

6

Multiprocessors and Thread-Level Parallelism

The turning away from the conventional organization came in the middle 1960s, when the law of diminishing returns began to take effect in the effort to increase the operational speed of a computer. ... Electronic circuits are ultimately limited in their speed of operation by the speed of light... and many of the circuits were already operating in the nanosecond range.

Bouknight et al., *The Illiac IV System* [1972]

... sequential computers are approaching a fundamental physical limit on their potential computational power. Such a limit is the speed of light...

A. L. DeCegama, *The Technology of Parallel Processing, Volume I* (1989)

... today's multiprocessors... are nearing an impasse as technologies approach the speed of light. Even if the components of a sequential processor could be made to work this fast, the best that could be expected is no more than a few million instructions per second.

Mitchell, *The Transputer: The Time Is Now* [1989]

6.1	Introduction	635
6.2	Characteristics of Application Domains	649
6.3	Symmetric Shared-Memory Architectures	658
6.4	Performance of Symmetric Shared-Memory Multiprocessors	670
6.5	Distributed Shared-Memory Architectures	687
6.6	Performance of Distributed Shared-Memory Multiprocessors	697
6.7	Synchronization	705
6.8	Models of Memory Consistency: An Introduction	719
6.9	Multithreading: Exploiting Thread-Level Parallelism within a Processor	723
6.10	Crosscutting Issues	728
6.11	Putting It All Together: Sun's Wildfire Prototype	735
6.13	Another View: Embedded Multiprocessors	751
6.14	Fallacies and Pitfalls	752
6.15	Concluding Remarks	758
6.16	Historical Perspective and References	765
	Exercises	780

Major changes

1. split up the longest sections
2. clearer discussion of the concept of thread and process
3. SMT and multithreading section
4. two another views
5. reordered the cross cutting issues--no big changes, just reordered

6.1 Introduction

As the quotations that open this chapter show, the view that advances in uni-processor architecture were nearing an end has been widely held at varying times. To counter this view, we observe that during the period 1985–2000, uni-

processor performance growth, driven by the microprocessor, was at its highest rate since the first transistorized computers in the late 1950s and early 1960s.

On balance, though, your authors believe that parallel processors will definitely have a bigger role in the future. This view is driven by three observations. First, since microprocessors are likely to remain the dominant uniprocessor technology, the logical way to improve performance beyond a single processor is by connecting multiple microprocessors together. This combination is likely to be more cost-effective than designing a custom processor. Second, it is unclear whether the pace of architectural innovation that has been based for more than fifteen years on increased exploitation of instruction-level parallelism can be sustained indefinitely. As we saw in Chapters 3 and 4, modern multiple-issue processors have become incredibly complex, and the increases achieved in performance for increasing complexity, increasing silicon, and increasing power seem to be diminishing. Third, there appears to be slow but steady progress on the major obstacle to widespread use of parallel processors, namely software. This progress is probably faster in the server and embedded markets, as we discussed in Chapter 3 and 4. Server and embedded applications exhibit natural parallelism that can be exploited without some of the burdens of rewriting a gigantic software base. This is more of a challenge in the desktop space.

Your authors, however, are extremely reluctant to predict the death of advances in uniprocessor architecture. Indeed, we believe that the rapid rate of performance growth will continue at least for the next five years. Whether this pace of innovation can be sustained longer is difficult to predict but hard to bet against. Nonetheless, if the pace of progress in uniprocessors does slow down, multiprocessor architectures will become increasingly attractive.

That said, we are left with two problems. First, multiprocessor architecture is a large and diverse field, and much of the field is in its youth, with ideas coming and going and, until very recently, more architectures failing than succeeding. Given that we are already on page 636, full coverage of the multiprocessor design space and its trade-offs would require another volume. (Indeed, Culler, Singh, and Gupta [1999] cover *only* multiprocessors in their 1000 page book!) Second, such coverage would necessarily entail discussing approaches that may not stand the test of time, something we have largely avoided to this point. For these reasons, we have chosen to focus on the mainstream of multiprocessor design: multiprocessors with small to medium numbers of processors (≤ 128). Such designs vastly dominate in terms of both units and dollars. We will pay only slight attention to the larger-scale multiprocessor design space (≥ 128 processors). At the present, the future architecture of such multiprocessors is unsettled and even the viability of that marketplace is in doubt. We will return to this topic briefly at the end of the chapter, in section 6.15.

A Taxonomy of Parallel Architectures

We begin this chapter with a taxonomy so that you can appreciate both the breadth of design alternatives for multiprocessors and the context that has led to the development of the dominant form of multiprocessors. We briefly describe the alternatives and the rationale behind them; a longer description of how these different models were born (and often died) can be found in the historical perspectives at the end of the chapter.

The idea of using multiple processors both to increase performance and to improve availability dates back to the earliest electronic computers. About 30 years ago, Flynn proposed a simple model of categorizing all computers that is still useful today. He looked at the parallelism in the instruction and data streams called for by the instructions at the most constrained component of the multiprocessor, and placed all computers in one of four categories:

1. *Single instruction stream, single data stream (SISD)*—This category is the uniprocessor.
2. *Single instruction stream, multiple data streams (SIMD)*—The same instruction is executed by multiple processors using different data streams. Each processor has its own data memory (hence multiple data), but there is a single instruction memory and control processor, which fetches and dispatches instructions. The multimedia extensions we considered in Chapter 2 are a limited form of SIMD parallelism. Vector architectures are the largest class of processors of this type.
3. *Multiple instruction streams, single data stream (MISD)*—No commercial multiprocessor of this type has been built to date, but may be in the future. Some special purpose stream processors approximate a limited form of this (there is only a single data stream that is operated on by successive functional units).
4. *Multiple instruction streams, multiple data streams (MIMD)*—Each processor fetches its own instructions and operates on its own data. The processors are often off-the-shelf microprocessors.

This is a coarse model, as some multiprocessors are hybrids of these categories. Nonetheless, it is useful to put a framework on the design space.

As discussed in the historical perspectives, many of the early multiprocessors were SIMD, and the SIMD model received renewed attention in the 1980s, and except for vector processors, was gone by the mid 1990s. MIMD has clearly emerged as the architecture of choice for general-purpose multiprocessors. Two factors are primarily responsible for the rise of the MIMD multiprocessors:

1. MIMDs offer flexibility. With the correct hardware and software support, MIMDs can function as single-user multiprocessors focusing on high performance for one application, as multiprogrammed multiprocessors running many tasks simultaneously, or as some combination of these functions.
2. MIMDs can build on the cost/performance advantages of off-the-shelf microprocessors. In fact, nearly all multiprocessors built today use the same microprocessors found in workstations and single-processor servers.

With an MIMD, each processor is executing its own instruction stream. In many cases, each processor executes a different process. Recall from the last chapter, that a process is an segment of code that may be run independently, and that the state of the process contains all the information necessary to execute that program on a processor. In a multiprogrammed environment, where the processors may be running independent tasks, each process is typically independent of the processes on other processors.

It is also useful to be able to have multiple processors executing a single program and sharing the code and most of their address space. When multiple processes share code and data in this way, they are often called *threads*. Today, the term thread is often used in a casual way to refer to multiple loci of execution that may run on different processors, even when they do not share an address space.

To take advantage of an MIMD multiprocessor with n processors, we must usually have at least n threads or processes to execute. The independent threads are typically identified by the programmer or created by the compiler. Since the parallelism in this situation is contained in the threads, it is called *thread-level parallelism*.

Threads may vary from large-scale, independent processes—for example, independent programs running in a multiprogrammed fashion on different processors—to parallel iterations of a loop, automatically generated by a compiler and each executing for perhaps less than a thousand instructions. Although the size of a thread is important in considering how to exploit thread-level parallelism efficiently, the important qualitative distinction is that such parallelism is identified at a high-level by the software system and that the threads consist of hundreds to millions of instructions that may be executed in parallel. In contrast, instruction-level parallelism is identified primarily by the hardware, though with software help in some cases, and is found and exploited one instruction at a time.

Existing MIMD multiprocessors fall into two classes, depending on the number of processors involved, which in turn dictate a memory organization and interconnect strategy. We refer to the multiprocessors by their memory organization, because what constitutes a small or large number of processors is likely to change over time.

The first group, which we call *centralized shared-memory architectures*, have at most a few dozen processors in 2000. For multiprocessors with small processor counts, it is possible for the processors to share a single centralized memory and to interconnect the processors and memory by a bus. With large caches, the bus and the single memory, possibly with multiple banks, can satisfy the memory demands of a small number of processors. By replacing a single bus with multiple buses, or even a switch, a centralized shared memory design can be scaled to a few dozen processors. Although scaling beyond that is technically conceivable, sharing a centralized memory, even organized as multiple banks, becomes less attractive as the number of processors sharing it increases.

Because there is a single main memory that has a symmetric relationship to all processors and a uniform access time from any processor, these multiprocessors are often called *symmetric (shared-memory) multiprocessors (SMPs)*, and this style of architecture is sometimes called *UMA* for *uniform memory access*. This type of centralized shared-memory architecture is currently by far the most popular organization. Figure 6.1 shows what these multiprocessors look like. The architecture of such multiprocessors is the topic of section 6.3.

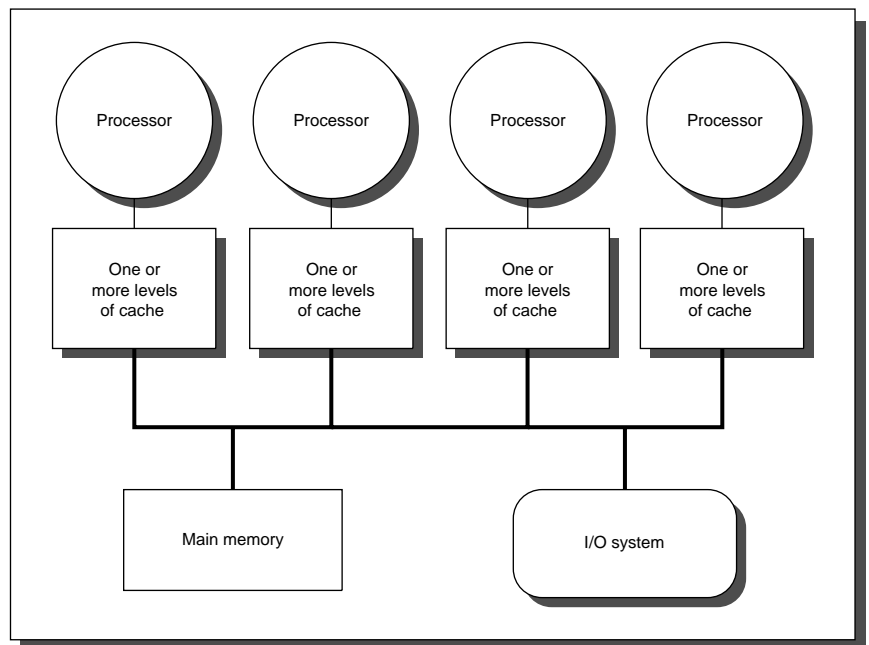


FIGURE 6.1 Basic structure of a centralized shared-memory multiprocessor. Multiple processor-cache subsystems share the same physical memory, typically connected by a bus. In larger designs, multiple buses, or even a switch may be used, but the key architectural property: uniform access time to all memory from all processors remains.

The second group consists of multiprocessors with physically distributed memory. To support larger processor counts, memory must be distributed among the processors rather than centralized; otherwise the memory system would not be able to support the bandwidth demands of a larger number of processors without incurring excessively long access latency. With the rapid increase in processor performance and the associated increase in a processor's memory bandwidth requirements, the scale of multiprocessor for which distributed memory is preferred over a single, centralized memory continues to decrease in number (which is another reason not to use small and large scale). Of course, the larger number of processors raises the need for a high bandwidth interconnect, of which we saw examples in Chapter 7. Both direct interconnection networks (i.e., switches) and indirect networks (typically multidimensional meshes) are used. Figure 6.2 shows what

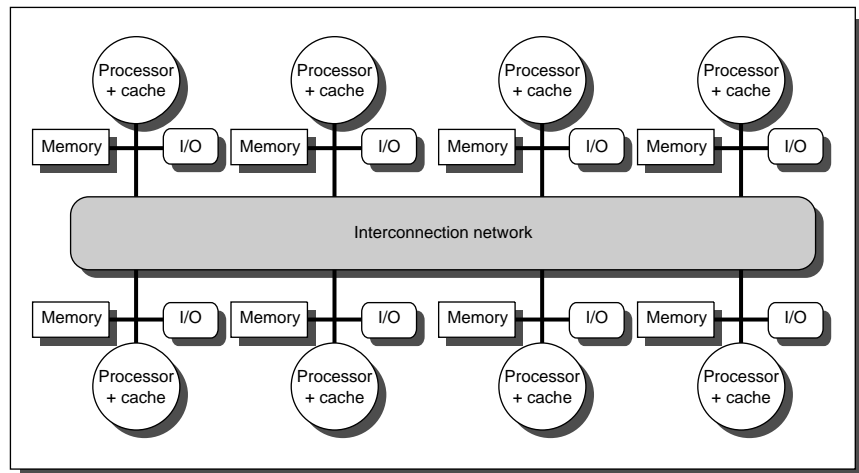


FIGURE 6.2 The basic architecture of a distributed-memory multiprocessor consists of individual nodes containing a processor, some memory, typically some I/O, and an interface to an interconnection network that connects all the nodes. Individual nodes may contain a small number of processors, which may be interconnected by a small bus or a different interconnection technology, which is less scalable than the global interconnection network.

these multiprocessors look like.

Distributing the memory among the nodes has two major benefits. First, it is a cost-effective way to scale the memory bandwidth, if most of the accesses are to the local memory in the node. Second, it reduces the latency for accesses to the local memory. These two advantages make distributed memory attractive at smaller processor counts as processors get ever faster and require more memory

bandwidth and lower memory latency. The key disadvantage for a distributed memory architecture is that communicating data between processors becomes somewhat more complex and has higher latency, at least when there is no contention, because the processors no longer share a single centralized memory. As we will see shortly, the use of distributed memory leads to two different paradigms for interprocessor communication.

Typically, I/O as well as memory is distributed among the nodes of the multiprocessor, and the nodes may be small SMPs (2–8 processors). Although the use of multiple processors in a node together with a memory and a network interface may be quite useful from a cost-efficiency viewpoint, it is not fundamental to how these multiprocessors work, and so we will focus on the one-processor-per-node design for most of this chapter.

Models for Communication and Memory Architecture

As discussed earlier, any large-scale multiprocessor must use multiple memories that are physically distributed with the processors. There are two alternative architectural approaches that differ in the method used for communicating data among processors.

In the first method, communication occurs through a shared address space. That is, the physically separate memories can be addressed as one logically shared address space, meaning that a memory reference can be made by any processor to any memory location, assuming it has the correct access rights. These multiprocessors are called *distributed shared-memory (DSM)* architectures. The term *shared memory* refers to the fact that the *address space* is shared; that is, the same physical address on two processors refers to the same location in memory. Shared memory does *not* mean that there is a single, centralized memory. In contrast to the symmetric shared-memory multiprocessors, also known as UMAs (uniform memory access), the DSM multiprocessors are also called *NUMAs*, *non-uniform memory access*, since the access time depends on the location of a data word in memory.

Alternatively, the address space can consist of multiple private address spaces that are logically disjoint and cannot be addressed by a remote processor. In such multiprocessors, the same physical address on two different processors refers to two different locations in two different memories. Each processor-memory module is essentially a separate computer; therefore these parallel processors have been called *multicomputers*. As pointed out in the previous chapter, a multicomputer can even consist of completely separate computers connected on a local area network, which, today, are popularly called *clusters*. For applications that require little or no communication and can make use of separate memories, such clusters of processors, whether using a standardized or customized interconnect, can form a very cost-effective approach (see Section 7.x).

With each of these organizations for the address space, there is an associated communication mechanism. For a multiprocessor with a shared address space, that address space can be used to communicate data implicitly via load and store operations; hence the name *shared memory* for such multiprocessors. For a multiprocessor with multiple address spaces, communication of data is done by explicitly passing messages among the processors. Therefore, these multiprocessors are often called *message passing multiprocessors*.

In message passing multiprocessors, communication occurs by sending messages that request action or deliver data just as with the network protocols discussed in section 7.2. For example, if one processor wants to access or operate on data in a remote memory, it can send a message to request the data or to perform some operation on the data. In such cases, the message can be thought of as a *remote procedure call (RPC)*. When the destination processor receives the message, either by polling for it or via an interrupt, it performs the operation or access on behalf of the remote processor and returns the result with a reply message. This type of message passing is also called *synchronous*, since the initiating processor sends a request and waits until the reply is returned before continuing. Software systems have been constructed to encapsulate the details of sending and receiving messages, including passing complex arguments or return values, presenting a clean RPC facility to the programmer.

Communication can also occur from the viewpoint of the writer of data rather than the reader, and this can be more efficient when the processor producing data knows which other processors will need the data. In such cases, the data can be sent directly to the consumer process without having to be requested first. It is often possible to perform such message sends asynchronously, allowing the sender process to continue immediately. Often the receiver will want to block if it tries to receive the message before it has arrived; in other cases, the reader may check whether a message is pending before actually trying to perform a blocking receive. Also the sender must be prepared to block if the receiver has not yet consumed an earlier message and no buffer space is available. The message passing facilities offered in different multiprocessors are fairly diverse. To ease program portability, standard message passing libraries (for example, message passing interface, or MPI) have been proposed. Such libraries sacrifice some performance to achieve a common interface.

Performance Metrics for Communication Mechanisms

Three performance metrics are critical in any communication mechanism:

1. *Communication bandwidth*—Ideally the communication bandwidth is limited by processor, memory, and interconnection bandwidths, rather than by some aspect of the communication mechanism. The bisection bandwidth (see Section 7.x) is determined by the interconnection network. The bandwidth in or

out of a single node, which is often as important as bisection bandwidth, is affected both by the architecture within the node and by the communication mechanism. How does the communication mechanism affect the communication bandwidth of a node? When communication occurs, resources within the nodes involved in the communication are tied up or occupied, preventing other outgoing or incoming communication. When this occupancy is incurred for each word of a message, it sets an absolute limit on the communication bandwidth. This limit is often lower than what the network or memory system can provide. Occupancy may also have a component that is incurred for each communication event, such as an incoming or outgoing request. In the latter case, the occupancy limits the communication rate, and the impact of the occupancy on overall communication bandwidth depends on the size of the messages.

2. *Communication latency*—Ideally the latency is as low as possible. As we will see in Chapter 8, communication latency is equal to

$$\text{Sender overhead} + \text{Time of flight} + \text{Transmission time} + \text{Receiver overhead}$$

Time of flight is fixed and transmission time is determined by the interconnection network. The software and hardware overheads in sending and receiving messages are largely determined by the communication mechanism and its implementation. Why is latency crucial? Latency affects both performance and how easy it is to program a multiprocessor. Unless latency is hidden, it directly affects performance either by tying up processor resources or by causing the processor to wait. Overhead and occupancy are closely related, since many forms of overhead also tie up some part of the node, incurring an occupancy cost, which in turn limits bandwidth. Key features of a communication mechanism may directly affect overhead and occupancy. For example, how is the destination address for a remote communication named, and how is protection implemented? When naming and protection mechanisms are provided by the processor, as in a shared address space, the additional overhead is small. Alternatively, if these mechanisms must be provided by the operating system for each communication, this increases the overhead and occupancy costs of communication, which in turn reduce bandwidth and increase latency.

3. *Communication latency hiding*—How well can the communication mechanism hide latency by overlapping communication with computation or with other communication? Although measuring this is not as simple as measuring the first two metrics, it is an important characteristic that can be quantified by measuring the running time on multiprocessors with the same communication latency but different support for latency hiding. We will see examples of latency hiding techniques for shared memory in sections 6.8 and 6.10. Although hiding latency is certainly a good idea, it poses an additional burden on the software system and ultimately on the programmer. Furthermore, the amount of latency that can be hidden is application dependent. Thus, it is usually best to reduce latency wherever possible.

Each of these performance measures is affected by the characteristics of the communications needed in the application. The size of the data items being communicated is the most obvious, since it affects both latency and bandwidth in a direct way, as well as affecting the efficacy of different latency hiding approaches. Similarly, the regularity in the communication patterns affects the cost of naming and protection, and hence the communication overhead. In general, mechanisms that perform well with smaller as well as larger data communication requests, and irregular as well as regular communication patterns, are more flexible and efficient for a wider class of applications. Of course, in considering any communication mechanism, designers must consider cost as well as performance.

Advantages of Different Communication Mechanisms

Each of these two primary communication mechanisms has its advantages. For shared-memory communication, advantages include

- n Compatibility with the well-understood mechanisms in use in centralized multiprocessors, which all use shared-memory communication.
- n Ease of programming when the communication patterns among processors are complex or vary dynamically during execution. Similar advantages simplify compiler design.
- n The ability to develop applications using the familiar shared-memory model, focusing attention only on those accesses that are performance critical.
- n Lower overhead for communication and better use of bandwidth when communicating small items. This arises from the implicit nature of communication and the use of memory mapping to implement protection in hardware, rather than through the I/O system.
- n The ability to use hardware-controlled caching to reduce the frequency of remote communication by supporting automatic caching of all data, both shared and private. As we will see, caching reduces both latency and contention for accessing shared data. This advantage also comes with a disadvantage, which we mention below.

The major advantages for message-passing communication include

- n The hardware can be simpler, especially by comparison with a scalable shared-memory implementation that supports coherent caching of remote data.
- n Communication is explicit, which means it is simpler to understand; in shared memory models, it can be difficult to know when communication is occurring and when it is not, as well as how costly the communication is.

- n Explicit communication focuses programmer attention on this costly aspect of parallel computation, sometimes leading to improved structure in a multiprocessor program.
- n Synchronization is naturally associated with sending messages, reducing the possibility for errors introduced by incorrect synchronization.
- n It makes it easier to use sender-initiated communication, which may have some advantages in performance,

Of course, the desired communication model can be created on top of a hardware model that supports either of these mechanisms. Supporting message passing on top of shared memory is considerably easier: Because messages essentially send data from one memory to another, sending a message can be implemented by doing a copy from one portion of the address space to another. The major difficulties arise from dealing with messages that may be misaligned and of arbitrary length in a memory system that is normally oriented toward transferring aligned blocks of data organized as cache blocks. These difficulties can be overcome either with small performance penalties in software or with essentially no penalties, using a small amount of hardware support.

Supporting shared memory efficiently on top of hardware for message passing is much more difficult. Without explicit hardware support for shared memory, all shared-memory references need to involve the operating system to provide address translation and memory protection, as well as to translate memory references into message sends and receives. Loads and stores usually move small amounts of data, so the high overhead of handling these communications in software severely limits the range of applications for which the performance of software-based shared memory is acceptable. An ongoing area of research is the exploration of when a software-based model is acceptable and whether a software-based mechanism is usable for the highest level of communication in a hierarchically structured system. One possible direction is the use of virtual memory mechanisms to share objects at the page level, a technique called *shared virtual memory*; we discuss this approach in section 6.10.

In distributed-memory multiprocessors, the memory model and communication mechanisms distinguish the multiprocessors. Originally, distributed-memory multiprocessors were built with message passing, since it was clearly simpler and many designers and researchers did not believe that a shared address space could be built with distributed memory. Shared-memory communication has been sup-

ported in virtually every multiprocessor designed since 1995. What hardware communication mechanisms will be supported in the very largest multiprocessors (called *massively parallel processors*, or *MPPs*), which typically have far more than 100 processors, is unclear; shared memory, message passing, and hybrid approaches are all contenders. Despite the symbolic importance of the MPPs, such multiprocessors are a small portion of the market and have little or no influence on the mainstream multiprocessors with tens of processors. We will return to a discussion of the possibilities and trends for MPPs in the concluding remarks and historical perspectives at the end of this chapter.

SMPs, which we focus on in Section 6.3, vastly dominate DSM multiprocessors in terms of market size (both units and dollars), and SMPs will probably be the architecture of choice for on-chip multiprocessors. For moderate scale multiprocessors (>8 processors) long-term technical trends favor distributing memory, which is also likely to be the dominant approach when on-chip SMPs are used as the building blocks in the future. These distributed shared-memory multiprocessors are a natural extension of the centralized multiprocessors that dominate the market, so we discuss these architectures in section 6.5. In contrast, multicomputers or message-passing multiprocessors build on advances in network technology, as we discussed in the last chapter. Since the technologies employed were well described in the last chapter, we focus our attention on shared-memory approaches in the rest of this chapter.

Challenges of Parallel Processing

Two important hurdles, both explainable with Amdahl's Law, make parallel processing challenging. The first has to do with the limited parallelism available in programs and the second arises from the relatively high cost of communications. Limitations in available parallelism make it difficult to achieve good speedups in any parallel processor, as our first Example shows.

EXAMPLE Suppose you want to achieve a speedup of 80 with 100 processors. What fraction of the original computation can be sequential?

ANSWER Amdahl's Law is

$$\text{Speedup} = \frac{1}{\frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}} + (1 - \text{Fraction}_{\text{enhanced}})}$$

For simplicity in this example, assume that the program operates in only two modes: parallel with all processors fully used, which is the enhanced

mode, or serial with only one processor in use. With this simplification, the speedup in enhanced mode is simply the number of processors, while the fraction of enhanced mode is the time spent in parallel mode. Substituting into the equation above:

$$80 = \frac{1}{\frac{\text{Fraction}_{\text{parallel}}}{100} + (1 - \text{Fraction}_{\text{parallel}})}$$

Simplifying this equation yields

$$\begin{aligned} 0.8 \times \text{Fraction}_{\text{parallel}} + 80 \times (1 - \text{Fraction}_{\text{parallel}}) &= 1 \\ 80 - 79.2 \times \text{Fraction}_{\text{parallel}} &= 1 \\ \text{Fraction}_{\text{parallel}} &= \frac{80 - 1}{79.2} \\ \text{Fraction}_{\text{parallel}} &= 0.9975 \end{aligned}$$

Thus to achieve a speedup of 80 with 100 processors, only 0.25% of original computation can be sequential. Of course, to achieve linear speedup (speedup of n with n processors), the entire program must usually be parallel with no serial portions. (One exception to this is *superlinear speedup* that occurs due to the increased memory and cache available when the processor count is increased. This effect is usually not very large and rarely scales linearly with processor count.) In practice, programs do not just operate in fully parallel or sequential mode, but often use less than the full complement of the processors when running in parallel mode. Exercise 6.2 asks you to extend Amdahl's Law to deal with such a case.

n

The second major challenge in parallel processing involves the large latency of remote access in a parallel processor. In existing shared-memory multiprocessors, communication of data between processors may cost anywhere from 100 clock cycles to over 1,000 clock cycles, depending on the communication mechanism, the type of interconnection network, and the scale of the multiprocessor. Figure 6.3 shows the typical round-trip delays to retrieve a word from a remote memory for several different shared-memory parallel processors.

The effect of long communication delays is clearly substantial. Let's consider a simple Example.

EXAMPLE Suppose we have an application running on a 32-processor multiprocessor, which has a 400 ns time to handle reference to a remote memory. For this application, assume that all the references except those involving communication hit in the local memory hierarchy, which may be only

Multiprocessor	Year shipped	SMP or NUMA	Max. processors	Interconnection network	Typical remote memory access time
Sun Starfire servers	1996	SMP	64	Multiple buses	500 ns
SGI Origin 3000	1999	NUMA	512	Fat hypercube	500 ns
Cray T3E	1996	NUMA	2,048	2-way 3D torus	300 ns
HP V series	1998	SMP	32	8x8 crossbar	1000 ns
Compaq Alphaserver GS	1999	SMP	32	Switched busses	400 ns

FIGURE 6.3 Typical remote access times to retrieve a word from a remote memory in shared-memory multiprocessors.

slightly pessimistic. Processors are stalled on a remote request, and the processor clock rate is 1GHz. If the base IPC (assuming that all references hit in the cache) is 2, how much faster is the multiprocessor if there is no communication versus if 0.2% of the instructions involve a remote communication reference?

ANSWER It is simpler to first calculate the CPI. The effective CPI for the multiprocessor with 0.2% remote references is

$$\begin{aligned} \text{CPI} &= \text{Base CPI} + \text{Remote request rate} \times \text{Remote request cost} \\ &= \frac{1}{\text{Base IPC}} + 0.2\% \times \text{Remote request cost} \\ &= 0.5 + 0.2\% \times \text{Remote request cost} \end{aligned}$$

The Remote request cost is

$$\frac{\text{Remote access cost}}{\text{Cycle time}} = \frac{400\text{ns}}{1 \text{ ns}} = 400 \text{ cycles}$$

Hence we can compute the CPI:

$$\text{CPI} = 0.5 + 0.8 = 1.3$$

The multiprocessor with all local references is 1.3/0.5 = 2.6 times faster. In practice, the performance analysis is much more complex, since some fraction of the noncommunication references will miss in the local hierarchy and the remote access time does not have a single constant value. For example, the cost of a remote reference could be quite a bit worse, since contention caused by many references trying to use the global interconnect can lead to increased delays.

These problems—insufficient parallelism and long latency remote communi-

cation—are the two biggest challenges in using multiprocessors. The problem of inadequate application parallelism must be attacked primarily in software with new algorithms that can have better parallel performance. Reducing the impact of long remote latency can be attacked both by the architecture and by the programmer. For example, we can reduce the frequency of remote accesses with either hardware mechanisms, such as caching shared data, or with software mechanisms, such as restructuring the data to make more accesses local. We can try to tolerate the latency by using prefetching or multithreading, which we examined in Chapters 4 and 5.

Much of this chapter focuses on techniques for reducing the impact of long remote communication latency. For example, sections 6.3 and 6.5 discuss how caching can be used to reduce remote access frequency, while maintaining a coherent view of memory. Section 6.7 discusses synchronization, which, because it inherently involves interprocessor communication, is an additional potential bottleneck. Section 6.8 talks about latency hiding techniques and memory consistency models for shared memory. Before we wade into these topics, it is helpful to have some understanding of the characteristics of parallel applications, both for better comprehension of the results we show using some of these applications and to gain a better understanding of the challenges in writing efficient parallel programs.

6.2 Characteristics of Application Domains

In earlier chapters, we examined the performance and characteristics of applications with only a small amount of insight into the structure of the applications. For understanding the key elements of uniprocessor performance, such as caches and pipelining, general knowledge of an application is often adequate, although we saw that deeper application knowledge was necessary to exploit higher levels of ILP.

In parallel processing, the additional performance-critical characteristics—such as load balance, synchronization, and sensitivity to memory latency—typically depend on high-level characteristics of the application. These characteristics include factors like how data is distributed, the structure of a parallel algorithm, and the spatial and temporal access patterns to data. Therefore at this point we take the time to examine the three different classes of workloads.

The three different domains of multiprocessor workloads we explore are a commercial workload, consisting of transaction processing, decision support, and web searching; a multiprogrammed workload with operating systems behavior included; and a workload consisting of individual parallel programs from the technical computing domain.

A Commercial Workload

Our commercial workload consists of three applications:

1. An online transaction processing workload (OLTP) modeled after TPC-B (which has similar memory behavior to its newer cousin TPC-C) and using Oracle 7.3.2 as the underlying database. The workload consists of a set of client processes that generate requests and a set of servers that handle them. The server processes consume 85% of the user time, with the remaining going to the clients. Although the I/O latency is hidden by careful tuning and enough requests to keep the CPU busy, the server processes typically block for I/O after about 25,000 instructions.
2. A decision support system (DSS) workload based on TPC-D and also using Oracle 7.3.2 as the underlying database. The workload includes only six of the 17 read queries in TPC-D, although the six queries examined in the benchmark span the range of activities in the entire benchmark. To hide the I/O latency, parallelism is exploited both within queries, where parallelism is detected during a query formulation process, and across queries. Blocking calls are much less frequent than in the OLTP benchmark; the six queries average about 1.5 million instructions before blocking.
3. A web index search (Altavista) benchmark based on a search of a memory mapped version of the Altavista database (200 GB). The inner loop is heavily optimized. Because the search structure is static, little synchronization is needed among the threads.

The fraction of time spent in user mode, in the kernel, and in the idle loop are shown in Figure 6.4. The frequency of I/O increases both the kernel time and the idle time (see the OLTP entry, which has the largest I/O to computation ratio). Altavista, which maps the entire search database into memory and has been extensively tuned, shows the least kernel or idle time.

Benchmark	% Time User Mode	% Time Kernel	% Time CPU Idle
OLTP	71%	18%	11%
DSS (range for the six queries)	82–94%	3–5%	4–13%
DSS (average across all queries)	87%	3.7%	9.3%
Altavista	> 98%	< 1%	<1%

FIGURE 6.4 The distribution of execution time in the commercial workloads. The OLTP benchmark has the largest fraction of both OS time and CPU idle time (which is I/O wait time). The DSS benchmark shows much less OS time, since it does less I/O, but still more than 9% idle time. The extensive tuning of the Altavista search engine is clear in these measurement. The data for this workload were collected by Barroso et. al. [1998] on a 4-processor Alphaserwer 4100.

Multiprogramming and OS Workload

For small-scale multiprocessors we will also look at a multiprogrammed workload consisting of both user activity and OS activity. The workload used is two independent copies of the compile phase of the Andrew benchmark. The compile phase consists of a parallel make using eight processors. The workload runs for 5.24 seconds on eight processors, creating 203 processes and performing 787 disk requests on three different file systems. The workload is run with 128 MB of memory, and no paging activity takes place.

The workload has three distinct phases: compiling the benchmarks, which involves substantial compute activity; installing the object files in a library; and removing the object files. The last phase is completely dominated by I/O and only two processes are active (one for each of the runs). In the middle phase, I/O also plays a major role and the CPU is largely idle.

Because both CPU idle time and instruction cache performance are important in this workload, we examine these two issues here, focusing on the data cache performance later in the chapter. For the workload measurements, we assume the following memory and I/O systems:

I/O system	Memory
Level 1 instruction cache	32K bytes, two-way set associative with a 64-byte block, one clock cycle hit time
Level 1 data cache	32K bytes, two-way set associative with a 32-byte block, one clock cycle hit time
Level 2 cache	1M bytes unified, two-way set associative with a 128-byte block, hit time 10 clock cycles
Main memory	Single memory on a bus with an access time of 100 clock cycles
Disk system	Fixed access latency of 3 ms (less than normal to reduce idle time)

Figure 6.5 shows how the execution time breaks down for the eight processors using the parameters just listed. Execution time is broken into four components: idle—execution in the kernel mode idle loop; user—execution in user code; synchronization—execution or waiting for synchronization variables; and kernel—execution in the OS that is neither idle nor in synchronization access.

Unlike the parallel scientific workload, this multiprogramming workload has a significant instruction cache performance loss, at least for the OS. The instruction cache miss rate in the OS for a 64-byte block size, two set-associative cache varies from 1.7% for a 32-KB cache to 0.2% for a 256-KB cache. User-level, instruction cache misses are roughly one-sixth of the OS rate, across the variety of cache sizes.

	User execution	Kernel execution	Synchronization wait	CPU Idle (waiting for I/O)
% instructions executed	27%	3%	1%	69%
% execution time	27%	7%	2%	64%

FIGURE 6.5 The distribution of execution time in the multiprogrammed parallel make workload. The high fraction of idle time is due to disk latency when only one of the eight processes is active. These data and the subsequent measurements for this workload were collected with the SimOS system [Rosenblum 1995]. The actual runs and data collection were done by M. Rosenblum, S. Herrod, and E. Bugnion of Stanford University, using the SimOS simulation system.

Scientific/Technical Applications

Our scientific/technical parallel workload consists of two applications and two computational kernels. The kernels are an FFT (fast Fourier transformation) and an LU decomposition, which were chosen because they represent commonly used techniques in a wide variety of applications and have performance characteristics typical of many parallel scientific applications. In addition, the kernels have small code segments whose behavior we can understand and directly track to specific architectural characteristics. Like many scientific application, I/O is essentially nonexistent in this workload.

The two applications that we use in this chapter are Barnes and Ocean, which represent two important but very different types of parallel computation. We briefly describe each of these applications and kernels and characterize their basic behavior in terms of parallelism and communication. We describe how the problem is decomposed for a distributed shared-memory multiprocessor; certain data decompositions that we describe are not necessary on multiprocessors that have a single centralized memory.

The FFT Kernel

The *fast Fourier transform* (FFT) is the key kernel in applications that use spectral methods, which arise in fields ranging from signal processing to fluid flow to climate modeling. The FFT application we study here is a one-dimensional version of a parallel algorithm for a complex-number FFT. It has a sequential execution time for n data points of $n \log n$. The algorithm uses a high radix (equal to \sqrt{n}) that minimizes communication. The measurements shown in this chapter are collected for a million-point input data set.

There are three primary data structures: the input and output arrays of the data being transformed and the roots of unity matrix, which is precomputed and only read during the execution. All arrays are organized as square matrices. The six steps in the algorithm are as follows:

1. Transpose data matrix.

2. Perform 1D FFT on each row of data matrix.
3. Multiply the roots of unity matrix by the data matrix and write the result in the data matrix.
4. Transpose data matrix.
5. Perform 1D FFT on each row of data matrix.
6. Transpose data matrix.

The data matrices and the roots of unity matrix are partitioned among processors in contiguous chunks of rows, so that each processor's partition falls in its own local memory. The first row of the roots of unity matrix is accessed heavily by all processors and is often replicated, as we do, during the first step of the algorithm just shown. The data transposes ensure good locality during the individual FFT steps, which would otherwise access nonlocal data.

The only communication is in the transpose phases, which require all-to-all communication of large amounts of data. Contiguous subcolumns in the rows assigned to a processor are grouped into blocks, which are transposed and placed into the proper location of the destination matrix. Every processor transposes one block locally and sends one block to each of the other processors in the system. Although there is no reuse of individual words in the transpose, with long cache blocks it makes sense to block the transpose to take advantage of the spatial locality afforded by long blocks in the source matrix.

The LU Kernel

LU is an LU factorization of a dense matrix and is representative of many dense linear algebra computations, such as QR factorization, Cholesky factorization, and eigenvalue methods. For a matrix of size $n \times n$ the running time is n^3 and the parallelism is proportional to n^2 . Dense LU factorization can be performed efficiently by blocking the algorithm, using the techniques in Chapter 5, which leads to highly efficient cache behavior and low communication. After blocking the algorithm, the dominant computation is a dense matrix multiply that occurs in the innermost loop. The block size is chosen to be small enough to keep the cache miss rate low, and large enough to reduce the time spent in the less parallel parts of the computation. Relatively small block sizes (8×8 or 16×16) tend to satisfy both criteria.

Two details are important for reducing interprocessor communication. First, the blocks of the matrix are assigned to processors using a 2D tiling: the $\frac{n}{B} \times \frac{n}{B}$ matrix of blocks is allocated by laying a grid of size $p \times p$ over the matrix of blocks in a cookie-cutter fashion until all the blocks are allocated to a processor. Second, the dense matrix multiplication is performed by the processor that owns the *destination* block. With this blocking and allocation scheme, communication during the reduction is both regular and predictable. For

the measurements in this chapter, the input is a 512×512 matrix and a block of 16×16 is used.

A natural way to code the blocked LU factorization of a 2D matrix in a shared address space is to use a 2D array to represent the matrix. Because blocks are allocated in a tiled decomposition, and a block is not contiguous in the address space in a 2D array, it is very difficult to allocate blocks in the local memories of the processors that own them. The solution is to ensure that blocks assigned to a processor are allocated locally and contiguously by using a 4D array (with the first two dimensions specifying the block number in the 2D grid of blocks, and the next two specifying the element in the block).

The Barnes Application

Barnes is an implementation of the Barnes-Hut n-body algorithm solving a problem in galaxy evolution. *N-body algorithms* simulate the interaction among a large number of bodies that have forces interacting among them. In this instance the bodies represent collections of stars and the force is gravity. To reduce the computational time required to model completely all the individual interactions among the bodies, which grow as n^2 , n-body algorithms take advantage of the fact that the forces drop off with distance. (Gravity, for example, drops off as $1/d^2$, where d is the distance between the two bodies.) The Barnes-Hut algorithm takes advantage of this property by treating a collection of bodies that are “far away” from another body as a single point at the center of mass of the collection and with mass equal to the collection. If the body is far enough from any body in the collection, then the error introduced will be negligible. The collections are structured in a hierarchical fashion, which can be represented in a tree. This algorithm yields an $n \log n$ running time with parallelism proportional to n .

The Barnes-Hut algorithm uses an octree (each node has up to eight children) to represent the eight cubes in a portion of space. Each node then represents the collection of bodies in the subtree rooted at that node, which we call a *cell*. Because the density of space varies and the leaves represent individual bodies, the depth of the tree varies. The tree is traversed once per body to compute the net force acting on that body. The force-calculation algorithm for a body starts at the root of the tree. For every node in the tree it visits, the algorithm determines if the center of mass of the cell represented by the subtree rooted at the node is “far enough away” from the body. If so, the entire subtree under that node is approximated by a single point at the center of mass of the cell, and the force this center of mass exerts on the body is computed. On the other hand, if the center of mass is not far enough away, the cell must be “opened” and each of its subtrees visited. The distance between the body and the cell, together with the error tolerances, determines which cells must be opened. This force calculation phase dominates the execution time. This chapter takes measurements using 16K bodies; the criterion for determining whether a cell needs to be opened is set to the middle of the range typically used in practice.

Obtaining effective parallel performance on Barnes-Hut is challenging because the distribution of bodies is nonuniform and changes over time, making partitioning the work among the processors and maintenance of good locality of reference difficult. We are helped by two properties: the system evolves slowly; and because gravitational forces fall off quickly, with high probability, each cell requires touching a small number of other cells, most of which were used on the last time step. The tree can be partitioned by allocating each processor a subtree. Many of the accesses needed to compute the force on a body in the subtree will be to other bodies in the subtree. Since the amount of work associated with a subtree varies (cells in dense portions of space will need to access more cells), the size of the subtree allocated to a processor is based on some measure of the work it has to do (e.g., how many other cells does it need to visit), rather than just on the number of nodes in the subtree. By partitioning the octree representation, we can obtain good load balance and good locality of reference, while keeping the partitioning cost low. Although this partitioning scheme results in good locality of reference, the resulting data references tend to be for small amounts of data and are unstructured. Thus this scheme requires an efficient implementation of shared-memory communication.

The Ocean Application

Ocean simulates the influence of eddy and boundary currents on large-scale flow in the ocean. It uses a restricted red-black Gauss-Seidel multigrid technique to solve a set of elliptical partial differential equations. *Red-black Gauss-Seidel* is an iteration technique that colors the points in the grid so as to consistently update each point based on previous values of the adjacent neighbors. *Multigrid methods* solve finite difference equations by iteration using hierarchical grids. Each grid in the hierarchy has fewer points than the grid below, and is an approximation to the lower grid. A finer grid increases accuracy and thus the rate of convergence, while requiring more execution time, since it has more data points. Whether to move up or down in the hierarchy of grids used for the next iteration is determined by the rate of change of the data values. The estimate of the error at every time-step is used to decide whether to stay at the same grid, move to a coarser grid, or move to a finer grid. When the iteration converges at the finest level, a solution has been reached. Each iteration has n^2 work for an $n \times n$ grid and the same amount of parallelism.

The arrays representing each grid are dynamically allocated and sized to the particular problem. The entire ocean basin is partitioned into square subgrids (as close as possible) that are allocated in the portion of the address space corresponding to the local memory of the individual processors, which are assigned responsibility for the subgrid. For the measurements in this chapter we use an input that has 130×130 grid points. There are five steps in a time iteration. Since data are exchanged between the steps, all the processors present synchronize at the end of each step before proceeding to the next. Communication occurs when the boundary points of a subgrid are accessed by the adjacent subgrid in nearest-neighbor fashion.

Computation/Communication for the Parallel Programs

A key characteristic in determining the performance of parallel programs is the ratio of computation to communication. If the ratio is high, it means the application has lots of computation for each datum communicated. As we saw in section 6.1, communication is the costly part of parallel computing; therefore high computation-to-communication ratios are very beneficial. In a parallel processing environment, we are concerned with how the ratio of computation to communication changes as we increase either the number of processors, the size of the problem, or both. Knowing how the ratio changes as we increase the processor count sheds light on how well the application can be sped up. Because we are often interested in running larger problems, it is vital to understand how changing the data set size affects this ratio.

To understand what happens quantitatively to the computation-to-communication ratio as we add processors, consider what happens separately to computation and to communication as we either add processors or increase problem size. Figure 6.6 shows that as we add processors, for these applications, the amount of computation per processor falls proportionately and the amount of communication per processor falls more slowly. As we increase the problem size, the computation scales as the $O(\)$ complexity of the algorithm dictates. Communication scaling is more complex and depends on details of the algorithm; we describe the basic phenomena for each application in the caption of Figure 6.6.

The overall computation-to-communication ratio is computed from the individual growth rate in computation and communication. In general, this ratio rises slowly with an increase in data set size and decreases as we add processors. This reminds us that performing a fixed-size problem with more processors leads to increasing inefficiencies because the amount of communication among processors grows. It also tells us how quickly we must scale data set size as we add processors, to keep the fraction of time in communication fixed. The following example illustrates this tradeoffs.

EXAMPLE Suppose we know that for a given multiprocessor the Ocean application spends 20% of its execution time waiting for communication when run on four processors. Assume that the cost of each communication event is independent on processor count, which is not true in general, since communication costs rise with processor count. How much faster might we expect Ocean to run on a 32-processor machine with the same problem size? What fraction of the execution time is spent on communication in this case? How much larger a problem should we run if we want the fraction of time spent communicating to be the same?

ANSWER The computation to communication ratio for Ocean is \sqrt{n}/\sqrt{p} , so if the problem size is the same, the communication frequency scales by \sqrt{p} .

Application	Scaling of computation	Scaling of communication	Scaling of computation-to-communication
FFT	$\frac{n \log n}{p}$	$\frac{n}{p}$	$\log n$
LU	$\frac{n}{p}$	$\frac{\sqrt{n}}{\sqrt{p}}$	$\frac{\sqrt{n}}{\sqrt{p}}$
Barnes	$\frac{n \log n}{p}$	Approximately $\frac{\sqrt{n}(\log n)}{\sqrt{p}}$	Approximately $\frac{\sqrt{n}}{\sqrt{p}}$
Ocean	$\frac{n}{p}$	$\frac{\sqrt{n}}{\sqrt{p}}$	$\frac{\sqrt{n}}{\sqrt{p}}$

FIGURE 6.6 Scaling of computation, of communication, and of the ratio are critical factors in determining performance on parallel multiprocessors. In this table p is the increased processor count and n is the increased data set size. Scaling is on a per-processor basis. The computation scales up with n at the rate given by $O(\)$ analysis and scales down linearly as p is increased. Communication scaling is more complex. In FFT all data points must interact, so communication increases with n and decreases with p . In LU and Ocean, communication is proportional to the boundary of a block, so it scales with data set size at a rate proportional to the side of a square with n points, namely \sqrt{n} ; for the same reason communication in these two applications scales inversely to \sqrt{p} . Barnes has the most complex scaling properties. Because of the fall-off of interaction between bodies, the basic number of interactions among bodies, which require communication, scales as \sqrt{n} . An additional factor of $\log n$ is needed to maintain the relationships among the bodies. As processor count is increased, communication scales inversely to \sqrt{p} .

This means that communication time increase by $\sqrt{8}$. We can use a variation on Amdahl's Law, recognizing that the computation is decreased but the communication time is increased. If T_4 is the total execution time for 4 processors, then the execution time for 32 processors is:

$$\begin{aligned}
 T_{32} &= \text{Compute time} + \text{Communication time} \\
 &= \frac{0.8 \times T_4}{8} + (0.2 \times T_4) \times \sqrt{8} \\
 &= 0.1 \times T_4 + 0.57 \times T_4 = 0.67 \times T_4
 \end{aligned}$$

Hence the speed-up is:

$$\text{Speedup} = \frac{T_4}{T_{32}} = \frac{T_4}{0.67 \times T_4} = 1.49$$

And the fraction of time spent in communication goes from 20% to $0.57/0.67 = 85\%$.

For the fraction of the communication time to remain the same, we must keep the computation to communication ratio the same, so the problem size must scale at the same rate as the processor count. Notice that because we have changed the problem size, we cannot measure of the scaled problem. We will return to the critical issue of scaling applications for multiprocessors in both in the Cross Cutting Issues and the Fallacies and Pitfalls.

n

6.3 Symmetric Shared-Memory Architectures

Multis are a new class of computers based on multiple microprocessors. The small size, low cost, and high performance of microprocessors allow design and construction of computer structures that offer significant advantages in manufacture, price-performance ratio, and reliability over traditional computer families.... Multis are likely to be the basis for the next, the fifth, generation of computers.
[p. 463]

Bell [1985]

As we saw in Chapter 5, the use of large, multilevel caches can substantially reduce the memory bandwidth demands of a processor. If the main memory bandwidth demands of a single processor are reduced, multiple processors may be able to share the same memory. Starting in the 1980s, this observation, combined with the emerging dominance of the microprocessor, motivated many designers to create small-scale multiprocessors where several processors shared a single physical memory connected by a shared bus. This type of design is called symmetric shared memory, because each processor has the same relationship to one single shared memory. Because of the small size of the processors and the significant reduction in the requirements for bus bandwidth achieved by large caches, such symmetric multiprocessors are extremely cost-effective, provided that a sufficient amount of memory bandwidth exists. Early designs of such multiprocessors were able to place an entire CPU and cache subsystem on a board, which plugged into the bus backplane. More recent designs have placed up to four processors per board; and a recent announcement by IBM includes 2 processors on the same die. Figure 6.1 on page 639 shows a simple diagram of such a multiprocessor.

Small-scale shared-memory machines usually support the caching of both shared and private data. *Private data* is used by a single processor, while *shared data* is used by multiple processors, essentially providing communication among the processors through reads and writes of the shared data. When a private item is cached, its location is migrated to the cache, reducing the average access time as well as the memory bandwidth required. Since no other processor uses the data,

the program behavior is identical to that in a uniprocessor. When shared data are cached, the shared value may be replicated in multiple caches. In addition to the reduction in access latency and required memory bandwidth, this replication also provides a reduction in contention that may exist for shared data items that are being read by multiple processors simultaneously. Caching of shared data, however, introduces a new problem: cache coherence.

What Is Multiprocessor Cache Coherence?

As we saw in Chapter 6, the introduction of caches caused a coherence problem for I/O operations, since the view of memory through the cache could be different from the view of memory obtained through the I/O subsystem. The same problem exists in the case of multiprocessors, because the view of memory held by two different processors is through their individual caches. Figure 6.7 illustrates the problem and shows how two different processors can have two different values for the same location. This difficulty is generally referred to as the *cache-coherence* problem.

Time	Event	Cache contents for CPU A	Cache contents for CPU B	Memory contents for location X
0				1
1	CPU A reads X	1		1
2	CPU B reads X	1	1	1
3	CPU A stores 0 into X	0	1	0

FIGURE 6.7 The cache-coherence problem for a single memory location (X), read and written by two processors (A and B). We initially assume that neither cache contains the variable and that X has the value 1. We also assume a write-through cache; a write-back cache adds some additional but similar complications. After the value of X has been written by A, A's cache and the memory both contain the new value, but B's cache does not, and if B reads the value of X, it will receive 1!

Informally, we could say that a memory system is coherent if any read of a data item returns the most recently written value of that data item. This definition, although intuitively appealing, is vague and simplistic; the reality is much more complex. This simple definition contains two different aspects of memory system behavior, both of which are critical to writing correct shared-memory programs. The first aspect, called *coherence*, defines what values can be returned by a read. The second aspect, called *consistency*, determines when a written value will be returned by a read. Let's look at coherence first.

A memory system is coherent if

1. A read by a processor, P, to a location X that follows a write by P to X, with no writes of X by another processor occurring between the write and the read by P, always returns the value written by P.
2. A read by a processor to location X that follows a write by another processor to X returns the written value if the read and write are sufficiently separated in time and no other writes to X occur between the two accesses.
3. Writes to the same location are *serialized*: that is, two writes to the same location by any two processors are seen in the same order by all processors. For example, if the values 1 and then 2 are written to a location, processors can never read the value of the location as 2 and then later read it as 1.

The first property simply preserves program order—we expect this property to be true even in uniprocessors. The second property defines the notion of what it means to have a coherent view of memory: If a processor could continuously read an old data value, we would clearly say that memory was incoherent.

The need for write serialization is more subtle, but equally important. Suppose we did not serialize writes, and processor P1 writes location X followed by P2 writing location X. Serializing the writes ensures that every processor will see the write done by P2 at some point. If we did not serialize the writes, it might be the case that some processor could see the write of P2 first and then see the write of P1, maintaining the value written by P1 indefinitely. The simplest way to avoid such difficulties is to serialize writes, so that all writes to the same location are seen in the same order; this property is called *write serialization*.

Although the three properties just described are sufficient to ensure coherence, the question of when a written value will be seen is also important. To see why, observe that we cannot require that a read of X instantaneously see the value written for X by some other processor. If, for example, a write of X on one processor precedes a read of X on another processor by a very small time, it may be impossible to ensure that the read returns the value of the data written, since the written data may not even have left the processor at that point. The issue of exactly *when* a written value must be seen by a reader is defined by a *memory consistency model*—a topic discussed in section 6.8.

Coherence and consistency are complementary: Coherence defines the behavior of reads and writes to the same memory location, while consistency defines the behavior of reads and writes with respect to accesses to other memory locations. For simplicity, and because we cannot explain the problem in full detail at this point, assume that we require that a write does not complete until all processors have seen the effect of the write and that the processor does not change the order of any write with any other memory access. This allows the processor to reorder reads, but forces the processor to finish a write in program order. We will rely on this assumption until we reach section 6.8, where we will see exactly the meaning of this definition, as well as the alternatives.

Basic Schemes for Enforcing Coherence

The coherence problem for multiprocessors and I/O, although similar in origin, has different characteristics that affect the appropriate solution. Unlike I/O, where multiple data copies are a rare event—one to be avoided whenever possible—a program running on multiple processors will normally have copies of the same data in several caches. In a coherent multiprocessor, the caches provide both *migration* and *replication* of shared data items.

Coherent caches provide migration, since a data item can be moved to a local cache and used there in a transparent fashion. This migration reduces both the latency to access a shared data item that is allocated remotely and the bandwidth demand on the shared memory.

Coherent caches also provide replication for shared data that is being simultaneously read, since the caches make a copy of the data item in the local cache. Replication reduces both latency of access and contention for a read shared data item. Supporting this migration and replication is critical to performance in accessing shared data. Thus, rather than trying to solve the problem by avoiding it in software, small-scale multiprocessors adopt a hardware solution by introducing a protocol to maintain coherent caches.

The protocols to maintain coherence for multiple processors are called *cache-coherence protocols*. Key to implementing a cache-coherence protocol is tracking the state of any sharing of a data block. There are two classes of protocols, which use different techniques to track the sharing status, in use:

- *Directory based*—The sharing status of a block of physical memory is kept in just one location, called the *directory*; we focus on this approach in section 6.5, when we discuss scalable shared-memory architecture.
- *Snooping*—Every cache that has a copy of the data from a block of physical memory also has a copy of the sharing status of the block, and no centralized state is kept. The caches are usually on a shared-memory bus, and all cache controllers monitor or *snoop* on the bus to determine whether or not they have a copy of a block that is requested on the bus. We focus on this approach in this section.

Snooping protocols became popular with multiprocessors using microprocessors and caches attached to a single shared memory because these protocols can use a preexisting physical connection—the bus to memory—to interrogate the status of the caches.

Snooping Protocols

There are two ways to maintain the coherence requirement described in the previous subsection. One method is to ensure that a processor has exclusive access to a data item before it writes that item. This style of protocol is called a *write invalidate protocol* because it invalidates other copies on a write. It is by far the most

common protocol, both for snooping and for directory schemes. Exclusive access ensures that no other readable or writable copies of an item exist when the write occurs: all other cached copies of the item are invalidated.

Figure 6.8 shows an example of an invalidation protocol for a snooping bus with write-back caches in action To see how this protocol ensures coherence, consider a write followed by a read by another processor: Since the write requires exclusive access, any copy held by the reading processor must be invalidated (hence the protocol name). Thus, when the read occurs, it misses in the cache and is forced to fetch a new copy of the data. For a write, we require that the writing processor have exclusive access, preventing any other processor from being able to write simultaneously. If two processors do attempt to write the same data simultaneously, one of them wins the race (we’ll see how we decide who wins shortly), causing the other processor’s copy to be invalidated. For the other processor to complete its write, it must obtain a new copy of the data, which must now contain the updated value. Therefore, this protocol enforces write serialization.

Processor activity	Bus activity	Contents of CPU A’s cache	Contents of CPU B’s cache	Contents of memory location X
				0
CPU A reads X	Cache miss for X	0		0
CPU B reads X	Cache miss for X	0	0	0
CPU A writes a 1 to X	Invalidation for X	1		0
CPU B reads X	Cache miss for X	1	1	1

FIGURE 6.8 An example of an invalidation protocol working on a snooping bus for a single cache block (X) with write-back caches. We assume that neither cache initially holds X and that the value of X in memory is 0. The CPU and memory contents show the value after the processor and bus activity have both completed. A blank indicates no activity or no copy cached. When the second miss by B occurs, CPU A responds with the value canceling the response from memory. In addition, both the contents of B’s cache and the memory contents of X are updated. This update of memory, which occurs when a block becomes shared, is typical in most protocols and simplifies the protocol, as we will see shortly.

The alternative to an invalidate protocol is to update all the cached copies of a data item when that item is written. This type of protocol is called a *write update* or *write broadcast* protocol. To keep the bandwidth requirements of this protocol under control it is useful to track whether or not a word in the cache is shared—that is, is contained in other caches. If it is not, then there is no need to broadcast or update any other caches. Figure 6.8 shows an example of a write update protocol in operation. In the decade since these protocols were developed, invalidate has emerged as the winner for the vast majority of designs. To understand why, let’s look at the qualitative performance differences.

The performance differences between write update and write invalidate protocols arise from three characteristics:

Processor activity	Bus activity	Contents of CPU A's cache	Contents of CPU B's cache	Contents of memory location X
				0
CPU A reads X	Cache miss for X	0		0
CPU B reads X	Cache miss for X	0	0	0
CPU A writes a 1 to X	Write broadcast of X	1	1	1
CPU B reads X		1	1	1

FIGURE 6.9 An example of a write update or broadcast protocol working on a snooping bus for a single cache block (X) with write-back caches. We assume that neither cache initially holds X and that the value of X in memory is 0. The CPU and memory contents show the value after the processor and bus activity have both completed. A blank indicates no activity or no copy cached. When CPU A broadcasts the write, both the cache in CPU B and the memory location of X are updated.

1. Multiple writes to the same word with no intervening reads require multiple write broadcasts in an update protocol, but only one initial invalidation in a write invalidate protocol.
2. With multiword cache blocks, each word written in a cache block requires a write broadcast in an update protocol, although only the first write to any word in the block needs to generate an invalidate in an invalidation protocol. An invalidation protocol works on cache blocks, while an update protocol must work on individual words (or bytes, when bytes are written). It is possible to try to merge writes in a write broadcast scheme, just as we did for write buffers in Chapter 5, but the basic difference remains.
3. The delay between writing a word in one processor and reading the written value in another processor is usually less in a write update scheme, since the written data are immediately updated in the reader's cache (assuming that the reading processor has a copy of the data). By comparison, in an invalidation protocol, the reader is invalidated first, then later reads the data and is stalled until a copy can be read and returned to the processor.

Because bus and memory bandwidth is usually the commodity most in demand in a bus-based multiprocessor and invalidation protocols generate less bus and memory traffic, invalidation has become the protocol of choice for almost all multiprocessors. Update protocols also cause problems for memory consistency models, reducing the potential performance gains of update, mentioned in point 3, even further. In designs with very small processor counts (2, or at most, 4) where the processors are tightly coupled (perhaps even on the same chip), the larger bandwidth demands of update may be acceptable. Nonetheless, given the trends in increasing processor performance and the related increase in bandwidth

demands, we can expect update schemes to be used very infrequently. For this reason, we will focus only on invalidate protocols for the rest of the chapter.

Basic Implementation Techniques

The key to implementing an invalidate protocol in a small-scale multiprocessor is the use of the bus to perform invalidates. To perform an invalidate the processor simply acquires bus access and broadcasts the address to be invalidated on the bus. All processors continuously snoop on the bus watching the addresses. The processors check whether the address on the bus is in their cache. If so, the corresponding data in the cache is invalidated.

The serialization of access enforced by the bus also forces serialization of writes, since when two processors compete to write to the same location, one must obtain bus access before the other. The first processor to obtain bus access will cause the other processor's copy to be invalidated, causing writes to be strictly serialized. One implication of this scheme is that a write to a shared data item cannot complete until it obtains bus access.

In addition to invalidating outstanding copies of a cache block that is being written into, we also need to locate a data item when a cache miss occurs. In a write-through cache, it is easy to find the recent value of a data item, since all written data are always sent to the memory, from which the most recent value of a data item can always be fetched. (Write buffers can lead to some additional complexities, which are discussed in section 6.8.)

For a write-back cache, however, the problem of finding the most recent data value is harder, since the most recent value of a data item can be in a cache rather than in memory. Happily, write-back caches can use the same snooping scheme both for cache misses and for writes: Each processor snoops every address placed on the bus. If a processor finds that it has a dirty copy of the requested cache block, it provides that cache block in response to the read request and causes the memory access to be aborted. Since write-back caches generate lower requirements for memory bandwidth, they are greatly preferable in a multiprocessor, despite the slight increase in complexity. Therefore, we focus on implementation with write-back caches.

The normal cache tags can be used to implement the process of snooping, and the valid bit for each block makes invalidation easy to implement. Read misses, whether generated by an invalidation or by some other event, are also straightforward since they simply rely on the snooping capability. For writes we'd like to know whether any other copies of the block are cached, because, if there are no other cached copies, then the write need not be placed on the bus in a write-back cache. Not sending the write reduces both the time taken by the write and the required bandwidth.

To track whether or not a cache block is shared we can add an extra state bit associated with each cache block, just as we have a valid bit and a dirty bit. By adding a bit indicating whether the block is shared, we can decide whether a write must generate an invalidate. When a write to a block in the shared state oc-

curs, the cache generates an invalidation on the bus and marks the block as private. No further invalidations will be sent by that processor for that block. The processor with the sole copy of a cache block is normally called the *owner* of the cache block.

When an invalidation is sent, the state of the owner's cache block is changed from shared to unshared (or exclusive). If another processor later requests this cache block, the state must be made shared again. Since our snooping cache also sees any misses, it knows when the exclusive cache block has been requested by another processor and the state should be made shared.

Every bus transaction must check the cache-address tags, which could potentially interfere with CPU cache accesses. This potential interference is reduced by one of two techniques: duplicating the tags or employing a multilevel cache with *inclusion*, whereby the levels closer to the CPU are a subset of those further away. If the tags are duplicated, then the CPU and the snooping activity may proceed in parallel. course, on a cache miss the processor needs to arbitrate for and update both sets of tags. Likewise, if the snoop finds a matching tag entry, it needs to arbitrate for and access both sets of cache tags (to perform an invalidate or to update the shared bit), as well as possibly the cache data array to retrieve a copy of a block. Thus with duplicate tags the processor only needs to be stalled when it does a cache access at the same time that a snoop has detected a copy in the cache. Furthermore, snooping activity is delayed only when the cache is dealing with a miss.

If the CPU uses a multilevel cache with the inclusion property, then every entry in the primary cache is required to be in the secondary cache. Thus the snoop activity can be directed to the second-level cache, while most of the processor's activity is directed to the primary cache. If the snoop gets a hit in the secondary cache, then it must arbitrate for the primary cache to update the state and possibly retrieve the data, which usually requires a stall of the processor. Since many multiprocessors use a multilevel cache to decrease the bandwidth demands of the individual processors, this solution has been adopted in many designs. Sometimes it may even be useful to duplicate the tags of the secondary cache to further decrease contention between the CPU and the snooping activity. We discuss the inclusion property in more detail in section 6.10 on page 728.

An Example Protocol

A bus-based coherence protocol is usually implemented by incorporating a finite state controller in each node. This controller responds to requests from the processor and from the bus, changing the state of the selected cache block, as well as using the bus to access data or to invalidate it. Figure 6.10 shows the requests generated by the processor-cache module in a node, in the top half of the table, as well as those coming from the bus, in the bottom half of the table. For simplicity,

the protocol we explain does not distinguish between a write hit and a write miss to a shared cache block: in both cases, we treat such an access as a write miss. When the write miss is placed on the bus, any processors with copies of the cache block invalidate it. In a write-back cache, if the block is exclusive in just one cache, that cache also writes back the block. Treating write hits to shared blocks as cache misses reduces the number of different bus transactions and simplifies the controller. In more sophisticated protocols, these “misses” are treated as upgrade requests that generate a bus transaction and an invalidate, but do not actually transfer the data, since the copy in the cache is up-to-date.

Request	Source	State of addressed cache block	Function and explanation
Read hit	Processor	Shared or Exclusive	Read data in cache
Read miss	Processor	Invalid	Place read miss on bus.
Read miss	Processor	Shared	Address conflict miss: place read miss on bus
Read miss	Processor	Exclusive	Address conflict miss: write back block, then place read miss on bus
Write hit	Processor	Exclusive	Write data in cache.
Write hit	Processor	Shared	Place write miss on bus.
Write miss	Processor	Invalid	Place write miss on bus.
Write miss	Processor	Shared	Address conflict miss: place write miss on bus
Write miss	Processor	Exclusive	Address conflict miss: write back block, then place write miss on bus
Read Miss	Bus	Shared	No action; allow memory to service read miss.
Read Miss	Bus	Exclusive	Attempt to share data: place cache block on bus and change state to Shared.
Write miss	Bus	Shared	Attempt to write shared block; invalidate the block.
Write miss	Bus	Exclusive	Attempt to write block that is exclusive elsewhere: write back the cache block and make its state Invalid.

FIGURE 6.10 The cache-coherence mechanism receives requests from both the processor and the bus and responds to these based on the type of request, whether it hits or misses in the cache, and the state of the cache block specified in the request. For read or write misses snooped from the bus, an action is required *only* if the read or write addresses matches a block in the cache and the block is valid. Placing a write miss on the bus when a write hits in the Shared state, ensures an exclusive copy, though the data need not actually be transferred. This is referred to as an upgrade, and some protocols distinguish it from a write miss to avoid the data transfer.

Figure 6.11 shows a finite-state transition diagram for a single cache block using a write-invalidation protocol and a write-back cache. For simplicity, the three states of the protocol are duplicated to represent transitions based on CPU requests (on the left, which corresponds to the top half of the table in Figure 6.11), as opposed to transitions based on bus requests (on the right, which corresponds

to the bottom half of the table in Figure 6.11). Boldface type is used to distinguish the bus actions, as opposed to the conditions on which a state transition depends. The state in each node represents the state of the selected cache block specified by the processor or bus request.

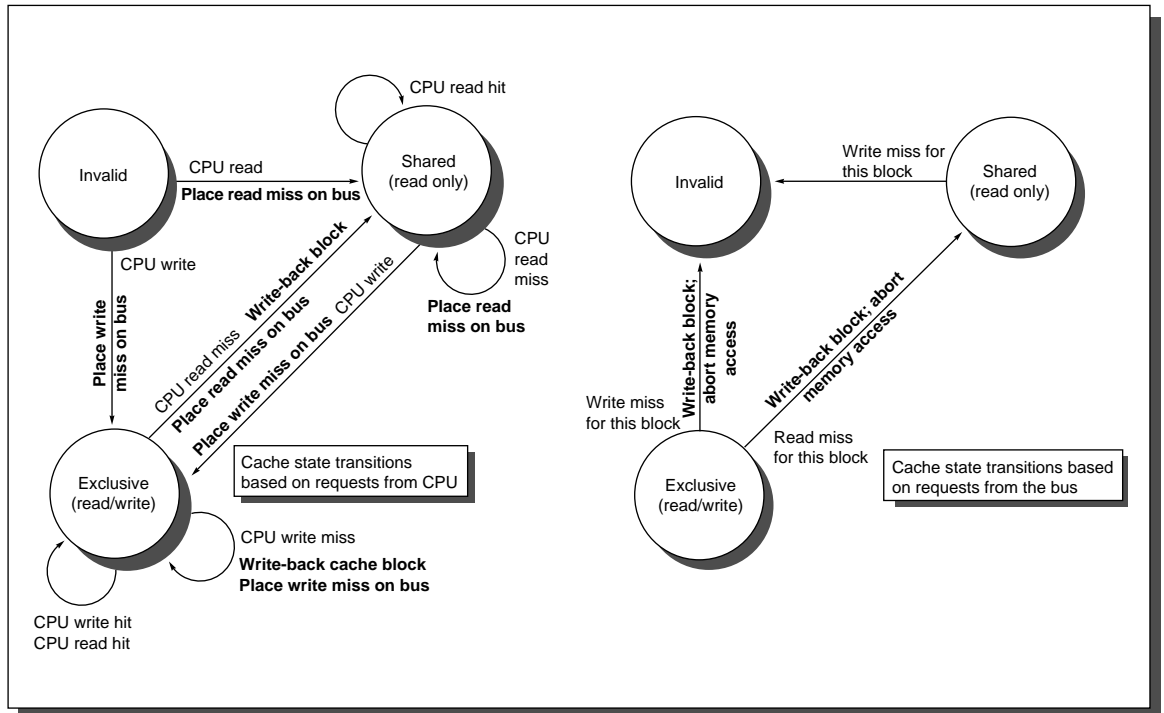


FIGURE 6.11 A write-invalidate, cache-coherence protocol for a write-back cache showing the states and state transitions for each block in the cache. The cache states are shown in circles with any access permitted by the CPU without a state transition shown in parenthesis under the name of the state. The stimulus causing a state change is shown on the transition arcs in regular type, and any bus actions generated as part of the state transition are shown on the transition arc in bold. The stimulus actions apply to a block in the cache, not to a specific address in the cache. Hence, a read miss to a block in the shared state is a miss for that cache block but for a different address. The left side of the diagram shows state transitions based on actions of the CPU associated with this cache; the right side shows transitions based on operations on the bus. A read miss in the exclusive or shared state and a write miss in the exclusive state occur when the address requested by the CPU does not match the address in the cache block. Such a miss is a standard cache replacement miss. An attempt to write a block in the shared state always generates a miss, even if the block is present in the cache, since the block must be made exclusive. Whenever a bus transaction occurs, all caches that contain the cache block specified in the bus transaction take the action dictated by the right half of the diagram. The protocol assumes that memory provides data on a read miss for a block that is clean in all caches. In actual implementations, these two sets of state diagrams are combined. This protocol is somewhat simpler than those in use in existing multiprocessors.

All of the states in this cache protocol would be needed in a uniprocessor cache, where they would correspond to the invalid, valid (and clean), and dirty

states. All of the state changes indicated by arcs in the left half of Figure 6.11 would be needed in a write-back uniprocessor cache; the only difference in a multiprocessor with coherence is that the controller must generate a write miss when the controller has a write hit for a cache block in the shared state. The state changes represented by the arcs in the right half of Figure 6.11 are needed only for coherence and would not appear at all in a uniprocessor cache controller.

In reality, there is only one finite-state machine per cache, with stimuli coming either from the attached CPU or from the bus. Figure 6.12 shows how the state

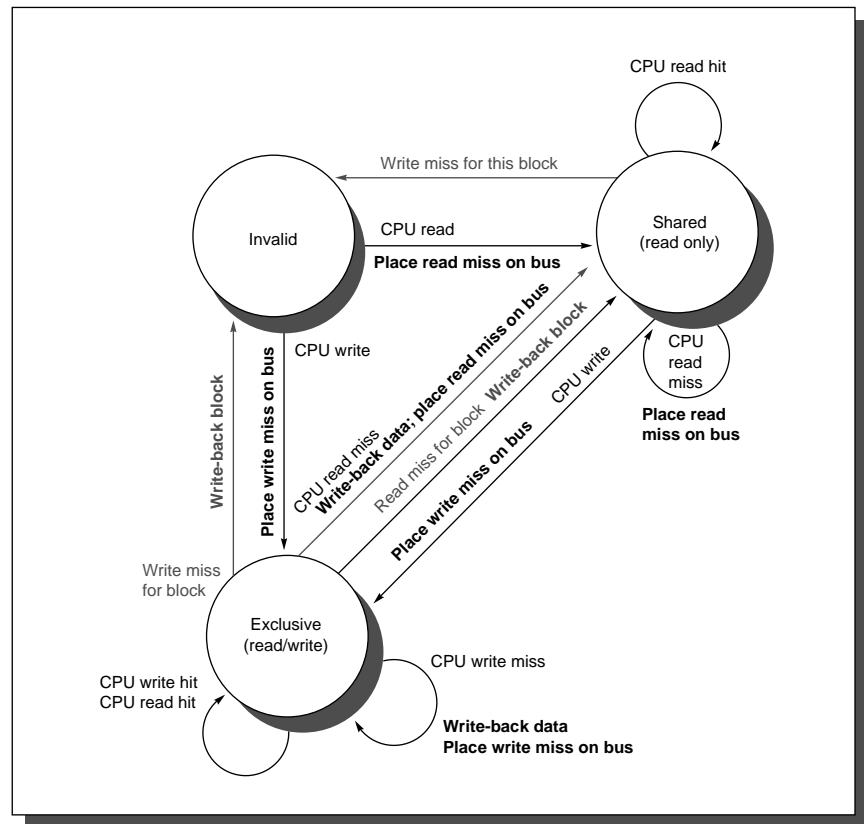


FIGURE 6.12 Cache-coherence state diagram with the state transitions induced by the local processor shown in black and by the bus activities shown in gray. As in Figure 6.11, the activities on a transition are shown in bold.

transitions in the right half of Figure 6.11 are combined with those in the left half of the figure to form a single state diagram for each cache block.

To understand why this protocol works, observe that any valid cache block is either in the shared state in multiple caches or in the exclusive state in exactly one cache. Any transition to the exclusive state (which is required for a processor to write to the block) requires a write miss to be placed on the bus, causing all caches to make the block invalid. In addition, if some other cache had the block in exclusive state, that cache generates a write back, which supplies the block containing the desired address. Finally, if a read miss occurs on the bus to a block in the exclusive state, the owning cache also makes its state shared, forcing a subsequent write to require exclusive ownership.

The actions in gray in Figure 6.12, which handle read and write misses on the bus, are essentially the snooping component of the protocol. One other property that is preserved in this protocol, and in most other protocols, is that any memory block in the shared state is always up to date in the memory. This simplifies the implementation, as we will see in detail in section 6.7.

Although our simple cache protocol is correct, it omits a number of complications that make the implementation much trickier. The most important of these is that the protocol assumes that operations are *atomic*—that is, an operation can be done in such a way that no intervening operation can occur. For example, the protocol described assumes that write misses can be detected, acquire the bus, and receive a response as a single atomic action. In reality this is not true. Similarly, if we used a split transaction bus (see Chapter 6, section 6.3), as most modern bus-based multiprocessors do, then even read misses would also not be atomic.

Nonatomic actions introduce the possibility that the protocol can *deadlock*, meaning that it reaches a state where it cannot continue. Appendix E deals with these complex issues, showing how the protocol can be modified to deal with nonatomic writes without introducing deadlock.

As stated earlier, this coherence protocol is actually simpler than those used in practice. There are two major simplifications. First, in this protocol all transitions to the exclusive state generate a write miss on the bus, and we assume that the requesting cache always fills the block with the contents returned. This simplifies the detailed implementation. Most real protocols distinguish between a write miss and a write hit, which can occur when the cache block is initially in the shared state. Such misses are called *ownership* or *upgrade* misses, since they involve changing the state of the block, but do not actually require a data fetch. To support such state changes, the protocol uses an *invalidate operation*, in addition to a write miss. With such operations, however, the actual implementation of the protocol becomes slightly more complex.

The second major simplification is that many multiprocessors distinguish between a cache block that is really shared and one that exists in the clean state in exactly one cache. This addition of a “clean and private” state eliminates the need to generate a bus transaction on a write to such a block. Another enhancement in wide use allows other caches to supply data on a miss to a shared block.

Constructing small (2-4) processor bus-based multiprocessors has become very easy. Many modern microprocessors provide basic support for cache coherency and also allow the construction of a shared memory bus by direct connection of the external memory bus of two processors. These capabilities reduce the number of chips required to build a small-scale multiprocessor. For example, the Intel Pentium III Xeon and Pentium 4 Xeon processors are designed for use in cache coherent multiprocessors and have an external interface that supports snooping and allows two processors to be directly connected. They also have larger on-chip caches to reduce bus utilization. A system chip set containing an external memory controller is used to connect the shared processor memory bus with a set of memory chips. The memory controller also implements the coherency protocol. Since different size multiprocessors generate different demands for bus bandwidth, Intel has two different system chip sets designed for dual processor systems and for midrange range systems (2-4 processors). A small-scale multiprocessor may be built with only two additional system chips: the memory controller memory controller mentioned above and an I/O hub chip that interfaces standard I/O buses, such as PCI, to the memory bus.

The next part of this section examines the performance of these protocols for our parallel and multiprogrammed workloads; the value of these extensions to a basic protocol will be clear when we examine the performance.

6.4 Performance of Symmetric Shared-Memory Multiprocessors

In a bus-based multiprocessor using an invalidation protocol, several different phenomena combine to determine performance. In particular, the overall cache performance is a combination of the behavior of uniprocessor cache miss traffic and the traffic caused by communication, which results in invalidations and subsequent cache misses. Changing the processor count, cache size, and block size can affect these two components of the miss rate in different ways, leading to overall system behavior that is a combination of the two effects.

In Chapter 5, we saw how breaking the uniprocessor miss rate into the 3C classification could provide insight into both application behavior and potential improvements to the cache design. Similarly, the misses that arise from interprocessor communication, which are often called *coherence misses*, can be broken into two separate sources.

The first source are the so-called *true sharing misses* that arise from the communication of data through the cache coherence mechanism. In an invalidation-based protocol, the first write by a processor to a shared cache block causes an invalidation to establish ownership of that block. Additionally, when another processor attempts to read a modified word in that cache block, a miss occurs and the

resultant block is transferred. Both these misses are classified as true sharing misses since they directly arise from the sharing of data among processors.

The second effect, called *false sharing*, arises from the use of an invalidation-based coherence algorithm with a single valid bit per cache block. False sharing occurs when a block is invalidated (and a subsequent reference causes a miss) because some word in the block, other than the one being read, is written into. If the word written into is actually used by the processor that received the invalidate, then the reference was a true sharing reference and would have caused a miss independent of the block size or position of words. If, however, the word being written and the word read are different and the invalidation does not cause a new value to be communicated, but only causes an extra cache miss, then it is a false sharing miss. In a false sharing miss, the block is shared, but no word in the cache is actually shared, and the miss would not occur if the block size were a single word. The following Example makes the sharing patterns clear.

EXAMPLE Assume that words x_1 and x_2 are in the same cache block, which is in the shared state in the caches of P1 and P2. Assuming the following sequence of events, identify each miss as a true sharing miss, a false sharing miss, or a hit. Any miss that would occur if the block size were one word is designated a true sharing miss.

Time	P1	P2
1	Write x_1	
2		Read x_2
3	Write x_1	
4		Write x_2
5	Read x_2	

ANSWER Here are classifications by time step:

1. This event is a true sharing miss, since x_1 was read by P2 and needs to be invalidated from P2.
2. This event is a false sharing miss, since x_2 was invalidated by the write of x_1 in P1, but that value of x_1 is not used in P2.
3. This event is a false sharing miss, since the block containing x_1 is marked shared due to the read in P2, but P2 did not read x_1 . The cache block containing x_1 will be in the shared state after the read by P2; a write miss is required to obtain exclusive access to the block. In some protocols this will be handled as an *upgrade request*, which

generates a bus invalidate, but does not transfer the cache block.

4. This event is a false sharing miss for the same reason as step 3.
5. This event is a true sharing miss, since the value being read was written by P2.

n

True sharing and false sharing miss rates can be affected by a variety of changes in the cache architecture. Thus, we will find it useful to decompose not only the uniprocessor and multiprocessor miss rates, but also the true-sharing and false-sharing miss rates.

Performance Measurements of the Commercial Workload

The performance measurements of the commercial workload, which we examine in this section, were taken either on a Alphaserer 4100, or using a configurable simulator modeled after the Alphaserer 4100. The Alphaserer 4100 used for these measurements has four processors, each of which is an Alpha 21164 running at 300 MHz. Each processor has a three-level cache hierarchy:

- n L1 consist of a pair of 8 KB direct-mapped on-chip caches, one for instruction and one for data. The block size is 32-bytes, and the data cache is write-through to L2, using a write buffer.
- n L2 is a 96 KB on-chip unified 3-way set associative cache with a 32-byte block size, using write-back.
- n L3 is an off-chip, combined, direct-mapped 2 MB caches with 64-byte blocks also using write-back.

The latency for an access to L2 is 7 cycles, to L3 it is 21 cycles, and to main memory it is 80 clock cycles (typical without contention). Cache to cache transfers, which occur on a miss to an exclusive block held in another cache, require 125 clock cycles. All the measurement shown in this section were collected by Barroso, Gharachorloo, and Bugnion [1998].

We start by looking at the overall CPU execution for these benchmarks on the 4-processor system; as discussed on page 650, these benchmarks include substantial I/O time, which is ignored in the CPU time measurements. We group the six DSS queries as a single benchmark, reporting the average behavior. The effective CPI varies widely for these benchmarks, from a CPI of 1.3 for the Altavista web search to an average CPI of 1.6 for the DSS workload, to 7.0 for the OLTP workload. Figure 6.13 shows how the execution time breaks down into instruction execution, cache and memory system access time, and other stalls (which are primarily pipeline resource stalls, but also include TLB and branch mispredict stalls). Although the performance of the DSS and Altavista workloads is reasonable, the performance of the OLTP workload is very poor, due to a poor performance of the memory hierarchy.

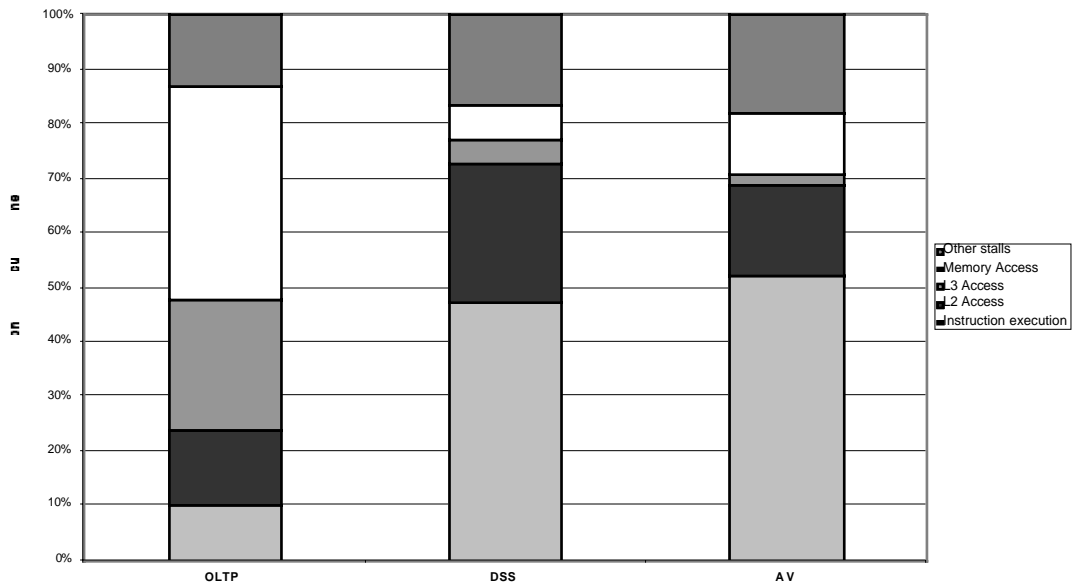


FIGURE 6.13 The execution time breakdown for the three programs (OLTP, DSS, and Altavista) in the commercial workload. The DSS numbers are the average across six different queries. The CPI varies widely from a low of 1.3 for Altavista, to 1.61 for the DSS queries, to 7.0 for Oracle. (Individually, the DSS queries show a CPI range of 1.3 to 1.9.) Other stalls includes: resource stalls (implemented with replay traps on the 21164), branch mispredict, memory barrier, and TLB misses. For these benchmarks resource-based pipeline stalls are the dominant factor. This data combines the behavior of user and kernel accesses. Only OLTP has a significant fraction of kernel accesses, and the kernel accesses tend to be better behaved than the user accesses!

Since the OLTP workload demands the most from the memory system with large numbers of expensive L3 misses, we focus on examining the impact of L3 cache size, processor count, and block size on the OLTP benchmark. Figure 6.14 shows the effect of increasing the cache size, using 2-way set associative caches, which reduces the large number of conflict misses. The execution time is improved as the L3 cache grows due to the reduction in L3 misses. The idle time also grows, reducing some of the performance gains. This growth occurs because with fewer memory system stalls, more server processes are needed to cover the I/O latency. The workload could be retuned to increase the computation/communication balance, holding the idle time in check.

To better understand how the L3 miss rate responds, we ask: What factors contribute to the L3 miss rate and how do they change as the L3 cache grows? Figure 6.15 shows this data, displaying the number of memory access cycles contributed per instruction from five sources. The two largest sources of memory access cycles (due to L3 misses) with a 1 MB L3 are instruction and capacity/conflict misses. With a larger L3 these two sources shrink to be minor contributors. Un-

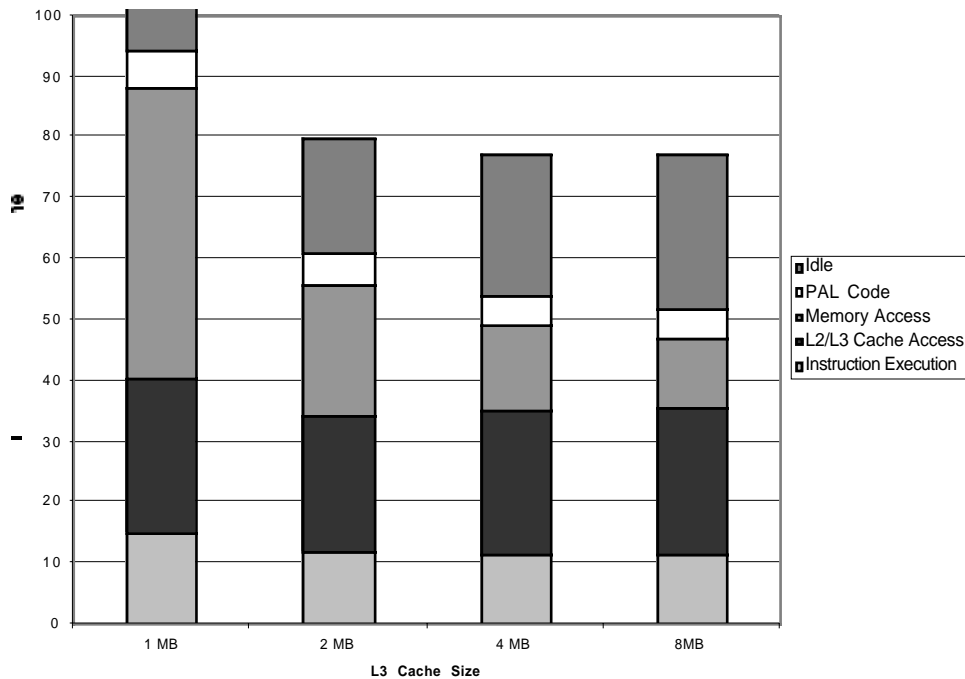


FIGURE 6.14 The relative performance of the OLTP workload as the size of the L3 cache, which is set as 2-way set associative, is grown from 1 MB to 8 MB. Interestingly, the performance of the 1 MB, 2-way set associative cache is very similar to the direct-mapped 2 MB cache that is used in the Alphaserver 4100.

fortunately, the cold, false sharing, and true sharing misses are unaffected by a larger L3. Thus, at 4 and 8 MB, the true sharing misses generate the dominant fraction of the misses.

Clearly, increasing the cache size eliminates most of the uniprocessor misses, while leaving the multiprocessor misses untouched. How does increasing the processor count affect different types of misses? Figure 6.16 shows this data assuming a base configuration with a 2 MB, 2-way set associative L3 cache. As we might expect, the increase in the true sharing miss rate, which is not compensated for by any decrease in the uniprocessor misses, leads to an overall increase in the memory access cycles per instruction.

The final question we examine is whether increasing the block size, which should decrease the instruction and cold miss rate and, within limits, also reduce the capacity/ conflict miss rate, is helpful for this workload. Figure 6.17 shows

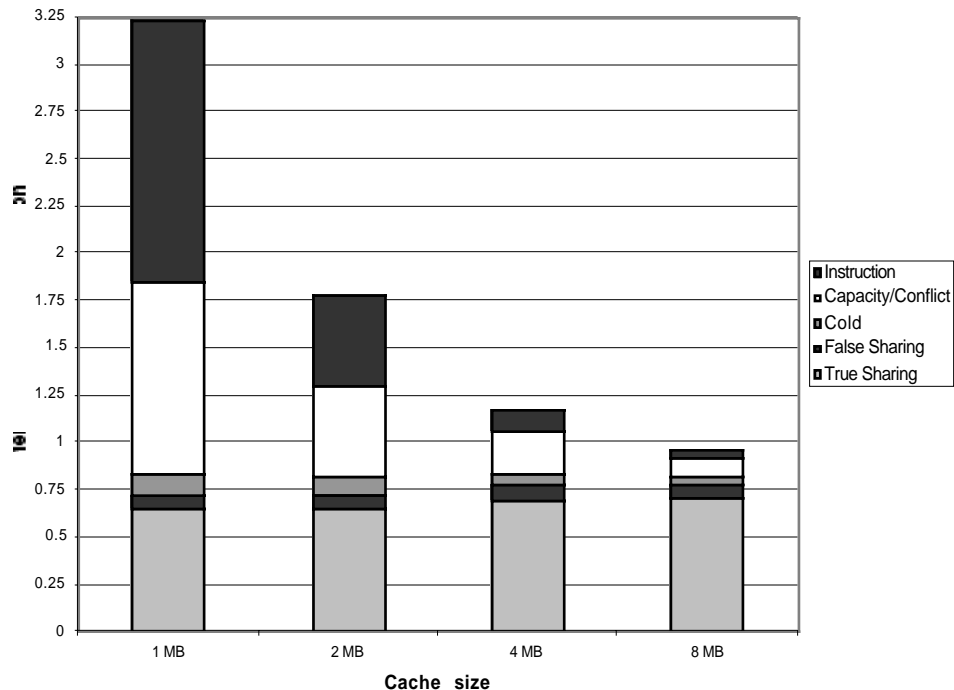


FIGURE 6.15 The contributing causes of memory access cycles shift as the cache size is increased. The L3 cache is simulated as 2-way set associative.

the number of misses per one-thousand instructions as the block size is increased from 32 to 256. Increasing the block size from 32 to 256 affects four of the miss rate components:

- n the true sharing miss rate decreases by more than a factor of 2, indicating locality in the true sharing patterns,
- n the cold start miss rate significantly decreases, as we would expect,
- n the conflict/capacity misses show a small decrease (a factor of 1.26 compared to a factor of 8 increase in block size), indicating that the spatial locality is not high in the uniprocessor misses, and
- n the false sharing miss rate, although small in absolute terms, nearly doubles.

The lack of a significant effect on the instruction miss rate is startling and clearly indicates that the large instruction footprint has very poor spatial locality! Over-

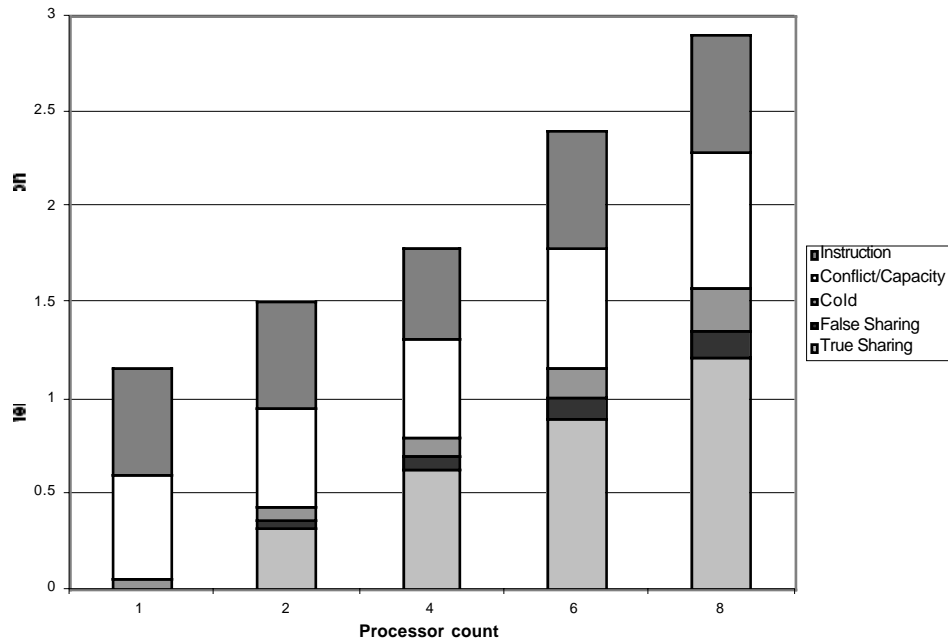


FIGURE 6.16 The contribution to memory access cycles increases as processor count increases primarily due to increased true sharing. The cold misses slightly increase since each processor must now handle more cold misses.

all, increasing the block size of the of the third-level cache to 128 or possibly 256 bytes seems appropriate.

Performance of the Multiprogramming and OS Workload

In this subsection we examine the cache performance of the multiprogrammed workload as the cache size and block size are changed. The workload remains the same as described in the previous section: two independent parallel makes, each using up to eight processors. Because of differences between the behavior of the kernel and that of the user processes, we keep these two components separate. Remember, though, that the user processes execute more than eight times as many instructions, so that the overall miss rate is determined primarily by the miss rate in user code, which, as we will see, is often one-fifth of the kernel miss rate.

Figure 6.18 shows the data miss rate versus data cache size for the kernel and user components. The misses can be broken into three significant classes:

- n Compulsory, or cold, misses represent the first access to this block by this pro-

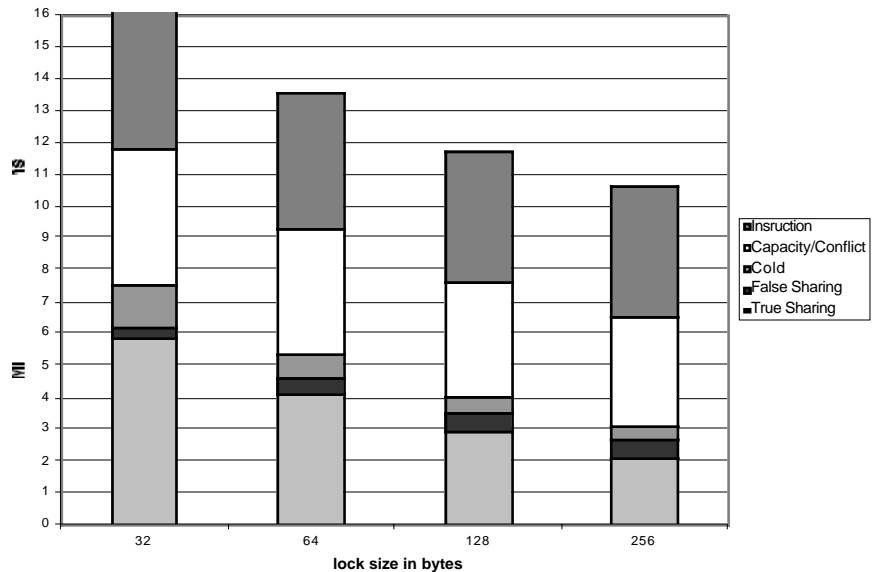


FIGURE 6.17 The number of misses per one-thousand instructions drops steadily as the block size of the L3 cache is increased making a good case for an L3 block size of at least 128 bytes. The L3 cache is a 2MB, 2-way set associative,

cessor and are significant in this workload.

- n Coherence misses represent misses due to invalidations.
- n Normal capacity misses include misses caused by interference between the OS and the user process and between multiple user processes. Conflict misses are included in this category.

For this workload the behavior of the operating system is more complex than the user processes. This is for two reasons. First, the kernel initializes all pages before allocating them to a user, which significantly increases the compulsory component of the kernel's miss rate. Second, the kernel actually shares data and thus has a nontrivial coherence miss rate. In contrast, user processes cause coherence misses only when the process is scheduled on a different processor; this component of the miss rate is small. Figure 6.19 shows the breakdown of the kernel miss rate as the cache size is increased.

Increasing the block size is likely to have beneficial effects for this workload, since a larger fraction of the misses arise from compulsory and capacity, both of which can be potentially improved with larger block sizes. Since coherence misses are relatively more rare, the negative effects of increasing block size should be

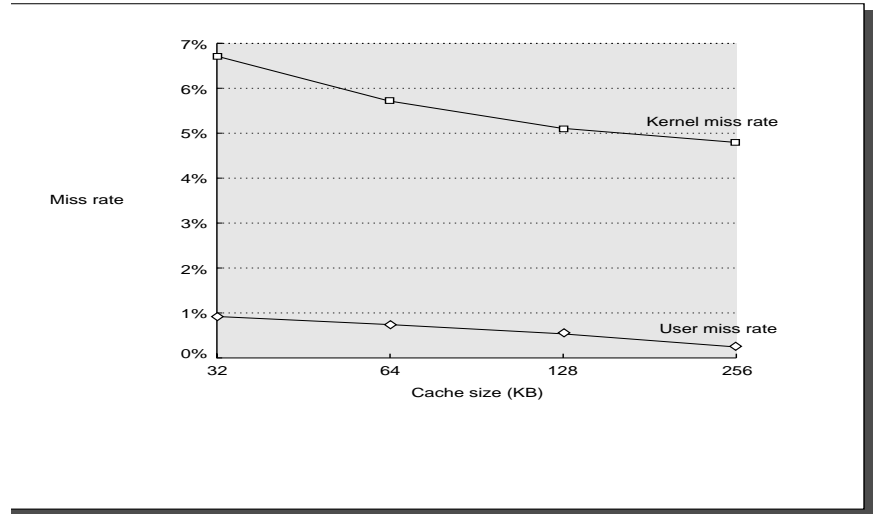


FIGURE 6.18 The data miss rate drops faster for the user code than for the kernel code as the data cache is increased from 32 KB to 256 KB with a 32-byte block. Although the user level miss rate drops by a factor of 3, the kernel level miss rate drops only by a factor of 1.3. As Figure 6.19 shows, this is due to a higher rate of compulsory misses and coherence misses. This multiprogramming workload is run on eight processors.

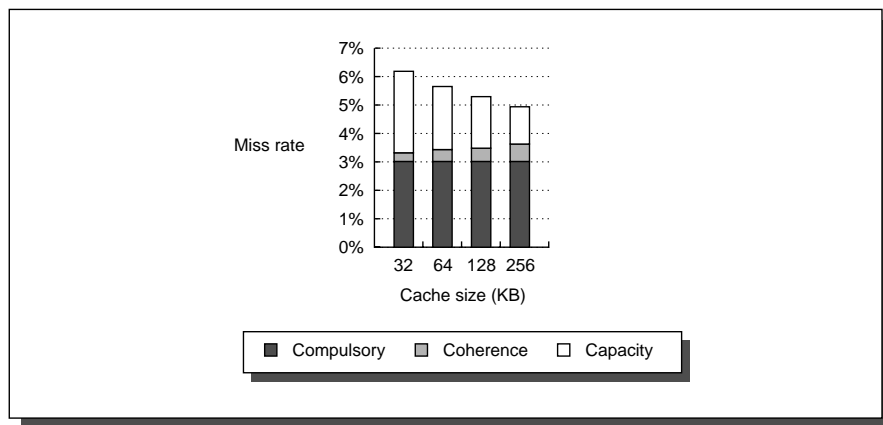


FIGURE 6.19 The components of the kernel data miss rate change as the data cache size is increased from 32KB to 256 KB, when the multiprogramming workload is run on eight processors. The compulsory miss rate component stays constant, since it is unaffected by cache size. The capacity component drops by more than a factor of two, while the coherence component nearly doubles. The increase in coherence misses occurs because the probability of a miss being caused by an invalidation increases with cache size, since fewer entries are bumped due to capacity.

small. Figure 6.20 shows how the miss rate for the kernel and user references changes as the block size is increased, assuming a 32 KB two-way set-associative data cache. Figure 6.21 confirms that, for the kernel references, the largest im-

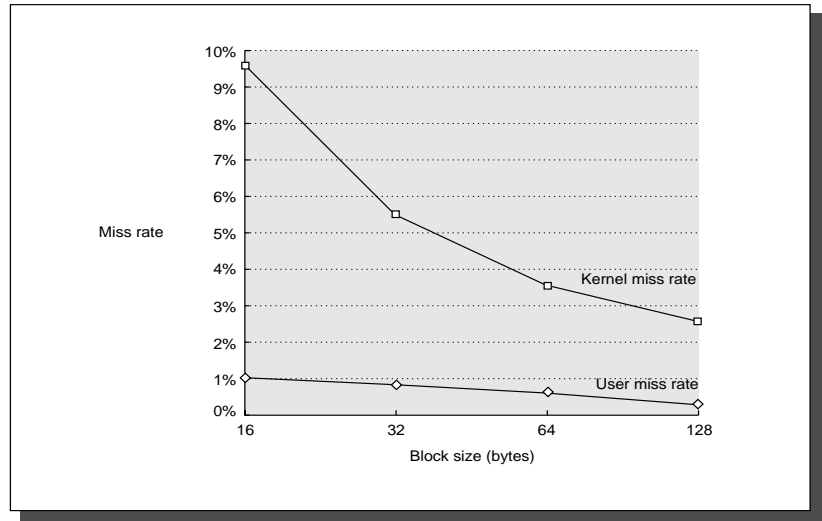


FIGURE 6.20 Miss rate for the multiprogramming workload drops steadily as the block size is increased for a 32-KB two-way set-associative data cache and an eight-CPU multiprocessor. As we might expect based on the higher compulsory component in the kernel, the improvement in miss rate for the kernel references is larger (almost a factor of 4 for the kernel references when going from 16-byte to 128-byte blocks versus just under a factor of 3 for the user references).

provement is the reduction of the compulsory miss rate. The absence of large increases in the coherence miss rate as block size is increased means that false sharing effects are insignificant.

If we examine the number of bytes needed per data reference, as in Figure 6.22, we see that the kernel has a higher traffic ratio that grows quickly with block size. This is despite the significant reduction in compulsory misses; the smaller reduction in capacity and coherence misses drives an increase in total traffic. The user program has a much smaller traffic ratio that grows very slowly.

For the multiprogrammed workload, the OS is a much more demanding user of the memory system. If more OS or OS-like activity is included in the workload, it will become very difficult to build a sufficiently capable memory system.

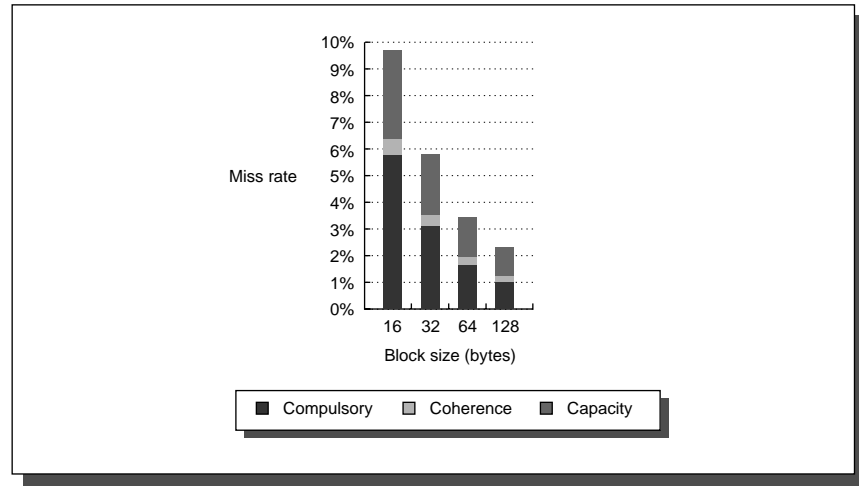


FIGURE 6.21 As we would expect, the increasing block size substantially reduces the compulsory miss rate in the kernel references. It also has a significant impact on the capacity miss rate, decreasing it by a factor of 2.4 over the range of block sizes. The increased block size has a small reduction in coherence traffic, which appears to stabilize at 64 bytes, with no change in the coherence miss rate in going to 128-byte lines. Because there are not significant reductions in the coherence miss rate as the block size increases, the fraction of the miss rate due to coherence grows from about 7% to about 15%.

Performance for the Scientific/Technical Workload

In this section, we use a simulator to study the performance of our four scientific parallel programs. For these measurements, the problem sizes are as follows:

- *Barnes-Hut*—16K bodies run for six time steps (the accuracy control is set to 1.0, a typical, realistic value);
- *FFT*—1 million complex data points
- *LU*—A 512×512 matrix is used with 16×16 blocks
- *Ocean*—A 130×130 grid with a typical error tolerance

In looking at the miss rates as we vary processor count, cache size, and block size, we decompose the total miss rate into *coherence misses* and normal uniprocessor misses. The normal uniprocessor misses consist of capacity, conflict, and compulsory misses. We label these misses as capacity misses, because that is the dominant cause for these benchmarks. For these measurements, we include as a coherence miss any write misses needed to upgrade a block from shared to exclusive, even though no one is sharing the cache block. This measurement reflects a protocol that does not distinguish between a private and shared cache block.

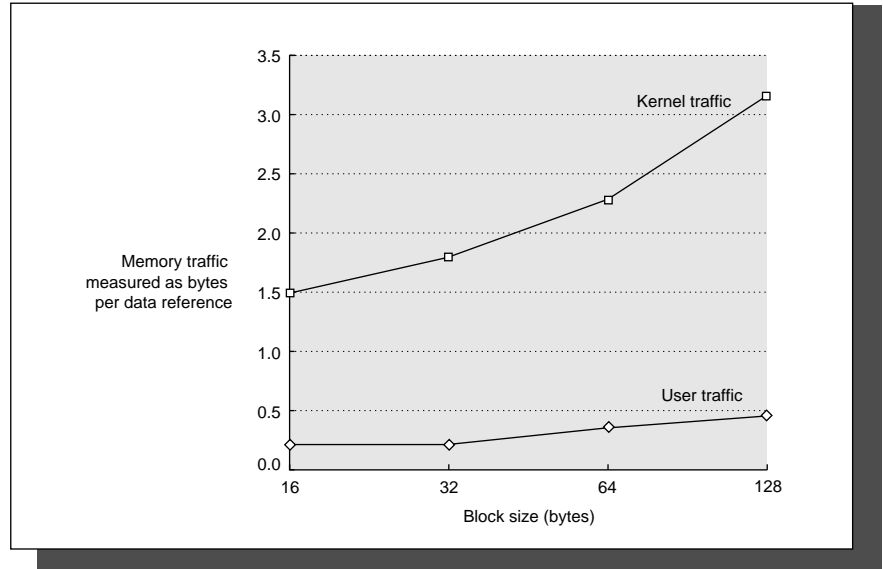


FIGURE 6.22 The number of bytes needed per data reference grows as block size is increased for both the kernel and user components. It is interesting to compare this chart against the same chart for the parallel program workload shown in Figure 6.26.

Figure 6.23 shows the data miss rates for our four applications, as we increase the number of processors from one to sixteen, while keeping the problem size constant. As we increase the number of processors, the total amount of cache increases, usually causing the capacity misses to drop. In contrast, increasing the processor count usually causes the amount of communication to increase, in turn causing the coherence misses to rise. The magnitude of these two effects differs by application.

In FFT, the capacity miss rate drops (from nearly 7% to just over 5%) but the coherence miss rate increases (from about 1% to about 2.7%), leading to a constant overall miss rate. Ocean shows a combination of effects, including some that relate to the partitioning of the grid and how grid boundaries map to cache blocks. For a typical 2D grid code the communication-generated misses are proportional to the boundary of each partition of the grid, while the capacity misses are proportional to the area of the grid. Therefore, increasing the total amount of cache while keeping the total problem size fixed will have a more significant effect on the capacity miss rate, at least until each subgrid fits within an individual processor's cache. The significant jump in miss rate between one and two proces-

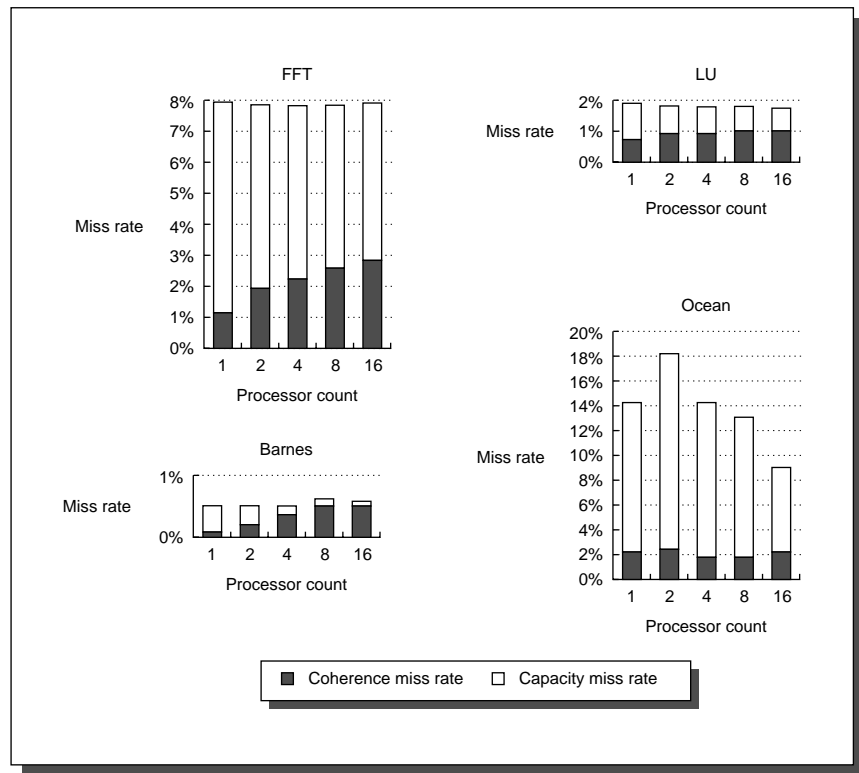


FIGURE 6.23 Data miss rates can vary in nonobvious ways as the processor count is increased from one to sixteen. The miss rates include both coherence and capacity miss rates. The compulsory misses in these benchmarks are all very small and are included in the capacity misses. Most of the misses in these applications are generated by accesses to data that is potentially shared, although in the applications with larger miss rates (FFT and Ocean), it is the capacity misses rather than the coherence misses that comprise the majority of the miss rate. Data is potentially shared if it is allocated in a portion of the address space used for shared data. In all except Ocean, the potentially shared data is heavily shared, while in Ocean only the boundaries of the subgrids are actually shared, although the entire grid is treated as a potentially shared data object. Of course, since the boundaries change as we increase the processor count (for a fixed-size problem), different amounts of the grid become shared. The anomalous increase in capacity miss rate for Ocean in moving from one to two processors arises because of conflict misses in accessing the subgrids. In all cases except Ocean, the fraction of the cache misses caused by coherence transactions rises when a fixed-size problem is run on an increasing number of processors. In Ocean, the coherence misses initially fall as we add processors due to a large number of misses that are write ownership misses to data that is potentially, but not actually, shared. As the subgrids begin to fit in the aggregate cache (around 16 processors), this effect lessens. The single processor numbers include write upgrade misses, which occur in this protocol even if the data is not actually shared, since it is in the shared state. For all these runs, the cache size is 64 KB, two-way set associative, with 32-byte blocks. Notice that the scale on the y-axis for each benchmark is different, so that the behavior of the individual benchmarks can be seen clearly.

sors occurs because of conflicts that arise from the way in which the multiple grids are mapped to the caches. This conflict is present for direct-mapped and two-way set associative caches, but fades at higher associativities. Such conflicts are not unusual in array-based applications, especially when there are multiple grids in use at once. In Barnes and LU the increase in processor count has little effect on the miss rate, sometimes causing a slight increase and sometimes causing a slight decrease.

Increasing the cache size usually has a beneficial effect on performance, since it reduces the frequency of costly cache misses. Figure 6.24 illustrates the change in miss rate as cache size is increased for 16 processors, showing the portion of the miss rate due to coherence misses and to uniprocessor capacity misses. Two effects can lead to a miss rate that does not decrease—at least not as quickly as we might expect—as cache size increases: inherent communication and plateaus in the miss rate. Inherent communication leads to a certain frequency of coherence misses that are not significantly affected by increasing cache size. Thus if the cache size is increased while maintaining a fixed problem size, the coherence miss rate eventually limits the decrease in cache miss rate. This effect is most obvious in Barnes, where the coherence miss rate essentially becomes the entire miss rate.

A less important effect is a temporary plateau in the capacity miss rate that arises when the application has some fraction of its data present in cache but some significant portion of the data set does not fit in the cache or in caches that are slightly bigger. In LU, a very small cache (about 4 KB) can capture the pair of 16×16 blocks used in the inner loop; beyond that the next big improvement in capacity miss rate occurs when both matrices fit in the caches, which occurs when the total cache size is between 4 MB and 8 MB. This effect, sometimes called a *working set effect*, is partly at work between 32 KB and 128 KB for FFT, where the capacity miss rate drops only 0.3%. Beyond that cache size, a faster decrease in the capacity miss rate is seen, as a major data structure begins to reside in the cache. These plateaus are common in programs that deal with large arrays in a structured fashion.

Increasing the block size is another way to change the miss rate in a cache. In uniprocessors, larger block sizes are often optimal with larger caches. In multiprocessors, two new effects come into play: a reduction in spatial locality for shared data and a potential increase in miss rate due to false sharing. Several studies have shown that shared data have lower spatial locality than unshared data. Poorer locality means that for shared data, fetching larger blocks is less effective than in a uniprocessor, because the probability is higher that the block will be replaced before all its contents are referenced. Likewise, increasing the basic size also increases the potential frequency of false sharing, increasing the miss rate.

Figure 6.25 shows the miss rates as the cache block size is increased for a 16-processor run with a 64-KB cache. The most interesting behavior is in Barnes, where the miss rate initially declines and then rises due to an increase in the num-

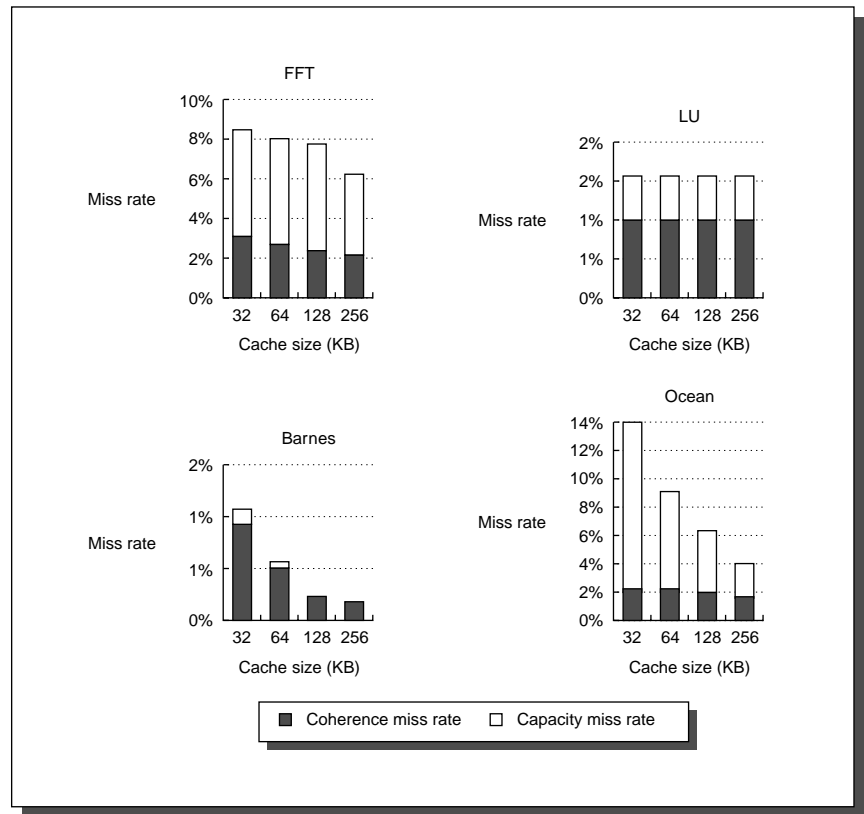


FIGURE 6.24 The miss rate usually drops as the cache size is increased, although coherence misses dampen the effect. The block size is 32 bytes and the cache is two-way set-associative. The processor count is fixed at 16 processors. Observe that the scale for each graph is different.

ber of coherence misses, which probably occurs because of false sharing. In the other benchmarks, increasing the block size decreases the overall miss rate. In Ocean and LU, the block size increase affects both the coherence and capacity miss rates about equally. In FFT, the coherence miss rate is actually decreased at a faster rate than the capacity miss rate. This reduction occurs because the communication in FFT is structured to be very efficient. In less optimized programs, we would expect more false sharing and less spatial locality for shared data, resulting in more behavior like that of Barnes.

Although the drop in miss rates with longer blocks may lead you to believe that choosing a longer block size is the best decision, the bottleneck in bus-based

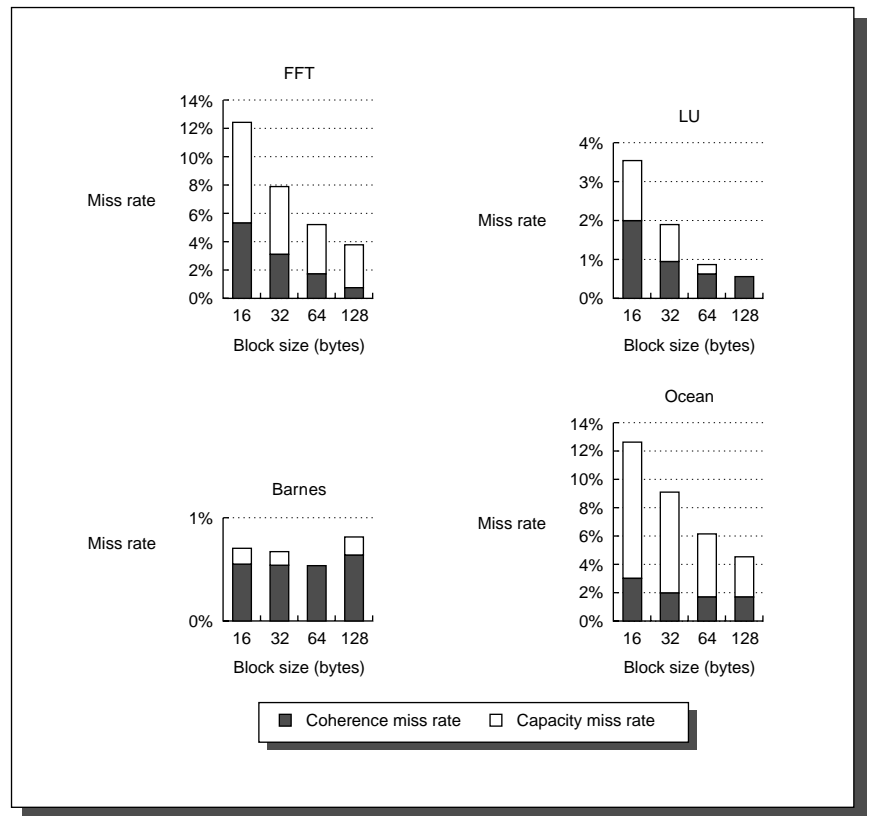


FIGURE 6.25 The data miss rate drops as the cache block size is increased. All these results are for a 16-processor run with a 64-KB cache and two-way set associativity. Once again we use different scales for each benchmark.

multiprocessors is often the limited memory and bus bandwidth. Larger blocks mean more bytes on the bus per miss. Figure 6.26 shows the growth in bus traffic as the block size is increased. This growth is most serious in the programs that have a high miss rate, especially Ocean. The growth in traffic can actually lead to performance slowdowns due both to longer miss penalties and to increased bus contention.

Summary: Performance of Snooping Cache Schemes

In this section we examined the cache performance of three very different workloads. We saw that the coherence traffic can introduce new behaviors in the memory system that do not respond as easily to changes in cache size or block size that are normally used to improve uniprocessor cache performance.

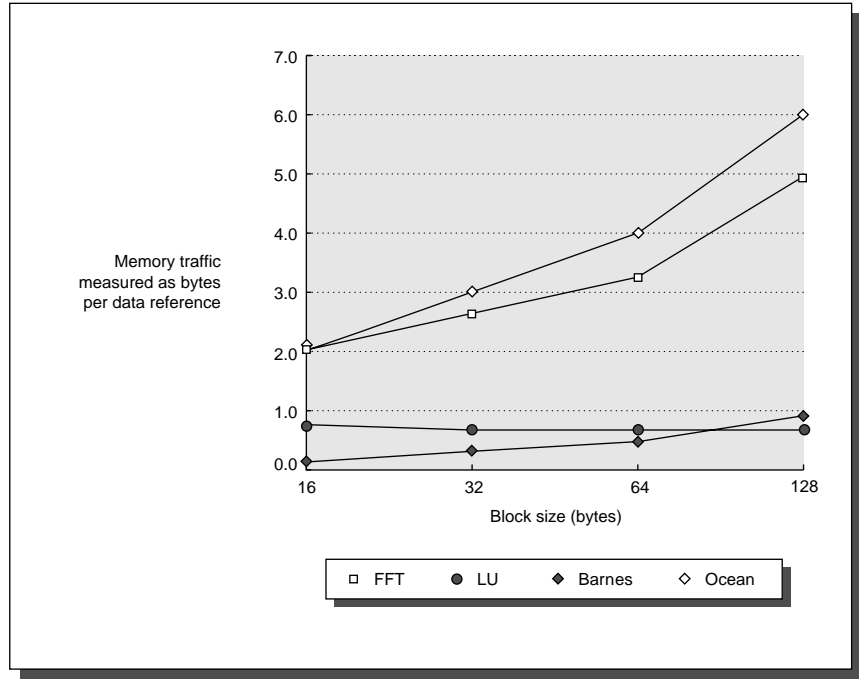


FIGURE 6.26 Bus traffic for data misses climbs steadily as the block size in the data cache is increased. The factor of 3 increase in traffic for Ocean is the best argument against larger block sizes. Remember that our protocol treats ownership or upgrade misses the same as other misses, slightly increasing the penalty for large cache blocks; in both Ocean and FFT this simplification accounts for less than 10% of the traffic.

In the commercial workload, the performance of the web searching and DSS benchmarks is reasonable (CPI of 1.3 and 1.6, respectively), while the OLTP benchmark is much worse (CPI=7.0). For OLTP, the large instruction working set demands a large cache to achieve acceptable performance. Increasing the cache size reduces the execution time, but is limited by the true and false sharing misses, which do not decrease as the cache grows. Similarly, increasing the processor counts increases true and false sharing, leading to an increase in memory access cycles. Fortunately, this workload responds favorably to an increase in block size, although the instruction miss rate remains similar. For these large workloads, it appears that very large (≥ 4 MB) off-chip caches with large block sizes (64-128 bytes) could work reasonably well.

In the multiprogrammed workload, the user and OS portions perform very differently, although neither has significant coherence traffic. In the OS portion, the compulsory and capacity contributions to the miss rate are much larger, leading to overall miss rates that are comparable to the worst programs in the parallel

program workload. User cache performance, on the other hand, is very good and compares to the best programs in the parallel program workload.

Coherence requests are a significant but not overwhelming component in the scientific processing workload. We can expect, however, that coherence requests will be more important in parallel programs that are less optimized.

The question of how these cache miss rates affect CPU performance depends on the rest of the memory system, including the latency and bandwidth of the interconnect and memory, a topic we return to in Section 6.11.

6.5 Distributed Shared-Memory Architectures

A scalable multiprocessor supporting shared memory could choose to exclude or include cache coherence. The simplest scheme for the hardware is to exclude cache coherence, focusing instead on a scalable memory system. Several companies have built this style of multiprocessor; the Cray T3D/E is best-known example. In such multiprocessors, memory is distributed among the nodes and all nodes are interconnected by a network. Access can be either local or remote—a controller inside each node decides, on the basis of the address, whether the data resides in the local memory or in a remote memory. In the latter case a message is sent to the controller in the remote memory to access the data.

These systems have caches, but to prevent coherence problems, shared data is marked as uncacheable and only private data is kept in the caches. Of course, software can still explicitly cache the value of shared data by copying the data from the shared portion of the address space to the local private portion of the address space that is cached. Coherence is then controlled by software. The advantage of such a mechanism is that little hardware support is required, although support for features such as block copy may be useful, since remote accesses fetch only single words (or double words) rather than cache blocks.

There are several disadvantages to this approach. First, compiler mechanisms for transparent software cache coherence are very limited. The techniques that currently exist apply primarily to programs with well-structured loop-level parallelism or a very strict form of object-oriented programming, and these techniques have significant overhead arising from explicitly copying data. For irregular problems or problems involving dynamic data structures and pointers (including operating systems, for example), compiler-based software cache coherence is currently impractical. The basic difficulty is that software-based coherence algorithms must be conservative: every block that *might* be shared must be treated as

if it *is* shared. Being conservative results in excess coherence overhead, because the compiler cannot predict the actual sharing accurately enough. Due to the complexity of the possible interactions, asking programmers to deal with coherence is unworkable.

Second, without cache coherence, the multiprocessor loses the advantage of being able to fetch and use multiple words in a single cache block for close to the cost of fetching one word. The benefits of spatial locality in shared data cannot be leveraged when single words are fetched from a remote memory for each reference. Support for a DMA mechanism among memories can help, but such mechanisms are often either costly to use (since they may require OS intervention) or expensive to implement since special-purpose hardware support and a buffer are needed. For message-passing programs, however, such mechanisms can be extremely useful, since programmers can overcome the usage penalties by using large messages.

Third, mechanisms for tolerating latency such as prefetch are more useful when they can fetch multiple words, such as a cache block, and where the fetched data remain coherent; we will examine this advantage in more detail later.

These disadvantages are magnified by the large latency of access to remote memory versus a local cache. For example, on the Cray T3E a local cache access has a latency of two cycles and is pipelined. A remote memory access takes up to 400 processor clock cycles for a remote memory using the 450 MHz Alpha processor in the T3E-900.

For these reasons, cache coherence is an accepted requirement in small-scale multiprocessors. For larger-scale architectures, there are new challenges to extending the cache-coherent shared-memory model. Although the bus can certainly be replaced with a more scalable interconnection network (the SUN Enterprise servers use up to four buses, e.g.), and we could certainly distribute the memory so that the memory bandwidth could also be scaled, the lack of scalability of the snooping coherence scheme needs to be addressed.

A snooping protocol requires communication with all caches on every cache miss, including writes of potentially shared data. The absence of any centralized data structure that tracks the state of the caches is both the fundamental advantage of a snooping-based scheme, since it allows it to be inexpensive, as well as its Achilles' heel when it comes to scalability. For example, with only 16 proces-

sors and a block size of 64 bytes and a 512 KB data cache, the total bus bandwidth demand (ignoring stall cycles) for the four programs in the scientific/technical workload ranges from about 1 GB/sec (for Barnes) to about 42 GB/sec (for FTT), assuming a processor that sustains a data reference every 1 ns, which is what a 2000 superscalar processor with nonblocking caches might generate. In comparison, the Sun Enterprise system with the Starfire interconnect, the highest bandwidth SMP system in 2000, can support about 12 GB/sec of random accesses for the 16x16 crossbar and has a maximum bandwidth of 21.3 GB/sec at the memory system. Although the cache size used in these simulations is moderate (though large enough to eliminate much of the uniprocessor miss traffic), so is the problem size.

Alternatively, we could build scalable shared-memory architectures that include cache coherency. The key is to find an alternative coherence protocol to the snooping protocol. One alternative protocol is a directory protocol. A *directory* keeps the state of every block that may be cached. Information in the directory includes which caches have copies of the block, whether it is dirty, and so on. (Section 6.11 on page 735 describes a hybrid approach that uses directories to extend a snooping protocol.)

Existing directory implementations associate an entry in the directory with each memory block. In typical protocols, the amount of information is proportional to the product of the number of memory blocks and the number of processors. This overhead is not a problem for multiprocessors with less than about two hundred processors, because the directory overhead will be tolerable. For larger multiprocessors, we need methods to allow the directory structure to be efficiently scaled. The methods that have been proposed either try to keep information for fewer blocks (e.g., only those in caches rather than all memory blocks) or try to keep fewer bits per entry.

To prevent the directory from becoming the bottleneck, directory entries can be distributed along with the memory, so that different directory accesses can go to different locations, just as different memory requests go to different memories. A distributed directory retains the characteristic that the sharing status of a block is always in a single known location. This property is what allows the coherence protocol to avoid broadcast. Figure 6.27 shows how our distributed-memory multiprocessor looks with the directories added to each node.

Directory-Based Cache-Coherence Protocols: The Basics

Just as with a snooping protocol, there are two primary operations that a directory protocol must implement: handling a read miss and handling a write to a shared,

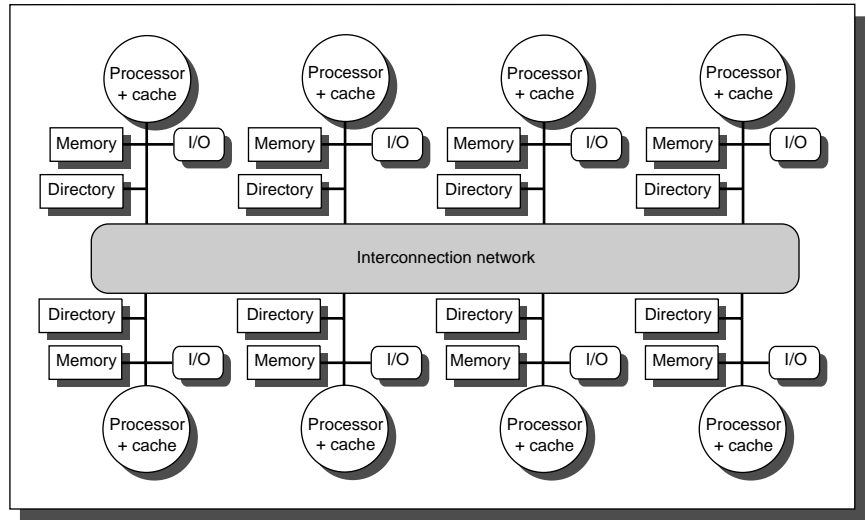


FIGURE 6.27 A directory is added to each node to implement cache coherence in a distributed-memory multiprocessor. Each directory is responsible for tracking the caches that share the memory addresses of the portion of memory in the node. The directory may communicate with the processor and memory over a common bus, as shown, or it may have a separate port to memory, or it may be part of a central node controller through which all intranode and internode communications pass.

clean cache block. (Handling a write miss to a shared block is a simple combination of these two.) To implement these operations, a directory must track the state of each cache block. In a simple protocol, these states could be the following:

- ▮ *Shared*—One or more processors have the block cached, and the value in memory is up to date (as well as in all the caches).
- ▮ *Uncached*—No processor has a copy of the cache block.
- ▮ *Exclusive*—Exactly one processor has a copy of the cache block and it has written the block, so the memory copy is out of date. The processor is called the *owner* of the block.

In addition to tracking the state of each cache block, we must track the processors that have copies of the block when it is shared, since they will need to be invalidated on a write. The simplest way to do this is to keep a bit vector for each memory block. When the block is shared, each bit of the vector indicates whether the corresponding processor has a copy of that block. We can also use the bit vec-

tor to keep track of the owner of the block when the block is in the exclusive state. For efficiency reasons, we also track the state of each cache block at the individual caches.

The states and transitions for the state machine at each cache are identical to what we used for the snooping cache, although the actions on a transition are slightly different. We make the same simplifying assumptions that we made in the case of the snooping cache: attempts to write data that is not exclusive in the writer's cache always generate write misses, and the processors block until an access completes. Since the interconnect is no longer a bus and since we want to avoid broadcast, there are two additional complications. First, we cannot use the interconnect as a single point of arbitration, a function the bus performed in the snooping case. Second, because the interconnect is message oriented (unlike the bus, which is transaction oriented), many messages must have explicit responses.

Before we see the protocol state diagrams, it is useful to examine a catalog of the message types that may be sent between the processors and the directories. Figure 6.28 shows the type of messages sent among nodes. The *local* node is the node where a request originates. The *home* node is the node where the memory location and the directory entry of an address reside. The physical address space is statically distributed, so the node that contains the memory and directory for a given physical address is known. For example, the high-order bits may provide the node number, while the low-order bits provide the offset within the memory on that node. The local node may also be the home node. The directory must be accessed when the home node is the local node, since copies may exist in yet a third node, called a remote node.

A *remote* node is the node that has a copy of a cache block, whether exclusive (in which case it is the only copy) or shared. A remote node may be the same as either the local node or the home node. In such cases, the basic protocol does not change, but interprocessor messages may be replaced with intraprocessor messages.

In this section, we assume a simple model of memory consistency. To minimize the type of messages and the complexity of the protocol, we make an assumption that messages will be received and acted upon in the same order they are sent. This assumption may not be true in practice, and can result in additional complications, some of which we address in section 6.8 when we discuss memory consistency models. In this section, we use this assumption to ensure that invalidates sent by a processor are honored immediately.

Message type	Source	Destination	Message contents	Function of this message
Read miss	Local cache	Home directory	P, A	Processor P has a read miss at address A; request data and make P a read sharer.
Write miss	Local cache	Home directory	P, A	Processor P has a write miss at address A; — request data and make P the exclusive owner.
Invalidate	Home directory	Remote cache	A	Invalidate a shared copy of data at address A.
Fetch	Home directory	Remote cache	A	Fetch the block at address A and send it to its home directory; change the state of A in the remote cache to shared.
Fetch/invalidate	Home directory	Remote cache	A	Fetch the block at address A and send it to its home directory; invalidate the block in the cache.
Data value reply	Home directory	Local cache	D	Return a data value from the home memory.
Data write back	Remote cache	Home directory	A, D	Write back a data value for address A.

FIGURE 6.28 The possible messages sent among nodes to maintain coherence are shown with the source and destination node, the contents (where P=requesting processor number), A=requested address, and D=data contents), and the function of the message. The first two messages are miss requests sent by the local cache to the home. The third through fifth messages are messages sent to a remote cache by the home when the home needs the data to satisfy a read or write miss request. Data value replies are used to send a value from the home node back to the requesting node. Data value write backs occur for two reasons: when a block is replaced in a cache and must be written back to its home memory, and also in reply to fetch or fetch/invalidate messages from the home. Writing back the data value whenever the block becomes shared simplifies the number of states in the protocol, since any dirty block must be exclusive and any shared block is always available in the home memory.

An Example Directory Protocol

The basic states of a cache block in a directory-based protocol are exactly like those in a snooping protocol, and the states in the directory are also analogous to those we showed earlier. Thus we can start with simple state diagrams that show

the state transitions for an individual cache block and then examine the state diagram for the directory entry corresponding to each block in memory. As in the snooping case, these state transition diagrams do not represent all the details of a coherence protocol; however, the actual controller is highly dependent on a number of details of the multiprocessor (message delivery properties, buffering structures, and so on). In this section we present the basic protocol state diagrams. The knotty issues involved in implementing these state transition diagrams are examined in Appendix E, along with similar problems that arise for snooping caches.

Figure 6.29 shows the protocol actions to which an individual cache responds.

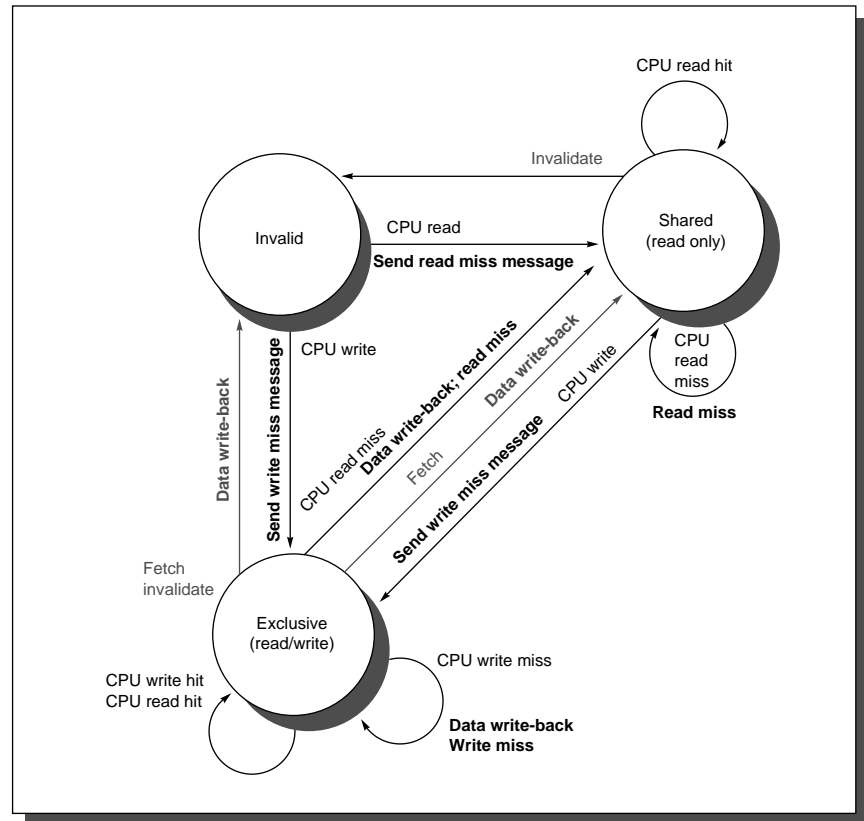


FIGURE 6.29 State transition diagram for an individual cache block in a directory-based system. Requests by the local processor are shown in black and those from the home directory are shown in gray. The states are identical to those in the snooping case, and the transactions are very similar, with explicit invalidate and write-back requests replacing the write misses that were formerly broadcast on the bus. As we did for the snooping controller, we assume that an attempt to write a shared cache block is treated as a miss; in practice, such a transaction can be treated as an ownership request or upgrade request and can deliver ownership without requiring that the cache block be fetched.

We use the same notation as in the last section, with requests coming from outside the node in gray and actions in bold. The state transitions for an individual cache are caused by read misses, write misses, invalidates, and data fetch requests; these operations are all shown in Figure 6.29. An individual cache also generates read and write miss messages that are sent to the home directory. Read

and write misses require data value replies, and these events wait for replies before changing state.

The operation of the state transition diagram for a cache block in Figure 6.29 is essentially the same as it is for the snooping case: the states are identical, and the stimulus is almost identical. The write miss operation, which was broadcast on the bus in the snooping scheme, is replaced by the data fetch and invalidate operations that are selectively sent by the directory controller. Like the snooping protocol, any cache block must be in the exclusive state when it is written and any shared block must be up to date in memory.

In a directory-based protocol, the directory implements the other half of the coherence protocol. A message sent to a directory causes two different types of actions: updates of the directory state, and sending additional messages to satisfy the request. The states in the directory represent the three standard states for a block; unlike in a snoopy scheme, however, the directory state indicates the state of all the cached copies of a memory block, rather than for a single cache block. The memory block may be uncached by any node, cached in multiple nodes and readable (shared), or cached exclusively and writable in exactly one node. In addition to the state of each block, the directory must track the set of processors that have a copy of a block; we use a set called *Sharers* to perform this function. In multiprocessors with less than 64 nodes (which may represent 2-4 times as many processors), this set is typically kept as a bit vector. In larger multiprocessors, other techniques, which we discuss in the Exercises, are needed. Directory requests need to update the set *Sharers* and also read the set to perform invalidations.

Figure 6.30 shows the actions taken at the directory in response to messages received. The directory receives three different requests: read miss, write miss, and data write back. The messages sent in response by the directory are shown in bold, while the updating of the set *Sharers* is shown in bold italics. Because all the stimulus messages are external, all actions are shown in gray. Our simplified protocol assumes that some actions are atomic, such as requesting a value and sending it to another node; a realistic implementation cannot use this assumption.

To understand these directory operations, let's examine the requests received and actions taken state by state. When a block is in the uncached state the copy in memory is the current value, so the only possible requests for that block are

- ⁿ *Read miss*—The requesting processor is sent the requested data from memory and the requestor is made the only sharing node. The state of the block is made shared.
- ⁿ *Write miss*—The requesting processor is sent the value and becomes the Sharing node. The block is made exclusive to indicate that the only valid copy is cached. *Sharers* indicates the identity of the owner.

When the block is in the shared state the memory value is up-to-date, so the same two requests can occur:

was the owner (since it still has a readable copy).

- ⁿ *Data write-back*—The owner processor is replacing the block and therefore must write it back. This write-back makes the memory copy up to date (the home directory essentially becomes the owner), the block is now uncached, and the sharer set is empty.
- ⁿ *Write miss*—The block has a new owner. A message is sent to the old owner causing the cache to invalidate the block and send the value to the directory, from which it is sent to the requesting processor, which becomes the new owner. Sharers is set to the identity of the new owner, and the state of the block remains exclusive.

This state transition diagram in Figure 6.30 is a simplification, just as it was in the snooping cache case. In the directory case it is a larger simplification, since our assumption that bus transactions are atomic no longer applies. Appendix E explores these issues in depth.

In addition, the directory protocols used in real multiprocessors contain additional optimizations. In particular, in this protocol when a read or write miss occurs for a block that is exclusive, the block is first sent to the directory at the home node. From there it is stored into the home memory and also sent to the original requesting node. Many of the protocols in use in commercial multiprocessors forward the data from the owner node to the requesting node directly (as well as performing the write back to the home). Such optimizations often add complexity by increasing the possibility of deadlock and by increasing the types of messages that must be handled.

6.6 Performance of Distributed Shared-Memory Multiprocessors

The performance of a directory-based multiprocessor depends on many of the same factors that influence the performance of bus-based multiprocessors (e.g., cache size, processor count, and block size), as well as the distribution of misses to various locations in the memory hierarchy. The location of a requested data item depends on both the initial allocation and the sharing patterns. We start by examining the basic cache performance of our scientific/technical workload and then look at the effect of different types of misses.

Because the multiprocessor is larger and has longer latencies than our snooping-based multiprocessor, we begin with a slightly larger cache (128 KB) and a larger block size of 64 bytes.

In distributed memory architectures, the distribution of memory requests between local and remote is key to performance, because it affects both the consumption of global bandwidth and the latency seen by requests. Therefore, for the figures in this section we separate the cache misses into local and remote requests. In looking at the figures, keep in mind that, for these applications, most of the remote misses that arise are coherence misses, although some capacity misses can also be remote, and in some applications with poor data distribution, such misses can be significant (see the Pitfall on page 758).

As Figure 6.31 shows, the miss rates with these cache sizes are not affected much by changes in processor count, with the exception of Ocean, where the miss rate rises at 64 processors. This rise results from two factors: an increase in mapping conflicts in the cache that occur when the grid becomes small, which leads to a rise in local misses, and an increase in the number of the coherence misses, which are all remote.

Figure 6.32 shows how the miss rates change as the cache size is increased, assuming a 64-processor execution and 64-byte blocks. These miss rates decrease at rates that we might expect, although the dampening effect caused by little or no reduction in coherence misses leads to a slower decrease in the remote misses than in the local misses. By the time we reach the largest cache size shown, 512 KB, the remote miss rate is equal to or greater than the local miss rate. Larger caches would amplify this trend.

We examine the effect of changing the block size in Figure 6.33. Because these applications have good spatial locality, increases in block size reduce the miss rate, even for large blocks, although the performance benefits for going to the largest blocks are small. Furthermore, most of the improvement in miss rate comes from a reduction in the local misses.

Rather than plot the memory traffic, Figure 6.34 plots the number of bytes required per data reference versus block size, breaking the requirement into local and global bandwidth. In the case of a bus, we can simply aggregate the demands of each processor to find the total demand for bus and memory bandwidth. For a scalable interconnect, we can use the data in Figure 6.34 to compute the required per-node global bandwidth and the estimated bisection bandwidth, as the next Example shows.

EXAMPLE Assume a 64-processor multiprocessor with 1GHz processors that sustain one memory reference per processor clock. For a 64-byte block size, the remote miss rate is 0.7%. Find the per-node and estimated bisection bandwidth for FFT. Assume that the processor does not stall for remote memory requests; this might be true if, for example, all remote data were prefetched. How do these bandwidth requirements compare to various interconnection technologies?

ANSWER The per-node bandwidth is simply the number of data bytes per reference

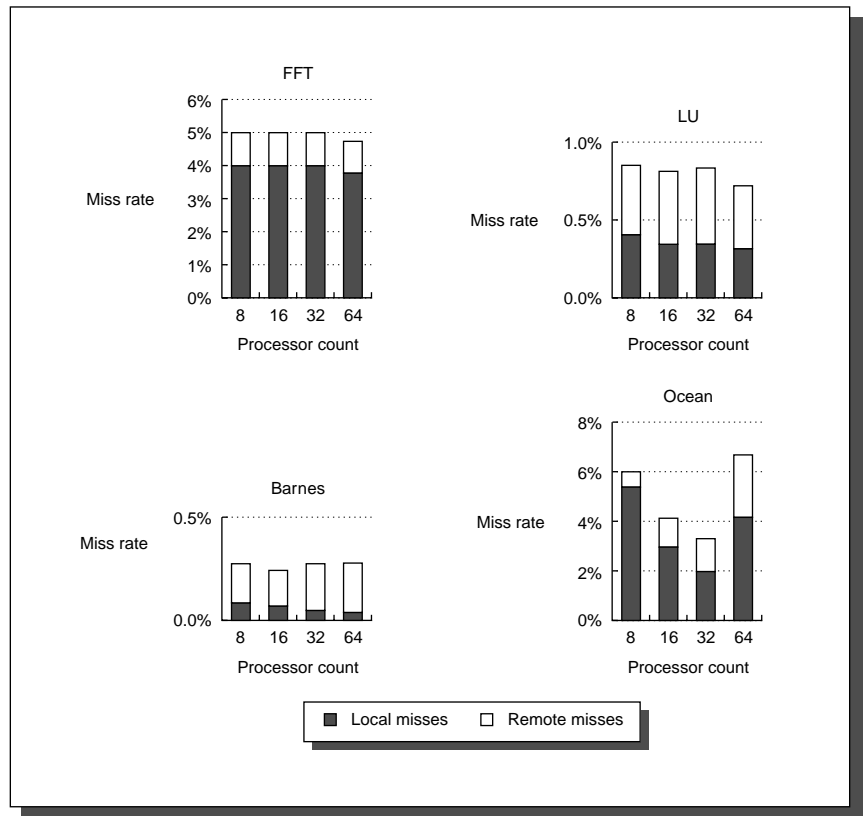


FIGURE 6.31 The data miss rate is often steady as processors are added for these benchmarks. Because of its grid structure, Ocean has an initially decreasing miss rate, which rises when there are 64 processors. For Ocean, the local miss rate drops from 5% at 8 processors to 2% at 32, before rising to 4% at 64. The remote miss rate in Ocean, driven primarily by communication, rises monotonically from 1% to 2.5%. Note that to show the detailed behavior of each benchmark, different scales are used on the y-axis. The cache for all these runs is 128 KB, two-way set associative, with 64-byte blocks. Remote misses include any misses that require communication with another node, whether to fetch the data or to deliver an invalidate. In particular, in this figure and other data in this section, the measurement of remote misses includes write upgrade misses where the data is up to date in the local memory but cached elsewhere and, therefore, requires invalidations to be sent. Such invalidations do indeed generate remote traffic, but may or may not delay the write, depending on the consistency model (see section 6.8).

times the reference rate: $0.7\% \times 1000 \times 64 = 448$ MB/sec. This rate is somewhat higher than the hardware sustainable transfer rate for the CrayT3E (using a block prefetch) and lower than that for an SGI Origin 3000 (1.6 GB/processor pair). The FFT per-node bandwidth demand ex-

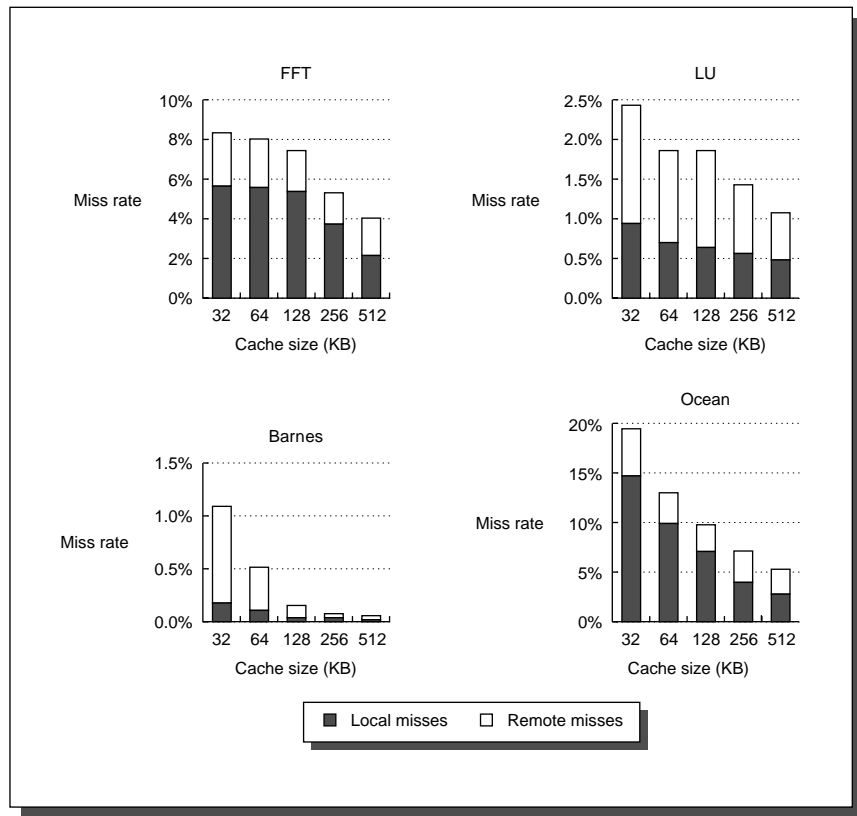


FIGURE 6.32 Miss rates decrease as cache sizes grow. Steady decreases are seen in the local miss rate, while the remote miss rate declines to varying degrees, depending on whether the remote miss rate had a large capacity component or was driven primarily by communication misses. In all cases, the decrease in the local miss rate is larger than the decrease in the remote miss rate. The plateau in the miss rate of FFT, which we mentioned in the last section, ends once the cache exceeds 128 KB. These runs were done with 64 processors and 64-byte cache blocks.

ceeds the bandwidth sustainable from the fastest standard networks by more than a factor of 5.

FFT performs all-to-all communication, so the bisection bandwidth is equal to the number of processors times the per-node bandwidth, or about $64 \times 448 \text{ MB/sec} = 28.7 \text{ GB/sec}$. The SGI Origin 3000 with 64-processors has a bisection bandwidth of about 50 GB/sec. No standard networking technology comes close.

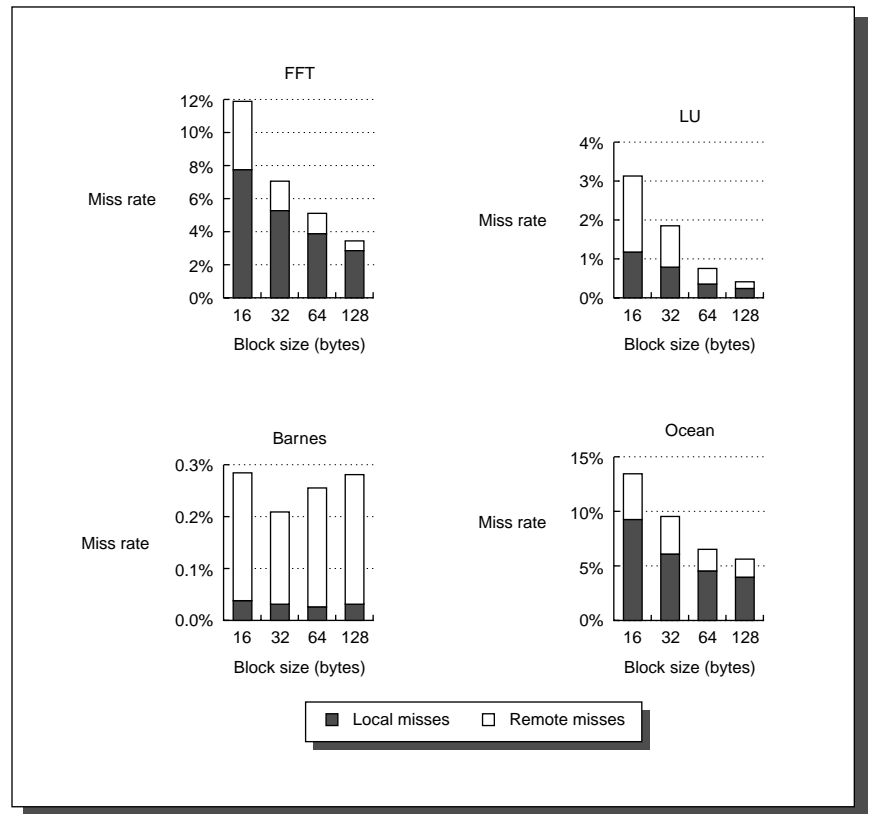


FIGURE 6.33 Data miss rate versus block size assuming a 128-KB cache and 64 processors in total. Although difficult to see, the coherence miss rate in Barnes actually rises for the largest block size, just as in the last section.

The previous Example looked at the bandwidth demands. The other key issue for a parallel program is remote memory access time, or latency. To get insight into this, we use a simple example of a directory-based multiprocessor. Figure 6.35 shows the parameters we assume for our simple multiprocessor model. It assumes that the time to first word for a local memory access is 85 processor cycles and that the path to local memory is 16 bytes wide, while the network interconnect is 4 bytes wide. This model ignores the effects of contention, which are probably not too serious in the parallel benchmarks we examine, with the possible exception of FFT, which uses all-to-all communication. Contention could have a serious performance impact in other workloads.

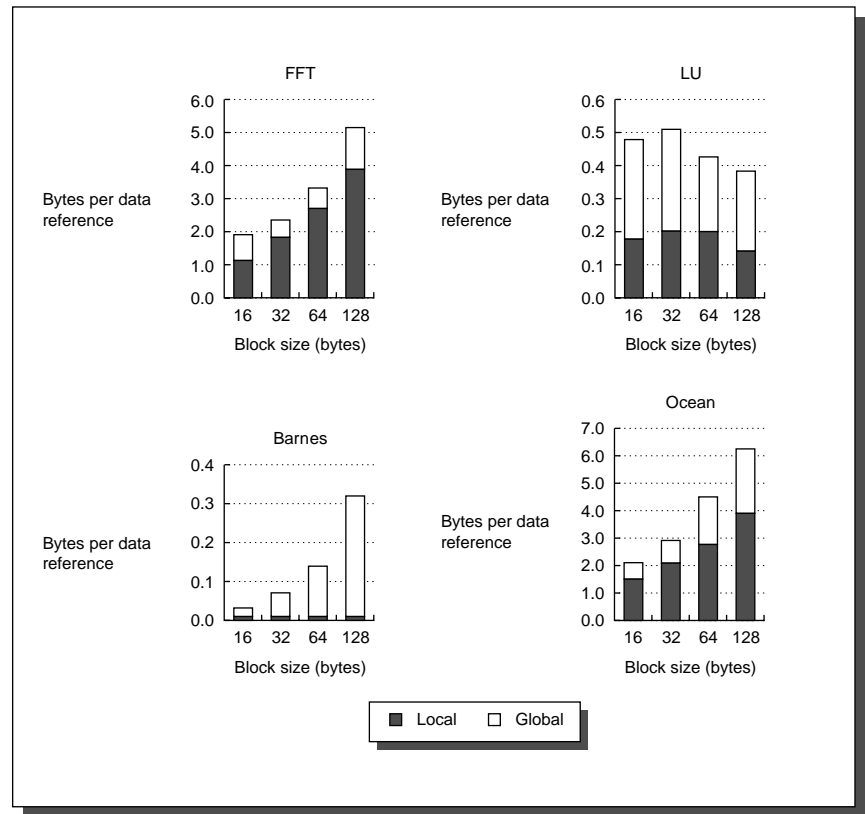


FIGURE 6.34 The number of bytes per data reference climbs steadily as block size is increased. These data can be used to determine the bandwidth required per node both internally and globally. The data assumes a 128-KB cache for each of 64 processors.

Figure 6.36 shows the cost in cycles for the average memory reference, assuming the parameters in Figure 6.35. Only the latencies for each reference type are counted. Each bar indicates the contribution from cache hits, local misses, remote misses, and 3-hop remote misses. The cost is influenced by the total frequency of cache misses and upgrades, as well as by the distribution of the location where the miss is satisfied. The cost for a remote memory reference is fairly steady as the processor count is increased, except for Ocean. The increasing miss rate in Ocean for 64 processors is clear in Figure 6.31. As the miss rate increases, we should expect the time spent on memory references to increase also.

Although Figure 6.36 shows the memory access cost, which is the dominant multiprocessor cost in these benchmarks, a complete performance model would

Characteristic	Processor clock cycles ≤ 16 processor	Processor clock cycles 17–64 processor
Cache hit	1	1
Cache miss to local memory	85	85
Cache miss to remote home directory	125	150
Cache miss to remotely cached data (3-hop miss)	140	170

FIGURE 6.35 Characteristics of the example directory-based multiprocessor. Misses can be serviced locally (including from the local directory), at a remote home node, or using the services of both the home node and another remote node that is caching an exclusive copy. This last case is called a 3-hop miss and has a higher cost because it requires interrogating both the home directory and a remote cache. Note that this simple model does not account for invalidation time, but does include some factor for increasing interconnect time. These remote access latencies are based on those in an SGI Origin 3000, the fastest scalable interconnect system in 2000, and assume a 500 MHz processor.

need to consider the effect of contention in the memory system, as well as the losses arising from synchronization delays.

The coherence protocols that we have discussed so far have made several simplifying assumptions. In practice, real protocols must deal with two realities: nonatomicity of operations and finite buffering. We have seen why certain operations (such as a write miss) cannot be atomic. In DSM multiprocessors the presence of only a finite number of buffers to hold message requests and replies introduces additional possibilities for deadlock. The challenge for the designer is to create a protocol that works correctly and without deadlock, using nonatomic actions and finite buffers as the building blocks. These factors are fundamental challenges in all parallel multiprocessors, and the solutions are applicable to a wide variety of protocol design environments, both in hardware and in software.

Because this material is extremely complex and not necessary to comprehend the rest of the chapter, we have placed it in Appendix E. For the interested reader, Appendix E shows how the specific problems in our coherence protocols are solved and illustrates the general principles that are more globally applicable. It describes the problems arising in snooping cache implementations, as well as the more complex problems that arise in more distributed systems using directories. If you want to understand how either state-of-the-art SMPs (which use split transactions buses and nonblocking memory accesses) or DSM multiprocessors really work and why designing them is such a challenge, go read Appendix E!

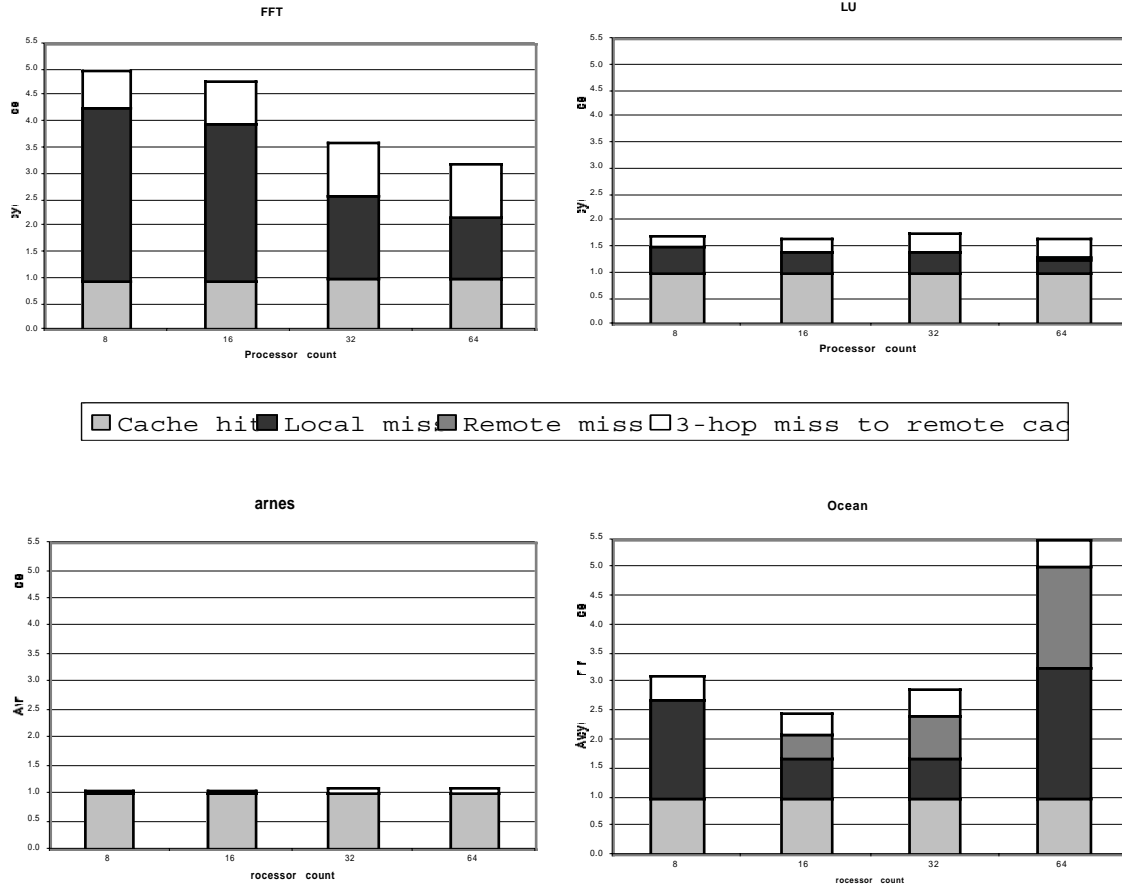


FIGURE 6.36 The effective latency of memory references in a DSM multiprocessor depends both on the relative frequency of cache misses and on the location of the memory where the accesses are served. These plots show the memory access cost (a metric called average memory access time in Chapter 5) for each of the benchmarks for 8, 16, 32, and 64 processors, assuming a 512KB data cache that is two-way set associative with 64-byte blocks. The average memory access cost is composed of four different types of accesses, with the cost of each type given in Figure 6.35. For the Barnes and LU benchmarks, the low miss rates lead to low overall access times. In FFT, the higher access cost is determined by a higher local miss rate (1-4%) and a significant 3-hop miss rate (1%). The improvement in FFT comes from the reduction in local miss rate from 4% to 1%, as the aggregate cache increases. Ocean shows the biggest change in the cost of memory accesses, and the highest overall cost at 64 processors. The high cost is driven primarily by a high local miss rate (average 1.6%). The memory access cost drops from 8 to 16 processors as the grids more easily fit in the individual caches. At 64 processors, the data set size is too small to map properly and both local misses and coherence misses rise, as we saw in Figure 6.31.

6.7 Synchronization

Synchronization mechanisms are typically built with user-level software routines that rely on hardware-supplied synchronization instructions. For smaller multiprocessors or low-contention situations, the key hardware capability is an uninterruptible instruction or instruction sequence capable of atomically retrieving and changing a value. Software synchronization mechanisms are then constructed using this capability. For example, we will see how very efficient spin locks can be built using a simple hardware synchronization instruction and the coherence mechanism. In larger-scale multiprocessors or high-contention situations, synchronization can become a performance bottleneck, because contention introduces additional delays and because latency is potentially greater in such a multiprocessor. We will see how contention can arise in implementing some common user-level synchronization operations and examine more powerful hardware-supported synchronization primitives that can reduce contention as well as latency.

We begin by examining the basic hardware primitives, then construct several well-known synchronization routines with the primitives, and then turn to performance problems in larger multiprocessors and solutions for those problems.

Basic Hardware Primitives

The key ability we require to implement synchronization in a multiprocessor is a set of hardware primitives with the ability to atomically read and modify a memory location. Without such a capability, the cost of building basic synchronization primitives will be too high and will increase as the processor count increases. There are a number of alternative formulations of the basic hardware primitives, all of which provide the ability to atomically read and modify a location, together with some way to tell if the read and write were performed atomically. These hardware primitives are the basic building blocks that are used to build a wide variety of user-level synchronization operations, including things such as locks and barriers. In general, architects do not expect users to employ the basic hardware primitives, but instead expect that the primitives will be used by system programmers to build a synchronization library, a process that is often complex and tricky. Let's start with one such hardware primitive and show how it can be used to build some basic synchronization operations.

One typical operation for building synchronization operations is the *atomic exchange*, which interchanges a value in a register for a value in memory. To see how to use this to build a basic synchronization operation, assume that we want to build a simple lock where the value 0 is used to indicate that the lock is free and a 1 is used to indicate that the lock is unavailable. A processor tries to set the lock by doing an exchange of 1, which is in a register, with the memory address corresponding to the lock. The value returned from the exchange instruction is 1 if some other processor had already claimed access and 0 otherwise. In the latter

case, the value is also changed to be 1, preventing any competing exchange from also retrieving a 0.

For example, consider two processors that each try to do the exchange simultaneously: This race is broken since exactly one of the processors will perform the exchange first, returning 0, and the second processor will return 1 when it does the exchange. The key to using the exchange (or swap) primitive to implement synchronization is that the operation is atomic: the exchange is indivisible and two simultaneous exchanges will be ordered by the write serialization mechanisms. It is impossible for two processors trying to set the synchronization variable in this manner to both think they have simultaneously set the variable.

There are a number of other atomic primitives that can be used to implement synchronization. They all have the key property that they read and update a memory value in such a manner that we can tell whether or not the two operations executed atomically. One operation, present in many older multiprocessors, is *test-and-set*, which tests a value and sets it if the value passes the test. For example, we could define an operation that tested for 0 and set the value to 1, which can be used in a fashion similar to how we used atomic exchange. Another atomic synchronization primitive is *fetch-and-increment*: it returns the value of a memory location and atomically increments it. By using the value 0 to indicate that the synchronization variable is unclaimed, we can use fetch-and-increment, just as we used exchange. There are other uses of operations like fetch-and-increment, which we will see shortly.

A slightly different approach to providing this atomic read-and-update operation has been used in some recent multiprocessors. Implementing a single atomic memory operation introduces some challenges, since it requires both a memory read and a write in a single, uninterruptible instruction. This requirement complicates the implementation of coherence, since the hardware cannot allow any other operations between the read and the write, and yet must not deadlock.

An alternative is to have a pair of instructions where the second instruction returns a value from which it can be deduced whether the pair of instructions was executed as if the instructions were atomic. The pair of instructions is effectively atomic if it appears as if all other operations executed by any processor occurred before or after the pair. Thus, when an instruction pair is effectively atomic, no other processor can change the value between the instruction pair.

The pair of instructions includes a special load called a *load linked* or *load locked* and a special store called a *store conditional*. These instructions are used in sequence: If the contents of the memory location specified by the load linked are changed before the store conditional to the same address occurs, then the store conditional fails. If the processor does a context switch between the two instructions, then the store conditional also fails. The store conditional is defined to return a value indicating whether or not the store was successful. Since the load linked returns the initial value and the store conditional returns 1 if it succeeds and 0 otherwise, the following sequence implements an atomic exchange on the memory location specified by the contents of R1:

```

try:  MOV    R3,R4,R0      ;mov exchange value
      LL     R2,0(R1)      ;load linked
      SC     R3,0(R1)      ;store conditional
      BEQZ   R3,try        ;branch store fails
      MOV    R4,R2        ;put load value in R4

```

At the end of this sequence the contents of R4 and the memory location specified by R1 have been atomically exchanged (ignoring any effect from delayed branches). Any time a processor intervenes and modifies the value in memory between the LL and SC instructions, the SC returns 0 in R3, causing the code sequence to try again.

An advantage of the load linked/store conditional mechanism is that it can be used to build other synchronization primitives. For example, here is an atomic fetch-and-increment:

```

try:  LL     R2,0(R1)      ;load linked
      DADDUI R3,R2,#1      ;increment
      SC     R3,0(R1)      ;store conditional
      BEQZ   R3,try        ;branch store fails

```

These instructions are typically implemented by keeping track of the address specified in the LL instruction in a register, often called the *link register*. If an interrupt occurs, or if the cache block matching the address in the link register is invalidated (for example, by another SC), the link register is cleared. The SC instruction simply checks that its address matches that in the link register; if so, the SC succeeds; otherwise, it fails. Since the store conditional will fail after either another attempted store to the load linked address or any exception, care must be taken in choosing what instructions are inserted between the two instructions. In particular, only register-register instructions can safely be permitted; otherwise, it is possible to create deadlock situations where the processor can never complete the SC. In addition, the number of instructions between the load linked and the store conditional should be small to minimize the probability that either an unrelated event or a competing processor causes the store conditional to fail frequently.

Implementing Locks Using Coherence

Once we have an atomic operation, we can use the coherence mechanisms of a multiprocessor to implement *spin locks*: locks that a processor continuously tries to acquire, spinning around a loop until it succeeds. Spin locks are used when a programmer expects the lock to be held for a very short amount of time and when

she wants the process of locking to be low latency when the lock is available. Because spin locks tie up the processor, waiting in a loop for the lock to become free, they are inappropriate in some circumstances.

The simplest implementation, which we would use if there were no cache coherence, would keep the lock variables in memory. A processor could continually try to acquire the lock using an atomic operation, say exchange, and test whether the exchange returned the lock as free. To release the lock, the processor simply stores the value 0 to the lock. Here is the code sequence to lock a spin lock whose address is in R1 using an atomic exchange:

```
                DADDUI R2,R0,#1
lockit:         EXCH   R2,0(R1)      ;atomic exchange
                BNEZ   R2,lockit    ;already locked?
```

If our multiprocessor supports cache coherence, we can cache the locks using the coherence mechanism to maintain the lock value coherently. Caching locks has two advantages. First, it allows an implementation where the process of “spinning” (trying to test and acquire the lock in a tight loop) could be done on a local cached copy rather than requiring a global memory access on each attempt to acquire the lock. The second advantage comes from the observation that there is often locality in lock accesses: that is, the processor that used the lock last will use it again in the near future. In such cases, the lock value may reside in the cache of that processor, greatly reducing the time to acquire the lock.

Obtaining the first advantage—being able to spin on a local cached copy rather than generating a memory request for each attempt to acquire the lock—requires a change in our simple spin procedure. Each attempt to exchange in the loop directly above requires a write operation. If multiple processors are attempting to get the lock, each will generate the write. Most of these writes will lead to write misses, since each processor is trying to obtain the lock variable in an exclusive state.

Thus we should modify our spin-lock procedure so that it spins by doing reads on a local copy of the lock until it successfully sees that the lock is available. Then it attempts to acquire the lock by doing a swap operation. A processor first reads the lock variable to test its state. A processor keeps reading and testing until the value of the read indicates that the lock is unlocked. The processor then races against all other processes that were similarly “spin waiting” to see who can lock the variable first. All processes use a swap instruction that reads the old value and stores a 1 into the lock variable. The single winner will see the 0, and the losers will see a 1 that was placed there by the winner. (The losers will continue to set the variable to the locked value, but that doesn’t matter.) The winning processor executes the code after the lock and, when finished, stores a 0 into the lock variable to release the lock, which starts the race all over again. Here is the code to perform this spin lock (remember that 0 is unlocked and 1 is locked):

```

lockit: LD      R2,0(R1)      ;load of lock
        BNEZ    R2,lockit    ;not available-spin
        DADDUI  R2,R0,#1     ;load locked value
        EXCH    R2,0(R1)     ;swap
        BNEZ    R2,lockit    ;branch if lock wasn't 0

```

Let's examine how this "spin-lock" scheme uses the cache-coherence mechanisms. Figure 6.37 shows the processor and bus or directory operations for multiple processes trying to lock a variable using an atomic swap. Once the processor with the lock stores a 0 into the lock, all other caches are invalidated and must fetch the new value to update their copy of the lock. One such cache gets the copy of the unlocked value (0) first and performs the swap. When the cache miss of other processors is satisfied, they find that the variable is already locked, so they must return to testing and spinning.

Step	Processor P0	Processor P1	Processor P2	Coherence state of lock	Bus/directory activity
1	Has lock	Spins, testing if lock = 0	Spins, testing if lock = 0	Shared	None
2	Set lock to 0	(Invalidate received)	(Invalidate received)	Exclusive	Write invalidate of lock variable from P0
3		Cache miss	Cache miss	Shared	Bus/directory services P2 cache miss; write back from P0
4		(Waits while bus/directory busy)	Lock = 0	Shared	Cache miss for P2 satisfied
5		Lock = 0	Executes swap, gets cache miss	Shared	Cache miss for P1 satisfied
6		Executes swap, gets cache miss	Completes swap: returns 0 and sets Lock = 1	Exclusive	Bus/directory services P2 cache miss; generates invalidate
7		Swap completes and returns 1	Enter critical section	Shared	Bus/directory services P1 cache miss; generates write back
8		Spins, testing if lock = 0			None

FIGURE 6.37 Cache-coherence steps and bus traffic for three processors, P0, P1, and P2. This figure assumes write-invalidate coherence. P0 starts with the lock (step 1). P0 exits and unlocks the lock (step 2). P1 and P2 race to see which reads the unlocked value during the swap (steps 3–5). P2 wins and enters the critical section (steps 6 and 7), while P1's attempt fails so it starts spin waiting (steps 7 and 8). In a real system, these events will take many more than eight clock ticks, since acquiring the bus and replying to misses takes much longer.

This example shows another advantage of the load-linked/store-conditional primitives: the read and write operation are explicitly separated. The load linked need not cause any bus traffic. This fact allows the following simple code sequence, which has the same characteristics as the optimized version using exchange (R1 has the address of the lock):

```
lockit:  LL      R2,0(R1)      ;load linked
         BNEZ    R2,lockit     ;not available-spin
         DADDUI  R2,R0,#1      ;locked value
         SC      R2,0(R1)      ;store
         BEQZ    R2,lockit     ;branch if store fails
```

The first branch forms the spinning loop; the second branch resolves races when two processors see the lock available simultaneously.

Although our spin lock scheme is simple and compelling, it has difficulty scaling up to handle many processors because of the communication traffic generated when the lock is released. The next section discusses these problems in more detail, as well as techniques to overcome these problems in larger multiprocessors.

Synchronization Performance Challenges

To understand why the simple spin-lock scheme of the previous section does not scale well, imagine a large multiprocessor with all processors contending for the same lock. The directory or bus acts as a point of serialization for all the processors, leading to lots of contention, as well as traffic. The following Example shows how bad things can be.

EXAMPLE Suppose there are 10 processors on a bus that each try to lock a variable simultaneously. Assume that each bus transaction (read miss or write miss) is 100 clock cycles long. You can ignore the time of the actual read or write of a lock held in the cache, as well as the time the lock is held (they won't matter much!). Determine the number of bus transactions required for all 10 processors to acquire the lock, assuming they are all spinning when the lock is released at time 0. About how long will it take to process the 10 requests? Assume that the bus is totally fair so that every pending request is serviced before a new request and that the processors are equally fast.

ANSWER Figure 6.38 shows the sequence of events from the time of the release to the time to the next release. Of course, the number of processors contending for the lock drops by one each time the lock is acquired, which reduces the average cost to 1550 cycles. Thus for 10 lock-unlock pairs it will take over 15,000 cycles for the processors to pass through the lock. Fur-

thermore, the average processor will spend half this time idle, simply trying to get the lock. The number of bus transactions involved is over 200!

Event	Duration
Read miss by all waiting processors to fetch lock (10×100)	1000
Write miss by releasing processor and invalidates	100
Read miss by all waiting processors (10×100)	1000
Write miss by all waiting processors, one successful lock (100), and invalidation of all lock copies (9×100)	1000
Total time for one processor to acquire and release lock	3100 clocks

FIGURE 6.38 The time to acquire and release a single lock when 10 processors contend for the lock, assuming each bus transaction takes 100 clock cycles. Because of fair bus arbitration, the releasing processor must wait for *all* other 9 processors to try to get the lock in vain!

n

The difficulty in this Example arises from contention for the lock and serialization of lock access, as well as the latency of the bus access. (The fairness property of the bus actually makes things worse, since it delays the processor that claims the lock from releasing it; unfortunately, for any bus arbitration scheme some worst-case scenario does exist.) The key advantages of spin locks, namely that they have low overhead in terms of bus or network cycles and offer good performance when locks are reused by the same processor, are both lost in this example. We will consider alternative implementations in the next section, but before we do that, let's consider the use of spin locks to implement another common high-level synchronization primitive.

Barrier Synchronization

One additional common synchronization operation in programs with parallel loops is a *barrier*. A barrier forces all processes to wait until all the processes reach the barrier and then releases all of the processes. A typical implementation of a barrier can be done with two spin locks: one used to protect a counter that tallies the processes arriving at the barrier and one used to hold the processes until the last process arrives at the barrier. To implement a barrier we usually use the ability to spin on a variable until it satisfies a test; we use the notation `spin(condition)` to indicate this. Figure 6.40 is a typical implementation, assuming that lock and unlock provide basic spin locks and `total` is the number of processes that must reach the barrier.

In practice, another complication makes barrier implementation slightly more complex. Frequently a barrier is used within a loop, so that processes released from the barrier would do some work and then reach the barrier again. Assume that one of the processes never actually leaves the barrier (it stays at the spin op-

```

lock (counterlock);/* ensure update atomic */
if (count==0) release=0;/*first=>reset release */
count = count +1;/* count arrivals */
unlock(counterlock);/* release lock */
if (count==total) {/* all arrived */
    count=0;/* reset counter */
    release=1;/* release processes */
}
else {/* more to come */

    spin (release==1);/* wait for arrivals */
}

```

FIGURE 6.39 Code for a simple barrier. The lock `counterlock` protects the counter so that it can be atomically incremented. The variable `count` keeps the tally of how many processes have reached the barrier. The variable `release` is used to hold the processes until the last one reaches the barrier. The operation `spin (release==1)` causes a process to wait until all processes reach the barrier.

eration), which could happen if the OS scheduled another process, for example. Now it is possible that one process races ahead and gets to the barrier again before the last process has left. The “fast” process then traps the remaining “slow” process in the barrier by resetting the flag `release`. Now all the processes will wait infinitely at the next instance of this barrier, because one process is trapped at the last instance, and the number of processes can never reach the value of `total`.

The important observation in this example is that the programmer did nothing wrong. Instead, the implementer of the barrier made some assumptions about forward progress that cannot be assumed. One obvious solution to this is to count the processes as they exit the barrier (just as we did on entry) and not to allow any process to reenter and reinitialize the barrier until all processes have left the prior instance of this barrier. This extra step would significantly increase the latency of the barrier and the contention, which as we will see shortly are already large. An alternative solution is a *sense-reversing barrier*, which makes use of a private per-process variable, `local_sense`, which is initialized to 1 for each process. Figure 6.40 shows the code for the sense-reversing barrier. This version of a barrier is

safely usable; as the next example shows, however, its performance can still be quite poor.

```

local_sense =! local_sense; /*toggle local_sense*/
lock (counterlock);/* ensure update atomic */
count=count+1;/* count arrivals */
unlock (counterlock);/* unlock */
if (count==total) {/* all arrived */
    count=0;/* reset counter */
    release=local_sense;/* release processes */
}
else {/* more to come */
    spin (release==local_sense);/*wait for signal*/
}

```

FIGURE 6.40 Code for a sense-reversing barrier. The key to making the barrier reusable is the use of an alternating pattern of values for the flag release, which controls the exit from the barrier. If a process races ahead to the next instance of this barrier while some other processes are still in the barrier, the fast process cannot trap the other processes, since it does not reset the value of release as it did in Figure 6.40.

EXAMPLE Suppose there are 10 processors on a bus that each try to execute a barrier simultaneously. Assume that each bus transaction is 100 clock cycles, as before. You can ignore the time of the actual read or write of a lock held in the cache as the time to execute other nonsynchronization operations in the barrier implementation. Determine the number of bus transactions required for all 10 processors to reach the barrier, be released from the barrier, and exit the barrier. Assume that the bus is totally fair, so that every pending request is serviced before a new request and that the processors are equally fast. Don't worry about counting the processors out of the barrier. How long will the entire process take?

ANSWER The following table shows the sequence of events for one processor to traverse the barrier, assuming that the first process to grab the bus does not have the lock.

Event	Duration in clocks for one processor	Duration in clocks for 10 processors
Time for each processor to grab lock, increment, release lock	3100	31,000
Time to execute release	100	100
Time for each processor to get the release flag	100	1000
Total	3300	31,100

Our barrier operation takes about as long as the 10-processor lock-unlock sequence we considered earlier. The total number of bus transactions is about 220.

n

As we can see from these examples, synchronization performance can be a real bottleneck when there is substantial contention among multiple processes. When there is little contention and synchronization operations are infrequent, we are primarily concerned about the latency of a synchronization primitive—that is, how long it takes an individual process to complete a synchronization operation. Our basic spin-lock operation can do this in two bus cycles: one to initially read the lock and one to write it. We could improve this to a single bus cycle by a variety of methods. For example, we could simply spin on the swap operation. If the lock were almost always free, this could be better, but if the lock were not free, it would lead to lots of bus traffic, since each attempt to lock the variable would lead to a bus cycle. In practice, the latency of our spin lock is not quite as bad as we have seen in this example, since the write miss for a data item present in the cache is treated as an upgrade and will be cheaper than a true read miss.

The more serious problem in these examples is the serialization of each process's attempt to complete the synchronization. This serialization is a problem when there is contention, because it greatly increases the time to complete the synchronization operation. For example, if the time to complete all 10 lock and unlock operations depended only on the latency in the uncontended case, then it would take 1000 rather than 15,000 cycles to complete the synchronization operations. The barrier situation is as bad, and in some ways worse, since it is highly likely to incur contention. The use of a bus interconnect exacerbates these problems, but serialization could be just as serious in a directory-based multiprocessor, where the latency would be large. The next section presents some solutions that are useful when either the contention is high or the processor count is large.

Synchronization Mechanisms for Larger-Scale Multiprocessors

What we would like are synchronization mechanisms that have low latency in uncontended cases and that minimize serialization in the case where contention is significant. We begin by showing how software implementations can improve the performance of locks and barriers when contention is high; we then explore two basic hardware primitives that reduce serialization while keeping latency low.

Software Implementations

The major difficulty with our spin-lock implementation is the delay due to contention when many processes are spinning on the lock. One solution is to artificially delay processes when they fail to acquire the lock. The best performance is obtained by increasing the delay exponentially whenever the attempt to acquire

the lock fails. Figure 6.41 shows how a spin lock with *exponential back-off* is implemented. Exponential back-off is a common technique for reducing contention in shared resources, including access to shared networks and buses (see section 7.7). This implementation still attempts to preserve low latency when contention is small by not delaying the initial spin loop. The result is that if many processes are waiting, the back-off does not affect the processes on their first attempt to acquire the lock. We could also delay that process, but the result would be poorer performance when the lock was in use by only two processes and the first one happened to find it locked.

```

                                ADDUI   R3,R0,#1      ;R3 = initial delay
lockit:  LL      R2,0(R1)        ;load linked
                                BNEZ    R2,lockit     ;not available-spin
                                DADDUI  R2,R2,#1      ;get locked value
                                SC      R2,0(R1)      ;store conditional
                                BNEZ    R2,gotit      ;branch if store succeeds
                                DSL     R3,R3,#1      ;increase delay by factor of 2
                                PAUSE   R3           ;delays by value in R3
                                J      lockit
gotit:   use data protected by lock

```

FIGURE 6.41 A spin lock with exponential back-off. When the store conditional fails, the process delays itself by the value in R3. The delay can be implemented by decrementing R3 until it reaches 0. The exact timing of the delay is multiprocessor dependent, although it should start with a value that is approximately the time to perform the critical section and release the lock. The statement `pause R3` should cause a delay of R3 of these time units. The value in R3 is increased by a factor of 2 every time the store conditional fails, which causes the process to wait twice as long before trying to acquire the lock again. The small variations in the rate at which competing processors execute instructions are usually sufficient to ensure that processes will not continually collide. If the natural perturbation in execution time was insufficient, R3 could be initialized with a small random value, increasing the variance in the successive delays and reducing the probability of successive collisions.

Another technique for implementing locks is to use queuing locks. Queuing locks work by constructing a queue of waiting processors; whenever a processor frees up the lock, it causes the next processor in the queue to attempt access. This eliminates contention for a lock when it is freed. We show how queuing locks operate in the next section using a hardware implementation, but software implementations using arrays can achieve most of the same benefits (see Exercise 6.25). Before we look at hardware primitives, let's look at a better mechanism for barriers.

Our barrier implementation suffers from contention both during the *gather* stage, when we must atomically update the count, and at the *release* stage, when

all the processes must read the release flag. The former is more serious because it requires exclusive access to the synchronization variable and thus creates much more serialization; in comparison, the latter generates only read contention. We can reduce the contention by using a *combining tree*, a structure where multiple requests are locally combined in tree fashion. The same combining tree can be used to implement the release process, reducing the contention there; we leave the last step for the Exercises.

Our combining tree barrier uses a predetermined n -ary tree structure. We use the variable k to stand for the fan-in; in practice $k = 4$ seems to work well. When the k th process arrives at a node in the tree, we signal the next level in the tree. When a process arrives at the root, we release all waiting processes. As in our earlier example, we use a sense-reversing technique. A tree-based barrier, as shown in Figure 6.42, uses a tree to combine the processes and a single signal to release the barrier; Exercises 6.23 and 6.24 ask you to analyze the time for the combining barrier versus the noncombining version. Some MPPs (e.g., the T3D and CM-5) have also included hardware support for barriers, but more recent machines have relied on software libraries for this support.

Hardware Primitives

In this section we look at two hardware synchronization primitives. The first primitive deals with locks, while the second is useful for barriers and a number of other user-level operations that require counting or supplying distinct indices. In both cases we can create a hardware primitive where latency is essentially identical to our earlier version, but with much less serialization, leading to better scaling when there is contention.

The major problem with our original lock implementation is that it introduces a large amount of unneeded contention. For example, when the lock is released all processors generate both a read and a write miss, although at most one processor can successfully get the lock in the unlocked state. This sequence happens on each of the 20 lock/unlock sequences, as we saw in the example on page 710.

We can improve this situation by explicitly handing the lock from one waiting processor to the next. Rather than simply allowing all processors to compete every time the lock is released, we keep a list of the waiting processors and hand the lock to one explicitly, when its turn comes. This sort of mechanism has been called a *queuing lock*. Queuing locks can be implemented either in hardware, which we describe here, or in software using an array to keep track of the waiting processes. The basic concepts are the same in either case. Our hardware implementation assumes a directory-based multiprocessor where the individual processor caches are addressable. In a bus-based multiprocessor, a software implementation would be more appropriate and would have each processor using

```

struct node{/* a node in the combining tree */
    int counterlock; /* lock for this node */
    int count; /* counter for this node */
    int parent; /* parent in the tree = 0..P-1 except for root */
};
struct node tree [0..P-1]; /* the tree of nodes */
int local_sense; /* private per processor */
int release; /* global release flag */

/* function to implement barrier */
barrier (int mynode) {
    lock (tree[mynode].counterlock); /* protect count */
    tree[mynode].count=tree[mynode].count+1;
    /* increment count */
    unlock (tree[mynode].counterlock); /* unlock */
    if (tree[mynode].count==k) { /* all arrived at mynode */
        if (tree[mynode].parent >=0) {
            barrier(tree[mynode].parent);
        } else{
            release = local_sense;
        };
        tree[mynode].count = 0; /* reset for the next time */
    } else{
        spin (release==local_sense); /* wait */
    };
};
/* code executed by a processor to join barrier */
local_sense =! local_sense;
barrier (mynode);

```

FIGURE 6.42 An implementation of a tree-based barrier reduces contention considerably. The tree is assumed to be prebuilt statically using the nodes in the array `tree`. Each node in the tree combines k processes and provides a separate counter and lock, so that at most k processes contend at each node. When the k th process reaches a node in the tree it goes up to the parent, incrementing the count at the parent. When the count in the parent node reaches k , the release flag is set. The count in each node is reset by the last process to arrive. Sense-reversing is used to avoid races as in the simple barrier.

a different address for the lock, permitting the explicit transfer of the lock from one process to another.

How does a queuing lock work? On the first miss to the lock variable, the miss is sent to a synchronization controller, which may be integrated with the memory controller (in a bus-based system) or with the directory controller. If the lock is free, it is simply returned to the processor. If the lock is unavailable, the controller creates a record of the node's request (such as a bit in a vector) and sends the

processor back a locked value for the variable, which the processor then spins on. When the lock is freed, the controller selects a processor to go ahead from the list of waiting processors. It can then either update the lock variable in the selected processor's cache or invalidate the copy, causing the processor to miss and fetch an available copy of the lock.

EXAMPLE How many bus transaction and how long does it take to have 10 processors lock and unlock the variable using a queuing lock that updates the lock on a miss? Make the other assumptions about the system the same as those in the earlier example on page 710.

ANSWER Each processor misses once on the lock initially and once to free the lock, so it takes only 20 bus cycles. The first 10 initial misses take 1000 cycles, followed by a 100-cycle delay for each of the 10 releases. This sequence yields a total of 2100 cycles—significantly better than the case with conventional coherence-based spin locks. n

There are a couple of key insights in implementing such a queuing lock capability. First, we need to be able to distinguish the initial access to the lock, so we can perform the queuing operation, and also the lock release, so we can provide the lock to another processor. The queue of waiting processes can be implemented by a variety of mechanisms. In a directory-based multiprocessor, this queue is akin to the sharing set, and similar hardware can be used to implement the directory and queuing lock operations. One complication is that the hardware must be prepared to reclaim such locks, since the process that requested the lock may have been context-switched and may not even be scheduled again on the same processor.

Queuing locks can be used to improve the performance of our barrier operation (see Exercise 6.17). Alternatively, we can introduce a primitive that reduces the amount of time needed to increment the barrier count, thus reducing the serialization at this bottleneck, which should yield comparable performance to using queuing locks. One primitive that has been introduced for this and for building other synchronization operations is *fetch-and-increment*, which atomically fetches a variable and increments its value. The returned value can be either the incremented value or the fetched value. Using fetch-and-increment we can dramatically improve our barrier implementation, compared to the simple code-sensing barrier.

EXAMPLE Write the code for the barrier using fetch-and-increment. Making the same assumptions as in our earlier example and also assuming that a

fetch-and-increment operation takes 100 clock cycles, determine the time for 10 processors to traverse the barrier. How many bus cycles are required?

ANSWER

Figure 6.40 shows the code for the barrier. This implementation requires 10 fetch-and-increment operations and 10 cache misses for the release operation for a total time of 2000 cycles and 20 bus/interconnect operations versus an earlier implementation that took over 15 times longer and 10 times more bus operations to complete the barrier. Of course, fetch-and-increment can also be used in implementing the combining tree barrier, reducing the serialization at each node in the tree.

```
local_sense =! local_sense; /*toggle local_sense*/
fetch_and_increment(count);/* atomic update*/
if (count==total) {/* all arrived */
    count=0;/* reset counter */
    release=local_sense;/* release processes */
}
else {/* more to come */
    spin (release==local_sense);/*wait for signal*/
}
```

FIGURE 6.43 Code for a sense-reversing barrier using fetch-and-increment to do the counting.

n

As we have seen, synchronization problems can become quite acute in larger-scale multiprocessors. When the challenges posed by synchronization are combined with the challenges posed by long memory latency and potential load imbalance in computations, we can see why getting efficient usage of large-scale parallel processors is very challenging.

6.8 Models of Memory Consistency: An Introduction

Cache coherence ensures that multiple processors see a consistent view of memory. It does not answer the question of *how* consistent the view of memory must be. By “how consistent” we mean, when must a processor see a value that has been updated by another processor? Since processors communicate through

shared variables (used both for data values and for synchronization), the question boils down to this: In what order must a processor observe the data writes of another processor? Since the only way to “observe the writes of another processor” is through reads, the question becomes, what properties must be enforced among reads and writes to different locations by different processors?

Although the question of how consistent memory be seems simple, it is remarkably complicated, as we can see with a simple example. Here are two code segments from processes P1 and P2, shown side by side:

P1:	A = 0;	P2:	B = 0;

	A = 1;		B = 1;
L1:	if (B == 0) ...	L2:	if (A == 0) ...

Assume that the processes are running on different processors, and that locations A and B are originally cached by both processors with the initial value of 0. If writes always take immediate effect and are immediately seen by other processors, it will be impossible for *both* if-statements (labeled L1 and L2) to evaluate their conditions as true, since reaching the if-statement means that either A or B must have been assigned the value 1. But suppose the write invalidate is delayed, and the processor is allowed to continue during this delay; then it is possible that both P1 and P2 have not seen the invalidations for B and A (respectively) *before* they attempt to read the values. The question is, should this behavior be allowed, and if so, under what conditions?

The most straightforward model for memory consistency is called *sequential consistency*. Sequential consistency requires that the result of any execution be the same as if the memory accesses executed by each processor were kept in order and the accesses among different processors were arbitrarily interleaved. Sequential consistency eliminates the possibility of some nonobvious execution in the previous example, because the assignments must be completed before the if statements are initiated.

The simplest way to implement sequential consistency is to require a processor to delay the completion of any memory access until all the invalidations caused by that access are completed. Of course, it is equally effective to delay the next memory access until the previous one is completed. Remember that memory consistency involves operations among different variables: the two accesses that must be ordered are actually to different memory locations. In our example, we must delay the read of A or B ($A=0$ or $B=0$) until the previous write has completed ($B=1$ or $A=1$). Under sequential consistency, we cannot, for example, simply place the write in a write buffer and continue with the read. Although sequential consistency presents a simple programming paradigm, it reduces potential

performance, especially in a multiprocessor with a large number of processors or long interconnect delays, as we can see in the following Example.

EXAMPLE Suppose we have a processor where a write miss takes 40 cycles to establish ownership, 10 cycles to issue each invalidate after ownership is established, and 50 cycles for an invalidate to complete and be acknowledged once it is issued. Assuming that four other processors share a cache block, how long does a write miss stall the writing processor if the processor is sequentially consistent? Assume that the invalidates must be explicitly acknowledged before the directory controller knows they are completed. Suppose we could continue executing after obtaining ownership for the write miss without waiting for the invalidates; how long would the write take?

ANSWER When we wait for invalidates, each write takes the sum of the ownership time plus the time to complete the invalidates. Since the invalidates can overlap, we need only worry about the last one, which starts $10 + 10 + 10 + 10 = 40$ cycles after ownership is established. Hence the total time for the write is $40 + 40 + 50 = 130$ cycles. In comparison, the ownership time is only 40 cycles. With appropriate write-buffer implementations it is even possible to continue before ownership is established. n

To provide better performance, researchers and architects have explored two different routes. First, they developed ambitious implementations that preserve sequential consistency but use latency hiding techniques to reduce the penalty; we discuss these in the section on cross-cutting issues (see page 731). Second, they developed less restrictive memory consistency models that allow for faster hardware. Such models can affect how the programmer sees the multiprocessor, so before we discuss these less restrictive models, let's look at what the programmer expects.

The Programmer's View

Although the sequential consistency model has a performance disadvantage, from the viewpoint of the programmer it has the advantage of simplicity. The challenge is to develop a programming model that is simple to explain and yet allows a high performance implementation.

One such programming model that allows us to have a more efficient implementation is to assume that programs are *synchronized*. A program is synchronized if all access to shared data is ordered by synchronization operations. A data reference is ordered by a synchronization operation if, in every possible execution, a write of a variable by one processor and an access (either a read or a write)

of that variable by another processor are separated by a pair of synchronization operations, one executed after the write by the writing processor and one executed before the access by the second processor. Cases where variables may be updated without ordering by synchronization are called *data races*, because the execution outcome depends on the relative speed of the processors, and like races in hardware design, the outcome is unpredictable, which leads to another name for synchronized programs: *data-race-free*.

As a simple example, consider a variable being read and updated by two different processors. Each processor surrounds the read and update with a lock and an unlock, both to ensure mutual exclusion for the update and to ensure that the read is consistent. Clearly, every write is now separated from a read by the other processor by a pair of synchronization operations: one unlock (after the write) and one lock (before the read). Of course, if two processors are writing a variable with no intervening reads, then the writes must also be separated by synchronization operations.

It is a broadly accepted observation that most programs are synchronized. This observation is true primarily because if the accesses were unsynchronized, the behavior of the program would be quite difficult to determine because the speed of execution would determine which processor won a data race and thus affect the results of the program. Even with sequential consistency, reasoning about such programs is very difficult. Programmers could attempt to guarantee ordering by constructing their own synchronization mechanisms, but this is extremely tricky, can lead to buggy programs, and may not be supported architecturally, meaning that they may not work in future generations of the multiprocessor. Instead, almost all programmers will choose to use synchronization libraries that are correct and optimized for the multiprocessor and the type of synchronization. Finally, the use of standard synchronization primitives ensures that even if the architecture implements a more relaxed consistency model than sequential consistency, a synchronized program will behave as if the hardware implemented sequential consistency.

Relaxed Consistency Models: The Basics

The key idea in relaxed consistency models is to allow reads and writes to complete out of order, but to use synchronization operations to enforce ordering, so that a synchronized program behaves as if the processor were sequentially consistent. There are a variety of relaxed models that are classified according to what orderings they relax. The three major sets of orderings that are relaxed are:

1. The $W \rightarrow R$ ordering: which yields a model known as total store ordering or processor consistency. Because this ordering retains ordering among writes, many programs that operate under sequential consistency operate under this model, without additional synchronization.

2. The $W \rightarrow W$ ordering: which yields a model known as partial store order.
3. The $R \rightarrow W$ and $R \rightarrow R$ orderings: which yields a variety of models including weak ordering, the Alpha consistency model, the PowerPC consistency model, and release consistency depending on the details of the ordering restrictions and how synchronization operations enforce ordering.

By relaxing these orderings, the processor can possibly obtain significant performance advantages. There are, however, many complexities in describing relaxed consistency models, including the advantages and complexities of relaxing different orders, defining precisely what it means for a write to complete, and deciding when processors can see values that the processor itself has written. These complexities, as well as an assessment of the performance of relaxed model and a discussion of the implementation issues, are described in more detail in Appendix F.

Final Remarks on Consistency Models

At the present time, many multiprocessors being built support some sort of relaxed consistency model, varying from processor consistency to release consistency. Since synchronization is highly multiprocessor specific and error prone, the expectation is that most programmers will use standard synchronization libraries and will write synchronized programs, making the choice of a weak consistency model invisible to the programmer and yielding higher performance.

An alternative viewpoint, which we discuss more extensively in the next section (specifically on page 731), argues that with speculation much of the performance advantage of relaxed consistency models can be obtained with sequential or processor consistency.

A key part of this argument in favor of relaxed consistency revolves the role of the compiler and its ability to optimize memory access to potentially shared variables. This topic is also discussed on page 731.

6.9 Multithreading: Exploiting Thread-Level Parallelism within a Processor

Multithreading allows multiple threads to share the functional units of a single processor in an overlapping fashion. To permit this sharing, the processor must duplicate the independent state of each thread. For example, a separate copy of the register file, a separate PC, and a separate page table are required for each thread. The memory itself can be shared through the virtual memory mechanisms, which already support multiprogramming. In addition, the hardware must

support the ability to change to a different thread relatively quickly; in particular, a thread switch should be much more efficient than a process switch, which typically requires hundreds to thousands of processor cycles.

There are two main approaches to multithreading. *Fine-grained multithreading* switches between threads on each instruction, causing the execution of multiples threads to be interleaved. This interleaving is often done in a round-robin fashion, skipping any threads that are stalled at that time. To make fine-grained multithreading practical, the CPU must be able to switch threads on every clock cycle. One key advantage of fine-grained multithreading is that it can hide the throughput losses that arise from both short and long stalls, since instructions from other threads can be executed when one thread stalls. The primary disadvantage of fine-grained multithreading is that it slows down the execution of the individual threads, since a thread that is ready to execute without stalls will be delayed by instructions from other threads.

Coarse-grained multithreading was invented as an alternative to fine-grained multithreading. Coarse-grained multithreading switches threads only on costly stalls, such as level two cache misses. This change relieves the need to have thread-switching be essentially free and is much less likely to slow the processor down, since instructions from other threads will only be issued, when a thread encounters a costly stall. Coarse-grained multithreading suffers, however, from a major drawback: it is limited in its ability to overcome throughput losses, especially from shorter stalls. This limitation arises from the pipeline start-up costs of coarse-grain multithreading. Because a CPU with coarse-grained multithreading issues instructions from a single thread, when a stall occurs, the pipeline must be emptied or frozen. The new thread that begins executing after the stall must fill the pipeline before instructions will be able to complete. Because of this start-up overhead, coarse-grained multithreading is much more useful for reducing the penalty of high cost stalls, where pipeline refill is negligible compared to the stall time.

The next section explores a variation on fine-grained multithreading that enables a superscalar processor to exploit ILP and multithreading in an integrated and efficient fashion. Section 6.12 examines a commercial processor using coarse-grained multithreading.

Simultaneous Multithreading: Converting Thread-Level Parallelism into Instruction-Level Parallelism

Simultaneous multithreading (SMT) is a variation on multithreading that uses the resources of a multiple-issue, dynamically-scheduled processor to exploit TLP at the same time it exploits ILP. The key insight that motivates SMT is that modern multiple-issue processors often have more functional unit parallelism available than a single thread can effectively use. Furthermore, with register renaming and dynamic scheduling, multiple instructions from independent threads can be is-

sued without regard to the dependences among them; the resolution of the dependences can be handled by the dynamic scheduling capability.

Figure 6.44 conceptually illustrates the differences in a processor's ability to exploit the resources of a superscalar for the following processor configurations:

- n a superscalar with no multithreading support,
- n a superscalar with coarse-grained multithreading,
- n a superscalar with fine-grained multithreading, and
- n a superscalar with simultaneous multithreading.

In the superscalar without multithreading support, the use of issue slots is limited by a lack of ILP, a topic we discussed extensively in Chapter 3. In addition, a major stall, such as an instruction cache miss, can leave the entire processor idle.

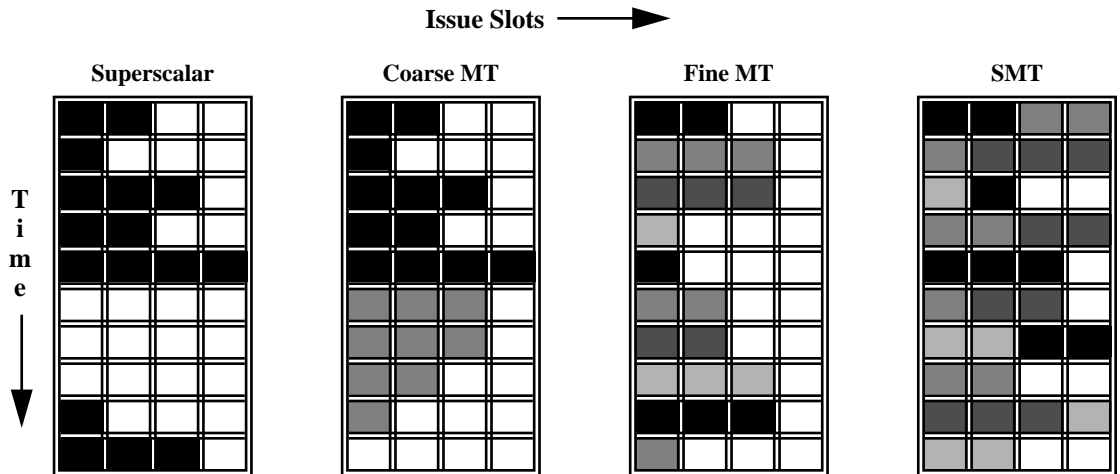


FIGURE 6.44 This illustration shows how these four different approaches use the issue slots of a superscalar processor. The horizontal dimension represents the instruction issue capability in each clock cycle. The vertical dimension represents a sequence of clock cycles. An empty (white) box indicates that the corresponding issue slot is unused in that clock cycle. The shades of grey and black correspond to four different threads in the multithreading processors. Black is also used to indicate the occupied issue slots in the case of the superscalar without multithreading support.

In the coarse-grained multithreaded superscalar, the long stalls are partially hidden by switching to another thread that uses the resources of the processor. Although this reduces the number of completely idle clock cycles, within each clock cycle, the ILP limitations still lead to idle cycles. Furthermore, in a coarse-grained multithreaded processor, since thread switching only occurs when there is a stall and the new thread has a start-up period, there are likely to be some fully idle cycles remaining.

In the fine-grained case, the interleaving of threads eliminates fully empty slots. Because only one thread issues instructions in a given clock cycle, however, ILP limitations still lead to a significant number of idle slots within individual clock cycles.

In the SMT case, thread-level parallelism (TLP) and instruction-level parallelism (ILP) are exploited simultaneously; with multiple threads using the issue slots in a single clock cycle. Ideally, the issue slot usage is limited by imbalances in the resource needs and resource availability over multiple threads. In practice, other factors—including how many active threads are considered, finite limitations on buffers, the ability to fetch enough instructions from multiple threads, and practical limitations of what instruction combinations can issue from one thread and from multiple threads—can also restrict how many slots are used. Although Figure 6.44 greatly simplifies the real operation of these processors it does illustrate the potential performance advantages of multithreading in general and SMT in particular.

As mentioned above, simultaneous multithreading uses the insight that a dynamically scheduled processor already has many of the hardware mechanisms needed to support the integrated exploitation of TLP through multithreading. In particular, dynamically scheduled superscalars have a large set of virtual registers that can be used to hold the register sets of independent threads (assuming separate renaming tables are kept for each thread). Because register renaming provides unique register identifiers, instructions from multiple threads can be mixed in the datapath without confusing sources and destinations across the threads. This observation leads to the insight that multithreading can be built on top of an out-of-order processor by adding a per thread renaming table, keeping separate PCs, and providing the capability for instructions from multiple threads to commit. There are complications in handling instruction commit, since we would like instructions from independent threads to be able to commit independently. The independent commitment of instructions from separate threads can be supported by logically keeping a separate reorder buffer for each thread.

Design Challenges in SMT processors

Because a dynamically scheduled superscalar processor is likely to have a deep pipeline, SMT will be unlikely to gain much in performance if it were coarse-grained. Since SMT will likely make sense only in a fine-grained implementation, we must worry about the impact of fine-grained scheduling on single thread performance. This effect can be minimized by having a preferred thread, which still permits multithreading to preserve some of its performance advantage with a smaller compromise in single thread performance. At first glance, it might appear that a preferred thread approach sacrifices neither throughput nor single-thread performance. Unfortunately, with a preferred thread, the processor is likely to sacrifice some throughput, when the preferred thread encounters a stall. The rea-

son is that the pipeline is less likely to have a mix of instructions from several threads, resulting in greater probability that either empty slots or a stall will occur. Throughput is maximized by having a sufficient number of independent threads to hide all stalls in any combination of threads.

Unfortunately, mixing many threads will inevitably compromise the execution time of individual threads. Similar problems exist in instruction fetch. To maximize single thread performance, we should fetch as far ahead as possible in that single thread and always have the fetch unit free when a branch is mispredicted and a miss occurs in the prefetch buffer. Unfortunately, this limits the number of instructions available for scheduling from other threads, reducing throughput. All multithreaded processor must seek to balance this tradeoff.

In practice, the problems of dividing resources and balancing single-thread and multiple-thread performance turn out not to be as challenging as they sound, at least for current superscalar back-ends. In particular, for current machines that issue four to eight instructions per cycle, it probably suffices to have a small number of active threads, and an even smaller number of “preferred” threads. Whenever possible, the processor acts on behalf of a preferred thread. This starts with prefetching instructions: whenever the prefetch buffers for the preferred threads are not full, instructions are fetched for those threads. Only when the preferred thread buffers are full is the instruction unit directed to prefetch for other threads. Note that having two preferred threads means that we are simultaneously prefetching for two instruction streams and this adds complexity to the instruction fetch unit and the instruction cache. Similarly, the instruction issue unit can direct its attention to the preferred threads, considering other threads only if the preferred threads are stalled and cannot issue. In practice, having four to eight threads and two to four preferred threads is likely to completely utilize the capability of a superscalar back-end that is roughly double the capability of those available in 2001.

There are a variety of other design challenges for an SMT processor, including:

- ▮ dealing with a larger register file needed to hold multiple contexts,
- ▮ maintaining low overhead on the clock cycle, particularly in critical steps such as instruction issue, where more candidate instructions need to be considered, and in instruction completion, where choosing what instructions to commit may be challenging, and
- ▮ ensuring that the cache conflicts generated by the simultaneous execution of multiple threads do not cause significant performance degradation.

In viewing these problems, two observations are important. In many cases, the potential performance overhead due to multithreading is small, and simple choices work well enough. Second, the efficiency of current superscalars is low enough that there is room for significant improvement, even at the cost of some overhead.

SMT appears to be the most promising way to achieve that improvement in throughput.

Because SMT exploits thread-level parallelism on a multiple-issue superscalar, it is most likely to be included in high-end processors targeted at server markets. In addition, it is likely that there will be some mode to restrict the multithreading, so as to maximize the performance of a single thread.

Prior to deciding to abandon the Alpha architecture in mid 2001, Compaq had announced that the Alpha 21364 would have SMT capability when it became available in 2002. In July 2001, Intel announced that a future processor based on the Pentium 4 microarchitecture and targeted at the server market, most likely Pentium 4 Xenon, would support SMT, initially with two-thread implementation. Intel claims a 30% improvement in throughput for server applications with this new support.

6.10 Crosscutting Issues

Because multiprocessors redefine many system characteristics (e.g., performance assessment, memory latency, and the importance of scalability), they introduce interesting design problems that cut across the spectrum, affecting both hardware and software. In this section we give several examples including: measuring and reporting the performance of multiprocessors, enhancing latency tolerance in memory systems, and a method for using virtual memory support to implement shared memory.

Memory System Issues

As we have seen in this chapter, memory system issues are at the core of the design of shared-memory multiprocessors. Indeed, multiprocessing introduces many new memory system complications that do not exist in uniprocessors. In this section we look at two implementation issues that have a significant impact on the design and implementation of a memory system in a multiprocessor context.

Inclusion and Its Implementation

Many multiprocessors use multilevel cache hierarchies to reduce both the demand on the global interconnect and the latency of cache misses. If the cache also provides *multilevel inclusion*—every level of cache hierarchy is a subset of the level further away from the processor—then we can use the multilevel structure to reduce the contention between coherence traffic and processor traffic, as explained earlier. Thus most multiprocessors with multilevel caches enforce the

inclusion property. This restriction is also called the *subset property*, because each cache is a subset of the cache below it in the hierarchy.

At first glance, preserving the multilevel inclusion property seems trivial. Consider a two-level example: any miss in L1 either hits in L2 or generates a miss in L2, causing it to be brought into both L1 and L2. Likewise, any invalidate that hits in L2 must be sent to L1, where it will cause the block to be invalidated, if it exists.

The catch is what happens when the block size of L1 and L2 are different. Choosing different block sizes is quite reasonable, since L2 will be much larger and have a much longer latency component in its miss penalty, and thus will want to use a larger block size. What happens to our “automatic” enforcement of inclusion when the block sizes differ? A block in L2 represents multiple blocks in L1, and a miss in L2 causes the replacement of data that is equivalent to multiple L1 blocks. For example, if the block size of L2 is four times that of L1, then a miss in L2 will replace the equivalent of four L1 blocks. Let’s consider a detailed example.

EXAMPLE Assume that L2 has a block size four times that of L1. Show how a miss for an address that causes a replacement in L1 and L2 can lead to violation of the inclusion property.

ANSWER Assume that L1 and L2 are direct mapped and that the block size of L1 is b bytes and the block size of L2 is $4b$ bytes. Suppose L1 contains two blocks with starting addresses x and $x + b$ and that $x \bmod 4b = 0$, meaning that x also is the starting address of a block in L2, then that single block in L2 contains the L1 blocks x , $x + b$, $x + 2b$, and $x + 3b$. Suppose the processor generates a reference to block y that maps to the block containing x in both caches and hence misses. Since L2 missed, it fetches $4b$ bytes and replaces the block containing x , $x + b$, $x + 2b$, and $x + 3b$, while L1 takes b bytes and replaces the block containing x . Since L1 still contains $x + b$, but L2 does not, the inclusion property no longer holds. □

To maintain inclusion with multiple block sizes, we must probe the higher levels of the hierarchy when a replacement is done at the lower level to ensure that any words replaced in the lower level are invalidated in the higher-level caches. Most systems chose this solution rather than the alternative of not relying on inclusion and snooping the higher-level caches. In the Exercises we explore inclusion further and show that similar problems exist if the associativity of the levels is different. Baer and Wang [1988] describe the advantages and challenges of inclusion in detail.

Nonblocking Caches and Latency Hiding

We saw the idea of nonblocking or lockup-free caches in Chapter 5, where the concept was used to reduce cache misses by overlapping them with execution and by pipelining misses. There are additional benefits in the multiprocessor case. The first is that the miss penalties are likely to be larger, meaning there is more latency to hide, and the opportunity for pipelining misses is also probably larger, since the memory and interconnect system can often handle multiple outstanding memory references also.

Second, a multiprocessor needs nonblocking caches to take advantage of weak consistency models. For example, to implement a model like processor consistency requires that writes be nonblocking with respect to reads so that a processor can continue either immediately, by buffering the write, or as soon as it establishes ownership of the block and updates the cache. Relaxed consistency models allow further reordering of misses, but nonblocking caches are needed to take full advantage of this flexibility.

Finally, nonblocking support is critical to implementing prefetching. Prefetching, which we also discussed in Chapter 5, is even more important in multiprocessors than in uniprocessors, again due to longer memory latencies. In Chapter 5 we described why it is important that prefetches not affect the semantics of the program, since this allows them to be inserted anywhere in the program without changing the results of the computation.

In a multiprocessor, maintaining the absence of any semantic impact from the use of prefetches requires that prefetched data be kept coherent. A prefetched value is kept coherent if, when the value is actually accessed by a load instruction, the most recently written value is returned, even if that value was written after the prefetch. This result is exactly the property that cache coherence gives us for other variables in memory. A prefetch that brings a data value closer, and guarantees that on the actual memory access to the data (a load of the prefetched value) the most recent value of the data item is obtained, is called *nonbinding*, since the data value is not bound to a local copy, which would be incoherent. By contrast, a prefetch that moves a data value into a general-purpose register is binding, since the register value is a new variable, as opposed to a cache block, which is a coherent copy of a variable. A nonbinding prefetch maintains the coherence properties of any other value in memory, while a binding prefetch appears more like a register load, since it removes the data from the coherent address space.

Why is nonbinding prefetch critical? Consider a simple but typical example: a data value written by one processor and used by another. In this case, the consumer would like to prefetch the value as early as possible; but suppose the producing process is delayed for some reason. Then the prefetch may fetch the old value of the data item. If the prefetch is nonbinding, the copy of the old data is invalidated when the value is written, maintaining coherence. If the prefetch is binding, however, then the old, incoherent value of the data is used by the prefetching process. Because of the long memory latencies, a prefetch may need to be placed a hundred or more instructions earlier than the data use, if we aim to

hide the entire latency. This requirement makes the nonbinding property vital to ensure coherent usage of the prefetch in multiprocessors.

Implementing prefetch requires the same sort of support that a lockup-free cache needs, since there are multiple outstanding memory accesses. This requirement causes several complications:

1. A local node will need to keep track of the multiple outstanding accesses, since the replies may return in a different order than they were sent. This accounting can be handled by adding tags to the requests, or by incorporating the address of the memory block in the reply.
2. Before issuing a request (either a normal fetch or a prefetch), the node must ensure that it has not already issued a request for the same block, since two write requests for the same block could lead to incorrect operation of the protocol. For example, if the node issues a write prefetch to a block, while it has a write miss or write prefetch outstanding, both our snooping protocol and directory protocol can fail to operate properly.
3. Our implementation of the directory and snooping controllers assumes that the processor stalls on a miss. Stalling allows the cache controller to simply wait for a reply when it has generated a request. With a nonblocking cache or with prefetching, a processor can generate additional requests while it is waiting for replies. This complicates the directory and snooping controllers; Appendix E shows how these issues can be addressed.

Compiler Optimization and the Consistency Model

Another reason for defining a model for memory consistency is to specify the range of legal compiler optimizations that can be performed on shared data. In explicitly parallel programs, unless the synchronization points are clearly defined and the programs are synchronized, the compiler could not interchange a read and a write of two different shared data items, because such transformations might affect the semantics of the program. This prevents even relatively simple optimizations, such as register allocation of shared data, because such a process usually interchanges reads and writes. In implicitly parallelized programs—for example, those written in High Performance FORTRAN (HPF)—programs must be synchronized and the synchronization points are known, so this issue does not arise.

Using Speculation to Hide Latency in Strict Consistency Models

As we saw in Chapters 4 and 5, speculation can be used to hide memory latency. It can also be used to hide latency arising from a strict consistency model, giving much of the benefit of a relaxed memory model. The key idea is for: the processor to use dynamic scheduling to reorder memory references, letting them possibly execute out-of-order. Executing the memory references out-of-order may generate violations of sequential consistency, which might affect the execution of the program. This possibility is avoided by using the delayed commit feature of a

speculative processor. Assume the coherency protocol is based on invalidation. If the processor receives an invalidation for a memory reference before the memory reference is committed, the processor uses speculation recovery to back-out the computation and restart with the memory reference whose address was invalidated.

If the reordering of memory requests by the processor yields an execution order that could result in an outcome that differs from what would have been seen under sequential consistency, the processor will redo the execution. The key to using this approach is that the processor need only guarantee that the result would be the same as if all access were completed in order, and it can achieve this by detecting when the results might differ. The approach is attractive because the speculative restart will rarely be triggered. It will only be triggered when there are unsynchronized access that actually cause a race. Gharachorloo, et. al. made this observation in a 1992 paper.

Hill in a 1998 paper advocates the combination of sequential or processor consistency together with speculative execution as the consistency model of choice. His argument has three parts. First, that an aggressive implementation of either sequential consistency or processor consistency will gain most of the advantage of a more relaxed model. Second, that such an implementation adds very little to the implementation cost of a speculative processor. Third, that such an approach allows the programmer to reason using the simpler programming models of either sequential or processor consistency.

The MIPS R10000 design team had this insight in the mid 1990s and used the R10000's out-of-order capability to support this type of aggressive implementation of sequential consistency. Hill's arguments are likely to motivate others to follow this approach.

One open question is how successful compiler technology will be in optimizing memory references to shared variables. The state of optimization technology and the fact that shared data is often accessed via pointers or array indexing has limited the use of such optimizations. If this technology became available and led to significant performance advantages, compiler writers would want to be able to take advantage of a more relaxed programming model.

Using Virtual Memory Support to Build Shared Memory

Suppose we wanted to support a shared address space among a group of workstations connected to a network. One approach is to use the virtual memory mechanism and operating system (OS) support to provide shared memory. This approach, which was first explored more than 10 years ago, has been called *distributed virtual memory (DVM)* or *shared virtual memory (SVM)*. The key observation that this idea builds on is that the virtual memory hardware has the ability to control access to portions of the address space for both reading and writing. By using the hardware to check and intercept accesses and the operating system to ensure coherence, we can create a coherent, shared address space across the distributed memory of multiple processors.

In SVM, pages become the units of coherence, rather than cache blocks. The OS can allow pages to be replicated in read-only fashion, using the virtual memory support to protect the pages from writing. When a process attempts to write such a page, it traps to the operating system. The operating system on that processor can then send messages to the OS on each node that shares the page, requesting that the page be invalidated. Just as in a directory system, each page has a home node, and the operating system running in that node is responsible for tracking who has copies of the page.

The mechanisms are quite similar to those at work in coherent shared memory. The key differences are that the unit of coherence is a page and that software is used to implement the coherence algorithms. It is exactly these two differences that lead to the major performance differences. A page is considerably bigger than a cache block, and the possibilities for poor usage of a page and for false sharing are very high. Such events can lead to much less stable performance and sometimes even lower performance than a uniprocessor. Because the coherence algorithms are implemented in software, they have much higher overhead.

The result of this combination is that shared virtual memory has become an acceptable substitute for loosely coupled message passing, since in both cases the frequency of communication must be low, and communication that is structured in larger blocks is favored. Distributed virtual memory is not currently competitive with schemes that have hardware-supported, coherent memory, such as the distributed shared-memory schemes we examined in section 6.5: Most programs written for coherent shared memory cannot be run efficiently on shared virtual memory without changes.

Several factors could change the attractiveness of shared virtual memory. Better implementation and small amounts of hardware support could reduce the overhead in the operating system. Compiler technology, as well as the use of smaller or multiple page sizes, could allow the system to reduce the disadvantages of coherence at a page-level granularity. The concept of software-supported

shared memory remains an active area of research, and such techniques may play an important role in connecting more loosely coupled multiprocessors, such as networks of workstations.

Performance Measurement of Parallel Processors

One of the most controversial issues in parallel processing has been how to measure the performance of parallel processors. Of course, the straightforward answer is to measure a benchmark as supplied and to examine wall-clock time. Measuring wall-clock time obviously makes sense; in a parallel processor, measuring CPU time can be misleading because the processors may be idle but unavailable for other uses.

Users and designers are often interested in knowing not just how well a multiprocessor performs with a certain fixed number of processors, but also how the performance scales as more processors are added. In many cases, it makes sense to scale the application or benchmark, since if the benchmark is unscaled, effects arising from limited parallelism and increases in communication can lead to results that are pessimistic when the expectation is that more processors will be used to solve larger problems. Thus, it is often useful to measure the speedup as processors are added both for a fixed-size problem and for a scaled version of the problem, providing an unscaled and a scaled version of the speedup curves. The choice of how to measure the uniprocessor algorithm is also important to avoid anomalous results, since using the parallel version of the benchmark may understate the uniprocessor performance and thus overstate the speedup, as discussed with an example in section 6.14.

Once we have decided to measure scaled speedup, the question is *how* to scale the application. Let's assume that we have determined that running a benchmark of size n on p processors makes sense. The question is how to scale the benchmark to run on $m \times p$ processors. There are two obvious ways to scale the problem: keeping the amount of memory used per processor constant; and keeping the total execution time, assuming perfect speedup, constant. The first method, called *memory-constrained scaling*, specifies running a problem of size $m \times n$ on $m \times p$ processors. The second method, called *time-constrained scaling*, requires that we know the relationship between the running time and the problem size, since the former is kept constant. For example, suppose the running time of the application with data size n on p processors is proportional to n^2/p . Then with time-constrained scaling, the problem to run is the problem whose ideal running time on $m \times p$ processors is still n^2/p . The problem with this ideal running time has size $\sqrt{m} \times n$.

EXAMPLE Suppose we have a problem whose execution time for a problem of size n is proportional to n^3 . Suppose the actual running time on a 10-processor multiprocessor is 1 hour. Under the time-constrained and memory-constrained scaling models, find the size of the problem to run and the effec-

tive running time for a 100-processor multiprocessor.

ANSWER For the time-constrained problem, the ideal running time is the same, 1 hour, so the problem size is $\sqrt[3]{10} \times n$ or 2.15 times larger than the original. For memory-constrained scaling, the size of the problem is $10n$ and the ideal execution time is $10^3/10$, or 100 hours! Since most users will be reluctant to run a problem on an order of magnitude more processors for 100 times longer, this size problem is probably unrealistic. n

In addition to the scaling methodology, there are questions as to how the program should be scaled when increasing the problem size affects the quality of the result. Often, we must change other application parameters to deal with this effect. As a simple example, consider the effect of time to convergence for solving a differential equation. This time typically increases as the problem size increases, since, for example, we often require more iterations for the larger problem. Thus when we increase the problem size, the total running time may scale faster than the basic algorithmic scaling would indicate.

For example, suppose that the number of iterations grows as the log of the problem size. Then for a problem whose algorithmic running time is linear in the size of the problem, the effective running time actually grows proportional to $n \log n$. If we scaled from a problem of size m on 10 processors, purely algorithmic scaling would allow us to run a problem of size $10m$ on 100 processors. Accounting for the increase in iterations means that a problem of size $k \times m$, where $k \log k = 10$, will have the same running time on 100 processors. This problem size yields a scaling of $5.72m$, rather than $10m$.

In practice, scaling to deal with error requires a good understanding of the application and may involve other factors, such as error tolerances (for example, it affects the cell-opening criteria in Barnes-Hut). In turn, such effects often significantly affect the communication or parallelism properties of the application as well as the choice of problem size.

Scaled speedup is not the same as unscaled (or true) speedup; confusing the two has led to erroneous claims (e.g., see the fallacy on page 753). Scaled speedup has an important role, but only when the scaling methodology is sound and the results are clearly reported as using a scaled version of the application. Singh et al.[1993] describe these issues in detail.

6.11 Putting It All Together: Sun's Wildfire Prototype

In Sections 6.3 and 6.5 we examined centralized memory architectures (also known as SMPs or symmetric multiprocessors) and distributed memory architectures (also known as DSMs or distributed shared memory multiprocessors). SMPs have the advantage of maintaining a single centralized memory with uni-

form access time, and although cache hit rates are crucial, memory placement is not. In comparison, DSMs have a nonuniform memory architecture and memory placement can be important; in return, they can achieve far greater scalability.

The question is whether there is a way to combine the advantages of the two approaches: maximizing the uniform memory access property while simultaneously allowing greater scalability. The answer is a partial yes, if we accept some compromises on the uniformity of the memory model and some limits of scalability. The machine we discuss in this section, an experimental prototype multiprocessor called Wildfire, built by Sun Microsystems, attempts to do exactly this.

One key motivation for an approach that maximizes the uniformity of memory access while accepting some limits on scalability is the rising importance of OLTP and web server markets for large-scale multiprocessors. In comparison to scientific applications, which played a key role in driving both SMP and DSM development, commercial server applications have both less predictable memory access patterns and less demand for scalability to hundreds or thousands of processors.

The Wildfire Architecture

Wildfire attempts to maximize the benefits of SMP, while allowing scalability by creating a DSM architecture using large SMPs as the nodes. The individual nodes in the Wildfire design are Sun E series multiprocessors (E6x00, E5x00, E4x00, or E3x00). Our measurements in this section are all done with E6000 multiprocessors as the nodes. An E6000 can accept up to 15 processor or I/O boards on the Gigaplane bus interconnect, which supports 50 million bus transactions per second, up to 112 outstanding transactions, and has a peak bandwidth of 3.2 GB/sec. Each processor board contains 2 UltraSPARC III processors.

Wildfire can connect 2 to 4 E6000 multiprocessors by replacing one dual processor (or I/O) board with a Wildfire Interface board (WFI), yielding up to 112 processors (4 x 28), as shown in Figure 6.45. The WFI board supports one coherent address space across all four multiprocessor nodes with the two high-order address bits used to designate which node contains a memory address. Hence, Wildfire is a cache-coherent nonuniform memory access architecture (cc-NU-MA) with large nodes. Within each node of 28 processors, memory is uniformly accessible, only processes that span nodes need to worry about the non uniformity in memory access times.

The WFI plugs into the bus and sees all memory requests; it implements the global coherence across up to four nodes. Each WFI has three ports that connect to up to three additional Wildfire nodes, each with a dual directional 800 MB/sec connection. WFI uses a simple directory scheme, similar to that discussed in Section 6.7. To keep the amount of directory state small, the directory is actually a cache, which is backed by the main memory in the node. When a miss occurs, the request is routed to the home node for the requested address. When the re-

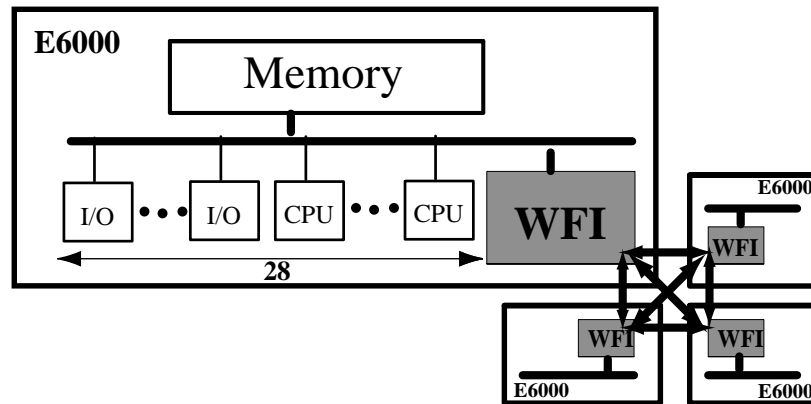


FIGURE 6.45 The Wildfire Architecture uses a bus-based SUN Enterprise server as its building blocks. The WFI replaces a processor or I/O board and provides internode coherency and communication, as well as hardware support for page replication.

quest arrives at the WFI of the home node, the WFI does a directory look-up. If the address is cached locally or clean in memory, a bus transaction is used to retrieve it. If the requested data is cached exclusively in a remote node, a request is sent to that remote node, where the WFI on that node generates a bus request to fetch the data. When the data is returned from either the remote owner or the home node, it is placed on the bus by the WFI and returned to the requesting processor.

We can see from this discussion one major disadvantage of this design: each remote request requires either four or five bus transactions. Two transactions are required at the local node and two or three others are required elsewhere, depending on where the data is cached:

- If the referenced data is cached only in the home node, then two additional bus transactions in the home are sufficient to retrieve the data.
- If the referenced data is cached exclusively in a third remote node, then two bus transactions are required at the remote node and one is required at the home node (to write the shared data back into the home memory).

These are one-way transactions and the E6000 bus is a split transaction bus, so even a normal memory access takes two bus transactions. Nonetheless, there is an increase in the required bus bandwidth of between 1.5 and 2.5. This increase

means that the processor count at which the buses within the nodes become saturated is lowered by a factor of at least 1.5, so that if a 28 processor design saturated the bus bandwidth of an E6000, a four-node Wildfire design could accommodate about 18 processors per node before saturating the bus bandwidth, assuming an even distribution of remote requests. A significant fraction of requests to data cached remotely from its home would further lower the useful processor count in each node. We will return to a further discussion of the advantages and disadvantages of this approach shortly, but first, let us look at how the Wildfire design reduces the fraction of costly remote memory accesses.

Using Page Replication and Migration to Reduce NUMA Effects

Wildfire uses special support, called CMR for *Coherent Memory Replication*, for page migration and replication. The idea is inspired by a more sophisticated hardware scheme for supporting migration and replication, called COMA for *Cache Only Memory Architecture*. COMA is an approach that treats all main memory as a cache allowing replication and migration of memory blocks. Full COMA implementations are quite complex, so a variety of simplifications have been proposed. CMR is based on one of these simplifications called S-COMA, for Simple COMA. S-COMA, like CMR, uses page-level mechanisms for migrating and replicating pages in memory, although coherence is still maintained at the cache-block level. We discuss the COMA ideas, as well as other approaches to migration and replication, in more detail in the historical perspectives and in the exercises.

To decide when to replicate or migrate pages, CMR uses a set of page counters that record the frequency of misses to remote pages. Migration is preferred when a page is primarily used by a node other than the one where the page is currently allocated. Replication is useful when multiple nodes share a page; the drawback of replication is that it requires extra memory. When the node sizes in a DSM are small, page migration and replication can lead to both excessive overhead for moving pages and excessive memory overhead from duplication of pages. With the large nodes in Wildfire, however, page-level migration and replication are much more attractive.

CMR, like S-COMA, maintains coherence at the unit of a cache-block, rather than at the page level. This choice is important for two reasons. First, maintaining coherence at the page level is likely to lead to a significant numbers of false sharing misses; we saw this increase in false sharing misses with increases in block size in Section 6.3. Second, the large size of a page means that even true sharing misses are likely to end up moving many bytes of data that are never used. These two drawbacks have limited the usefulness of the Shared Virtual Memory approach, which we discussed on page 733. CMR avoids these problems by making the unit of coherence a cache block and by selectively migrating and replicating some pages, while leaving others as standard NUMA pages that are accessed remotely when a cache miss occurs.

In addition to the page counters that the operating system uses to decide when to migrate or replicate a page, CMR requires special support to map between physical and virtual addresses of replicated pages. First, when a page is replicated the page tables are changed to refer to the local physical memory address of the duplicated page. To maintain coherence, however, a miss to this page must be sent to the home node to check the directory entry in that node. Thus, the WFI maintains a structure that maps the address of a replicated page (the local physical address) to its original physical address (called the global address) and generates the appropriate remote memory request, just as if the page were never replicated. When a write-back request or invalidation request is received, the global address must be translated to the local address, and the WFI maintains such a mapping for all pages that have been replicated. By maintaining these two maps, pages can be replicated while maintaining coherence at the unit of a cache block, which increases the usefulness of page replication.

Performance of Wildfire

In this section we look at the performance of the Wildfire prototype starting first with basic performance measures such as latency for memory accesses and bandwidth and then turning to application performance. Since Wildfire is a research prototype, rather than a product, its performance evaluation for applications is limited, but some interesting experiments that evaluate the use of page migration and replication are available.

Basic Performance Measures: Latency and Bandwidth

To better understand the design trade-offs between DSM architectures with nodes that have small, medium, and large processor counts, we compare the latency and bandwidth measurements of two different machines: the Sun Wildfire and the SGI Origin 2000.

The SGI Origin 2000 is a highly scalable cc-NUMA architecture capable of accommodating up to 2,048 processors. Each node consists of a pair of MIPS R1000 processors sharing a single memory module. An interface processor called the Hub (see Figure 6.46) provides an interface to the memory and directory in each node and implements the coherence protocol. The Hub interfaces directly to the routing chip, which provides a hypercube interconnection network that maintains a bisection bandwidth of 200 MB/sec per processor. The high dimension of the router also reduces hop counts leading to a lower ratio of remote to local access.

The Origin and Wildfire designs have significantly different motivations, so a comparison of the design trade-offs must acknowledge this fact. Among the most important differences are:

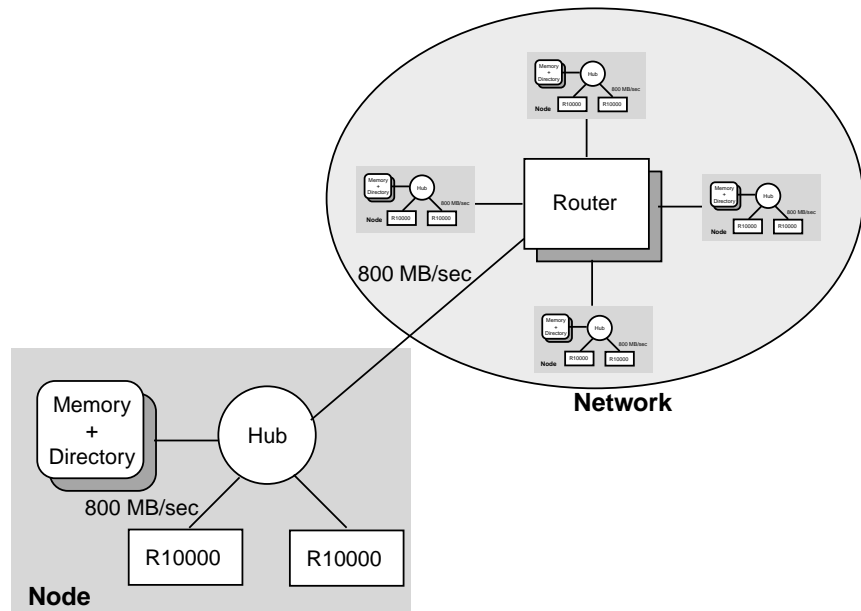


FIGURE 6.46 The SGI Origin 2000 uses an architecture that contains two processors per node and a scalable interconnection network that can handle up to 2,048 processors. A higher dimension network leads to scalable bisection bandwidth and a low ratio per out-of-node and in-node references.

- n The range of scalability: Origin can scale to thousands of processors, while the Wildfire design can scale to 112. Practically, the Wildfire design limit is likely to be closer to 64 to 80 processors, since bus bandwidth limits and the need for I/O boards will reduce the effective size of each node.
- n The Origin is designed primarily, though not exclusively, for scientific computation and the Wildfire design is oriented primarily for commercial processing. For the Origin design, this means that scalable bandwidth is crucial, and for the Wildfire design, it means that hiding more of the NUMA-ness is crucial.
- n The processors are also different in ways that affect both the bandwidth and latency of the nodes, including the block sizes of the L2 caches. We try to reduce this artifact by supplying multiple comparison numbers (e.g., latency to restart and back-to-back worst-case latency).

In Figure 6.47 we compare a variety of latency measurements for the two machines showing the variation arising both from local versus remote accesses and the variation arising from the cache organization. The first portion of the table concentrates on local memory accesses, which remain within one node. We compare both the restart latency, which is the time from miss detection to pipeline re-

start, and a worst-case, back-to-back measurement, which is measured by a sequence of dependent loads. The performance differences arise from the cache architecture (including a factor of two difference in block size), the pipeline architecture, and the main memory access time. Local memory latency also depends on the state of the cache block. We show three cases:

Characteristic	How measured?	Target status?	Sun Wildfire	SGI Origin 2000
Local memory latency	Restart	Unowned	342	338
Local memory latency	Back-to-back	Unowned	330	472
Local memory latency	Restart	Exclusive	362	656
Local memory latency	Back-to-back	Exclusive	350	707
Local memory latency	Restart	Dirty	482	892
Local memory latency	Back-to-back	Dirty	470	1036
Remote memory latency to nearest node	Restart	Unowned	1774	570
Remote memory latency to nearest node	Restart	Dirty	2162	1128
Remote memory latency to furthest node (< 128)	Restart	Unowned	1774	1219
Remote memory latency to furthest node (< 128)	Restart	Dirty	2162	1787
Avg. remote memory latency processors (< 128)	Restart	Unowned	1774	973
Avg. remote memory latency: processors (< 128)	Restart	Dirty	2162	1531
Average memory latency all processors (< 128)	Restart	Unowned	1416	963
Average memory latency all processors (< 128)	Restart	Dirty	1742	1520
Three hop miss to nearest node	Restart	Dirty	2550	953
Three hop miss to furthest node (worst case)	Restart	Dirty	2550	1967
Average three hop miss	Restart	Dirty	2453	1582

FIGURE 6.47 A comparison of memory access latencies (in ns) between the Sun Wildfire prototype (using E6000 nodes) and a SGI Origin 2000 shows significant differences in both local and remote access times. This table has four parts corresponding to local memory accesses (which are within the node), remote memory access involving only the requesting and home node, a third section that compares the average memory latency for the combination of local and remote (but not 3-hop) misses, and a final section showing the 3-hop latencies. The second column describes whether the latency is measured by time to restart the pipeline or by the back-to-back miss cost. For local accesses we show both; for remote accesses, we show the restart latency, which is the more likely case. The third column indicates the state of the remote data. Unowned means that it is in the shared or invalid state in the other caches. Exclusive means exclusive but clean, which requires an intervention to be completed before the memory access can complete, so that write serialization may be maintained. Dirty indicates that the data is exclusive and has been updated; an access, therefore, requires retrieving the data from the cache. In the local case, we show all three possibilities, to show the effect of the processor architecture (e.g., intervention cost and cache block size both affect the access times), while for remote accesses we show the unowned and dirty case, which are likely to be the most frequent.

1. the accessed block is unowned or it is in the shared state
2. the accessed block is owned exclusively but clean, which requires that the block be invalidated,
3. the accessed block is owned and dirty, which requires that the block be retrieved from the cache to satisfy the miss.

These 6 combinations (3 possible states of the target block x 2 possible miss timings) are the most likely cases of a local miss, though there are several other possibilities. These latencies are primarily dominated by choices in the microprocessor design (such as minimizing time to restart or minimizing total miss time) as well as in the local memory system and coherence implementation. These choices increase the difficulty in comparing memory latency for a multiprocessor, since some of these design choices affect the remote latencies as well.

The second section of the table compares the remote access times under a variety of different circumstances but all assuming that the home address is in a different node and that any cached copies are in the home node. For these numbers we use restart latency and consider the two most probable coherence states for a remotely accessed datum: unowned and dirty. The first two entries describe the time to access a datum whose home is in the nearest node; for the Wildfire system all remote nodes are equidistant, while for the Origin, the nearest node is one router hop away. The second pair of numbers deals with the latency when the home is as far away as possible for Origin. Finally, the third and last pair provide the average latency for a uniform distribution of the home address across a multiprocessor with 128 processors.

The fourth set of numbers deals with 3-hop misses, assuming that the owner is in a different node from either the home or the originating node. Here the most likely case is that the data is Dirty, and we show the restart latency for this case under the best, worst, and average assumptions.

From these measurements, we can see several of the trade-offs at work in a design that uses large nodes versus one that uses smaller nodes. Large nodes increase the number of processors reachable with a local access, but also typically have a longer remote access time. The latter is driven primarily by the higher overhead of acquiring access to the bus either for the directory or to access a remote cached copy. Of course, access latency is only part of the picture, bandwidth is also affected by these design decisions.

As Figure 6.48 shows, the pipeline memory bandwidth can be measured in many different ways. The Origin design supports greater memory bandwidth by every measure except local bandwidth to dirty data. Local bandwidth and bisection bandwidth are almost three times higher on a per processor basis for Origin.

Application performance of Wildfire

In this section, we examine the performance of Wildfire, first on an OLTP application and then on a scientific application. We look at both the basic performance

Characteristic	Sun Wildfire MB/sec	SGI Origin 2000 MB/sec
Pipelined local memory bandwidth: unowned data	312	554
Pipelined local memory bandwidth: exclusive data	266	340
Pipelined local memory bandwidth: dirty data	246	182
Total local memory bandwidth (per node)	2,700	631
Local memory bandwidth per processor	96	315
Aggregate local memory bandwidth (all nodes, 112 processors)	10,800	39,088
Remote memory bandwidth, unowned data		508
Remote 3-hop bandwidth, dirty data		238
Total bisection bandwidth (112 processors)	9,600	25,600
Bisection bandwidth per processor (112 processors)	86	229

FIGURE 6.48 A comparison of memory bandwidth measurements (in MB/sec) between the Sun Wildfire prototype (using E6000 nodes) and a SGI Origin 2000 shows significant differences in both local and remote memory bandwidth. The first section of the table compares pipelined local memory bandwidth, which is defined as the sustainable bandwidth for independent accesses generated by a single processor; like restart latency, this measure depends on the state of the addressed cache block. The second section of the table compares the total local memory bandwidth (i.e., within a node) on a per processor basis and system wide. The third section compares the memory bandwidth for remote accesses, both a two-hop access to an unowned cache block and a three-hop access to a dirty cache block. The final section compares the total bisection bandwidth for the entire system and on a per processor basis.

of the architecture versus alternatives such as a strict SMP or a small-node NUMA and then consider the effect of Wildfire's support for replication and migration.

Performance of the OLTP Workload

In this study an OLTP application supporting 900 warehouses was run on a 16-processor E6000 and on a two-node, 16-processor Wildfire configuration. I/O was supplied by 240 disks connected by fiber-channel. To examine the performance of Wildfire and the effect of its support for replication and migration, we consider six system alternatives:

1. Ideal SMP: a 16-processor SMP design, modeled using the E6000.
2. Wildfire with CMR and locality scheduling: a 2-node, 16-processor Wildfire with replication and migration enabled and using the locality scheduling in the OS.
3. Wildfire with CMR only.
4. Wildfire base with neither CMR nor locality scheduling.
5. Unoptimized Wildfire with poor data placement: Wildfire with poor data

placement and unintelligent scheduling. Poor data placement is modeled by assuming that 50% of the cache misses are remote, which in practice is unrealistic.

6. Unoptimized Wildfire with thin nodes (2 processors per node) and poor data placement. This system assumes Wildfires interconnection characteristics, but with eight two-processor nodes. Poor data placement is modeled by assuming that 87.5% (i.e., 14/16) of the cache misses are remote, which in practice is unrealistic.

To examine performance we first look at the fraction of cache misses satisfied within a node. Figure 6.49 shows the fraction of local accesses for each of these configurations. For this OLTP application the Wildfire optimizations improve fraction of local accesses by a factor of 1.23 over unoptimized Wildfire, bringing the fraction of local accesses to 87%.

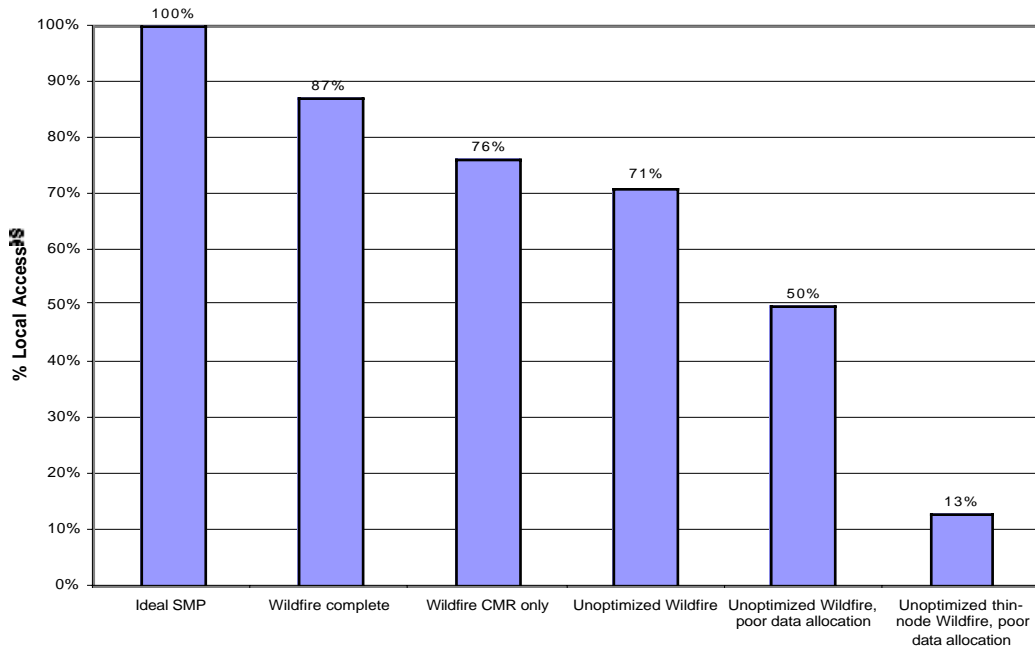


FIGURE 6.49 The fraction of local accesses (defined as within the node) is shown for six different configurations, ranging from an ideal SMP (with only one node and 16 processors) to four configurations with 8-processor nodes, to a configuration with thin, 2-processor nodes. The fraction of remote accesses is set as a parameter for the two right-most data points, while the other numbers are measured.

Figure 6.50 shows how these changes in local versus remote access fractions translate to performance for this OLTP application. The performance of each system in Figure 6.50 is relative to the E6000; however, as we will see when we examine a scientific application, the E6000 can encounter performance losses from bus contention at 16 processors, so that, in fact, the performance of the E6000 does not represent an upper bound for a multiprocessor using sixteen of the same processors. The E6000 performance is probably within 10-20% of contention-free performance for this benchmark. As we can see from the data the penalty for off-node accesses translates directly to reduced performance. The next section examine how Wildfire performs for a scientific application.

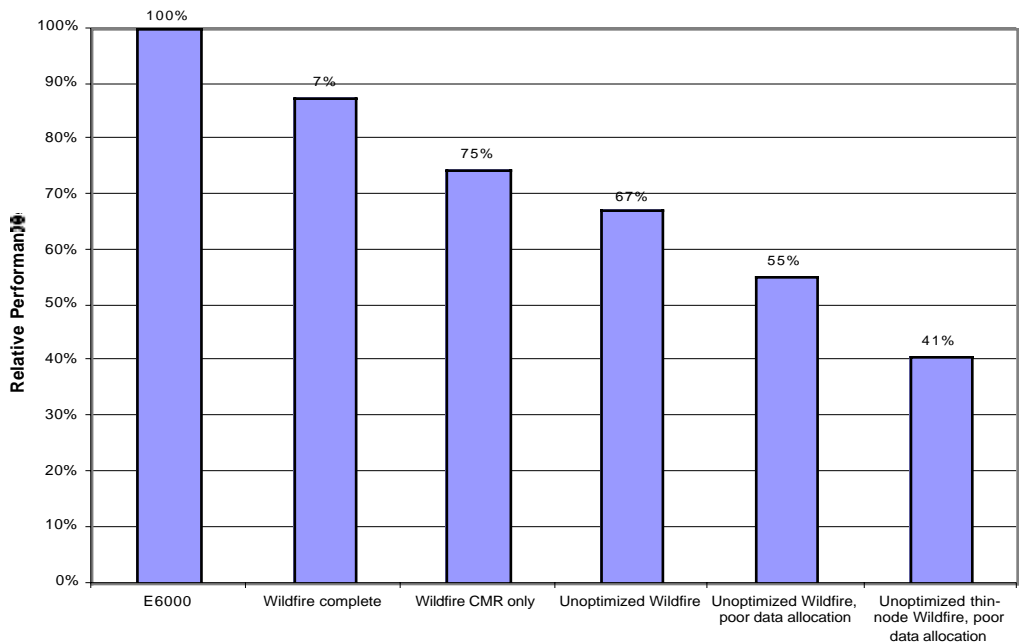


FIGURE 6.50 The performance of the OLTP application using 16 processors is highest for the E6000, and drops off as remote memory accesses become a major performance loss.

Performance of Wildfire on a Scientific Application

In this section we examine a performance study of Wildfire using a Red-Black finite difference solver to solve a 2-dimensional Poisson equation for a square grid. In this implementation, each 2x2 block of grid points is assigned either a red or

black color, so that the overall grid looks like a checkerboard. Red data points are updated based on values of black data points and vice versa, which allows all red points to be updated in parallel and all black points to be updated in parallel. A point is updated by accessing the four neighboring points (all of which are a different color). This data access pattern is common in two-dimensional solvers.

Our first performance comparisons examine the performance of Wildfire versus the E6000 and the E10000. The E 10000 uses a two-level interconnect. Four processors are connected with a 4x4 cross-bar to four memory modules, creating a 4-processor SMP. Up to 16 of these 4-processor nodes can be connected with the Starfire interconnect, which uses a 16x16 cross-bar. Coherence is maintained by a global broadcast scheme.

Figure 6.51 shows the performance of the generalized red-black (GRB) solver for six different configurations. The performance is given in terms of iterations per second with more iterations being better. The leftmost group of columns compares 24-processor measurements on an E6000, E10000, and Wildfire; while the rightmost bars compare 36-processor runs on two different Wildfire configurations and an E10000. Both the 24 and 36 processor runs use the same processor (250 MHZ UltraSPARC II) with 4MB secondary caches.

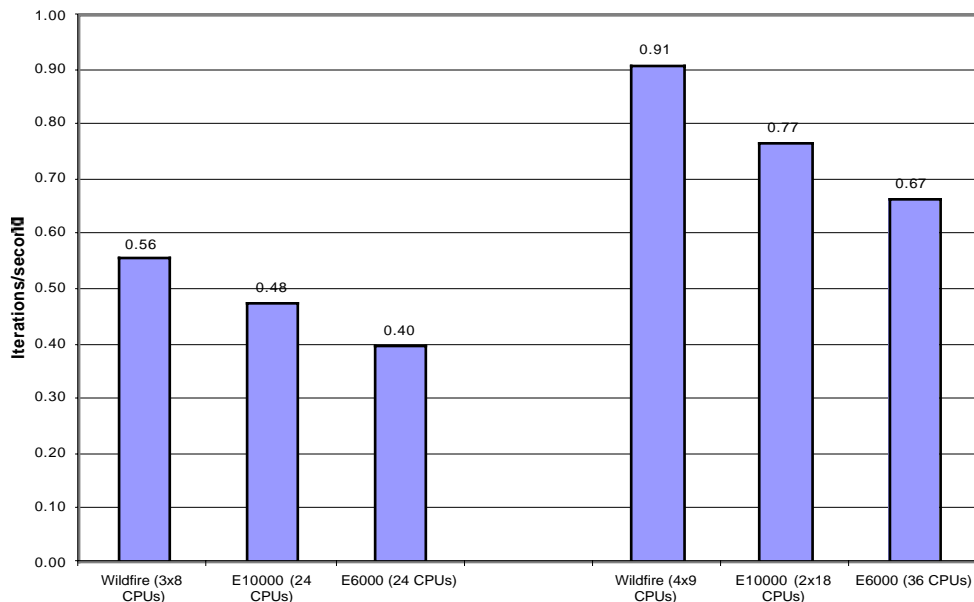


FIGURE 6.51 Wildfire performance for the Red-Black solver measured as iterations per second shows the performance for three different 24-processor and three different 36-processor machines. Iterations per second is directly proportional to performance.

The 24-processor runs include a 3-node Wildfire configuration (with an 8-processor E6000 in each Wildfire node), a 6-node E10000 and a 24-processor E6000. The performance differences among the 24-processor runs on Wildfire, the E1000, and the E6000 arise primarily from bus and interconnect differences. The global broadcast of the E10000 has nontrivial overhead. Thus, despite the fact that the E10000 interconnect has performance equal to that of Wildfire, the performance of Wildfire is about 1.17 times better. For the E6000, the measured bus usage for the 24-processor runs is between 90% and 100%, leading to a significant bottleneck and lengthened memory access time. Overall, Wildfire has a performance advantage of about 1.19 versus the E6000. Equally importantly, these measurements tell us that configurations of Wildfire with larger processor counts per node will not have good performance, at least for applications with behavior similar to this solver. The 36 processor runs confirm this view.

The 36-processor runs compare three alternatives: a 9-node E10000, a 2x18 configure of Wildfire (each Wildfire node is an 18-processor E6000) and a 4x9 configuration of Wildfire (each Wildfire node is an 9-processor E6000). The most interesting comparison here involve the 36-processor versus 24-processor results. The E10000 shows a faster than linear speedup (1.67 in runtime versus 1.5 in processor count); this probably results from improved cache behavior due to the smaller data set that each processor must access in the 36-processor case. The Wildfire results are even more interesting; the 4x9 configuration also shows faster than linear speed-up versus the 24-processor result. The 2x18 configuration, however, shows speed-up that is slower than linear (1.38 vs. 1.5), most probably because the bus has become a major bottleneck.

How well do the migration and replication capabilities of Wildfire work for scientific applications? To examine this question, this solver was executed starting with a memory allocation that placed all the data on a single node. Wildfire's migration and replication capabilities were used to allow data to migrate and replicate to one of the other nodes. Figure 6.52 shows the performance in iterations per second over time for a 1, 2, 3, and 4 node Wildfire, each with 24 processors/node. As shown, the 2, 3, and 4 node runs converge to stable and best performance after somewhere between 120 and 180 seconds. Since during the initial time period, the application averages about 0.2 iterations per second, it requires between 600 and 900 iterations to reach the stable performance levels.

Although the eventual convergence to a good operating point from an initial pathological memory allocation is impressive, the number of iterations required is rather large, and leaves open the question of how well the migration and replication strategies might work in problems where the memory allocation continued to change over time.

A key question is what the relative benefits of migration and replication are? Figure 6.53 examines this question by showing the iteration rate and time to reach that rate. We also show the number of replications and migrations. The pri-

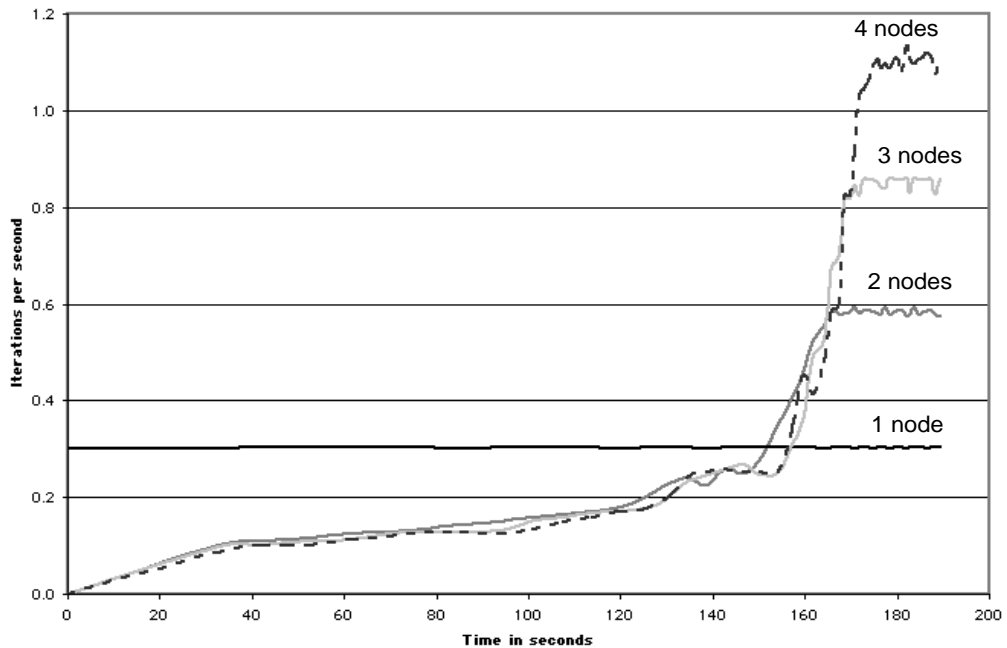


FIGURE 6.52 The replication and migration support of Wildfire allows an application to start with a pathological memory allocation (all memory on one node) and converge to a stable allocation that gives nearly linear speed-up. The final iterations/second number shows that the 96-processor, 4-node version achieves 90% of linear speedup. As expected, the two node runs converge slightly faster than three or four node runs.

mary conclusion we can draw from the performance of these three cases is that the stable performance level for migration is competitive with the combination of migration and replication. Since supporting migration had much lower hardware costs than supporting replication (because the reverse memory maps are not needed), a design that supports migration may be equally or more cost effective than supporting both migration and replication. The large data set coupled with well

defined access patterns by the “owner” of each portion of the grid means that replication buys little over only migration.

Policy	Iterations per second	Iterations needed to reach stability	# Migrations	# Replications
No migration or replication	0.10	0	0	0
Migration only	1.06	154 sec.	99,251	
Replication only	1.15	61 sec.		98,545
Migration + replication	1.09	151 sec.	98,543	85

FIGURE 6.53 Migration only, replication only, and the combination of all three achieve about the same performance given enough execution time and that number is roughly 10 times the performance achieved with the initial data allocation and no replication or migration. For this experiment, which used a 96-processor, 4-node Wildfire, the pages were allocated in a cyclic fashion, meaning that roughly 25% of were allocated to the correct location initially. A large data set size (16K x 8K) that exceeds the capacity of the secondary caches, leads to a high miss rate, which requires migration, replication, or careful initial data placement to reduce the miss penalty.

Concluding Remarks on Wildfire

Wildfire represents an alternative to thin-node NUMAs with 2 to 4 processors per node, while permitting greater scalability than strict SMP designs. The shift in market interest from scientific and supercomputing applications to large-scale servers for database and web applications may favor a fat-node design with 8 to 16 processors per node. The two primary reasons for this are:

1. Although a moderate range of scalability, up to a few hundred processors may be of interest, the “sweet spot” of the server market is likely to be tens of processors. Few, if any, customers will express interest in the thousand processor machines that are a key part of the supercomputer marketplace.
2. The memory access patterns of commercial applications tend to have less sharing and less predictable sharing and data access. The lower rates of sharing are key because a fat node design will tend to have lower bisection bandwidth per processor than a thin-node design. Since a fat-node design has somewhat less dependence on exact memory allocation and data placement, it is likely to perform better for applications with irregular or changing data access patterns. Furthermore, fat-nodes make it easier for migration and replication to work well.

The drawbacks of a fat node design are essentially the dual of its advantages. These include: less scalability, lower bisection bandwidth per processor, and higher internode latencies. For applications that require significant amounts of internode communication even with fat nodes, a fat-node design will face a more challenging programming and optimization task, since the ratio of local to remote accesses times is likely to be quite a bit larger. To read more on Wildfire see:

Hagersten and Koster [1998] and Noordergraaf and van der Pas [1999], which are also the sources for the data in this section.

Considering the growing significance of the commercial server market with its less predictable memory access patterns, its reduced emphasis on ultimate scalability, and its lower interprocess communication requirements, it is likely the “plump” node designs will become more attractive. Growing processor demands and avoidance of bus limits, is likely to lead to designs with 4-8 processors per node rather than the 16-24 limit in Wildfire. Although fatter nodes are likely to be beneficial, the nonuniform access time to memory cannot be ignored when the local node provide SMP-style access to only 3-7 other nodes.

6.12 Another View: Multithreading in a Commercial Server

As we have seen, dynamic scheduling can be used to make a single program run faster, as we saw in the Pentium III. Alternatively, multithreading can use a different form of dynamic scheduling (scheduling across multiple threads) to increase the throughput of multiple simultaneously executing programs. This is the approach used in the IBM RS64 III.

The IBM RS64 III processor, also called Pulsar, is a PowerPC microprocessor that supports two different IBM product lines: the RS/6000 series, where it is called the RS64 III processor, and the AS/400 series, where it is called the A50. Both product lines are aimed at commercial servers and focus on throughput in common commercial applications.

Motivated by the observation that such applications encounter high cache and TLB miss rates and thus degraded CPI, the designers decided to include a multithreading capability to enhance throughput and make use of the processor during long TLB or cache-miss stalls. In deciding how to support multithreading, the designers considered three facts:

1. The Pulsar processor, which was based on the earlier Northstar, is a statically scheduled processor.
2. The performance penalty for multithreading must be small both in silicon area and in clock rate.
3. Single thread performance on Pulsar must not suffer.

This combination of considerations led to a multithreading architecture with the following characteristics:

1. Pulsar supports precisely two threads: this minimizes both the incremental silicon area and the potential clock rate impact.

2. The multithreading is coarsely scheduled; that is, threads are not interleaved, instead a thread switch occurs only when a long latency stall is encountered. Coarse multithreading was chosen to maximize single thread performance and make use of the statically scheduled pipeline structure, which makes SMT an impractical choice.

To implement the multithreading architecture, Pulsar includes two copies of the register files and PC register, which resulted in relatively minor silicon overhead ($< 10\%$). In addition, a special register that determines the maximum number of cycles between a thread switch ensures that no thread is ever completely starved for cycles. The overall architecture provides a significant improvement in multithreaded throughput, a key metric for the commercial server workloads. The Pulsar microprocessor is the first widely available, mainline microprocessor to support multithreading; it is likely that future microprocessors will include such a capability either a coarse-grained form or using the SMT approach.

6.13 Another View: Embedded Multiprocessors

Multiprocessors are now common in server environments, and several desktop multiprocessors are available from vendors, such as Sun, Compaq, and Apple. In the embedded space, a number of special-purpose designs have used customized multiprocessors, including the Sony Playstation described in Chapters 2 and 5. Many special-purpose embedded designs consist of a general-purpose programmable processor with special purpose finite-state machines that are used for stream-oriented I/O. In applications ranging from computer graphics and media processing to telecommunications, this style of special-purpose multiprocessor is becoming common. Although the interprocessor interactions in such designs is highly regimented and relatively simple—consisting primarily of a simple communication channel—because much of the design is committed to silicon, ensuring that the communication protocols among the input/output processors and the general-purpose processor are correct is a major challenge in such designs.

More recently, we have seen the first appearance, in the embedded space, of embedded multiprocessors built from several general-purpose processors. These multiprocessors have been focused primarily on the high-end telecommunications and networking market, where scalability is critical. An example of such a design is the MXP processor designed by empowerTel Networks for use in voice over IP systems. The MXP processor consists of four main components:

1. An interface to serial voice streams, including support for handling jitter.
2. Support for fast packet routing and channel lookup.

3. A complete Ethernet interface, including the MAC layer.
4. Four MIPS32 R4000-class processors each with its own caches (a total of 48 KB or 12 KB per processor).

The MIPS processors are used to run the code responsible for maintaining the voice over IP channels, including the assurance of quality of service, echo cancellation, simple compression, and packet encoding. Since the goal is to run as many independent voice streams as possible, a multiprocessor is an ideal solution.

Because of the small size of the MIPS cores, the entire chip takes only 13.5M transistors. Future generations of the chip are expected to handle more voice channels, as well as do more sophisticated echo cancellation, voice activity detection, and more sophisticated compression.

Your authors expect that multiprocessing will become widespread in the embedded computing arena in the future for two primary reasons. First, the issues of binary software compatibility, which plague desktop and server systems, are less relevant in the embedded space. Often software in an embedded application is written from scratch for an application or significantly modified. Second, the applications often have natural parallelism, especially at the high-end of the embedded space. Examples of this natural parallelism abound in applications such as a set-top box, a network switch, or a game system. The lower barriers to use of thread-level parallelism together with the greater sensitivity to die cost (and hence efficient use of silicon) will likely lead to more ready adoption of multiprocessing in the embedded space, as the application needs grow to demand more performance.

6.14 Fallacies and Pitfalls

Given the lack of maturity in our understanding of parallel computing, there are many hidden pitfalls that will be uncovered either by careful designers or by unfortunate ones. Given the large amount of hype that has surrounded multiprocessors, especially at the high end, common fallacies abound. We have included a selection of these.

Pitfall: Measuring performance of multiprocessors by linear speedup versus execution time.

“Mortar shot” graphs—plotting performance versus number of processors showing linear speedup, a plateau, and then a falling off—have long been used to judge the success of parallel processors. Although speedup is one facet of a parallel program, it is not a direct measure of performance. The first question is the power of the processors being scaled: A program that linearly improves performance to equal 100 Intel 486s may be slower than the sequential version on a workstation. Be especially careful of floating-point-intensive programs; process-

ing elements without hardware assist may scale wonderfully but have poor collective performance.

Comparing execution times is fair only if you are comparing the best algorithms on each computer. Comparing the identical code on two processors may seem fair, but it is not; the parallel program may be slower on a uniprocessor than a sequential version. Developing a parallel program will sometimes lead to algorithmic improvements, so that comparing the previously best-known sequential program with the parallel code—which seems fair—will not compare equivalent algorithms. To reflect this issue, the terms *relative speedup* (same program) and *true speedup* (best program) are sometimes used.

Results that suggest *super-linear* performance, when a program on n processors is more than n times faster than the equivalent uniprocessor, may indicate that the comparison is unfair, although there are instances where “real” superlinear speedups have been encountered. For example, when Ocean is run on two processors, the combined cache produces a small superlinear speedup (2.1 vs. 2.0).

In summary, comparing performance by comparing speedups is at best tricky and at worst misleading. Comparing the speedups for two different multiprocessors does not necessarily tell us anything about the relative performance of the multiprocessors. Even comparing two different algorithms on the same multiprocessor is tricky, since we must use true speedup, rather than relative speedup, to obtain a valid comparison.

Fallacy: Amdahl's Law doesn't apply to parallel computers.

In 1987, the head of a research organization claimed that Amdahl's Law (see section 1.6) had been broken by an MIMD multiprocessor. This statement hardly meant, however, that the law has been overturned for parallel computers; the neglected portion of the program will still limit performance. To understand the basis of the media reports, let's see what Amdahl [1967] originally said:

A fairly obvious conclusion which can be drawn at this point is that the effort expended on achieving high parallel processing rates is wasted unless it is accompanied by achievements in sequential processing rates of very nearly the same magnitude. [p. 483]

One interpretation of the law was that since portions of every program must be sequential, there is a limit to the useful economic number of processors—say 100. By showing linear speedup with 1000 processors, this interpretation of Amdahl's Law was disproved.

The basis for the statement that Amdahl's Law had been “overcome” was the use of scaled speedup. The researchers scaled the benchmark to have a data set size that is 1000 times larger and compared the uniprocessor and parallel execution times of the scaled benchmark. For this particular algorithm the sequential

portion of the program was constant independent of the size of the input, and the rest was fully parallel—hence, linear speedup with 1000 processors.

We have already described the dangers of relating scaled speedup as true speedup. Additional problems with this sort of scaling methodology, which can result in unrealistic running times, were examined in section 6.10.

Fallacy: Linear speedups are needed to make multiprocessors cost-effective.

It is widely recognized that one of the major benefits of parallel computing is to offer a “shorter time to solution” than the fastest uniprocessor. Many people, however, also hold the view that parallel processors cannot be as cost-effective as uniprocessors unless they can achieve perfect linear speedup. This argument says that because the cost of the multiprocessor is a linear function of the number of processors, anything less than linear speedup means that the ratio of performance/cost decreases, making a parallel processor less cost-effective than using a uniprocessor.

The problem with this argument is that cost is not only a function of processor count, but also depends on memory and I/O. The effect of including memory in the system cost was pointed out by Wood and Hill [1995], and we use an example from their article to demonstrate the effect of looking at a complete system. They compare a uniprocessor server, the Challenge DM (a desktide unit with one processor and up to 6 GB of memory), against a multiprocessor Challenge XL, a rack-mounted, bus-based multiprocessor holding up to 32-processors. (The XL also has faster processors than those of the Challenge DM—150 MHz versus 100 MHz—but we will ignore this difference.)

First, Wood and Hill introduce a cost function: $cost(p, m)$, which equals the list price of a multiprocessor with p processors and m megabytes of memory. For the Challenge DM:

$$cost(1, m) = \$38,400 + \$100 \times m$$

For the Challenge XL:

$$cost(p, m) = \$81,600 + \$20,000 \times p + \$100 \times m$$

Suppose our computation requires 1 GB of memory on either multiprocessor. Then the cost of the DM is \$138,400, while the cost of the Challenge XL is $\$181,600 + \$20,000 \times p$.

For different numbers of processors, we can compute what speedups are necessary to make the use of parallel processing on the XL *more* cost effective than that of the uniprocessor. For example, the cost of an 8-processor XL is \$341,600, which is about 2.5 times higher than the DM, so if we have a speedup on 8 processors of more than 2.5, the multiprocessor is actually *more* cost effective than the uniprocessor. If we are able to achieve linear speedup, the 8-processor XL

system is actually more than *three times* more cost effective! Things get better with more processors: On 16 processors, we need to achieve a speedup of only 3.6, or less than 25% parallel efficiency, to make the multiprocessor as cost effective as the uniprocessor.

The use of a multiprocessor may involve some additional memory overhead, although this number is likely to be small for a shared-memory architecture. If we assume an extremely conservative number of 100% overhead (i.e., double the memory is required on the multiprocessor), the 8-processor multiprocessor needs to achieve a speedup of 3.2 to break even, and the 16-processor multiprocessor needs to achieve a speedup of 4.3 to break even.

Surprisingly, the XL can even be cost effective when compared against a headless workstation used as a server. For example, the cost function for a Challenge S, which can have at most 256 MB of memory, is

$$\text{cost}(1, m) = \$16,600 + \$100 \times m$$

For problems small enough to fit in 256 MB of memory on both multiprocessors, the XL breaks even with a speedup of 6.3 on 8 processors and 10.1 on 16 processors.

In comparing the cost/performance of two computers, we must be sure to include accurate assessments of both total system cost and what performance is achievable. For many applications with larger memory demands, such a comparison can dramatically increase the attractiveness of using a multiprocessor.

Fallacy: Multiprocessors are “free.”

This fallacy has two different interpretations, and both are erroneous. The first is, given that modern microprocessors contain support for snooping caches, we can build small-scale, bus-based multiprocessors for no additional cost in dollars (other than the microprocessor cost) or sacrifice of performance. Many designers believed this to be true and have even tried to build multiprocessors to prove it.

To understand why this doesn't work, you need to compare a design with no multiprocessing extensibility against a design that allows for a moderate level of multiprocessing (say 2–4 processors). The 2–4 processor design requires some sort of bus and a coherence controller that is more complicated than the simple memory controller required for the uniprocessor design. Furthermore, the memory access time is almost always faster in the uniprocessor case, since the processor can be directly connected to memory with only a simple single-master bus. Thus the strictly uniprocessor solution typically has better performance and lower cost than the 1-processor configuration of even a very small multiprocessor.

It also became popular in the 1980s to believe that the multiprocessor design was free in the sense that an MP could be quickly constructed from state-of-the-art microprocessors and then quickly updated using newer processors as they

became available. This viewpoint ignores the complexity of cache coherence and the challenge of designing high-bandwidth, low-latency memory systems, which for modern processors is extremely difficult. Moreover, there is additional software effort: compilers, operating systems, and debuggers all must be adapted for a parallel system. The next two fallacies are closely related to this one.

Fallacy: Scalability is almost free.

The goal of scalable parallel computing was a focus of much of the research and a significant segment of the high-end multiprocessor development from the mid-1980s through the late 1990s. In the first half of that period, it was widely held that you could build scalability into a multiprocessor and then simply offer the multiprocessor at any point on the scale from a small to large number of processors without sacrificing cost effectiveness. The difficulty with this view is that multiprocessors that scale to larger processor counts require substantially more investment (in both dollars and design time) in the interprocessor communication network, as well as in aspects such as operating system support, reliability, and reconfigurability.

As an example, consider the Cray T3E, which uses 3D torus capable of scaling to 2,048 processors as an interconnection network. At 128 processors, it delivers a peak bisection bandwidth of 38.4 GB/s, or 300 MB/s per processor. But for smaller configurations, the Compaq Alphaserp ES40 can accept up to 4 processors and has 5.6 GB/s of interconnect bandwidth, or almost four times the bandwidth per processor. Furthermore, the cost per CPU in a Cray T3E is several times higher than the cost in the ES40.

The cost of scalability can be seen even in more limited design ranges, such as the Sun Enterprise server line that all use the same basic Ultraport interconnect, scaling the amount of interconnect for different systems. For example, the 4 processor Enterprise 450 places all four processors on a single board and uses an on-board crossbar. The midrange system, designed to support 6 to 30 processors, uses a single address bus and a 32-byte wide data bus to connect the processors. The Enterprise 10000 series uses four addresses buses (memory address interleaved) and a 16x16 crossbar to connect the processors. Although the solution gives better scalability across the product range than forcing the low-end systems to accommodate four address buses and a multiboard crossbar, the cost of the interconnect system grows faster than linear as the number of processors grows, leading to a higher per processor cost for the 6000 series versus the 450 and for the 10000 series versus the 6000 series.

Scalability is also not free in software: To build software applications that scale requires significantly more attention to load balance, locality, potential contention for shared resources, and the serial (or partly parallel) portions of the program. Obtaining scalability for real applications, as opposed to toys or small kernels, across factors of more than 10 in processor count, is a *major* challenge.

In the future, better compiler technology and performance analysis tools may help with this critical problem.

Pitfall: Not developing the software to take advantage of, or optimize for, a multiprocessor architecture.

There is a long history of software lagging behind on massively parallel processors, possibly because the software problems are much harder. Two examples from mainstream, bus-based multiprocessors illustrate the difficulty of developing software for new multiprocessors. The first has to do with not being able to take advantage of a potential architectural capability, and the second arises from the need to optimize the software for a multiprocessor.

The SUN SPARCCenter was an earlier bus-based multiprocessor with one or two buses. Memory is distributed on the boards with the processors to create a simple building block consisting of processor, cache, and memory. With this structure, the multiprocessor could also have a fast local access and use the bus only to access remote memory. The SUN operating system, however, was not able to deal with the NUMA (non-uniform memory access) aspect of memory, including such issues as controlling where memory was allocated (local versus global). If memory pages were allocated randomly, then successive runs of the same application could have substantially different performance, and the benefits of fast local access might be small or nonexistent. In addition, providing both a remote and a local access path to memory slightly complicated the design because of timing. Since neither the system software nor the application software would not have been able to take advantage of faster local memory and the design was believed to be more complicated, the designers decided to require all requests to go over the bus.

Our second example shows the subtle kinds of problems that can arise when software designed for a uniprocessor is adapted to a multiprocessor environment. The SGI operating system protects the page table data structure with a single lock, assuming that page allocation is infrequent. In a uniprocessor this does not represent a performance problem. In a multiprocessor situation, it can become a major performance bottleneck for some programs. Consider a program that uses a large number of pages that are initialized at start-up, which UNIX does for statically allocated pages. Suppose the program is parallelized so that multiple processes allocate the pages. Because page allocation requires the use of the page table data structure, which is locked whenever it is in use, even an OS kernel that allows multiple threads in the OS will be serialized if the processes all try to allocate their pages at once (which is exactly what we might expect at initialization time!).

This page table serialization eliminates parallelism in initialization and has significant impact on overall parallel performance. This performance bottleneck persists even under multiprogramming. For example, suppose we split the parallel program apart into separate processes and run them, one process per proces-

sor, so that there is no sharing between the processes. (This is exactly what one user did, since he reasonably believed that the performance problem was due to unintended sharing or interference in his application.) Unfortunately, the lock still serializes all the processes—so even the multiprogramming performance is poor. This pitfall indicates the kind of subtle but significant performance bugs that can arise when software runs on multiprocessors. Like many other key software components, the OS algorithms and data structures must be rethought in a multiprocessor context. Placing locks on smaller portions of the page table effectively eliminates the problem.

Pitfall: Neglecting data distribution in a distributed shared-memory multiprocessor.

Consider the Ocean benchmark running on a 32-processor DSM architecture. As Figure 6.31 (page 699) shows, the miss rate is 3.1% for a 64KB cache. Because the grid used for the calculation is allocated in a tiled fashion (as described on page 658), 2.5% of the accesses are local capacity misses and 0.6% are remote communication misses needed to access data at the boundary of each grid. Assuming a 50-cycle local memory access cost and a 150-cycle remote memory access cost, the average miss has a cost of 69.3 cycles.

If the grid was allocated in a straightforward fashion by round-robin allocation of the pages, we could expect 1/32 of the misses to be local and the rest to be remote, which would lead to local miss rate of $3.1\% \times 1/32 = 0.1\%$ and a remote miss rate of 3.0%, for an average miss cost of 146.7 cycles. If the average CPI without cache misses is 0.6, and 45% of the instructions are data references, the version with tiled allocation is

$$\frac{0.6 + 45\% \times 3.1\% \times 146.7}{0.6 + 45\% \times 3.1\% \times 69.3} = \frac{0.6 + 2.05}{0.6 + 0.97} = \frac{2.65}{1.57} = 1.69 \text{ times faster}$$

This analysis only considers latency, and assumes that contention effects do not lead to increased latency, which is very optimistic. Round-robin is also not the worst possible data allocation: if the grid fit in a subset of the memory and was allocated to only a subset of the nodes, contention for memory at those nodes could easily lead to a difference in performance of more than a factor of 2.

6.15 Concluding Remarks

For over a decade prophets have voiced the contention that the organization of a single computer has reached its limits and that truly significant advances can be made only by interconnection of a multiplicity of computers in such a manner as to permit cooperative solution. ... Demonstration is made of the continued validity of the single processor approach. ... [p. 483]

Amdahl [1967]

The dream of building computers by simply aggregating processors has been around since the earliest days of computing. Progress in building and using effective and efficient parallel processors, however, has been slow. This rate of progress has been limited by difficult software problems as well as by a long process of evolving architecture of multiprocessors to enhance usability and improve efficiency. We have discussed many of the software challenges in this chapter, including the difficulty of writing programs that obtain good speedup due to Amdahl's law, dealing with long remote access or communication latencies, and minimizing the impact of synchronization. The wide variety of different architectural approaches and the limited success and short life of many of the architectures to date has compounded the software difficulties. We discuss the history of the development of these multiprocessors in section 6.16.

Despite this long and checkered past, progress in the last fifteen years leads to some reasons to be optimistic about the future of parallel processing and multiprocessors. This optimism is based on a number of observations about this progress and the long-term technology directions:

1. The use of parallel processing in some domains is beginning to be understood. Probably first among these is the domain of scientific and engineering computation. This application domain has an almost limitless thirst for more computation. It also has many applications that have lots of natural parallelism. Nonetheless, it has not been easy: programming parallel processors even for these applications remains very challenging. Another important, and much larger (in terms of market size), application area is large-scale data base and transaction processing systems. This application domain also has extensive natural parallelism available through parallel processing of independent requests, but its needs for large-scale computation, as opposed to purely access to large-scale storage systems, are less well understood. There are also several contending architectural approaches that may be viable—a point we discuss shortly.
2. It is now widely held that the most effective way to build a computer that offers more performance than that achieved with a single-chip microprocessor is by building a multiprocessor or a cluster that leverages the significant price/performance advantages of mass-produced microprocessors.
3. Multiprocessors are highly effective for multiprogrammed workloads, which are often the dominant use of mainframes and large servers, as well as for file servers or web servers, which are effectively a restricted type of parallel workload. In the future, such workloads may well constitute a large portion of the market for higher-performance multiprocessors. When a workload wants to share resources, such as file storage, or can efficiently timeshare a resource, such as a large memory, a multiprocessor can be a very efficient host. Further-

more, the OS software needed to efficiently execute multiprogrammed workloads is commonplace.

4. More recently, multiprocessors have proved very effective for certain intensive commercial workloads, such as OLTP (assuming the system supports enough I/O to be CPU-limited), DSS applications (where query optimization is critical), and large-scale, web searching applications. For commercial applications with undemanding communication requirements, little need for very large memories (typically used to cache databases), or limited demand for computation, clusters are likely to be more cost-effective than multiprocessors. The commercial space is currently a mix of clusters of basic PCs, SMPs, and clustered SMPs with different architectural styles appearing to hold some lead in different application spaces.
5. On-chip multiprocessing appears to be growing in importance for two reasons. First, in the embedded market where natural parallelism often exists, such approaches are an obvious alternative to faster, and possibly less silicon efficient, processors. Second, diminishing returns in high-end microprocessor design will encourage designers to pursue on-chip multiprocessing as a potentially more cost-effective direction. We explore the challenges to this direction at the end of this section.

Although there is reason to be optimistic about the growing importance of multiprocessors, many areas of parallel architecture remain unclear. Two particularly important questions are, How will the largest-scale multiprocessors (the massively parallel processors, or MPPs) be built? and What is the role of multiprocessing as a long-term alternative to higher-performance uniprocessors?

The Future of MPP Architecture

Hennessy and Patterson should move MPPs to Chapter 11.

Jim Gray, when asked about coverage of MPPs
in the second edition of this book, alludes to
Chapter 11 bankruptcy protection in U.S. law (1995)

Small-scale multiprocessors built using snooping-bus schemes are extremely cost-effective. Microprocessors traditionally have even included much of the logic for cache coherence in the processor chip, and several allow the buses of two or more processors to be directly connected—implementing a coherent bus with no additional logic. With modern integration levels, multiple processors can be placed on a board, on a single multi-chip module (MCM), or even within a single die (as we saw in Section 6.13) resulting in a highly cost-effective multiprocessor. Recent microprocessors have been including support for DSM approaches,

making it possible to connect small to moderate numbers of processors with little overhead. It is premature to predict that such architectures will dominate the middle range of processor counts (16–64), but it appears at the present that this approach is the most attractive.

What is totally unclear at the present is how the very largest multiprocessors will be constructed. The difficulties that designers face include the relatively small market for very large multiprocessors (> 64 nodes and often > \$5 million) and the need for multiprocessors that scale to larger processor counts to be extremely cost-effective at the lower processor counts where most of the multiprocessors will be sold. At the present there appear to be four slightly different alternatives for large-scale multiprocessors:

1. Large-scale multiprocessors that simply scale up naturally, using proprietary interconnect and communications controller technology. This approach has been followed in multiprocessors like the Cray T3E and the SGI Origin. There are two primary difficulties with such designs. First, the multiprocessors are not cost-effective at small scales, where the cost of scalability is not valued. Second, these multiprocessors have programming models that are incompatible, in varying degrees, with the mainstream of smaller and midrange multiprocessors.
2. Large-scale multiprocessors constructed from clusters of midrange multiprocessors with combinations of proprietary and standard technologies to interconnect such multiprocessors. The Wildfire design is just such a system. This cluster approach gets its cost-effectiveness through the use of cost-optimized building blocks. In some approaches, the basic architectural model (e.g., coherent shared memory) is extended. Many companies offer a high-end version of such a machine including HP, Sun, and SGI. Due to the two-level nature of the design, the programming model sometimes must be changed from shared memory to message passing or to a different variation on shared memory, among clusters. The migration and replication features in Wildfire offer a way to minimize this disadvantage. This class of machines has made important inroads, especially in commercial applications.
3. Designing clustered multicomputers that use off-the-shelf uniprocessor nodes and a custom interconnect. The advantage of such a design is the cost-effectiveness of the standard uniprocessor node, which is often a repackaged workstation; the disadvantage is that the programming model will probably need to be message passing even at very small node counts. In some application environments where little or no sharing occurs, this may be acceptable. In addition, the cost of the interconnect, because it is custom, can be significant, making the multiprocessor costly, especially at small node counts. The IBM SP-2 is the best example of this approach today.
4. Designing a cluster using *all* off-the-shelf components, which promises the

lowest cost. The leverage in this approach lies in the use of commodity technology everywhere: in the processors (PC or workstation nodes), in the interconnect (high-speed local area network technology, such as ATM or Gigabit Ethernet), and in the software (standard operating systems and programming languages). Of course, such multiprocessors will use message passing, and communication is likely to have higher latency and lower bandwidth than in the alternative designs. Like the previous class of designs, for applications that do not need high bandwidth or low-latency communication, this approach can be extremely cost-effective. Web servers, for example, may be a good match to these multicomputers, as we saw for the Google cluster in Chapter 8.

Each of these approaches has advantages and disadvantages, and the importance of the shortcomings of any one approach are dependent on the application class. In 2000 it is unclear which if any of these models will win out for larger-scale multiprocessors, although the growth of the market for web servers has made “racks of PCs” the dominant form at least by processor count. It is likely that the current bifurcation by market and scale will continue for some time, although in some area a hybridization of these ideas may emerge, given the similarity in several of the approaches.

The Future of Microprocessor Architecture

As we saw in Chapters 3 and 4, architects are using ever more complex techniques to try to exploit more instruction-level parallelism. As we also saw in that chapter, the prospects for finding ever-increasing amounts of instruction-level parallelism in a manner that is efficient to exploit are somewhat limited. Likewise, there are increasingly difficult problems to be overcome in building memory hierarchies for high-performance processors. Of course, continued technology improvements will allow us to continue to advance clock rate. But the use of technology improvements that allow a faster gate speed alone is not sufficient to maintain the incredible growth of performance that the industry has experienced for over 15 years. Maintaining a rapid rate of performance growth will depend to an increasing extent on exploiting the dramatic growth in effective silicon area, which will continue to grow much faster than the basic speed of the process technology.

Unfortunately, for almost ten years, increases in performance have come at the cost of ever-increasing inefficiencies in the use of silicon area, external connections, and power. This diminishing-returns phenomenon has only recently (as of 2001) appeared to have slowed the rate of performance growth. Whether or not this is slowdown temporary is unclear. What is clear, is that we cannot sustain the rapid rate of performance improvements without significant new innovations in computer architecture.

Unlike the prophets quoted at the beginning of the chapter, your authors do not believe that we are about to “hit a brick wall” in our attempts to improve single-

processor performance. Instead, we may see a gradual slowdown in performance growth, especially for integer performance, with the eventual growth being limited primarily by improvements in the speed of the technology. When these limitation will become serious is hard to say, but possibly as early as 2005 and likely by 2010. Even if such a slowdown were to occur, performance might well be expected to grow at the annual rate of 1.35 that we saw prior to 1985 at least until fundamental limitations in silicon are become serious in the 2015 time frame.

Furthermore, we do not want to rule out the possibility of a breakthrough in uniprocessor design. In the early 1980s, many people predicted the end of growth in uniprocessor performance, only to see the arrival of RISC technology and an unprecedented 15-year growth in performance averaging 1.5 times per year!

With this in mind, we cautiously ask whether the long-term direction will be to use increased silicon to build multiple processors on a single chip. Such a direction is appealing from the architecture viewpoint—it offers a way to scale performance without increasing hardware complexity. It also offers an approach to easing some of the challenges in memory-system design, since a distributed memory can be used to scale bandwidth while maintaining low latency for local accesses. The challenge lies in software and in what architecture innovations may be used to make the software easier.

In 2000, IBM announced the first commercial chips with two general-purpose processors on a single die, the Power4 processor. Each Power4 contains two Power3 microprocessors, a shared secondary cache, an interface to an off-chip tertiary cache or main memory, and chip-to-chip communication system, which allows a four processor cross-bar connected module to be built with no additional logic. Using 4 Power4 chips and the appropriate DRAMS, an eight-processor system can be integrated onto a board about 8 inches on a side. The board would contain 700 million transistors, not including the third level cache or main memory, and would have a peak instruction execution rate of 32 billion instructions per second!

Evolution Versus Revolution and the Challenges to Paradigm Shifts in the Computer Industry

Figure 6.54 shows what we mean by the *evolution-revolution spectrum* of computer architecture innovation. To the left are ideas that are invisible to the user (presumably excepting better cost, better performance, or both) and are at the evolutionary end of the spectrum. At the other end are revolutionary architecture ideas. These are the ideas that require new applications from programmers who must learn new programming languages and models of computation, and must invent new data structures and algorithms.

Revolutionary ideas are easier to get excited about than evolutionary ideas, but to be adopted they must have a much higher payoff. Caches are an example of an evolutionary improvement. Within 5 years after the first publication about caches, almost every computer company was designing a computer with a cache. The

RISC ideas were nearer to the middle of the spectrum, for it took more than eight years for most companies to have a RISC product and more than fifteen years for the last hold out to announce their product. Most multiprocessors have tended to the revolutionary end of the spectrum, with the largest-scale multiprocessors (MPPs) being more revolutionary than others. Most programs written to use multiprocessors as parallel engines have been written especially for that class of multiprocessors, if not for the specific architecture.

The challenge for both hardware and software designers that would propose that multiprocessors and parallel processing become the norm, rather than the exception, is the disruption to the established base of programs. There are two possible ways this paradigm shift could be facilitated: if parallel processing offers the only alternative to enhance performance, and if advances in hardware and software technology can construct a gentle ramp that allows the movement to parallel processing, at least with small numbers of processors, to be more evolutionary.

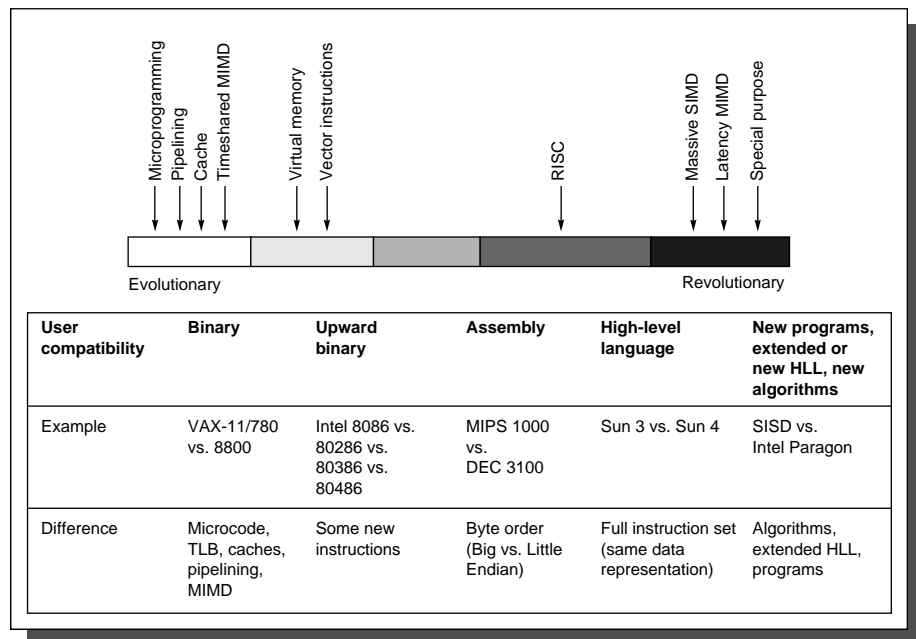


FIGURE 6.54 The evolution-revolution spectrum of computer architecture. The second through fourth columns are distinguished from the final column in that applications and operating systems can be ported from other computers rather than written from scratch. For example, RISC is listed in the middle of the spectrum because user compatibility is only at the level of high-level languages, while microprogramming allows binary compatibility, and latency-oriented MIMDs require changes to algorithms and extending HLLs. Timeshared MIMD means MIMDs justified by running many independent programs at once, while latency MIMD means MIMDs intended to run a single program faster.

6.16 Historical Perspective and References

There is a tremendous amount of history in parallel processing; in this section we divide our discussion by both time period and architecture. We start with the SIMD approach and the Illiac IV. We then turn to a short discussion of some other early experimental multiprocessors and progress to a discussion of some of the great debates in parallel processing. Next we discuss the historical roots of the present multiprocessors and conclude by discussing recent advances.

SIMD Computers: Several Attempts, No Lasting Successes

The cost of a general multiprocessor is, however, very high and further design options were considered which would decrease the cost without seriously degrading the power or efficiency of the system. The options consist of recentralizing one of the three major components.... Centralizing the [control unit] gives rise to the basic organization of [an]... array processor such as the Illiac IV.

Bouknight et al. [1972]

The SIMD model was one of the earliest models of parallel computing, dating back to the first large-scale multiprocessor, the Illiac IV. The key idea in that multiprocessor, as in more recent SIMD multiprocessors, is to have a single instruction that operates on many data items at once, using many functional units.

The earliest ideas on SIMD-style computers are from Unger [1958] and Slotnick, Borck, and McReynolds [1962]. Slotnick's Solomon design formed the basis of the Illiac IV, perhaps the most infamous of the supercomputer projects. Although successful in pushing several technologies that proved useful in later projects, it failed as a computer. Costs escalated from the \$8 million estimate in 1966 to \$31 million by 1972, despite construction of only a quarter of the planned multiprocessor. Actual performance was at best 15 MFLOPS, versus initial predictions of 1000 MFLOPS for the full system [Hord 1982]. Delivered to NASA Ames Research in 1972, the computer took three more years of engineering before it was usable. These events slowed investigation of SIMD, with Danny Hillis [1985] resuscitating this style in the Connection Machine, which had 65,636 1-bit processors.

Real SIMD computers need to have a mixture of SISD and SIMD instructions. There is an SISD host computer to perform operations such as branches and address calculations that do not need parallel operation. The SIMD instructions are broadcast to all the execution units, each of which has its own set of registers. For flexibility, individual execution units can be disabled during a SIMD instruction.

In addition, massively parallel SIMD multiprocessors rely on interconnection or communication networks to exchange data between processing elements.

SIMD works best in dealing with arrays in for-loops. Hence, to have the opportunity for massive parallelism in SIMD there must be massive amounts of data, or *data parallelism*. SIMD is at its weakest in case statements, where each execution unit must perform a different operation on its data, depending on what data it has. The execution units with the wrong data are disabled so that the proper units can continue. Such situations essentially run at $1/n$ th performance, where n is the number of cases.

The basic trade-off in SIMD multiprocessors is performance of a processor versus number of processors. Recent multiprocessors emphasize a large degree of parallelism over performance of the individual processors. The Connection Multiprocessor 2, for example, offered 65,536 single bit-wide processors, while the Illiac IV had 64 64-bit processors.

After being resurrected in the 1980s, first by Thinking Machines and then by MasPar, the SIMD model has once again been put to bed as a general-purpose multiprocessor architecture, for two main reasons. First, it is too inflexible. A number of important problems cannot use such a style of multiprocessor, and the architecture does not scale down in a competitive fashion; that is, small-scale SIMD multiprocessors often have worse cost/performance compared with that of the alternatives. Second, SIMD cannot take advantage of the tremendous performance and cost advantages of microprocessor technology. Instead of leveraging this low-cost technology, designers of SIMD multiprocessors must build custom processors for their multiprocessors.

Although SIMD computers have departed from the scene as general-purpose alternatives, this style of architecture will continue to have a role in special-purpose designs. Many special-purpose tasks are highly data parallel and require a limited set of functional units. Thus designers can build in support for certain operations, as well as hardwire interconnection paths among functional units. Such organizations are often called *array processors*, and they are useful for tasks like image and signal processing.

Other Early Experiments

It is difficult to distinguish the first MIMD multiprocessor. Surprisingly, the first computer from the Eckert-Mauchly Corporation, for example, had duplicate units to improve availability. Holland [1959] gave early arguments for multiple processors.

Two of the best-documented multiprocessor projects were undertaken in the 1970s at Carnegie Mellon University. The first of these was C.mmp [Wulf and Bell 1972; Wulf and Harbison 1978], which consisted of 16 PDP-11s connected by a crossbar switch to 16 memory units. It was among the first multiprocessors with more than a few processors, and it had a shared-memory programming model. Much of the focus of the research in the C.mmp project was on software, espe-

cially in the OS area. A later multiprocessor, Cm* [Swan et al. 1977], was a cluster-based multiprocessor with a distributed memory and a nonuniform access time. The absence of caches and a long remote access latency made data placement critical. This multiprocessor and a number of application experiments are well described by Gehringer, Siewiorek, and Segall [1987]. Many of the ideas in these multiprocessors would be reused in the 1980s when the microprocessor made it much cheaper to build multiprocessors.

Great Debates in Parallel Processing

The quotes at the beginning of this chapter give the classic arguments for abandoning the current form of computing, and Amdahl [1967] gave the classic reply in support of continued focus on the IBM 370 architecture. Arguments for the advantages of parallel execution can be traced back to the 19th century [Menabrea 1842]! Yet the effectiveness of the multiprocessor for reducing latency of individual important programs is still being explored. Aside from these debates about the advantages and limitations of parallelism, several hot debates have focused on how to build multiprocessors.

Predictions of the Future

It's hard to predict the future, yet in 1989 Gordon Bell made two predictions for 1995. We included these predictions in the first edition of the book, when the outcome was completely unclear. We discuss them in this section, together with an assessment of the accuracy of the prediction.

The first is that a computer capable of sustaining a teraFLOPS—one million MFLOPS—will be constructed by 1995, either using a multicomputer with 4K to 32K nodes or a Connection Multiprocessor with several million processing elements [Bell 1989]. To put this prediction in perspective, each year the Gordon Bell Prize acknowledges advances in parallelism, including the fastest real program (highest MFLOPS). In 1989 the winner used an eight-processor Cray Y-MP to run at 1680 MFLOPS. On the basis of these numbers, multiprocessors and programs would have to have improved by a factor of 3.6 each year for the fastest program to achieve 1 TFLOPS in 1995. In 1999, the first Gordon Bell prize winner crossed the 1 TF bar, using a 5,832 processor IBM RS/6000 SST system designed specially for Livermore Laboratories, they achieved 1.18 Teraflops on a shock-wave simulation. This ratio represents a year-to-year improvement of 1.93, which is still quite impressive.

What has become recognized since 1989 is that although we may have the technology to build a teraFLOPS multiprocessor, it is not clear that the machine is cost-effective, except perhaps for a few very specialized and critically important application related to national security. Your authors estimated in 1990 that to achieve 1 TF would require a machine with about 5,000 processors and would cost about \$100 million. The 5,832 processor IBM system at Livermore cost

\$110 million. As might be expected, improvements in the performance of individual microprocessors both in cost and performance directly affect the cost and performance of large-scale multiprocessors, but a 5000 processor system will cost more than 5000 times the price of a desktop system using the same processor.

The second Bell prediction concerned the number of data streams in supercomputers shipped in 1995. Danny Hillis believed that although supercomputers with a small number of data streams may be the best sellers, the biggest multiprocessors will be multiprocessors with many data streams, and these will perform the bulk of the computations. Bell bet Hillis that in the last quarter of calendar year 1995 more sustained MFLOPS will be shipped in multiprocessors using few data streams (≤ 100) rather than many data streams (≥ 1000). This bet concerned only supercomputers, defined as multiprocessors costing more than \$1 million and used for scientific applications. Sustained MFLOPS was defined for this bet as the number of floating-point operations per *month*, so availability of multiprocessors affects their rating.

In 1989, when this bet was made, it was totally unclear who would win. In 1995, a survey of the current publicly known supercomputers showed only six multiprocessors in existence in the world with more than 1000 data streams, so Bell's prediction was a clear winner. In fact, in 1995, much smaller microprocessor-based multiprocessors (≤ 20 processors) were becoming dominant. In 1995, a survey of the 500 highest-performance multiprocessors in use (based on Linpack ratings), called the Top 500, showed that the largest number of multiprocessors were bus-based shared-memory multiprocessors! By 2000, the picture had become less clear: the top four vendors were IBM (144 SP systems), Sun (121 Enterprise systems), SGI (62 Origin systems), and Cray (54 T3E systems). Although IBM holds the largest number of spots, almost all the other systems on the TOP 500 list are shared-memory systems or clusters of such systems.

More Recent Advances and Developments

With the primary exception of the parallel vector multiprocessors (see Appendix B), all other recent MIMD computers have been built from off-the-shelf microprocessors using a bus and logically central memory or an interconnection network and a distributed memory. A number of experimental multiprocessors built in the 1980s further refined and enhanced the concepts that form the basis for many of today's multiprocessors.

The Development of Bus-Based Coherent Multiprocessors

Although very large mainframes were built with multiple processors in the 1970s, multiprocessors did not become highly successful until the 1980s. Bell [1985] suggests the key was that the smaller size of the microprocessor allowed the memory bus to replace the interconnection network hardware, and that porta-

ble operating systems meant that multiprocessor projects no longer required the invention of a new operating system. In this paper, Bell defines the terms *multi-processor* and *multicomputer* and sets the stage for two different approaches to building larger-scale multiprocessors.

The first bus-based multiprocessor with snooping caches was the Synapse N+1 described by Frank [1984]. Goodman [1983] wrote one of the first papers to describe snooping caches. The late 1980s saw the introduction of many commercial bus-based, snooping-cache architectures, including the Silicon Graphics 4D/240 [Baskett et al. 1988], the Encore Multimax [Wilson 1987], and the Sequent Symmetry [Lovett and Thakkar 1988]. The mid 1980s saw an explosion in the development of alternative coherence protocols, and Archibald and Baer [1986] provide a good survey and analysis, as well as references to the original papers. Figure 6.55 summarizes several snooping cache-coherence protocols and shows some multiprocessors that have used or are using that protocol.

Name	Protocol type	Memory-write policy	Unique feature	Multiprocessors using
Write Once	Write invalidate	Write back after first write	First snooping protocol described in literature	
Synapse N+1	Write invalidate	Write back	Explicit state where memory is the owner	Synapse multiprocessors; first cache-coherent multiprocessors available
Berkeley (MOESI)	Write invalidate	Write back	Owned shared state	Berkeley SPUR multiprocessor; SUN Enterprise servers
Illinois (MESI)	Write invalidate	Write back	Clean private state; can supply data from any cache with a clean copy	SGI Power and Challenge series
“Firefly”	Write broadcast	Write back when private, write through when shared	Memory updated on broadcast	No current multiprocessors; SPARCCenter 2000 closest.

FIGURE 6.55 Five snooping protocols summarized. Archibald and Baer [1986] use these names to describe the five protocols, and Eggers [1989] summarizes the similarities and differences as shown in this figure. The Firefly protocol was named for the experimental DEC Firefly multiprocessor, in which it appeared. The alternative names for protocols are based on the states they support: M=Modified, E=Exclusive (shared clean), S=Shared, I=Invalid, O=Owner (shared dirty).

The early 1990s saw the beginning of an expansion of such systems with the use of very wide, high speed buses (the SGI Challenge system used a 256-bit, packet-oriented bus supporting up to 8 processor boards and 32 processors) and later, the use of multiple buses and crossbar interconnects, e.g. in the SUN SPARCCenter and Enterprise systems (Charlesworth [1998] discusses the interconnect architecture of these multiprocessors). In 2001, the Sun Enterprise serv-

ers represent the primary example of large-scale (> 16 processors), symmetric multiprocessors in active use.

Toward Large-Scale Multiprocessors

In the effort to build large-scale multiprocessors, two different directions were explored: message passing multicomputers and scalable shared-memory multiprocessors. Although there had been many attempts to build mesh and hypercube-connected multiprocessors, one of the first multiprocessors to successfully bring together all the pieces was the Cosmic Cube built at Caltech [Seitz 1985]. It introduced important advances in routing and interconnect technology and substantially reduced the cost of the interconnect, which helped make the multicomputer viable. The Intel iPSC 860, a hypercube-connected collection of i860s, was based on these ideas. More recent multiprocessors, such as the Intel Paragon, have used networks with lower dimensionality and higher individual links. The Paragon also employed a separate i860 as a communications controller in each node, although a number of users have found it better to use both i860 processors for computation as well as communication. The Thinking Multiprocessors CM-5 made use of off-the-shelf microprocessors and a fat tree interconnect (see Chapter 7). It provided user-level access to the communication channel, thus significantly improving communication latency. In 1995, these two multiprocessors represent the state of the art in message-passing multicomputers.

Early attempts at building a scalable shared-memory multiprocessor include the IBM RP3 [Pfister et al. 1985], the NYU Ultracomputer [Schwartz 1980; Elder et al. 1985], the University of Illinois Cedar project [Gajski et al. 1983], and the BBN Butterfly and Monarch [BBN Laboratories 1986; Rettberg et al. 1990]. These multiprocessors all provided variations on a nonuniform distributed-memory model (and hence are distributed shared memory or DSM multiprocessors), but did not support cache coherence, which substantially complicated programming. The RP3 and Ultracomputer projects both explored new ideas in synchronization (fetch-and-operate) as well as the idea of combining references in the network. In all four multiprocessors, the interconnect networks turned out to be more costly than the processing nodes, raising problems for smaller versions of the multiprocessor. The Cray T3D/E (see Arpaci et. al. [1995] for an evaluation of the T3D and Scott [1996] for a description of the T3E enhancements) builds on these ideas, using a noncoherent shared address space but building on the advances in interconnect technology developed in the multicomputer domain (see Scott and Thorson [1996]).

Extending the shared-memory model with scalable cache coherence was done by combining a number of ideas. Directory-based techniques for cache coherence were actually known before snooping cache techniques. In fact, the first cache-coherence protocols actually used directories, as described by Tang [1976] and implemented in the IBM 3081. Censier and Feautrier [1978] described a directory coherence scheme with tags in memory. The idea of distributing directories

with the memories to obtain a scalable implementation of cache coherence was first described by Agarwal et al. [1988] and served as the basis for the Stanford DASH multiprocessor (see Lenoski et al. [1990, 1992]), which was the first operational cache-coherent DSM multiprocessor. DASH was a “plump” node cc-NUMA machine that used 4-processor SMPs as its nodes; interconnecting them in a style similar to that of Wildfire but using a more scalable 2-dimensional grid rather than a crossbar for the interconnect.

The Kendall Square Research KSR-1 [Burkhardt et al. 1992] was the first commercial implementation of scalable coherent shared memory. It extended the basic DSM approach to implement a concept called *COMA* (*cache-only memory architecture*), which makes the main memory a cache. Like the Wildfire CMR scheme, in the KSR-1 memory blocks could be replicated in the main memories of each node with hardware support to handle the additional coherence requirements for these replicated blocks. (The KSR-1 was not strictly a pure COMA because it did not migrate the home location of a data item, but always kept a copy at home. Essentially, it implemented only replication.)

In parallel, researchers at the Swedish Institute for Computer Science [Hagersten et al. 1992.] developed a concept called DDM (for Data Diffusion Machine) which is a true COMA, since all memory operates as a cache, and a memory block does not exist in a predefined node. The absence of a designated home for a memory block significantly complicates the protocols, since it means that there is no static look-up scheme to find the location and status of a block. Furthermore, a true COMA must contend with the problem of finding a place to move a memory block when it conflicts with another block for the same location in memory (which happens because the memory is a cache with a limited associativity). In the event that the displaced block is the last copy of a memory block, which in itself may be difficult to know precisely, the displaced block must be migrated to some other memory location, since it cannot be destroyed (as it is the only copy of the data). This migration process can be very complex requiring a potentially unbounded number of memory blocks to be displaced!

Although no pure COMA machines were ever built, the COMA idea has inspired many variations. COMA-F, or FLAT COMA was proposed by Stenström, Joe, and Gupta in 1992 as a simpler alternative to the original COMA proposals. By allocating a home location COMA-F eliminated the need for multilevel hierarchical look-ups and possible displacement misses, since the block status could always be looked up in the home and the home location always had space for the block. In 1995, Saulsbury et al. proposed Simple COMA (S-COMA), which implemented COMA using the virtual memory mechanisms for replication and migration, rather than hardware support at the cache-level. Reactive NUMA [Falsafi and Wood 1997] is a proposal to develop a protocol that merges the best of CC-NUMA protocols with S-COMA protocols. At the same time, several groups (see Chandra et al. 1994 and Soundararajan 1996) explored the use of page-level replication and migration, both to assist in reducing remote misses and as an alternative to other schemes such as strict COMA or remote access caches. Wildfire

builds on many of these ideas to create a blend of hardware and software mechanisms.

The Convex Exemplar implemented scalable coherent shared memory using a two-level architecture: at the lowest level eight-processor modules are built using a crossbar. A ring can then connect up to 32 of these modules, for a total of 256 processors (see Thekkath et. al. [1997] for an evaluation). Lenoski and Laudon [1997] describe the SGI Origin, which was first delivered in 1996 and is closely based on the original Stanford DASH machine, though including a number of innovations for scalability and ease of programming. Origin uses a bit-vector for the directory structure, which is either 16 or 32 bits long. Each bit represents a node, which consists of two processors; a coarse bit vector representation allows each bit to represent up to 8 nodes for a total of 1,024 processors. As Galles [1996] describes, a high performance fat hypercube is used for the global interconnect. Hristea et. al [1997] is a thorough evaluation of the performance of the Origin memory system.

More recent research has focused on enhanced scalability for cache-coherent designs, flexible and adaptable techniques for implementing coherency, and approaches that merge hardware and software schemes. The MIT Alewife machine [Agarwal et. al. 1995] incorporated several innovations including processor support for multithreading and the use of cooperative mechanisms for handling coherence. The Stanford FLASH multiprocessor [Kuskin et. al. 1994, Gibson et. al. 2000] makes use of a programmable processor that implements the coherence scheme, as well as alternative schemes for message-passing, synchronization primitives, or performance instrumentation. Reinhardt and his colleagues at the University of Wisconsin [1994] explored an alternative for a combination of user and-base software and hardware support for coherent shared-memory. The Star-T [Nikhil et. al 1992] and Star-T Voyager [Ang, et. al. 1998] projects at MIT explored the use of multithreading and combining customized and commodity approaches to building scalable multiprocessors.

Developments in Synchronization and Consistency Models

A wide variety of synchronization primitives have been proposed for shared-memory multiprocessors. Mellor-Crummey and Scott [1991] provide an overview of the issues as well as efficient implementations of important primitives, such as locks and barriers. An extensive bibliography supplies references to other important contributions, including developments in spin locks, queuing locks, and barriers.

Lamport [1979] introduced the concept of sequential consistency and what correct execution of parallel programs means. Dubois, Scheurich, and Briggs [1988] introduced the idea of weak ordering (originally in 1986). In 1990, Adve and Hill provided a better definition of weak ordering and also defined the concept of data-race-free; at the same conference, Gharachorloo [1990] and his colleagues introduced release consistency and provided the first data on the

performance of relaxed consistency models. More relaxed consistency models have been widely adopted in microprocessor architectures, including the Sun SPARC, Alpha, and IA-64. Adve and Gharachorloo [1996] is an excellent tutorial on memory consistency and the differences among these models.

Other References

The concept of using virtual memory to implement a shared address space among distinct machines was pioneered in Kai Li's Ivy system in 1988. There have been subsequent papers exploring both hardware support issues, software mechanisms, and programming issues. Amza et. al. [1996] describe a system built on workstations using a new consistency model, L. Kontothanassis, et. al. [1997] describe a software shared memory scheme using remote writes, and Erlichson et. al. [1996] describe the use of shared virtual memory to build large-scale multiprocessors using SMPs as nodes.

There is an almost unbounded amount of information on multiprocessors and multicomputers: Conferences, journal papers, and even books seem to appear faster than any single person can absorb the ideas. No doubt many of these papers will go unnoticed—not unlike the past. Most of the major architecture conferences contain papers on multiprocessors. An annual conference, *Supercomputing XY* (where X and Y are the last two digits of the year), brings together users, architects, software developers, and vendors and publishes the proceedings in book, CD-ROM, and online (see www.scXY.org) form. Two major journals, *Journal of Parallel and Distributed Computing* and the *IEEE Transactions on Parallel and Distributed Systems*, contain papers on all aspects of parallel processing. Several books focusing on parallel processing are included in the following references with Culler, Singh, and Gupta [1999] being the most recent, large-scale effort. For years, Eugene Miya of NASA Ames has collected an online bibliography of parallel-processing papers. The bibliography, which now contains that contains more than 35,000 entries, is available online at:

<http://liinwww.ira.uka.de/bibliography/Parallel/Eugene/index.html>.

In addition to documenting the discovery of concepts now used in practice, these references also provide descriptions of many ideas that have been explored and found wanting, as well as ideas whose time has just not yet come.

Multithreading and Simultaneous Multithreading

The concept of multithreading dates back to one of the earliest transistorized computers, the TX-2. TX-2 was one of the earliest transistorized computers and is also famous for being the computer on which Ivan Sutherland created Sketchpad, the first computer graphics system. TX-2 was built at MIT's Lincoln Laboratory and became operational in 1959. It used multiple threads to support fast context switching to handle I/O functions. Clark [1957] describes the basic architecture and Forgie [1957] describes the I/O architecture. Multithreading was also

used in the CDC 6600, where a fine-grained multithreading scheme with interleaved scheduling among threads was used as the architecture of the I/O processors. The HEP processor, a pipelined multiprocessor, designed by Denelcor. and shipped in 1982 used fine-grained multithreading to hide the pipeline latency as well as to hide the latency to a large memory shared among all the processors. Because the HEP had no cache, this hiding of memory latency was critical. Burton Smith, one the primary architects, describes the HEP architecture in a 1978 paper and Jordan [1983] published a performance evaluation. The Tera processor extends the multithreading ideas and is described by Alverson et. al. in a 1992 paper.

In the late 1980s and early 1990s, researchers explored the concept of coarse-grained (also called block multithreading), as a way to tolerate latency, especially in multiprocessor environments. The SPARCLE processor in the Alewife system used such a scheme, switching threads whenever a high latency exceptional event, such as a long cache miss, occurred. Agarwal et. al. describe SPARCLE in a 1993 paper. The IBM Pulsar processor uses similar ideas.

By the early 1990s, several research groups had arrived at two key insights. First, they realized that fine-grained multithreading was needed to get the maximum performance benefit, since in a coarse-grained approach, the overhead of thread switching and thread start-up (e.g., filling the pipeline from the new thread) negated much of the performance advantage (see Laudon et. al. 1994). Second, several groups realized that to effectively use large numbers of functional units would require both ILP and thread-level parallelism (TLP). These insights led to several architectures that used combinations of multithreading and multiple issue. Wolfe and Shen [1991] describe an architecture called XIMD that statically interleaves threads scheduled for a VLIW processor. Mirata et. al. (1992) describe a proposed processor for media use that combines a static superscalar pipeline with support for multithreading; they report speed-ups from combining both forms of parallelism. Keckler and Dally [1992] combine static scheduling of ILP and dynamic scheduling of threads combining the two forms for a processor with multiple functional units. The question of how to balance the allocation of functional units between ILP and TLP and how to schedule the two forms of parallelism remained open.

When it became clear in the middle of the 1990s that dynamically-scheduled superscalars would be delivered shortly, several research groups proposed using the dynamic scheduling capability to mix instructions from several threads on the fly. Yamamoto, Searing, Talcott, Wood, and Nemirosky [1994] appears to be the first such proposal, though the simulation results for their multithreaded superscalar architecture use simplistic assumptions. This work was quickly followed by Tullsen, Eggers, and Levy [1995], which was the first realistic simulation assessment and coined the name simultaneous multithreading. Subsequent work by the same group together with industrial coauthors addressed many of the open questions about SMT. For example, Tullsen et. al. [1996] addressed questions about the challenges of scheduling ILP vs. TLP. Lo et. al. [1997] is an extensive

discussion of the SMT concept and an evaluation of its performance potential, and Lo et. al. [1998] evaluates database performance on an SMT processor.

References

- A. AGARWAL, A., KUBIATOWICZ, J., KRANZ, D., LIM, B.-H., YEUNG, D., D'SOUZA, G. AND M. PARKIN [1993], "Sparcle: An evolutionary processor design for large-scale multiprocessors," IEEE Micro 13 (June), pp. 48--61.
- ALVERSON, G. ALVERSON, R., CALLAHAN, D. , KOBLLENZ, B., PORTERFIELD, A. AND B. SMITH [1992]. "Exploiting heterogeneous parallelism on a multithreaded multiprocessor," Proc. 1992 International Conf. on Supercomputing (November) , pp. 188--197.
- ADVE, S. V. AND K. GHARACHORLOO [1996]. "Shared Memory Consistency Models: A Tutorial," IEEE Computer 29:12 (December), 66--76.
- ADVE, S. V. AND M. D. HILL [1990]. "Weak ordering—A new definition," Proc. 17th Int'l Symposium on Computer Architecture (June), Seattle, 2--14.
- AGARWAL, A., BIANCHINI, R., CHAIKEN, D., JOHNSON, K., AND D. KRANZ [1995]. "THE MIT ALE-WIFE MACHINE: ARCHITECTURE AND PERFORMANCE", INTERNATIONAL SYMPOSIUM ON COMPUTER ARCHITECTURE, DENVER, JUNE, 2--13.
- AGARWAL, A., J. L. HENNESSY, R. SIMONI, AND M.A. HOROWITZ [1988]. "An evaluation of directory schemes for cache coherence," Proc. 15th Int'l Symposium on Computer Architecture (June), 280--289.
- ALMASI, G. S. AND A. GOTTLIEB [1989]. *Highly Parallel Computing*, Benjamin/Cummings, Redwood City, Calif.
- AMDAHL, G. M. [1967]. "Validity of the single processor approach to achieving large scale computing capabilities," Proc. AFIPS Spring Joint Computer Conf. 30, Atlantic City, N.J. (April), 483--485.
- AMZA C., COX, A. L., DWARKADAS, S., KELEHER, P., LU, H., RAJAMONY, R., YU, W. AND W. ZWAENEPOEL. [1996]. "TREADMARKS: SHARED MEMORY COMPUTING ON NETWORKS OF WORKSTATIONS". IEEE COMPUTER, 29(2) (FEBRUARY), 18--28.
- ANG, B., CHIOU, D., ROSENBAUM, D., EHRLICH, M., AND RUDOLPH, L., AND ARVIND [1998]. "START-VOYAGER: A FLEXIBLE PLATFORM FOR EXPLORING SCALABLE SMP ISSUES", PROCEEDINGS OF SC'98, ORLANDO, FLORIDA, NOV.
- ARCHIBALD, J. AND J.-L. BAER [1986]. "Cache coherence protocols: Evaluation using a multiprocessor simulation model," ACM Trans. on Computer Systems 4:4 (November), 273--298.
- ARPACI, R.H., CULLER, D.E., KRISHNAMURTHY, A., STEINBERG, S.G. AND K. YELICK [1995]. "Empirical evaluation of the CRAY-T3D: A compiler perspective," Proceedings of the International Symposium on Computer Architecture, Denver (June), pages 320-331.
- BAER J.-L. AND W.-H. WANG [1988]. "On the Inclusion Properties for Multi-Level Cache Hierarchies." In Proceedings of the 15th Annual International Symposium on Computer Architecture, Honolulu, June, 73--80.
- BARROSO, L.A., GHARACHORLOO, K. AND E. BUGNION [1998]. "Memory System Characterization of Commercial Workloads," Proceedings 25th International Symposium on Computer Architecture, Barcelona (July), 3-14.
- BASKETT, F., T. JERMOLUK, AND D. SOLOMON [1988]. "The 4D-MP graphics superworkstation: Computing + graphics = 40 MIPS + 40 MFLOPS and 10,000 lighted polygons per second," Proc. COMPCON Spring, San Francisco, 468--471.
- BBN LABORATORIES [1986]. "Butterfly parallel processor overview," Tech. Rep. 6148, BBN Labo-

- ratories, Cambridge, Mass.
- BELL, C. G. [1985]. "Multis: A new class of multiprocessor computers," *Science* 228 (April 26), 462–467.
- BELL, C. G. [1989]. "The future of high performance computers in science and engineering," *Comm. ACM* 32:9 (September), 1091–1101.
- BOUKNIGHT, W. J. S. A. DENEGER, D. E. MCINTYRE, J. M. RANDALL, A. H. SAMEH, AND D. L. SLOTNICK [1972]. "The Illiac IV system," *Proc. IEEE* 60:4, 369–379. Also appears in D. P. Siewiorek, C. G. Bell, and A. Newell, *Computer Structures: Principles and Examples*, McGraw-Hill, New York (1982), 306–316.
- BURKHARDT, H. III, S. FRANK, B. KNOBE, AND J. ROTHNIE [1992]. "Overview of the KSR1 computer system," Tech. Rep. KSR-TR-9202001, Kendall Square Research, Boston (February).
- CENSIER, L. AND P. FEAUTRIER [1978]. "A new solution to coherence problems in multicache systems," *IEEE Trans. on Computers* C-27:12 (December), 1112–1118.
- CHANDRA, R., DEVINE, S., VERGHESE, B., GUPTA, A. AND MENDEL ROSENBLUM [1994]. "Scheduling and Page Migration for Multiprocessor Compute Servers." In Sixth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-VI). ACM, Santa Clara, CA. October, 12–24. .
- CHARLESWORTH, A [1998]. "STARFIRE: EXTENDING THE SMP ENVELOPE," *IEEE MICRO* 18:1 (JAN/FEB), P 39–49.
- CLARK, W.A. [1957]. "The Lincoln TX-2 Computer Development." Proceedings of the Western Joint Computer Conference (February), Institute of Radio Engineers, Los Angeles, 143–145.
- CULLER, D. E., SINGH, J. P., AND A. GUPTA [1999]. *Parallel Computer Architecture A Hardware/Software Approach*. Morgan Kaufmann Publishers,
- 1 EDITION, 1999.
- DUBOIS, M., C. SCHEURICH, AND F. BRIGGS [1988]. "Synchronization, coherence, and event ordering," *IEEE Computer* 9-21 (February).
- EGGERS, S. [1989]. *Simulation Analysis of Data Sharing in Shared Memory Multiprocessors*, Ph.D. Thesis, Univ. of California, Berkeley. Computer Science Division Tech. Rep. UCB/CSD 89/501 (April).
- ELDER, J., A. GOTTLIEB, C. K. KRUSKAL, K. P. MCAULIFFE, L. RANDOLPH, M. SNIR, P. TELLER, AND J. WILSON [1985]. "Issues related to MIMD shared-memory computers: The NYU Ultracomputer approach," *Proc. 12th Int'l Symposium on Computer Architecture* (June), Boston, 126–135.
- ERLICHSON, A., NUCKOLLS, N., CHESSON, G. AND J. L. HENNESSY [1996]. "SoftFLASH: Analyzing the performance of clustered distributed virtual shared memory." In Proc. of the 7th Symp.on Architectural Support for Programming Languages and Operating Systems (ASPLOS-VII), pages 210–220, October.
- FLYNN, M. J. [1966]. "Very high-speed computing systems," *Proc. IEEE* 54:12 (December), 1901–1909.
- FALSAFI, B. AND WOOD, D.A. [1997]. "Reactive NUMA: a design for unifying S-COMA and CC-NUMA," Proceedings of the 24th international symposium on Computer architecture, June, Denver, CO, 229–240.
- FORGIE, J.W [1957]. "The Lincoln TX-2 Input-Output System," Proceedings of the Western Joint Computer Conference (February), Institute of Radio Engineers, Los Angeles, 156–160.
- FRANK, S. J. [1984] "TIGHTLY COUPLED MULTIPROCESSOR SYSTEMS SPEED MEMORY ACCESS TIME," *ELECTRONICS* 57:1 (JANUARY), 164–169.
- GALLES, M. [1996]. "Scalable Pipelined Interconnect for Distributed Endpoint Routing: The SGI SPIDER chip" . Proceedings Hot Interconnects '96, Stanford University, August.
- GAJSKI, D., D. KUCK, D. LAWRIE, AND A. SAMEH [1983]. "CEDAR—A large scale multiprocessor," *Proc. Int'l Conf. on Parallel Processing* (August), 524–529.
- GEHRINGER, E. F., D. P. SIEWIOREK, AND Z. SEGALL [1987]. *Parallel Processing: The Cm* Experience*, Digital Press, Bedford, Mass.

- GHARACHORLOO, K., GUPTA, A., AND J.L. HENNESSY [1992]. "Hiding memory latency using dynamic scheduling in shared-memory multiprocessors." In Proc. of the 19th Annual Int. Symp. on Computer Architecture, FGold Coast, Australia, June.
- GHARACHORLOO, K., D. LENOSKI, J. LAUDON, P. GIBBONS, A. GUPTA, AND J. L. HENNESSY [1990]. "Memory consistency and event ordering in scalable shared-memory multiprocessors," *Proc. 17th Int'l Symposium on Computer Architecture* (June), Seattle, 15–26.
- GIBSON, J, KUNZ, R, OFELT, D, HOROWITZ, M, HENNESSY, J, AND M. HEINRICH [2000]. "FLASH vs. (Simulated) FLASH: Closing the Simulation Loop". Proc. of the 9th Conference on Architectural Support for Programming Languages and Operating Systems (November), San Jose, 49–58.
2000. GOODMAN, J. R. [1983]. "Using cache memory to reduce processor memory traffic," *Proc. 10th Int'l Symposium on Computer Architecture* (June), Stockholm, Sweden, 124–131.
- HAGERSTEN E. AND M. KOSTER [1998]. "WILDFIRE: A SCALABLE PATH FOR SMPs," ROCEEDINGS OF THE THE FIFTH INTERNATIONAL SYMPOSIUM ON HIGH PERFORMANCE COMPUTER ARCHITECTURE, 1998.
- HAGERSTEN, E., LANDIN, A. AND S. HARIDI. DDM --- A Cache-Only Memory Architecture. *IEEE Computer*, 25(9):44–54, September, 1992.
- HILL, M.D. [1998]. "Multiprocessors should support simple memory consistency models," *IEEE Computer*, 31:8 (August), 28–34.
- HILLIS, W. D. [1985]. *The Connection Multiprocessor*, MIT Press, Cambridge, Mass.
- HIRATA, H., KIMURA, K., NAGAMINE, S., MOCHIZUKI, Y., NISHIMURA, A., NAKASE, Y., AND NISHIZAWA, T. [1992]. "An elementary processor architecture with simultaneous instruction issuing from multiple threads," Proc. 19th Annual International Symposium on Computer Architecture (May). 136–145.
- HOCKNEY, R. W. AND C. R. JESSHOPE [1988]. *Parallel Computers-2, Architectures, Programming and Algorithms*, Adam Hilger Ltd., Bristol, England.
- HOLLAND, J. H. [1959]. "A universal computer capable of executing an arbitrary number of subprograms simultaneously," *Proc. East Joint Computer Conf.* 16, 108–113.
- HORD, R. M. [1982]. *The Illiac-IV, The First Supercomputer*, Computer Science Press, Rockville, Md.
- HRISTEA, C., LENOSKI, D., AND J. KEEN [1997]. Measuring Memory Hierarchy Performance of Cache-Coherent Multiprocessors Using Micro Benchmarks, Proc. Supercomputing 97, San Jose, CA, November.
- HWANG, K. [1993]. *Advanced Computer Architecture and Parallel Programming*, McGraw-Hill, New York.
- KECKLER, S.W. AND DALLY, W. J. [1992]. "Processor coupling: Integrating compile time and runtime scheduling for parallelism," Proc. 19th Annual International Symposium on Computer Architecture (May). 202–213.
- KONTOTHANASSIS, L., HUNT, G., STETS, R., HARDAVELLAS, N., CIERNIAK, M., PARTHASARATHY, S., MEIRA, W., DWARKADAS, S. AND M. SCOTT [1997]. "VM-based shared memory on low-latency, remote-memory-access networks", . Proc., 24th Annual Int'l. Symp. on Computer Architecture, June, Denver.
- KUSKIN, J., OFELT, D., HEINRICH, M., HEINLEIN, J., SIMONI, R., GHARACHORLOO, K., CHAPIN, J., NAKAHIRA, D., BAXTER, J., HOROWITZ, M., GUPTA, A., ROSENBLUM, M., AND J.L. HENNESSY [1994]. "The Stanford FLASH Multiprocessor", Proceedings of the 21th International Symposium on Computer Architecture, Chicago, April.
- LAMPORT, L. [1979]. "How to make a multiprocessor computer that correctly executes multiprocess programs," *IEEE Trans. on Computers* C-28:9 (September), 241–248.
- LAUDON, J., GUPTA, A., AND M. HOROWITZ [1994]. "Interleaving: A multithreading technique target-

- ing multiprocessors and work-stations.," Proc Sixth International Conference on Architectural Support for Programming Languages and Operating Systems (October), Boston, 308–318.
- LAUDON J. AND D. LENOSKI [1997]. "THE SGI ORIGIN: A CCNUMA HIGHLY SCALABLE SERVER," Proceedings of the 24th international symposium on Computer architecture , June, Denver, p 241–251
- LENOSKI, D., J. LAUDON, K. GHARACHORLOO, A. GUPTA, AND J. L. HENNESSY [1990]. "The Stanford DASH multiprocessor," *Proc. 17th Int'l Symposium on Computer Architecture* (June), Seattle, 148–159.
- LENOSKI, D., J. LAUDON, K. GHARACHORLOO, W.-D. WEBER, A. GUPTA, J. L. HENNESSY, M. A. HOROWITZ, AND M. LAM [1992]. "The Stanford DASH multiprocessor," *IEEE Computer* 25:3 (March).
- LI, K., [1988] "IVY: A Shared Virtual Memory System for Parallel Computing," Proceedings of the 1988 International Conference on Parallel Processing, Pennsylvania State University Press.
- LO, J., EGGERS, S., EMER, J., LEVY, H., STAMM, R., AND D. TULLSEN [1997]. "Converting Thread-Level Parallelism Into Instruction-Level Parallelism via Simultaneous Multithreading," *ACM Transactions on Computer Systems* 15:2 (August), 322–354.
- LO, J., BARROSO, L., EGGERS, S., GHARACHORLOO, K., LEVY, H., AND S. PAREKH [1998]. "An Analysis of Database Workload Performance on Simultaneous Multithreaded Processors," *Proceedings of the 25th International Symposium on Computer Architecture* (June), 39–50.
- LOVETT, T. AND S. THAKKAR [1988]. "The Symmetry multiprocessor system," *Proc. 1988 Int'l Conf. of Parallel Processing*, University Park, Penn., 303–310.
- MELLOR-CRUMMEY, J. M. AND M. L. SCOTT [1991]. "Algorithms for scalable synchronization on shared-memory multiprocessors," *ACM Trans. on Computer Systems* 9:1 (February), 21–65.
- MENABREA, L. F. [1842]. "Sketch of the analytical engine invented by Charles Babbage," *Bibliothèque Universelle de Genève* (October).
- MITCHELL, D. [1989]. "The Transputer: The time is now," *Computer Design* (RISC supplement), 40–41.
- MIYA, E. N. [1985]. "Multiprocessor/distributed processing bibliography," *Computer Architecture News* (ACM SIGARCH) 13:1, 27–29.
- NIKHIL, R.S., PAPADOPOULOS, G.M. AND ARVIND [1992]. "T: A Multithreaded Massively Parallel Architecture." In Proceedings of the 19th International Symposium on Computer Architecture, Gold Coast, Australia, May, 156–167.
- NOORDERGRAAF, L. AND R. VAN DER PAS [1999]. "Performance Experiences on Sun's WildFire Prototype," *Proc. Supercomputing 99*, Portland, Oregon, November.
- PFISTER, G. F., W. C. BRANTLEY, D. A. GEORGE, S. L. HARVEY, W. J. KLEINFEDER, K. P. MCAULIFFE, E. A. MELTON, V. A. NORTON, AND J. WEISS [1985]. "The IBM research parallel processor prototype (RP3): Introduction and architecture," *Proc. 12th Int'l Symposium on Computer Architecture* (June), Boston, 764–771.
- REINHARDT, S.K., LARUS, J.R., AND D. A. WOOD [1994]. "Tempest and Typhoon: User-Level Shared Memory." In Proceedings of the 21st Annual International Symposium on Computer Architecture, . Chicago, April, 325–336.
- RETTBERG, R. D., W. R. CROWTHER, P. P. CARVEY, AND R. S. TOWLINSON [1990]. "The Monarch parallel processor hardware design," *IEEE Computer* 23:4 (April).
- ROSENBLUM, M., S. A. HERROD, E. WITCHEL, AND A. GUTPA [1995]. "Complete computer simulation: The SimOS approach," to appear in *IEEE Parallel and Distributed Technology* 3:4 (fall).
- SAULSBURY, A., WILKINSON, T., CARTER, J. AND A. LANDIN [1995]. "An Argument for Simple CO-MA," *Proc. First Conf. on High Performance Computer Architectures* (January), Raleigh, N. Carolina,, 276–285

- SCHWARTZ, J. T. [1980]. "Ultracomputers," *ACM Trans. on Programming Languages and Systems* 4:2, 484–521.
- SCOTT S. L. [1996] "SYNCHRONIZATION AND COMMUNICATION IN THE T3E MULTIPROCESSOR," "Proceeding Architectural Support for Programming Languages and Operating Systems (ASPLOS-VII), Cambridge, Massachusetts, October, pp. 26--36.
- SCOTT S. L. AND G. M. THORSON. "The Cray T3E Network: Adaptive Routing in a High Performance 3D Torus," In Proceedings of the Symposium on High Performance Interconnects (Hot Interconnects 4), Stanford University, August, pages 14-156.
- SEITZ, C. [1985]. "The Cosmic Cube," *Comm. ACM* 28:1 (January), 22–31.
- SINGH, J. P., HENNESSY, J. L. AND A. GUPTA., "Scaling Parallel Programs for Multiprocessors: Methodology and Examples," *Computer* 26: 7 (July), 22–33.
- SLOTNICK, D. L., W. C. BORCK, AND R. C. McREYNOLDS [1962]. "The Solomon computer," *Proc. Fall Joint Computer Conf.* (December), Philadelphia, 97–107.
- SMITH, B.J. [1978] "A pipelined, shared resource MIMD computer," *Proc. 1978 ICPP* (August) pp. 6-8.
- SOUNDARARAJAN, V., HEINRICH, M., VERGHESE, B., GHARACHORLOO, K., GUPTA, A., AND J.L. HENNESSY [1998]. "FLEXIBLE USE OF MEMORY FOR REPLICATION/MIGRATION IN CACHE-COHERENT DSM MULTIPROCESSORS," . *Proc. 25th Int'l Symposium on Computer Architecture* (June), Barcelona, Spain, 342-355.
- STENSTRÖM, P., JOE, T. AND A. GUPTA [1992]. "Comparative performance evaluation of cache-coherent NUMA and COMA architectures." Proceedings of the 19th annual international symposium on Computer architecture, May, Queensland Australia, 80-91.
- STONE, H. [1991]. *High Performance Computers*, Addison-Wesley, New York.
- SWAN, R. J., A. BECHTOLSHEIM, K. W. LAI, AND J. K. OUSTERHOUT [1977]. "The implementation of the Cm* multi-microprocessor," *Proc. AFIPS National Computing Conf.*, 645–654.
- SWAN, R. J., S. H. FULLER, AND D. P. SIEWIOREK [1977]. "Cm*—A modular, multi-microprocessor," *Proc. AFIPS National Computer Conf.* 46, 637–644.
- TANG, C. K. [1976]. "Cache design in the tightly coupled multiprocessor system," *Proc. AFIPS National Computer Conf.*, New York (June), 749–753.
- THEKKATH, R. SINGH, A.P. SINGH, J.P., JOHN, S. AND J.L. HENNESSY [1997]. "An Evaluation of a Commercial CC-NUMA Architecture---The CONVEX Exemplar SPP1200," Proceedings of the 11th International Parallel Processing Symposium (IPPS '97), Geneva, Switzerland, April.
- TULLSEN, D.M., EGGERS, S.J., EMER, J.S., LEVY, H.M., LO, J.L. AND R.L. STAMM [1996]. "Exploiting choice: Instruction fetch and issue on an implementable simultaneous multithreading processor." Proceedings of the 23rd Annual International Symposium on Computer Architecture (May), pages 191--202.
- TULLSEN, D.M., EGGERS, S.J., AND H.M. LEVY [1995], "Simultaneous multithreading: Maximizing on-chip parallelism," *Proc. 22nd International Symposium on Computer Architecture* (June), pp.392-403.
- UNGER, S. H. [1958]. "A computer oriented towards spatial problems," *Proc. Institute of Radio Engineers* 46:10 (October), 1744–1750.
- WILSON, A. W., JR. [1987]. "Hierarchical cache/bus architecture for shared-memory multiprocessors," *Proc. 14th Int'l Symposium on Computer Architecture* (June), Pittsburgh, 244–252.
- WOOD, D. A. AND M. D. HILL [1995]. "Cost-effective parallel computing," *IEEE Computer* 28:2 (February).
- WOLFE, A. AND J. P. SHEN [1991]. "A variable instruction stream extension to the VLIW architecture." *Proc. of the Fourth Conference on Architectural Support for Programming Languages and*

Operating Systems (April), Santa Clara, 2-14.

WULF, W. AND C. G. BELL [1972]. "C.mmp—A multi-mini-processor," *Proc. AFIPS Fall Joint Computing Conf.* 41, part 2, 765-777.

WULF, W. AND S. P. HARBISON [1978]. "Reflections in a pool of processors—An experience report on C.mmp/Hydra," *Proc. AFIPS 1978 National Computing Conf.* 48 (June), Anaheim, Calif., 939-951.

Yamamoto, W., Serrano, M.J., Talcott, A.R., Wood, R.C., and M. Nemirosky [1992]. "Performance estimation of multistreamed, superscalar processors," *Proc. Twenty-Seventh Hawaii International Conference on System Sciences* (January), pages I:195-204.

EXERCISES

6.1 [10] <6.1> Suppose we have an application that runs in three modes: all processors used, half the processors in use, and serial mode. Assume that 0.02% of the time is serial mode, and there are 100 processors in total. Find the maximum time that can be spent in the mode when half the processors are used, if our goal is a speedup of 80.

6.2 [15] <6.1> Assume that we have a function for an application of the form $F(i,p)$, which gives the fraction of time that exactly i processors are usable given that a total of p processors are available. This means that

$$\sum_{i=1}^p F(i,p) = 1$$

Assume that when i processors are in use, the application runs i times faster. Rewrite Amdahl's Law so that it gives the speedup as a function of p for some application.

6.3 [10] <6.1, 6.2> The Transaction Processing Council (TPC) has several different benchmarks. Visit their website at www.tpc.org and look at the top 10 performers in each benchmark class. Determine whether each of the top 10 configurations is a multiprocessor or if so what types (SMP, NUMA, cluster, e.g.). Does the ordering look different if price-performance is used as the metric?

6.4 [10] <6.1, 6.2> The Top 500 list categorizes the fastest scientific machines in the world according to their performance on the Linpack benchmark. Visit their website at www.top500.org and look at the top 100 performers (there are many repeats of a particular vendor product, since individual supercomputer sites rather than a product are counted). Determine how many different supercomputer products occur among the top 100 configurations and what type (SMP, NUMA, cluster, e.g.) each different supercomputer is. Try to obtain cost information and see how the data changes when cost-performance is considered.

6.5 [15] <6.3> In small bus-based multiprocessors, write-through caches are sometimes used. One reason is that a write-through cache has a slightly simpler coherence protocol. Show how the basic snooping cache coherence protocol of Figure 6.12 on page 668 can be changed for a write-through cache. From the viewpoint of an implementor, what is the major hardware functionality that is not needed with a write-through cache compared with a

write-back cache?

6.6 [20] <6.3> Add a clean private state to the basic snooping cache-coherence protocol (Figure 6.12 on page 668). Show the protocol in the format of Figure 6.12.

6.7 [15] <6.3> One proposed solution for the problem of false sharing is to add a valid bit per word (or even for each byte). This would allow the protocol to invalidate a word without removing the entire block, allowing a cache to keep a portion of a block in its cache while another processor wrote a different portion of the block. What extra complications are introduced into the basic snooping cache coherency protocol (Figure 6.12) if this capability is included? Remember to consider all possible protocol actions.

6.8 [12/10/15] <6.3> The performance differences for write invalidate and write update schemes can arise from both bandwidth consumption and latency. Assume a memory system with 64-byte cache blocks. Ignore the effects of contention.

- a. [12] <6.3> Write two parallel code sequences to illustrate the bandwidth differences between invalidate and update schemes. One sequence should make update look much better and the other should make invalidate look much better.
- b. [10] <6.3> Write a parallel code sequence to illustrate the latency advantage of an update scheme versus an invalidate scheme.
- c. [15] <6.3> Show, by example, that when contention is included, the latency of update may actually be worse. Assume a bus-based multiprocessor with 50-cycle memory and snoop transactions.

6.9 Use the data on miss rates versus block size for the scientific applications in Section 6.3 to compute AMAT and bus bandwidth making some assumptions about memory access time based on block size.

6.10 [15/15] <6.3–6.5> Restructure this exercise to use timing from E6000 series.

One possible approach to achieving the scalability of distributed shared memory and the cost-effectiveness of a bus design is to combine the two approaches, using a set of processors with memories attached directly to the processors, and interconnected with a bus. The argument in favor of such a design is that the use of local memories and a coherence scheme with limited broadcast results in a reduction in bus traffic, allowing the bus to be used for a larger number of processors. For these Exercises, assume the same parameters as for the Challenge bus. Assume that remote snoops and memory accesses take the same number of cycles as a memory access on the Challenge bus. Ignore the directory processing time for these Exercises. Assume that the coherency scheme works as follows on a miss: If the data are up-to-date in the local memory, it is used there. Otherwise, the bus is used to snoop for the data. Assume that local misses take 25 bus clocks.

- a. [15] <6.3–6.5> Find the time for a read or write miss to data that are remote.
- b. [15] <6.3–6.5> Ignoring contention and using the data from the Ocean benchmark run on 16 processors for the frequency of local and remote misses (Figure 6.31 on page 699), estimate the average memory access time versus that for a Challenge using the same total miss rate.

6.11 [12/15] <6.3,6.5,6.11> Restructure this exercise using the data comparing Origin to

E6000.

Although it is widely believed that buses are the ideal interconnect for small-scale multiprocessors, this may not always be the case. For example, increases in processor performance are lowering the processor count at which a more distributed implementation becomes attractive. Because a standard bus-based implementation uses the bus both for access to memory and for interprocessor coherency traffic, it has a uniform memory access time for both. In comparison, a distributed memory implementation may sacrifice on remote memory access, but it can have a much better local memory access time.

Consider the design of a DSM multiprocessor with 16 processors. Assume the R4400 cache miss overheads shown for the Challenge design (see pages 730–731). Assume that a memory access takes 150 ns from the time the address is available from either the local processor or a remote processor until the first word is delivered.

- a. [12] <6.3,6.5,6.11> How much faster is a local access than on the Challenge?
- b. [15] <6.3,6.5,6.11> Assume that the interconnect is a 2D grid with links that are 16 bits wide and clocked at 100 MHz, with a start-up time of five cycles for a message. Assume one clock cycle between nodes in the network, and ignore overhead in the messages and contention (i.e., assume that the network bandwidth is not the limit). Find the average remote memory access time, assuming a uniform distribution of remote requests. How does this compare to the Challenge case? What is the largest fraction of remote misses for which the DSM multiprocessor will have a lower average memory access time than that of the Challenge multiprocessor?

6.12 [20/15/30] <6.5> One downside of a straightforward implementation of directories using fully populated bit vectors is that the total size of the directory information scales as the product: Processor count \times Memory blocks. If memory is grown linearly with processor count, then the total size of the directory grows quadratically in the processor count. In practice, because the directory needs only 1 bit per memory block (which is typically 32 to 128 bytes), this problem is not serious for small to moderate processor counts. For example, assuming a 128-byte block, the amount of directory storage compared to main memory is Processor count/1024, or about 10% additional storage with 100 processors. This problem can be avoided by observing that we only need to keep an amount of information that is proportional to the cache size of each processor. We explore some solutions in these Exercises.

- a. [20] <6.5> One method to obtain a scalable directory protocol is to organize the multiprocessor as a logical hierarchy with the processors at the leaves of the hierarchy and directories positioned at the root of each subtree. The directory at each subtree root records which descendants cache which memory blocks, as well as which memory blocks with a home in that subtree are cached outside of the subtree. Compute the amount of storage needed to record the processor information for the directories, assuming that each directory is fully associative. Your answer should incorporate both the number of nodes at each level of the hierarchy as well as the total number of nodes.
- b. [15] <6.5> Assume that each level of the hierarchy in part (a) has a lookup cost of 50 cycles plus a cost to access the data or cache of 50 cycles, when the point is reached. We want to compute the AMAT (average memory access time—see Chapter 5) for a 64-processor multiprocessor with four-node subtrees. Use the data from the Ocean benchmark run on 64 processors (Figure 6.31) and assume that all noncoherence miss-

es occur within a subtree node and that coherence misses are uniformly distributed across the multiprocessor. Find the AMAT for this multiprocessor. What does this say about hierarchies?

- c. [30] <6.5> An alternative approach to implementing directory schemes is to implement bit vectors that are not dense. There are two such strategies: one reduces the number of bit vectors needed and the other reduces the number of bits per vector. Using traces, you can compare these schemes. First, implement the directory as a four-way set-associative cache storing full bit vectors, but only for the blocks that are cached outside of the home node. If a directory cache miss occurs, choose a directory entry and invalidate the entry. Second, implement the directory so that every entry has 8 bits. If a block is cached in only one node outside of its home, this field contains the node number. If the block is cached in more than one node outside its home, this field is a bit vector with each bit indicating a group of eight processors, at least one of which caches the block. Using traces of 64-processor execution, simulate the behavior of these two schemes. Assume a perfect cache for nonshared references, so as to focus on coherency behavior. Determine the number of extraneous invalidations as the directory cache size is increased.

6.13 [25/40] <6.10> Prefetching and relaxed consistency models are two methods of tolerating the latency of longer access in multiprocessors. Another scheme, originally used in the HEP multiprocessor and incorporated in the MIT Alewife multiprocessor, is to switch to another activity when a long-latency event occurs. This idea, called *multiple context* or *multithreading*, works as follows:

- n The processor has several register files and maintains several PCs (and related program states). Each register file and PC holds the program state for a separate parallel thread.
- n When a long-latency event occurs, such as a cache miss, the processor switches to another thread, executing instructions from that thread while the miss is being handled.
- a. [25] <6.10> Using the data for the Ocean benchmark running on 64 processors (Figure 6.31), determine how many contexts are needed to hide all the latency of remote accesses. Assume that local cache misses take 40 cycles and that remote misses take 120 cycles. Assume that the increased demands due to a higher request rate do not affect either the latency or the bandwidth of communications.
- b. [40] <6.10> Implement a simulator for a multiple-context directory-based multiprocessor. Use the simulator to evaluate the performance gains from multiple context. How significant are contention and the added bandwidth demands in limiting the gains?

6.14 [25] <6.10> Prove that in a two-level cache hierarchy, where L1 is closer to the processor, inclusion is maintained with no extra action if L2 has at least as much associativity as L1, both caches use LRU replacement, and both caches have the same block size.

6.15 [20] <6.5,6.11> As we saw in the *Putting it All Together* and in *Fallacies and Pitfalls*, data distribution can be important when an application has a nontrivial private data miss rate caused by capacity misses. This problem can be attacked with compiler technology (distributing the data in blocks) or through architectural support, as we saw in the descrip-

tion of CMR on Wildfire.

Assume that we have two DSM multiprocessors: one with CMR support and one without such support. Both multiprocessors have one processor per node and remote coherence misses, which are uniformly distributed, take 1 μ S. Assume that all capacity misses on the CMR multiprocessor hit in the local memory and require 250 ns. Assume that capacity misses take 200 ns when they are local on the DSM multiprocessor without CMR and 800 ns, otherwise. Using the Ocean data for 32 processors (Figure 6.23), find what fraction of the capacity misses on the DSM multiprocessor must be local if the performance of the two multiprocessors is identical.

6.16 [15] <6.7> Some multiprocessors have implemented a special broadcast coherence protocol just for locks, sometimes even using a different bus. Evaluate the performance of the spin lock in the Example on page 710 assuming a write broadcast protocol.

6.17 [15] <6.7> Implement the barrier in Figure 6.40 on page 713, using queuing locks. Compare the performance to the spin-lock barrier.

6.18 [15] <6.7> Implement the barrier in Figure 6.40 on page 713, using fetch-and-increment. Compare the performance to the spin-lock barrier.

6.19 [15] <6.7> Implement the barrier on page 717, so that barrier release is also done with a combining tree.

6.20 [30] <6.3–6.7,6.11> Using an available shared-memory multiprocessor, see if you can determine the organization and latencies of its memory hierarchy. For each level of the hierarchy, you can look at the total size, block size, and associativity, as well as the latency of each level of the hierarchy. If the multiprocessor uses a nonbus interconnection network, see if you can discover the topology and latency characteristics of the network. Try to make a table like that in Figure 6.47 for the machine. The Imbench (www.bitmover.com/Imbench/) and stream (<http://www.cs.virginia.edu/stream/>) benchmark may prove useful in this exercise.

6.21 [30] <6.3–6.7,6.11> Perform exercise 6.20 but looking at the bandwidth characteristics rather than latency. See if you can prepare a table like that in Figure 6.48. Extend the table by looking at the effect of strided accesses, as well as sequential and unrelated accesses.

6.22 [20] <6.5> As we discussed earlier, the directory controller can send invalidates for lines that have been replaced by the local cache controller. To avoid such messages, and to keep the directory consistent, replacement hints are used. Such messages tell the controller that a block has been replaced. Modify the directory coherence protocol of section 6.5 to use such replacement hints.

6.23 [15] <6.7> Find the time for n processes to synchronize using a standard barrier. Assume that the time for a single process to update the count and release the lock is c .

6.24 [15] <6.7> Find the time for n processes to synchronize using a combining tree barrier. Assume that the time for a single process to update the count and release the lock is c .

6.25 [25] <6.7> Implement a software version of the queuing lock for a bus-based system. Using the model in the Example on page 710, how long does it take for 20 processors to acquire and release the lock? You need only count bus cycles.

6.26 [20/30] <6.2–6.7> Both researchers and industry designers have explored the idea of having the capability to explicitly transfer data between memories. The argument in favor of such facilities is that the programmer can achieve better overlap of computation and communication by explicitly moving data when it is available. The first part of this exercise explores the potential on paper; the second explores the use of such facilities on real multiprocessors.

- a. [20] <6.2–6.7> Assume that cache misses stall the processor, and that block transfer occurs into the local memory of a DSM node. Assume that remote misses cost 100 cycles and that local misses cost 40 cycles. Assume that each DMA transfer has an overhead of 10 cycles. Assuming that all the coherence traffic can be replaced with DMA into main memory followed by a cache miss, find the potential improvement for Ocean running on 64 processors (Figure 6.31).
- b. [30] <6.2–6.7> Find a multiprocessor that implements both shared memory (coherent or incoherent) and a simple DMA facility. Implement a blocked matrix multiply using only shared memory and using the DMA facilities with shared memory. Is the latter faster? How much? What factors make the use of a block data transfer facility attractive?

6.27 [Discussion] <6.11> Construct a scenario whereby a truly revolutionary architecture—pick your favorite candidate—will play a significant role. *Significant* is defined as 10% of the computers sold, 10% of the users, 10% of the money spent on computers, or 10% of some other figure of merit.

6.28 [40] <6.2,6.10,6.14> A multiprocessor or cluster is typically marketed using programs that can scale performance linearly with the number of processors. The project here is to port programs written for one multiprocessor to the others and to measure their absolute performance and how it changes as you change the number of processors. What changes need to be made to improve performance of the ported programs on each multiprocessor? What is the ratio of processor performance according to each program?

6.29 [35] <6.2,6.10,6.14> Instead of trying to create fair benchmarks, invent programs that make one multiprocessor or cluster look terrible compared with the others, and also programs that always make one look better than the others. It would be an interesting result if you couldn't find a program that made one multiprocessor or cluster look worse than the others. What are the key performance characteristics of each organization?

6.30 [40] <6.2,6.10,6.14> Multiprocessors and cluster usually show performance increases as you increase the number of processors, with the ideal being n times speedup for n processors. The goal of this biased benchmark is to make a program that gets worse performance as you add processors. For example, this means that one processor on the multiprocessor or cluster runs the program fastest, two are slower, four are slower than two, and so on. What are the key performance characteristics for each organization that give inverse linear speedup?

6.31 [50] <6.2,6.10,6.14> Networked workstations can be considered multicomputers or clusters, albeit with somewhat slower, though perhaps cheaper, communication relative to computation. Port some cluster benchmarks to a network using remote procedure calls for communication. How well do the benchmarks scale on the network versus the cluster? What are the practical differences between networked workstations and a commercial clus-

ter, such as the IBM-SP series?

