size is set to $8000 * 4000$ and 40 iterations are performed.

The Integer Sort (IS) algorithm used in this experiment is based on the NPB version [34]. This is a counting-sort algorithm. In this experiment, the problem size is 2^{26} integers with a B_{max} of 2^{15} and 40 repetitions are performed.

The Gaussian Elimination (GE) implementation from [38] is used in this experiment and the matrix size is set to $4000 * 4000$.

The parallel Neural Network (NN) algorithm is based on [27]. This algorithm trains a back-propagation neural network in parallel using a training data set. In this experiment, the size of the neural network is set to $9 * 40 * 1$ and the number of epochs is set to 200.

The Mandelbrot algorithm is embarrassingly-parallel. However, the workload of pixels is extremely uneven, and thus requires a load-balancing mechanism to prevent process starvation [13, 36]. In this experiment, the size of the screen is set to $500 * 500$, the maximum number of iterations is set to 500 and each pixel is calculated 5000 times. The maximum number of processes / threads is set to eight for this experiment because hyperthreading relies on memory latency. Since this application has very few memory accesses, there is little speedup when more processes / threads than the number of CPU cores are used (The UltraSparc T1 has eight cores).

The parallel Mergesort algorithm is recursive [20, 32] and is implemented verbatim in Cilk and OpenMP to test performance of the newly-available task-parallelism feature in OpenMP [1]. The array consists of 200 million integers. This algorithm is converted to the iterative version for VOPP. The iterative version requires the number of processes to be a power of 2. This version first divides the array equally between the processes and each process sorts its own subarray. Then the merge procedure largely models the recursive version of the parallel merge algorithm.

Since the UltraSPARC T1 has only one floating-point unit, all floating-point calculations in the above algorithms are converted to integer calculation to avoid the bottleneck at the floating-point unit. Removal of floating point calculations is done in all implementations and does not affect the scalability of the algorithm or the fairness of the comparison.

4.1 Experimental Results

The experimental results are illustrated with speedup curves. For each application, we give its speedup curves on Maotai 2.0, Cilk, and OpenMP. In the discussion below, n refers to the number of processes / threads.

Speedup is calculated by:

$$speedup = \frac{time_{serialimplementation}}{time_{parallelimplementation}} \quad (1)$$

To ensure fair comparison, the same serial implementation of each benchmark application is used as a baseline for calculating speedups of all parallel programming models.

For SOR (Figure 3), Maotai 2.0 has the best performance. At $n = 32$, Maotai 2.0 is 13.6% better than Cilk and 17.9% better than OpenMP.

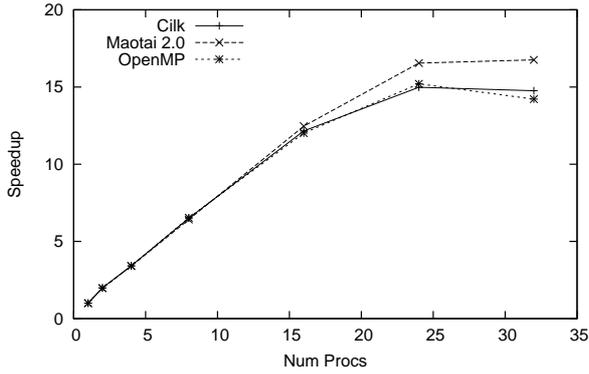


Figure 3: Speedup of SOR

For GE (Figure 4), Maotai 2.0 again has the highest speedup. At $n = 32$, Maotai 2.0 is 7.4% better than Cilk and 33% better than OpenMP.

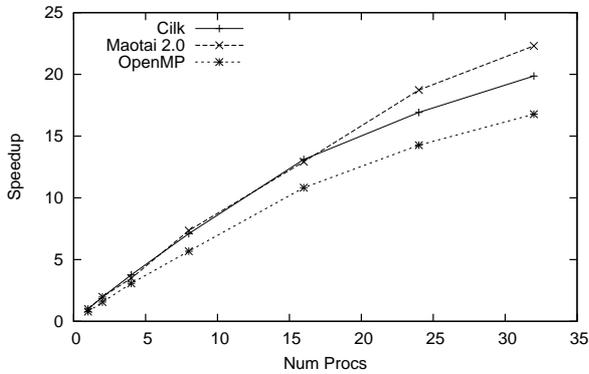


Figure 4: Speedup of GE

In IS (Figure 5), there are less variations in speedups in different models. However at $n = 32$, Maotai 2.0 is 5% faster than Cilk and 15% faster than OpenMP.

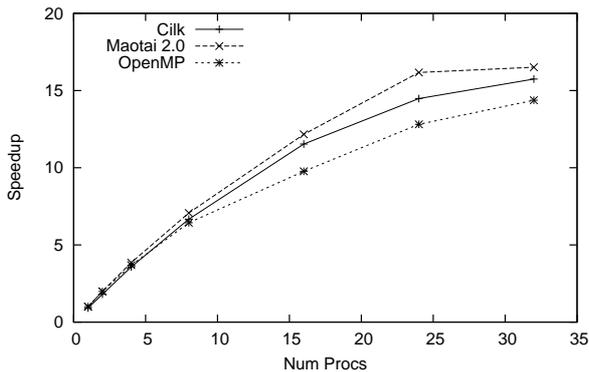


Figure 5: Speedup of IS

In NN (Figure 6), all models have similar speedups. Maotai 2.0 is 3.1% faster than OpenMP, but it is 1.8% slower than Cilk.

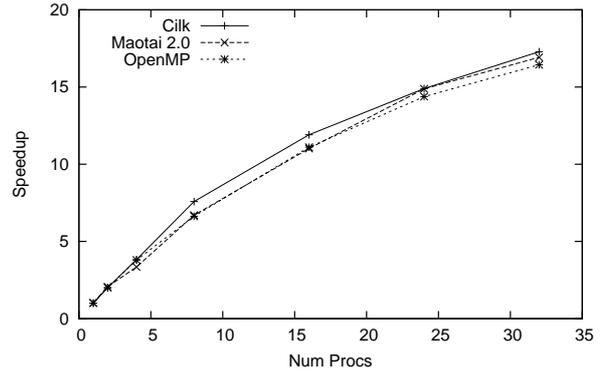


Figure 6: Speedup of NN

In Mandelbrot (Figure 7), there are relatively little differences between speedups of different models. At $n = 8$, Maotai 2.0 is 0.8% faster than Cilk and 7.2% faster than OpenMP.

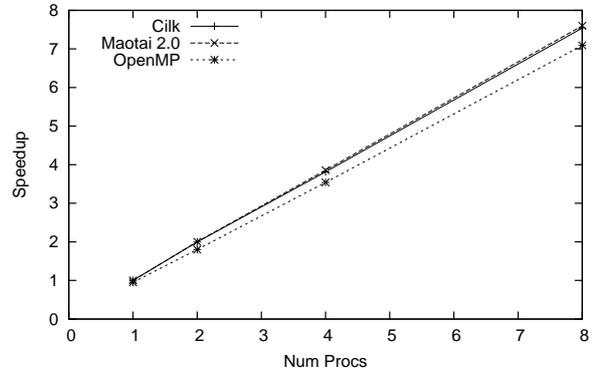


Figure 7: Speedup of Mandelbrot

For Mergesort, Figure 8 shows speedup of Maotai 2.0 is relatively slower. We will address this issue in Section 4.2.

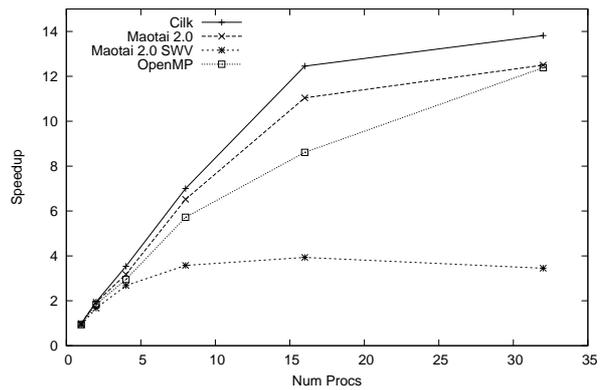


Figure 8: Speedup of Mergesort

Note that, in the above collected results, the standard deviations of the elapsed time at $n = 32$ for Maotai 2.0 and Cilk cases are less than 0.1s, but the standard deviations of the elapsed time for OpenMP are between 0.2 to 0.5s,

which may be due to the random nature of the OpenMP task scheduler.

Table 4 presents the startup and finalization time of each system. As expected, startup and finalization costs for thread-based models including Cilk and OpenMP are lower than process-based system like Maotai 2.0.

Table 4: Combined startup and finalization time (in ms) for different number of processes/threads on a Sun T2000 server

	1	2	4	8	16	24	32
Cilk	2	2	2	2	2	2	2
OpenMP	2	2	2	2	2	2	2
Maotai 2.0	9	10	11	13	15	19	22
Serial	2						

All thread-based models have the same combined startup and finalization time as the serial version regardless of the number of threads. Maotai 2.0 has a startup/finalization cost of 9ms (at $n = 1$) and the cost grows to 22ms at $n = 32$, almost linear to the number of processes. Despite Maotai 2.0 having a larger startup/finalization overhead, the 22ms is still negligible compared to the time consumed in $n = 32$ cases, which is at least 10 seconds. Also the startup/finalization time in Maotai 2.0 is only a one time event, therefore, this overhead should have negligible effect on the speedup curves.

4.2 Discussion

The following is an analysis on why Maotai 2.0 performs better or worse than other systems.

As we mentioned before, the producer/consumer view (PCV) in Maotai 2.0 enhances both programmability and performance of SOR and GE. In SOR, PCV is used to pass boundary rows to neighbour processes, thus allowing the natural expression of the message-passing relationship without the use of barrier, which would hold up irrelevant processes. Apart from programmability, the resultant performance gain is reflected in Figure 1, where the PCV VOPP version is 11.2% faster than the barrier-based SOR version.

Similarly in GE, PCV is used to *broadcast* the pivot row and the swap index, which improves programmability by mimicking the broadcasting semantics in the parallel algorithm. Also the removal of barriers by PCV improves the VOPP performance by 4.2% (Figure 2). Time is saved by replacing lock and barrier primitives with a PCV primitive.

Multiple-Program Multiple-Data (MPMD) models such as Cilk/Cilk++ and OpenMP do *not* have barriers because in this case, the parallel calculation part is conveniently expressed by parallel for-loop (or in case of Cilk, spawn recursive task decomposition threads and sync at end of parallel calculation) and the pivot part is run serially. Synchronization is implicit in the parallel for-loop construct, where tasks are forked at the beginning of the loop and joined at the end of the loop, therefore these fork-join actions are essentially barriers and have the simi-

lar overhead to the barriers in VOPP. In multiple-iterative cases such as GE and SOR, the cumulative task scheduling and synchronization overheads can be considerable. Therefore, the Maotai model would be more suited for these problems.

Mandelbrot is an embarrassingly-parallel algorithm. This application demonstrates the slight performance advantage of the VOPP SPMD model, in which a task queue is used to balance workload in the program, instead of using general runtime schedulers as in OpenMP and Cilk. This result has also demonstrated that our implementation of the system queue is efficient.

In IS, the performance advantage seen in Maotai 2.0 over other models can be attributed to the split of global keyden array into $nproc$ views. In the global keyden construction step, each process updates *all* global keyden parts in the round-robin fashion, starting from the $proc_{id}^{th}$ part. Here, the SWMR view access pattern removes the need for barriers for preventing data race due to multiple processes updating an element simultaneously. This removal of barriers can contribute to the performance gain by the VOPP program.

In NN, since multiple items are updated by multiple processes at the end of the iteration, barriers are still used in the VOPP program. Therefore, it has the same synchronization overhead of other models. However the performance of Maotai 2.0 is still comparable to other models, which shows that being data race free has little impact on performance.

However, the SWMR model in VOPP does have its limitations in cases where the access pattern changes in every iteration. In those cases, view data must be copied to a local buffer of a process, where the process works on the data. After the data is processed, the view is acquired again by the process and the results copied back to the view. In our applications, Mergesort is such an example. In Mergesort, the resultant excessive memory-copying renders the implementation unscalable (Refer to VOPP-SWV in Figure 8). For this application, VOPP does trade off some programming convenience and performance for data race prevention. However, Maotai 2.0 has provided a Multiple Writer View (MWV) to offer the programming convenience for experienced programmers. A MWV is a view that can be accessed at *different locations* simultaneously by multiple processes. Therefore, it is up to the programmer to make sure there is no data race in a MWV. In contrast to other programming models, the data races of a MWV are *confined* inside the view should they occur. This alternative MWV implementation allows multiple processes to work directly on the view and avoid memory copying. With MWV, the speedup of Mergesort in Maotai 2.0 is comparable to other shared-memory models. Figure 8 shows, at $n = 32$, Maotai 2.0 (refer to VOPP-MWV) is 1% faster than OpenMP, but 9 % slower than Cilk.

Cilk performs very well in cases like Mergesort and NN. This can be attributed to its recursive task decomposition that ensures cache locality [20].

The parallel for-loop in OpenMP allows easy specification of data-parallelism. However, it would introduce a task-scheduling cost, especially when the workload is fixed and no load-balancing is required. The lower speedups of GE, SOR and NN of OpenMP can be attributed to this parallel for-loop overhead. Although Cilk++ cannot be benchmarked in this experiment because it does not support sparc64-smp, its equivalent construct *cilk_for* can also have the similar task-scheduling overhead.

5 Conclusions and Future Work

Our data race prevention scheme based on views proves to be efficient and adds little extra overhead to parallel programming systems. Though there is some memory wastage due to page alignment in the implementation, architectural support for variable-size pages will significantly reduce the wastage. Even with a fixed page size, view constructs are useful to remove data races.

With the advanced features in Maotai 2.0, the performance and programmability of VOPP are enhanced. Though strict SWV views are rigid for some applications like Mergesort, Maotai 2.0 offers MWV views to allow programmers to fall back to traditional shared memory programming, with the risk of data races that are confined in a single MWV view.

Performance results demonstrate that Maotai 2.0 is very competent among modern parallel programming models, even with the unique data race prevention scheme.

In the near future, we will investigate if this data race prevention mechanism can be used for data race detection as well. We would like to provide programmers an alternative debugging mode in which the memory protection mechanism is used to detect data races. At runtime, the programmers can choose to disable the data race prevention mechanism if they are sure there is no data race, so that the memory wastage in the prevention scheme can be avoided at runtime.

We will also support more advanced views in VOPP. One of them could be Transactional Memory View (TMV). TMV is similar to a piece of transactional memory, except that the roll-back behavior is only applied to the involved TMV, not to other memory objects. TMV can help programmers to avoid deadlocks and locking problems, in addition to the features of VOPP such as data race free.

References

- [1] Eduard Ayguadé, Nawai Copty, Alejandro Duran, Jay Hoeflinger, Yuan Lin, Federico Massaioli, Xavier Teruel, Priya Unnikrishnan, and Guansong Zhang. The design of OpenMP tasks. *IEEE Transactions on Parallel And Distributed Systems*, 20(3):404–418, March 2009.
- [2] Michael A. Bender, Jeremy T. Fineman, Seth Gilbert, and Charles E. Leiserson. On-the-fly maintenance of series-parallel relationships in fork-join multithreaded programs. In *SPAA04*, June 2004.
- [3] Bershada, B. N., Zekauskas, and M. J. Midway: Shared memory parallel programming with entry consistency for distributed memory multiprocessors. Technical Report CMU-CS-91-170, Carnegie-Mellon University, 1991.
- [4] Robert L. Bocchino, Vikram S. Adve, Sarita V. Adve, and Marc Snir. Parallel programming must be deterministic by default. In *First USENIX Workshop on Hot Topics in Parallelism (HotPar)*, 2009.
- [5] Robert L. Bocchino, Vikram S. Adve, Danny Dig, Sarita Adve, Stephen Heumann, Rakesh Komuravelli, Jeffery Overbey, Patrick Simmons, Hyojin Sung, and Mohsen Vakilian. A type and effect system for deterministic parallel java. Technical Report UIUCDCS-R-2009-3032, University of Illinois at Urbana-Champaign, 2009.
- [6] Mihai Burcea, J. Gregory Steffan, and Cristiana Amza. The potential for variable-granularity access tracking for optimistic parallelism. In *The ACM SIGPLAN Workshop on Memory Systems Performance and Correctness (MSPC'08)*, 2008.
- [7] Luis Ceze, Pablo Montesinos, Christoph von Praun, and Josep Torrellas. Colorama: Architectural support for data-centric synchronization. In *HPCA'07*, 2007.
- [8] J-D Choi, B. P. Miller, and R. H. B. Netzer. Techniques for debugging parallel programs with flowback analysis. *ACM Transactions on Programming Languages and Systems*, 13:491–530, 1991.
- [9] Cilk Arts Inc. *Cilk++ User Guide*.
- [10] Anne Dining and Edith Schonberg. An empirical comparison of monitoring algorithms for access anomaly detection. In *The 2nd ACM SIGPLAN Symposium on Principles & Practice of Parallel Programming (PPoPP)*, pages 1–10, 1991.
- [11] Perry A. Emrath, Sanjoy Ghosh, and David A. Padua. Event synchronization analysis for debugging parallel programs. In *Supercomputing'91*, pages 580–588, 1991.
- [12] Dave Field, Deron Johnson, Don Mize, and Robert Stober. Scheduling to overcome the multi-core memory bandwidth bottleneck. <http://www.platform.com/>, 2007.
- [13] William Gropp, Ewing Lusk, and Rajeev Thakur. *Using MPI : Portable Parallel Programming with the Message-Passing Interface*. MIT Press, 2nd edition, 1999.
- [14] David P. Helmbold, Charles E. McDowell, and Jian-Zhong Wang. Analyzing traces with anonymous synchronization. In *The 19th International Conference on Parallel Processing (ICPP'90)*, pages 1170–1177, 1990.
- [15] Z. Huang and W. Chen. Revisit of View-Oriented Parallel Programming. In *Proc. of the Seventh IEEE Inter. Symp. on Cluster Computing and the Grid*, pages 801–810, 2007.
- [16] Zhiyi Huang, Wenguang Chen, Martin Purvis, and Weimin Zheng. VODCA: View-oriented, distributed, cluster-based approach to parallel computing. In *The 2006 International Workshop on DSM (DSM2006)*, In (CD-ROM) *Proc. of the IEEE/ACM Symposium on Cluster Computing and Grid 2006 (CCGrid06)*, Singapore, 2006. IEEE Computer Society.

- [17] Zhiyi Huang, Martin Purvis, and Paul Werstein. Performance evaluation of view-oriented parallel programming. In *In Proceedings of the 2005 International Conference on Parallel Processing (ICPP05)*, pages 251–258, Oslo, June 2005. IEEE Computer Society.
- [18] C. Jung, D. Lim, L. Lee, and Y. Solihin. Helper thread prefetching for loosely-coupled multiprocessor systems. In *Proc. of 20th IEEE Inter. Parallel & Distributed Processing Symp.*, 2006.
- [19] Tushara C. Karunaratna. Nondeterminator-3: A provably good data-race detector that runs in parallel. Master’s thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, 2005.
- [20] Charles E. Leiserson. A minicourse in multithreaded programming. Technical report, MIT Laboratory for Computer Science, 1998.
- [21] J. Lu et al. Dynamic helper threaded prefetching on the Sun UltraSPARC CMP processor. In *Proc. of the 38th Annual IEEE/ACM Inter. Symp. on Microarchitecture*, pages 93–104, 2004.
- [22] John Mellor-Crummey. On-the-fly detection of data races for programs with nested fork-join parallelism. In *Supercomputing’91*, pages 24–33, 1991.
- [23] Barton P. Miller and Jong-Deok Choi. A mechanism for efficient debugging of parallel programs. In *Proceedings of the SIGPLAN ’88 conference on Programming language design and implementation*, volume 23, pages 135–144, July 1988.
- [24] Robert H.B. Netzer and Sanjoy Ghosh. Efficient race condition detection for shared-memory programs with post/wait synchronization. In *The 21st International Conference on Parallel Processing (ICPP’92)*, 1992.
- [25] Bradford Nichols, Dick Buttlar, and Jacqueline Proulx Farrell. *PThreads Programming*. O’Reilly, 1996.
- [26] OpenMP Architecture Review Board. *OpenMP Application Program Interface Version 3.0*, May 2008.
- [27] Mark Pethick, Michael Liddle, Paul Werstein, and Zhiyi Huang. Parallelization of a backpropagation neural network on a cluster computer. In *International Conference on Parallel and Distributed Computing and Systems (PDCS 2003)*, 2003.
- [28] James Reinders. *Intel Threading Building Blocks : Outfitting C++ for Multi-core Processor Parallelism*. O’Reilly, 2007.
- [29] Stefan Savage, Michael Burrows, Greg Nelson, Patric Sobalvarro, and Thomas Anderson. Eraser: A dynamic race detector for multithreaded programs. In *The 16th ACM Symposium on Operating Systems Principles (SOSP’97)*, 1997.
- [30] L. Spracklen and S. G. Abraham. Chip multithreading: Opportunities and challenges. In *Proc. of Inter. Symp. on High-Performance Computer Architecture*, pages 248–252, 2005.
- [31] Sun Microsystems. *OpenSPARC T1 Microarchitecture Specification*, 2006.
- [32] Supercomputing Technologies Group, MIT Laboratory for Computer Science. *Cilk 5.4.6 Reference Manual*, 1998.
- [33] UPC Consortium. UPC specification, v1.2. Technical Report LBNL-59208, Lawrence Berkeley National Lab, 2005.
- [34] Rob F. van der Wijngaart and Michael Frumkin. NAS grid benchmarks version 1.0. Technical report, NASA Advanced Supercomputing Division, NASA Ames Research Center, 2002.
- [35] Mandana Vaziri, Frank Tip, and Julian Dolby. Associating synchronization constraints with data in an object-oriented language. In *The 33rd Annual ACM SIGPLAN Symposium on Principles of Programming Languages (POPL’06)*, 2006.
- [36] Barry Wilkinson and Michael Allen. *Parallel Programming Techniques and Applications Using Networked Workstations and Parallel Computers*. Prentice Hall, 2nd edition, 2005.
- [37] Emmett Witchel, Josh Cates, and Krste Asanović. Mondrian memory protection. In *ASPLOS-X 2002*, 2002.
- [38] Jiaqi Zhang, Zhiyi Huang, Wenguang Chen, Qihang Huang, and Weimin Zheng. Maotai: View-Oriented Parallel Programming on CMT processors. In *The 37th International Conference on Parallel Processing (ICPP’08)*, 2008.