

Delayed Consistency and Its Effects on the Miss Rate of Parallel Programs

Michel Dubois, Jin Chin Wang*, Luiz A. Barroso, Kangwoo Lee and Yung-Syau Chen

Department of EE-Systems
University of Southern California
Los Angeles, CA 90089-2562

*Tandem Computers Incorporated
M.P. Division
Austin, TX 78728

Abstract

In cache based multiprocessors a protocol must maintain coherence among replicated copies of shared writable data. In delayed consistency protocols the effect of out-going and in-coming invalidations or updates are delayed. Delayed coherence can reduce processor blocking time as well as the effects of false sharing. In this paper, we introduce several implementations of delayed consistency for cache-based systems in the framework of a weakly-ordered consistency model. A performance comparison of the delayed protocols with the corresponding On-the-Fly (non-delayed) consistency protocol is made, through execution-driven simulations of four parallel algorithms. The results show that, for parallel programs in which false sharing is a problem, significant reductions in the data miss rate of parallel programs can be obtained with just a small increase in the cost and complexity of the cache system.

1.0 Introduction

The design of shared memory multiprocessors that can scale up to large number of processors is a current topic of active research. It has been argued that technological constraints will ultimately put a limit on the processing rate of uniprocessors. Therefore, future systems will need to incorporate some form of parallelism. The shared-memory model appears at this point in time to be the choice parallel architecture for general-purpose computing.

These systems must efficiently support parallel multithreaded applications, as well as single thread processes, for three reasons. First, users need to run programs on a single processor. Second, some applications have very limited parallelism. Third, serial bottlenecks in the code set an upper limit on achievable speedups for large number of processors [3]; it is therefore critical that single threads run at peak efficiency.

This work is funded by NSF under Grant No. CCR-870997.

There are two major problems in shared memory multiprocessors: shared memory bandwidth and shared-memory access latency. Both problems can be addressed by private caches associated with each processor [24]. Most processor accesses are satisfied by the cache, at processor speed. However, coherence must be maintained among caches. Every time a miss occurs in a cache or every time a processor needs to modify a cache block present in several caches, the cache controller must access the global memory through an interconnection. During each such access, the processor and the cache are usually blocked. The time during which the processor is blocked will be referred to as a *penalty*. Penalties can be very high, especially for fast uniprocessors in large-scale multiprocessor configurations. The miss penalty in current systems is of the order of 10 to 20 processor cycles. In the future we can expect miss penalties in excess of 100 cycles [17]. Therefore, there is a need to reduce these penalties.

In some parallel programs a large fraction of the misses are caused by false sharing. False sharing is the sharing of cache blocks without actual sharing of data. It occurs because cache blocks contain more than one data item. False sharing results in non-optimum protocols. In the case of a write-invalidate protocol, such as the Illinois protocol [19], more invalidations are sent than strictly needed by the parallel application and its data-sharing requirements. Invalidations create traffic and delays in the processor issuing them; moreover they increase the miss rate, because an invalidated block must be reloaded if it is accessed again. The situation is similar in write-broadcast protocols, such as the Firefly protocol [27]. In these protocols sharing is detected dynamically and multiple copies of the same block can be modified at the same time by different processors, provided modifications are broadcast to all processors with a copy. The update traffic should be limited to data elements that are actually shared; however, because of false sharing, a lot of redundant updates corresponding to different data elements in the same block are propagated.

To reduce the effects of penalties and of false sharing we propose to delay consistency, under a weakly ordered

model [2] called release consistency [15]. In general, in a cache-based system, a Store may cause an invalidation (write-invalidate protocol) or updates to other caches (write-broadcast protocol). To keep this discussion general, we will refer to both as *coherence updates*. In present systems, when a coherence update must be sent, the cache controller is blocked while the update signals propagate and are acknowledged. If these update signals are buffered so that the cache is not blocked during the propagation of the signals, we say that coherence is *send delayed*. Note that this buffering is different from the usual Store buffer needed in systems with write-through caches (see Figure 1). Coherence can also be *receive delayed*. Namely, if a coherence update signal reaches a cache, the effect of the update can be temporarily buffered.

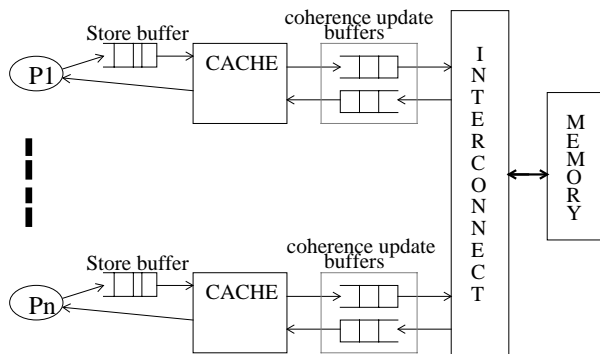


Figure 1: System structure with Store buffers and coherence update buffers

By delaying the sending of updates, update propagation can be overlapped with cache activity. In the protocol that we will describe, a processor or a cache never blocks on a Store hit, even if it hits on a non-unique, non-exclusive copy. Moreover, in a proposed variant of the protocol, the processor cache never blocks on a Store miss as well. Therefore, latency of Stores should be reduced. Also, the delaying of update propagation allows multiple processors to have dirty copies, and increases the concurrency of accesses to shared modifiable blocks.

The major contributions of this paper are the specifications of two delayed protocols derived from Censier and Feautrier's directory scheme [8], the hardware implementation details of both protocols, as well as simulation results showing the reduction in false sharing misses. In the following, we first present some background and the false sharing problem. In Section 4.0 we describe the protocols and their implementation. Performance results derived using execution-driven simulations are shown in Section 5.0, followed by some concluding remarks. Readers who are not familiar with the multicache consistency problem should consult survey papers [24] before proceeding.

2.0 False Sharing

In parallel applications, shared data structures are partitioned statically or dynamically and different processes work on different partitions of the structures. In general, partition boundaries do not coincide with cache block boundaries. As a result, cache blocks are shared while no data is actually shared. This gives rise to false sharing transitions [26], which create coherence or miss activity which would not happen if each cache block contained a single data item.

To demonstrate occurrences of false sharing we will show two simple examples. The first example is an algorithm with static partitioning of the data, the S.O.R. iterative algorithm to solve Poisson's equation on a square domain [28]. In this algorithm, an array (grid) of iterate components is updated iteratively by a linear combination of the iterate and its four neighbors in the 2-D grid. In the example of Figure 2.a, the grid has been partitioned among four processors. There are private iterate components and shared iterate components, as indicated in the Figure. In a shared memory organized as a linear address space, the array will be stored row-wise or column-wise. Assume that it is stored row-wise (i.e., first row 1, then row 2, and so on), and assume that the row size is not a multiple of the block size. Then false sharing (and true sharing) occurs also for blocks such as block 1.

The second example, in Figure 2.b, is an algorithm with dynamic partitioning of the shared data structure, the dynamic quicksort algorithm [22]. In this algorithm, a processor acquires exclusive access to a subfile, estimates a "pivot" element, and splits the subfile in two around the pivot. False sharing occurs at the boundaries between adjacent subfiles. The boundary between two subfiles cannot be predicted at compile time.

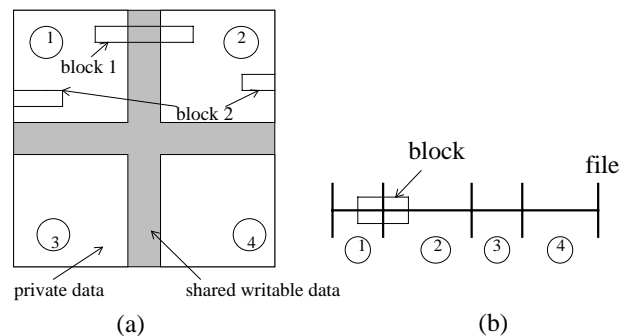


Figure 2: Illustration of false sharing in the S.O.R. (a) and Quicksort (b) algorithms

Because of false sharing a block may "ping-pong" several times between two processors, even if they reference

different data elements in the block. For a given number of processors, the effect of the block size is very similar to the effect of the block size in uniprocessor systems, but for different reasons. As the block size increases, the miss rate curve first decreases because of spatial locality, then the trend is reversed and the miss rate increases for larger block sizes. This behavior is observed even in caches of infinite sizes and is due to the increase in false sharing, which quickly offsets the gains due to spatial locality.

the effect of false sharing as the block size increases is small. In the second curve (worst case - dotted line), processor 2 is slightly slower so that it reaches a block such as block 2 (i.e., a block at the end of a row) at the same time as processor 1 (and similarly for processors 3 and 4); here the effect of false sharing is maximum. Whether the best or the worst case happens in an actual implementation depends on the relative speed of the processors and on the order in which the processors reach and execute the barrier synchronization.

The plots for the quicksort are shown in Figure 3.b; the number of processors is varied from 2 to 32 and the file to sort is made of 32K random integers drawn from a uniform distribution; each point is the average of the number of misses for 10 independent files. For 32 processors, the miss rate curve bottoms out for block sizes of $8 \times 4 = 32$ bytes.

In a system where we want to support efficiently both single and multiple thread programs, we are left with a *dilemma*. Namely, single thread programs benefit from bigger block sizes, while bigger blocks are detrimental to the performance of parallel threads.

3.0 Prior Related Work

In the IBM 3033, a multiprocessor with write-through caches, a mechanism called the *BIAS filter* was implemented [6]. In this buffer, invalidations with the same block address and coming from other processors are filtered before they reach the cache, in order to reduce the number of cache cycles needed by invalidations. In [11], buffering of memory accesses at both the sending and the receiving processors was studied in the context of *sequential consistency* [14] and weakly ordered protocols. For sequential consistency to hold it was required that buffered received invalidations have higher priority than accesses from local processors.

In [20], the notion of *isolated caches* was introduced, and it was informally shown that received and buffered invalidations may have lower priority than local processor accesses, in the context of sequential consistency. In [1], a theoretical model called *lazy caching* is developed, and it is proved that consistency can be send delayed and receive delayed in sequentially consistent systems. In this scheme, the sending of a coherence update can be delayed until the next Read (hit or miss) in the local processor. Moreover, the receiving of a coherence update in a cache may be delayed until the next miss in that cache. These conditions are very restrictive, especially for the sending of updates.

Under weak ordering [2][11], protocols can be delayed further, and Distributed Shared Memory system proposals try to take advantage of this. In *Munin*, a distributed shared memory system under development at Rice University [4], delayed consistency is advocated at the

Figure 3: Total number of misses for S.O.R (a) and quicksort (b) algorithms.

These effects are clear from the curves of Figure 3, which show the effect of false sharing on the total number of misses for executions of the S.O.R. algorithm (Figure 3.a) and the quicksort (Figure 3.b). The results in these Figures were obtained through execution-driven simulations [9]. In these simulations, all caches have infinite sizes and each simulated processor executes in turn until it accesses a shared data or executes a synchronization primitive; at that point, the simulator simulates a different processor. This is done in a round-robin fashion. In Figure 3.a we have plotted the total number of shared-data misses for the S.O.R. algorithm with four processors, a grid size of 128×128 and 100 iterations (the data structure size corresponding to this grid is actually 130×130 because of boundary conditions). Two curves are shown: in one curve, it is assumed that all processors are working at the same speed and start each iteration at the same time (best case - plain line); in this case

software level. Namely, updates by a thread of replicated objects are delayed until such a time that they can be detected by another thread, under a weakly-ordered model of consistency called *loose coherence*. A *delayed update queue* is maintained for each object and updates on an object are propagated at synchronizations. In [7], a delayed consistency protocol implemented on pages at the software level is outlined. In these papers delayed consistency is linked to the problem of false sharing for the first time.

4.0 Delayed Protocols

In a system with delayed consistency, synchronization variables must be stored in different regions of shared memory than other shared data. Accesses to the region of memory reserved for synchronization variables are not subject to delays, so that an On-the-Fly protocol is enforced on these variables. Therefore, the processor must have the ability to distinguish between synchronization variables and other variables. This is usually easy to do.

In the following we specify three protocols. The specification of each protocol is in three parts. The first part is a specification of block states in the caches and the system directory. The second part is a description of transactions between caches and memory to implement the protocol. The third part is the algorithm for the control of the cache.

4.1 On-the-Fly Protocol

This is the non-delayed protocol. Coherence actions are taken immediately and, during their propagation, the cache controller does not accept any request from the processor. This protocol was first introduced by Censier and Feautrier [8]. We reproduce the specification of this protocol for future reference in the paper.

Block states

(a) Cache states.

If a cache block is Valid, then the cache may be an Owner (i.e. the copy is unique in the system and it is dirty) or a Keeper (i.e. there may be copies in other caches and all copies are clean).

(b) System Directory states.

There is one Modified bit and P Presence bits (one per processor) per block in memory. A Presence bit is set only if the corresponding cache is an Owner or a Keeper of the block. If a cache is the Owner of the block (in which case only one P bit is set), then the Modified bit is also set.

Memory commands

(a) Issued by a memory controller to the caches.

- **Inv (Invalidate):** the receiving cache is either a Keeper (the cache controller must invalidate its copy of the block), or the Owner (the cache controller must invalidate its copy of the block and send it to memory).
- **UpdM (Update Memory):** the receiving cache must be an Owner. The copy of the block is sent to memory and the cache becomes a Keeper of the block.

(b) Issued by a cache to the memory controller.

- **ReqO (Request Ownership):** the memory controller sends an Inv command to all caches with a copy and the requesting cache becomes the Owner.
- **ReqOC (Request Owner Copy):** same as ReqO, but the memory copy of the block is also sent to the requesting cache.
- **ReqKC (Request Keeper Copy):** if there is an Owner, the memory controller sends an UpdM command to the Owner. The memory copy of the block is then sent to the requesting cache (which becomes a Keeper).
- **WB:** the block is written back to memory.

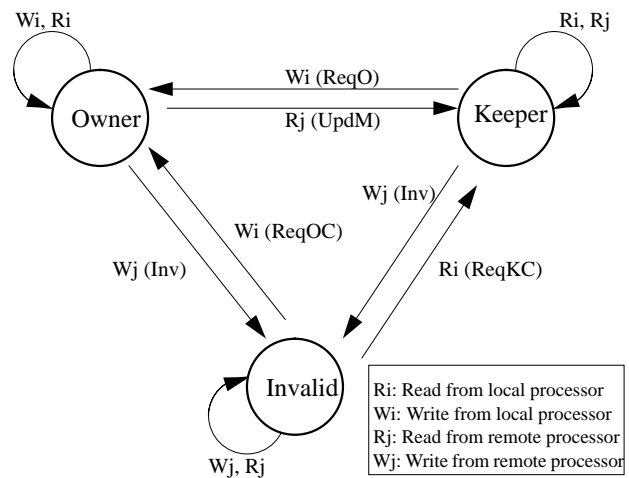


Figure 4: Cache states transition diagram for the On-the-Fly protocol

Cache algorithm

For the various types of cache accesses, the cache controller takes the following actions (Figure 4).

- **Read hit:** no action is taken.
- **Write hit:** if the cache is Keeper, a ReqO command is sent to the memory controller, otherwise no action is taken.
- **Read miss:** a ReqKC command is sent to the memory controller.

- Write miss: a ReqOC command is sent to the memory controller.
- Replacement: if the cache is the Owner, a WB command is issued to memory.

4.2 Receive Delayed Protocol

When an Inv signal is received by a cache, the invalidation does not need to reach the cache until the next Lock instruction executed by the local processor. The behavior is still correct because the programming model in weakly-ordered systems forbids accesses to a shared writable data outside a critical or semi-critical section. Therefore, between the times when the Inv is received and the next Lock instruction is executed in the local processor, any Read to the block must be for a different word or byte in the block than the one modified. We say that the block copy is *Stale*. Right after the processor acquires a Lock, all the Stale blocks in the cache must be invalidated; otherwise, Stale blocks are treated the same way as Fresh (non-Stale) blocks in the cache, and the protocol is very similar to the above On-the-Fly protocol.

Block states

(a) *Cache states.*

As for the On-the-Fly protocol, a cache may be a Keeper or the Owner of a block. However, in addition, a Valid block may also be Stale (i.e. Valid locally but Invalid in the system directory).

(b) *System Directory states.*

Same as for the On-the-Fly protocol.

Memory commands

(a) *Issued by a memory controller to the caches.*

The commands are Inv and UpdM as in the On-the-Fly protocol. However, the cached copy becomes Stale instead of Invalid when an Inv command is received by a cache.

(b) *Issued by a cache to the memory controller.*

The commands are ReqO, ReqOC, ReqKC, and WB as in the On-the-Fly protocol. The only difference is that the memory controller sends Inv and UpdM commands only to copies that are Fresh (non-Stale).

Cache algorithm

Similar to the On-the-Fly algorithm. However a Write hit on a Stale copy triggers a ReqOC command to memory and a block reload. Also, upon execution in the processor of a Lock instruction, all Stale blocks become

Invalid (see Figure 5).

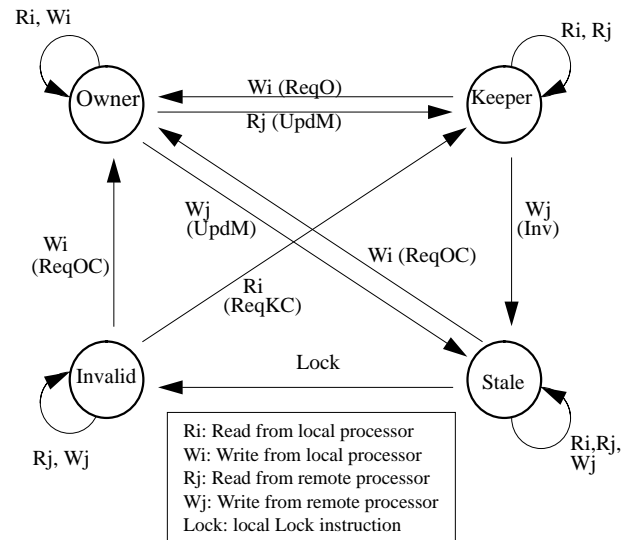


Figure 5: Cache state diagram for the Receive Delayed protocol.

4.3 Send-and-Receive Delayed Protocol

Receive Delayed protocols are effective at reducing the number of false sharing misses. They work well for blocks that are read by several processors, and modified by a small subset of these processors: the modifying processors experience false sharing transitions on each Write, but the reading processors do not (between the execution of two Lock instructions). Usually, however, between two Lock instructions, processors may read and write different words of the same block, or even sometimes only write into the same block. If two Writes to the same block can occur simultaneously, then they must be for different parts of a block. Therefore, such Writes can be executed without acquiring a unique copy. Writes executed on a non-unique copy must be propagated at the next execution of an Unlock instruction [15].

We have to keep track of all partial modifications of a non-owned block copy so that they can be written back to memory at the next Unlock instruction. Copies of these modifications must be kept in an Invalidation Send Buffer (ISB). At most this buffer must have as many entries as there are block frames in the cache. However, the optimum size of this buffer is probably much less than that. First of all, if the buffer is too large, Unlock instructions are very costly because a large number of memory updates must be propagated. Second, because of the locality of accesses to shared blocks, we can expect that accesses to a shared block by a given processor occur in runs or bursts. Therefore, we

will observe diminishing returns in the miss rate curves as the buffer size is increased. The ISB should be a small, fully associative buffer, capable of containing a few blocks.

Besides being receive delayed, the following protocol is also send delayed in the sense that ownership is acquired either on a Write miss or when a modified block is removed from the ISB. While the block is in the ISB, multiple Writes can be done locally and other processors can also read and write the block.

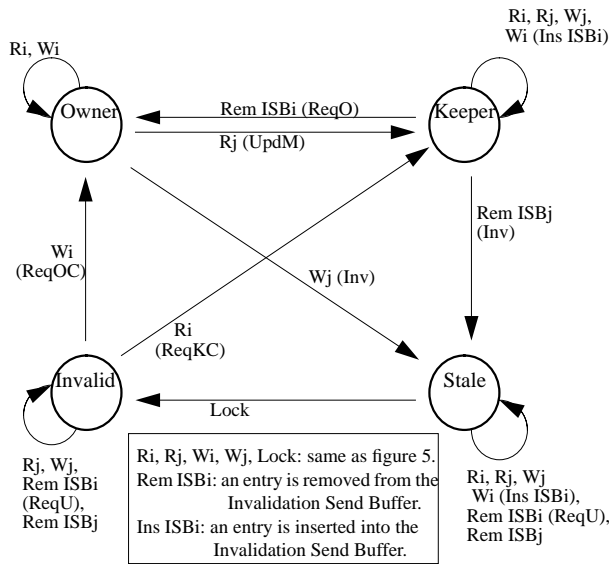


Figure 6: Cache state diagram for the Send-and-Receive Delayed protocol.

Block states

(a) Cache states.

Similar to the states for the Receive Delayed protocol. The difference is that a Keeper copy or a Stale copy may be Modified locally. In these cases, there is a copy of the modifications to the block in the ISB.

(b) System Directory states.

Same as for the On-the-Fly protocol.

Memory commands

(a) Issued by a memory controller to the caches.

Same as for the Receive Delayed protocol.

(b) Issued by a cache to the memory controller.

Commands ReqO, ReqOC, ReqKC and WB are

needed and trigger the same actions as in the Receive Delayed protocol. However, a new command must be introduced to notify the possible Owner or Keepers in case a Modified Keeper or Stale copy is removed from the ISB or is replaced in the cache. This new command is called ReqU (Request Update). Besides a partial update of memory (based on the words modified in the ISB), Inv signals are sent to all Fresh copies in the system.

Cache algorithm

For different types of cache accesses, the cache controller takes the following actions (Figure 6).

- Read hit: no action is taken.
- Write hit: if the copy is Owned, no action is taken.

If the cache is a Keeper or if the copy is Stale, then an entry is allocated to the block in the ISB (unless an entry is already present); the new values are stored in the buffer. Therefore, a Write hit is always a local operation.

• Read miss: a ReqKC command is sent to the memory controller.

• Write miss: a ReqOC is sent to the memory controller.

• Removing an entry from ISB: the Keepers or the Owner (if any) must be notified; if the block was Invalid¹ or Stale in the cache then a ReqU request is sent to memory and the block stays Invalid or Stale in the cache. If the cache was a Keeper then a ReqO command is sent to the memory controller (in this case, the part of the local copy that has not been modified is consistent with memory).

• Executing a Lock instruction: all Stale blocks become Invalid right after the successful acquisition of the Lock.

• Executing an Unlock instruction: all entries must be removed from the ISB right before releasing the Lock.

• Replacement: if the cache is the Owner, then a WB command is sent to memory. A ReqU command is sent to the memory if the copy is a Stale or a Keeper copy and it has been Modified; in this case, the modifications are stored in the ISB.

4.4 Hardware Implementation

At the system level, the Send-and-Receive Delayed protocol is very similar to the On-the-Fly protocol. The only major difference is the ReqU command, which uses controls needed for the ReqO command. Added complexity is in the implementations of the Stale state in the cache, of the Invalidation Send Buffer, and of the partial updates of memory.

1. Since we allow that modifications of a block may be in the ISB even if the block is not in cache, a miss in the cache must first check the ISB.

Each block frame in the cache requires a Valid bit (V) and an Ownership bit (O). A Stale bit (S) is also needed to distinguish between a Stale and a Fresh block. When an invalidation reaches the cache, the S-bit is set and the V-bit is reset to mark the block as Stale. When a Stale block is accessed, the reset V-bit is masked by the S-bit, i.e. if the S-bit is set then the value of the V-bit is ignored. After the successful execution of a Lock instruction, all S-bits must be reset in the cache (at this point all Stale blocks become Invalid). A simple way to do this is to store the S-bits in a clearable SRAM chip (for example, part No. SN 74ACT2154 in [25]).

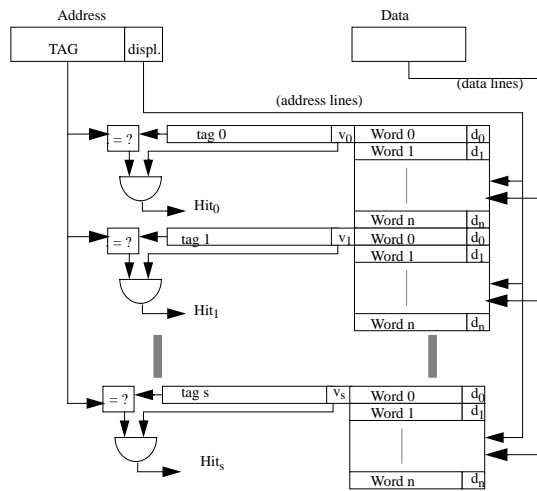


Figure 7: Block diagram of an ISB

The ISB should contain no more than a few entries (2, 4 or 8). Each entry contains the block address, one Valid bit, and a Dirty bit for each word of the block (see Figure 7). The buffer is accessed associatively with the block address when a Store is done in the cache. If the copy in cache is a Keeper or a Stale copy, as indicated by the S and O bits, the ISB is consulted. If an entry already exists for that block the word is updated in the buffer and its dirty bit is set. Otherwise an entry must be allocated in the buffer. It is important to manage the ISB efficiently to avoid slowing down the processor on a Write hit. In particular, there should always be at least one free entry in the ISB.

When a non-owned block copy updates memory the block in the ISB (containing both modified and “empty” words) is sent to the memory controller². The memory controller can use the dirty bits to enable/disable the Store of each word of the block.

2. Alternatively, in order to reduce the traffic, only the modified bytes could be sent to memory. However this optimization is probably not cost-effective.

4.5 Alternate Designs

We propose some refinements which could further improve the performance of the Send-and-Receive Delayed protocol.

In the present protocol, when a Keeper copy is removed from the ISB, the cache gains ownership for the block, in order to reduce the overhead for private, non-shared blocks. On the other hand, when a Stale block is removed from the ISB, all Fresh copies are notified and the block remains Stale. Besides leaving it Stale, we could envision two other strategies: invalidating the block, or acquiring ownership for the block. It may be useful to have different strategies, depending on *when* the block is removed from the ISB.

The ISB must be managed in such a way that Write hits are never slowed down, and synchronization points are executed efficiently. We propose that, in normal operation, one or two free entries be always maintained in the ISB, and that blocks be removed when the number of free entries falls below one or two. However, before an Unlock instruction, a special instruction “Prepare to Synchronize” could be issued by the compiler to remove all entries in the ISB.

In the delayed protocol, the cache never blocks on a Write hit. A simple extension is to avoid blocking on Write misses as well. A Write miss simply fills a word in an entry of the ISB (such a copy which exists in the ISB but not in the cache should be considered as Invalid¹). The block is only loaded at the first following Read miss.

Table 1: Characteristics of the Applications

| Application | # of Data Accesses | # of Reads | # of Writes |
|-------------|--------------------|------------|-------------|
| SOR | 9,829,600 | 8,192,000 | 1,637,600 |
| QSORT | 969,631 | 715,351 | 254,280 |
| FLOYD | 354,913 | 350,278 | 4,635 |
| INTERPOLATE | 17,028 | 8,836 | 8,192 |

5.0 Effect on False Sharing Misses

All simulation results reported in this Section were derived for infinite cache sizes and for an ISB of size two blocks (no gain was observed beyond an ISB of size 2 blocks). We only show data miss rate figures. The miss rates on data for the three protocols can be seen in Tables 2 through 6, for the set of programs analyzed. In the discussion we focus on block sizes of 16 words (64 bytes), which is the block size adopted in the SCI (Scalable Coherence Interface) protocol [12]. It is also a good choice for uniprocessor

caches.

The data access patterns of the four parallel programs used in our performance studies are summarized in Table 1. The programs are further described below.

1) SOR: 100 iterations of the Successive Over Relaxation iterative algorithm to solve Poisson's equation on a square domain [28] with single precision (32 bit) floating point numbers. The grid has size 128x128 (actually 130x130 because of boundary conditions) and is partitioned in four quadrants. Each quadrant is allocated to one processor. Partitioning of data is static. The size of one word is 32 bits. The miss rate due to false sharing is very sensitive to the actual timing of processor execution, as discussed in Section 2.0. Results are reported for the worst and best cases of false sharing.

Table 2: SOR (best case) - miss rates (%)

| Block size (bytes) | 4 | 8 | 16 | 32 | 64 | 128 |
|--------------------------|------|------|------|------|------|------|
| On-the-Fly | 0.69 | 0.60 | 0.69 | 0.87 | 0.96 | 1.00 |
| Receive Delayed | 0.69 | 0.54 | 0.60 | 0.78 | 0.92 | 0.98 |
| Send-and-Receive Delayed | 0.69 | 0.54 | 0.60 | 0.75 | 0.86 | 0.92 |

Table 3: SOR (worst case) - miss rates (%)

| Block size (bytes) | 4 | 8 | 16 | 32 | 64 | 128 |
|--------------------------|------|------|------|------|------|------|
| On-the-Fly | 0.69 | 0.60 | 0.69 | 1.03 | 1.61 | 2.71 |
| Receive Delayed | 0.69 | 0.54 | 0.60 | 0.85 | 1.17 | 1.75 |
| Send-and-Receive Delayed | 0.69 | 0.54 | 0.60 | 0.78 | 0.88 | 0.92 |

2) QSORT: Quicksort [22] of a 32K file of 32-bit integers. Results are the average of the results of 10 runs for 10 different random files. Partitioning of the data is dynamic. The number of processors is 16. The size of one word is 32 bits.

Table 4: QSORT - miss rates (%)

| Block size (bytes) | 4 | 8 | 16 | 32 | 64 | 128 |
|--------------------------|-------|------|------|------|------|------|
| On-the-Fly | 18.62 | 9.89 | 5.75 | 4.12 | 4.10 | 5.15 |
| Receive Delayed | 18.62 | 9.86 | 5.67 | 3.86 | 3.36 | 3.56 |
| Send-and-Receive Delayed | 18.62 | 9.86 | 5.49 | 3.35 | 2.36 | 1.91 |

3) FLOYD: Single source shortest path problem using the Floyd-Warshall algorithm [10]. There is a *path* and a *cost* array. Each graph is a random graph of 128 nodes with maximum connectivity of 96. Each result is an average over 10 runs of the algorithm on all nodes of 10 different random

graphs. The two arrays are frequently read but rarely modified. Data partitioning is done dynamically. The number of processors is 16. All data are 4 bytes, which is also the size of a word.

Table 5: FLOYD - miss rates (%)

| Block size (bytes) | 4 | 8 | 16 | 32 | 64 | 128 |
|--------------------------|-------|-------|-------|------|------|-------|
| On-the-Fly | 28.38 | 17.91 | 12.22 | 9.46 | 9.47 | 10.05 |
| Receive Delayed | 28.38 | 17.77 | 12.06 | 9.22 | 8.98 | 9.27 |
| Send-and-Receive Delayed | 28.38 | 17.76 | 12.05 | 9.20 | 8.96 | 9.23 |

4) INTERPOLATE: Picture interpolation program. Only one out of every 9 pixels in a picture has initial values. Based on a linear interpolation the missing pixel values are computed. The size of the picture after interpolation is 96x96. The picture is divided into 8 rectangles and each rectangle is assigned to one processor. Partitioning is static. Data is either write-only or read-only. All data are 8-bit pixels, and the size of a word is also 8 bits.

Of all these programs FLOYD exhibits the least number of false sharing transitions. The reason is that the two arrays are accessed mostly randomly and are read most of the time. INTERPOLATE exhibits large false sharing effects. It is not iterative like the other three algorithms. Processes read the known pixel values, compute each unknown value and store it. There is no Read/Write sharing of data, and consequently there is no locking in the whole program, except at the beginning (to fork the processes) and at the end (to terminate). There is some Read sharing, but all the coherence activity is due to false sharing transitions caused by Stores. This type of sharing is very frequent. It occurs for example in a Doall loop where successive iterations of a loop compute the components of a vector or an array, and are allocated to different processors. A classical example is the Doall loop computing the product of two matrices A and B and storing the result in a matrix C.

Table 6: INTERPOLATE - miss rates (%)

| Block size (bytes) | 4 | 8 | 16 | 32 | 64 | 128 |
|--------------------------|-------|-------|-------|-------|-------|-------|
| On-the-Fly | 16.16 | 10.15 | 32.28 | 52.20 | 53.89 | 53.61 |
| Receive Delayed | 16.16 | 10.15 | 7.14 | 5.64 | 5.07 | 3.38 |
| Send-and-Receive Delayed | 16.16 | 10.15 | 7.14 | 5.64 | 5.07 | 3.38 |

The results above show that delayed consistency is effective at reducing the number of misses when the effect of false sharing transitions is large. FLOYD has very few false sharing transitions and therefore it shows little gain. The best cases of SOR have few false sharing transitions, and most of the false sharing transitions occur across barrier

synchronizations, so that the Stale bit and the ISB are not very effective (a reduction of misses of only 10% is observed for a block size of 64 bytes); for the worst case of SOR, delaying consistency is extremely effective in the case of a block size of 64 bytes, because false sharing effects are intense and occur mostly between barrier synchronizations (within the same sweep). The gains in QSORT (P=16) only occur for block sizes larger than 16 bytes, because the false sharing effects are very small for smaller blocks. For a block size of 64 bytes (file size of 32K), the reduction in the number of misses approaches 42% for a Send-and-Receive Delayed protocol. Finally, as expected, INTERPOLATION benefits the most from delayed consistency (90% reduction in the number of misses for block size 64), because all misses (except initial loading misses) are due to false sharing and because the algorithm does not require any synchronization during most of its execution.

6.0 Other Approaches

6.1 False Sharing

There have been several other proposed solutions to the problem of false sharing and the relatively low spatial locality of shared data accesses [16][26]. The first one consists of not caching shared writable data. This solution usually requires a T.L.B. (Translation Lookaside Buffer) to discriminate dynamically between cacheable and non-cacheable blocks. The problem with this solution is that entire data structures must be deemed non-cacheable if any item in the structure is shared and can be modified. Another proposition would be to use two data caches: one for shared data, and one for private data. These two caches could have different block sizes. Theoretically, caches could have two block sizes: one for shared data and one for private data, but the complexity of the implementation of such proposal has never been investigated. Bitar and Despain [5] have proposed to allocate one cache block per shared data item. This scheme would probably have to rely on the compiler to expand the shared data structure with dummy items. This will result in significant waste of memory and cache.

The compiler and programmer can try to reduce false sharing through wise data placements. For example, in the S.O.R. algorithm above the compiler could try to allocate an integer number of blocks per row; this would remove the problem for blocks such as block 2 in Figure 2. However, for blocks such as block 1, it will be difficult to achieve this in general, without wasting a lot of cache space or complicating drastically the addressing to contiguous array components. Real applications, while exhibiting the type of sharing present in the S.O.R. Poisson solver, are seldom as simple [28]. Usually, the boundary calculations are complex, and the region is not square but is irregularly shaped and has internal cavities; sometimes the grid mesh size is different for different areas of the grid. Compilers may have problems dealing with false sharing in those cases, even for blocks such as block 2.

Reduction of false sharing by the compiler or the programmer is also difficult in the case of dynamic data partitioning, such as in quicksort. In this case, the user would have to take into account the block size in the computation of the pivot, and this will result in complex algorithms, optimized at the source code level for one machine but not another. Some simple compiling techniques, which in some cases can reduce the effect of false sharing, are proposed in [26] and evaluations based on 16 and 32 processor systems look encouraging.

6.2 Latency Tolerance

Other cache mechanisms to tolerate latencies can be used besides delayed consistency. One such approach is non-blocking or *lockup-free* caches [13][21]. When a miss occurs in a non-blocking cache, the cache does not block the processor and services the miss concurrently with processor accesses. Several misses can be pending. Non-blocking caches are especially effective for second-level shared caches, for superscalar processors [23], and for programs in which in-cache prefetching can be done effectively [18]. Penalties due to Stores in systems with write-through caches can also be reduced effectively with a Store buffer between the processor and the cache.

Finally, the DASH multiprocessor implements a form of weak ordering protocol based on *Release Consistency* [15] and attempts to overlap Stores with computation in a system with write-through caches, mostly by using a Store buffer. However, the second level cache locks out the processor when a Store requires ownership. The DASH protocol is neither send delayed nor receive delayed, in the sense that we mean it here: the cache locks out the processor on a coherence update and there is no provision for delaying the reception of invalidations. In the DASH multiprocessor as long as the first level cache [15] hits on Read accesses while the second level cache gains ownership for a previous Store, consistency is in fact delayed. Some of the benefits of delayed consistency may be obtained through this two-level organization. Since the DASH protocol relies mostly on a Store buffer associated with the processor, the propagation of invalidations are time-critical and therefore results in complex deadlock-prone sequences of transfers between memory and caches.

7.0 Conclusion

In this paper, we have introduced two write-invalidate delayed consistency protocols, built as extensions of an existing protocol. Delays are obtained through one buffer and the addition of the Stale bit in the cache state. Implementation of an ISB allows the overlapping of cache activity and sending of invalidations.

Delayed consistency also reduces the number of false sharing transitions. This is important in systems supporting efficiently both parallel and single thread processes.

Significant reductions in the miss rate on data can be obtained with the addition of a stale bit, and further reduction was observed by adding a small ISB. These reductions are obtained with no assistance from the programmer or the compiler, which makes delayed consistency particularly useful for general-purpose multiprocessors. The only restriction is that parallel programs must access shared-writable data in critical or semi-critical sections.

Acknowledgments

The idea of the stale state in the cache to implement partial invalidations of blocks is due to Andrew Glew.

References

- [1] Y. Afek, G. Brown, and M. Merritt, "A Lazy Cache algorithm," *Proc. of the 1st ACM Symp. on Parallel Algorithms and Architectures*, pp. 209-223, Jun 1989.
- [2] S.V. Adve and M. D. Hill, "Weak Ordering-A New Definition," *Proc. of the 17th Int. Symp. on Computer Architecture*, pp.2-14, 1990.
- [3] G.M. Amdahl, "Validity of the single processor approach to achieving large-scale computing capabilities," *Proc AFIPS*, Vol. 30, pp.483-465, 1967.
- [4] J.K. Bennett, J.B. Carter, and W. Zwaenepoel, "Munin: Distributed Shared Memory Based on Type-Specific Memory Coherence," *Proc. of the 2nd ACM Symposium on Principles and Practice of Parallel Programming, SIGPLAN Notices* 25:3, Mar 1990, pp. 168-176
- [5] P. Bitar and A. Despain, "Multiprocessor Cache Synchronization: Issues, Innovations, Evolution," *Proc. of the 13th Annual Int. Symp. on Comp. Architecture*, June 1986, pp. 424-433.
- [6] B.M. Bean et al., "Bias Filter Memory for Filtering Out Unnecessary Interrogations of Cache Directories in a Multiprocessor System," U.S. Patent 4,142,234, Feb 1979.
- [7] L. Borrman and M. Herdieckerhoff, "A Coherency Model for Virtual Shared Memory," *Proc. of Int. Conf. on Parallel Processing*, Vol.2, pp.252-257, June 1990.
- [8] L.M. Censier and P. Feautrier, "A New Solution to Coherence Problems in Multicache Systems," *IEEE Trans. on Computers*, Vol. C-27, No. 12, pp. 1112-1118, Dec. 1978.
- [9] M. Dubois, F.A. Briggs, I. Patil, and M. Balakrishnan, "Trace-driven Simulations of Parallel and Distributed Algorithms in Multiprocessors," *Proc. Int. Conf. on Parallel Processing*, pp.909-916, Aug. 1976.
- [10] N. Deo, C.Y. Pang, and R.E. Lord, "Two Parallel Algorithms for Shortest Path Problems," *Proc. of the Int. Conf. on Parallel Proc.*, pp. 244-253, Aug. 1980.
- [11] M. Dubois and C.Scheurich, "Memory Access Dependencies in Shared Memory Multiprocessors," *IEEE Transactions on Software Eng.*, 16(6), pp. 660-674, June 1990.
- [12] D.V. James et al., "Scalable Coherent Interface," *IEEE Computer*, Vol. 23, No. 6, pp. 74-77, June 1990.
- [13] D. Kroft, "Lockup-free Instruction Fetch/Prefetch Cache Organization," *Proc. of the 8th Ann. Int. Symp. on Comp. Arch.*, pp.81-87, 1981.
- [14] L. Lamport, "How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs," *IEEE Transactions on Computers*, Vol. C-28, No.9, pp.690-691, Sept. 1979.
- [15] D. Lenoski et al., "The Directory-Based Cache Coherence Protocol for the DASH Multiprocessor," *Proc of the 17th Annual Int Symposium on Comp Arch*, pp. 148-159, Jun 1990.
- [16] R.L. Lee, P.C. Yew, and D.H. Lawrie, "Multiprocessor Cache Design Considerations," *Proc. of the 14th Int. Symp. on Computer Architecture*, pp. 253-262, June 1987.
- [17] Mudge, T. N. et al., "The Design of a Microsupercomputer", *IEEE Computer*, January 1991, pp.57-64.
- [18] Mowry, T. and Gupta, A., "Tolerating Latency through Software-Controlled Prefetching in Scalable Shared-Memory Multiprocessors", to appear in the *Journal of Parallel and Distributed Computing*, 1991.
- [19] M. Papamarcos and J. Patel, "A Low Overhead Coherence Solution for Multiprocessors with Private Cache Memories," *Proc. of the 11th Int. Symp. on Computer Architecture*, pp. 348-354, June 1984.
- [20] C. Scheurich, "Access Ordering and Coherence in Shared-Memory Multiprocessors," PhD thesis, University of Southern California, May 1989.
- [21] C. Scheurich and M. Dubois, "Lockup-free Caches in High-Performance Multiprocessors," *Journal of Parallel and Distributed Computing*, January 1991.
- [22] R. Sedgewick, "Quicksort", New York: Garland Publishing, Inc., 1980.
- [23] Sohi, G. and Franklin, M., "High-Bandwidth Data Memory Systems for Superscalar Processors", *APLOS IV*, 1991.
- [24] P. Stenstrom, "A Survey of Cache Coherence Scheme for Multiprocessors," *IEEE Computer*, Vol. 23, No. 6, Jun 1990.
- [25] Texas Instruments MOS Memory Data Book, pp. 7-135 to 7-147, 1989.
- [26] J. Torrellas, M.S. Lam, and J.L. Hennessy, "Shared Data Placement Optimizations to Reduce Multiprocessor Cache Misses," *Proc. of the 1990 Int. Conf. on Parallel Proc.*, Aug 1990, pp. 266-270.
- [27] C.P. Thacker, L.C. Stewart, and E.H. Satterthwaite, "Firefly: A Multiprocessor Workstation," *IEEE Trans. on Computers*, Vol. 37, No. 8, pp. 909-920, Aug.1988.
- [28] D. Young, "Iterative Solution of Large Linear Systems", Academic Press: New York, 1971.