

# The M-Machine Multicomputer

Marco Fillo,<sup>1</sup> Stephen W. Keckler,<sup>1</sup> William J. Dally,<sup>1</sup>  
Nicholas P. Carter,<sup>1</sup> Andrew Chang,<sup>1</sup> Yevgeny Gurevich,<sup>1</sup>  
and Whay S. Lee<sup>1</sup>

---

The M-Machine is an experimental multicomputer being developed to test architectural concepts motivated by the constraints of modern semiconductor technology and the demands of programming systems. The M-Machine computing nodes are connected with a 3-D mesh network; each node is a multithreaded processor incorporating 9 function units, on-chip cache, and local memory. The multiple function units are used to exploit both instruction-level and thread-level parallelism. A user accessible message passing system yields fast communication and synchronization between nodes. Rapid access to remote memory is provided transparently to the user with a combination of hardware and software mechanisms. This paper presents the architecture of the M-Machine and describes how its mechanisms attempt to maximize both single thread performance and overall system throughput. The architecture is complete and the MAP chip, which will serve as the M-Machine processing node, is currently being implemented.

---

**KEY WORDS:** Computer architecture; parallelism; multithreading; coherence.

## 1. INTRODUCTION

Because of the increasing density of VLSI integrated circuits, most of the chip area of modern computers is now occupied by memory and not by processing resources. Nearly all of this memory is located far from the processor, resulting in long latency and limited bandwidth access to it. It is clear that the computer systems of the future must address the latency and bandwidth limitations of these technology trends. The M-Machine

---

<sup>1</sup> Artificial Intelligence Laboratory, Laboratory for Computer Science, Massachusetts Institute of Technology, 545 Technology Square, Cambridge, Massachusetts 02139. E-mail: fillo{skeckler,billd,npccarter,achang,yev,wslee}@ai.mit.edu.

multicomputer is an experimental machine being developed to test architectural concepts motivated by the constraints of semiconductor technology and the demands of programming systems, such as faster execution of fixed sized problems and easier programmability of parallel computers.

The normalized area (in  $\lambda^2$ ) of a VLSI chip<sup>2</sup> increasing by 50% per year, while gate speed and communication bandwidth are increasing by 20% per year.<sup>(1)</sup> As a result, a 64-bit processor with a pipelined FPU ( $400 M\lambda^2$ ) is only 8% of a  $5G\lambda^2$  1996  $0.35 \mu\text{m}$  chip.<sup>3</sup> In a system with 256 MBytes of DRAM, the processor accounts for 0.13% of the silicon area in the system. The memory system, cache, TLB, controllers, and DRAM account for most of the remaining area. Technology scaling has made the memory, rather than the processor, the most area-consuming resource in a computer system.

The processor-memory imbalance in conventional computer organization exacerbates both the latency and bandwidth mismatch between execution units and main memory. Current systems use multiple levels of caches to manage the increasing size of main memory. While this does reduce latency to access that data contained in the cache closest to the processor, the latency to access uncached data increases dramatically. Furthermore, since chip packaging and offchip interconnect technologies are not advancing as fast as VLSI, the bandwidth between the processor and its memory is not scaling with processor speed or memory capacity. Both increasing latency and slow improvement of bandwidth to main memory will only continue to become more critical bottlenecks.

The M-Machine addresses both latency and bandwidth bottlenecks by increasing the number of processors per unit memory. This allows any given word of memory to be accessed by some processor very quickly and enables the aggregate bandwidth between processors and memory to scale with the size of the memory. In the M-Machine, instead of a multi-level cache hierarchy, the communication network determines the latency to access memory on another processor. Both multiple arithmetic units and fine-grained multithreading are implemented so that a processor may continue to perform useful work, and keep the memory pins busy, during long latency operations. In addition, the high arithmetic execution rate on-chip coupled with multithreading will keep the memory pins busy even during periods of high computation in which the data resides completely within the on-chip cache and the processor registers. That some of the on-chip

<sup>2</sup> The parameter  $\lambda$  is a normalized, process independent unit of length equivalent to one half of the gate length.<sup>(2)</sup> For a  $0.5 \mu\text{m}$  process,  $\lambda$  is  $0.25 \mu\text{m}$ .

<sup>3</sup> Area was determined by measuring the processing components of various chips, in particular the R4600 described in Ref. 3.

arithmetic resources go idle is inconsequential relative to the importance of using the memory bandwidth effectively.

An M-Machine multi-ALU processor (MAP) chip contains three 64-bit three-issue *clusters* that comprise 46% of the  $5G\lambda^2$  chip and 16% of an 8 MByte (six-chip) node. The high ratio of arithmetic bandwidth to memory bandwidth (9 operations/word) allows the MAP to saturate the costly DRAM bandwidth even on code with high cache-hit ratios. A 32-node M-Machine system with 256 Mbytes of memory has 96 times the peak performance of a 1996 uniprocessor with the same memory capacity at 1.5 times the area, a 64:1 improvement in peak performance/area. Even at a small fraction of this peak performance, such a machine allows the expensive memory bandwidth to handle more problems per unit time, resulting in more cost-effective computing.

The M-Machine is intended to extract more parallelism from problems of a fixed size than traditional computers, rather than requiring enormous problems to achieve peak performance. To do this, nodes are designed to manage parallelism at a range of granularities, from the instruction level to the process level. The 9 function units in an M-Machine node are controlled using a form of Processor Coupling<sup>(4)</sup> to exploit instruction level parallelism by executing 9 operations per cycle from the same thread, or to exploit thread-level parallelism by executing operations from up to six different threads. The fast internode communication allows collaborating threads to reside on different nodes.

The M-Machine also addresses the demand for easier programmability by providing an incremental path for increasing parallelism and performance. An unmodified sequential program can run on a single M-Machine node, accessing both local and remote memory. This code can be incrementally parallelized by identifying tasks, such as loop iterations, that can be distributed both across nodes and within each node to run in parallel. A flat, shared address space simplifies naming and communication. Software support for caching of remote data in local DRAM will automatically migrate a task's data to exploit locality.

Previous publications have introduced some of the mechanisms used in the M-Machine. The first description of Processor Coupling, a method for exploiting instruction level parallelism, appeared in Ref. 4. The novel capability-based memory protection system of the M-Machine was described in Ref. 5. This paper describes the M-Machine's other features which include an improved form of Processor Coupling as well as communication and global addressing mechanisms. The M-Machine architectural design is complete and the MAP chip, which will serve as the M-Machine processing node, is currently being implemented.

Section 2 gives an overview of the machine architecture, including the physical resources of the M-Machine. Section 3 describes the updated

version of Processor Coupling that simplifies hardware implementation and is expected to improve performance. Instead of lock-step execution of the wide instruction words across all of the function units, an instruction stream is partitioned by the compiler into horizontal threads (*H-Threads*), which run concurrently on different execution clusters to exploit instruction level parallelism. Several mechanisms for synchronizing the clusters are provided, including a cluster barrier instruction, broadcast of single bit condition values, and instructions that write register files in remote clusters. In addition, the function units are time-shared among vertical threads (*V-Threads*) which exploit runtime parallelism and mask pipeline, memory, and communication latencies. Events are handled asynchronously in a dedicated *V-Thread* so that event handling may proceed in parallel with user program execution and the issued instructions of the thread that caused the event need not be cancelled. Section 4 discusses inter-node communication including user-level communication primitives, global mapping of virtual addresses to physical memory and remote processors, and how these mechanisms are used to provide global coherent memory access. Finally, Section 5 describes the M-Machine software effort, including a brief overview of the compiler and runtime system.

## 2. M-MACHINE ARCHITECTURE

The M-Machine consists of a collection of computing nodes interconnected by a bidirectional 3-D mesh network. Each six-chip node consists of a multi-ALU (MAP) chip and 1 MW (8 Mbytes) of synchronous DRAM (SDRAM) with ECC. The MAP chip includes the network interface and router, and provides bandwidth of 800 Mbytes/s to the local SDRAMs and to each network channel. Each node contains a dedicated I/O bus; I/O devices may be connected to either every node or a subset of nodes, for example, all nodes on a face of the mesh. The target clock rate for the MAP is 100 MHz.

As shown in Fig. 1, a MAP contains three execution clusters, a unified cache which is divided into four banks, an external memory interface, and a communication subsystem consisting of the network interfaces and the router. Two crossbar switches interconnect these components. Clusters make memory requests to the appropriate bank of the interleaved cache over the 142-bit wide (51 address bits, 66 data bits, 25 control bits)  $3 \times 4$  M-Switch. The 88-bit wide (66 data bits, 22 control bits)  $9 \times 3$  C-Switch is used for inter-cluster communication and to return data from the memory system. Both switches support up to three transfers per cycle; each cluster may send and receive one transfer per cycle.

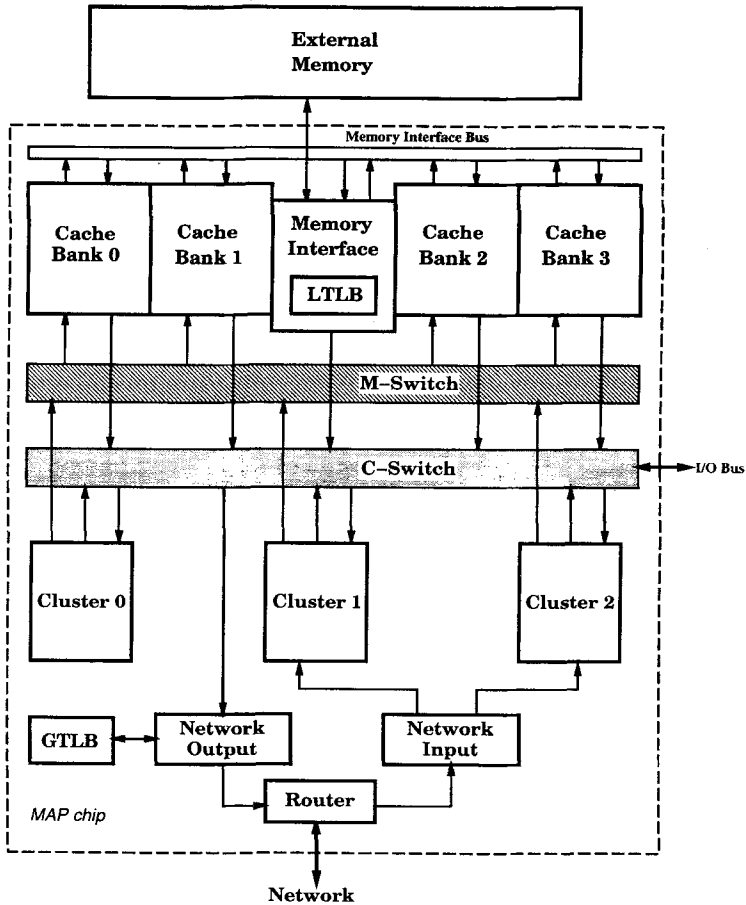


Fig. 1. The MAP architecture.

### 2.1. Map Execution Clusters

Each of the three MAP clusters is a 64-bit, three-issue, pipelined processor consisting of two integer ALUs, a floating-point ALU, associated register files, and a 1 KW (8 KB) instruction cache, as shown in Fig. 2. One of the integer ALUs in each cluster, termed the memory unit, is the interface to the memory system. Each MAP instruction contains 1, 2, or 3 operations. All operations in a single instruction issue together but may complete out of order. Every operation may be conditionally executed depending on the one-bit value of one of the condition code registers.

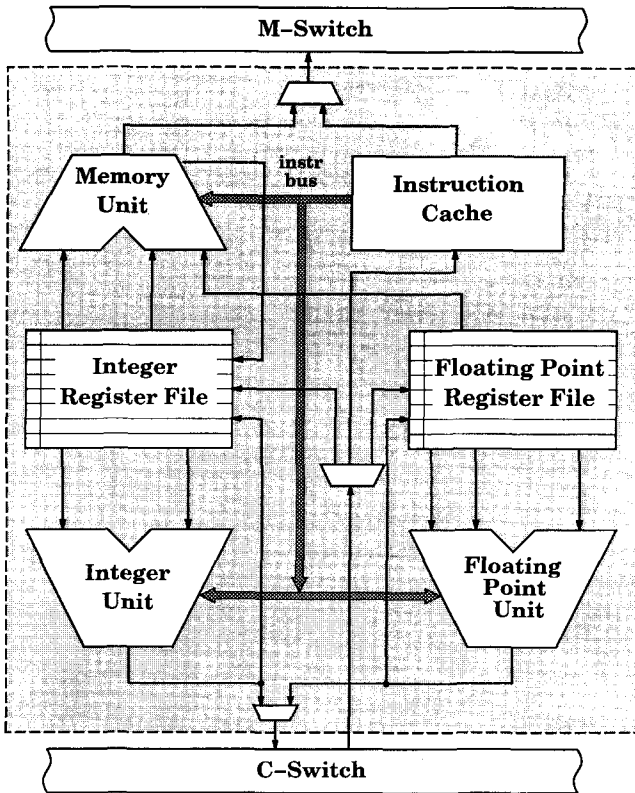


Fig. 2. A map cluster consists of 3 execution units, 2 register files, an instruction cache and ports onto the memory and cluster switches.

## 2.2. Memory System

As illustrated in Fig. 1, the 64 KB unified on-chip cache is organized as four 2 KW (16 KB) banks that are word-interleaved to permit accesses to consecutive addresses to proceed in parallel. The cache is virtually addressed and tagged. The cache banks are pipelined with a three-cycle read latency, including switch traversal. Each cluster has its own 8 KB instruction cache which fetches instructions from the unified cache when instruction cache misses occur.

The external memory interface consists of the SDRAM controller and a local translation lookaside buffer (LTLB) used to cache local page table (LPT) entries. Pages are 512 words (64 8-word cache blocks). The SDRAM controller exploits the pipeline and page modes of the external

SDRAM and performs single error correction and double error detection on the data transferred from external memory.

Each MAP word in memory is composed of a 64-bit data value, one synchronization bit and one pointer bit. A pair of special load and store operations specify a precondition and a postcondition on the synchronization bit and are used as atomic read-modify-write memory operations. The M-Machine supports a single global virtual address space. A light-weight capability system implements protection through the pointer bit and guarded pointers,<sup>(5)</sup> while paging is used to manage the relocation of data in physical memory within the virtual address space. The segmentation and paging mechanisms are independent so that protection may be preserved on variable-sized segments of memory. The memory subsystem is integrated with the communication system and can be used to access memory on remote nodes, as described in Section 4.2.

### 2.3. Communication Subsystem

Messages are composed in the general registers of a cluster and launched atomically using a user-level SEND instruction. To provide protection, messages must be sent to virtual addresses which are automatically translated into physical node identifiers via a global translation lookaside buffer (GTLB). The GTLB caches entries of a software global destination table (GDT), much like a TLB caches page table entries. Arriving messages are queued in a register-mapped hardware FIFO readable by a system-level message handler. Two network priorities are provided, one each for requests and replies. Messages are routed in dimension order using up to four virtual channels.

## 3. INTRA-NODE CONCURRENCY MECHANISMS

The amount and granularity of available parallelism varies enormously across application programs and even during different phases of the same program. Some phases have an abundance of instruction level parallelism that can be extracted at compile time. Others have data dependent parallelism that can only be exploited using multiple threads; thus the task size to achieve maximum concurrency may vary widely.

The M-Machine is designed to efficiently execute programs with either compiler or runtime scheduled parallelism and with a range of granularities. The M-Machine architecture contains two mechanisms for intra-node concurrency: Vertical Threads (V-Threads) and Horizontal Threads (H-Threads). A *V-Thread* is similar to a *process*; it has its own protection domain and may be swapped in and out of the processor by the system

software. A V-Thread is allocated one set of processor and pipeline registers on each cluster where its subthreads (or H-Threads) execute. Each *H-Thread* is a 3-wide instruction stream which is statically scheduled and executes on a single MAP cluster. The H-Threads of a V-Thread can be either independently scheduled or scheduled together by the compiler to achieve 9-wide instruction level parallelism (ILP).

The MAP has sufficient hardware resources to support up to six resident V-Threads; on each cluster, the H-Threads residing there are interleaved on a cycle-by-cycle basis over the shared execution resources. Consecutive instructions executed by a given cluster may be from distinct V-Threads and instructions executed at the same time on different clusters may be from different V-Threads. This flexible interleaving allows the MAP to exploit thread-level parallelism and to mask variable pipeline, memory, and communication delays.

The arrangement of V-Threads, H-Threads, instructions, and operations is summarized in Fig. 3. The contexts of six V-Threads are resident in the clusters' register files. Each V-Thread includes three H-Threads, one on each cluster. Each H-Thread consists of a sequence of 3-wide instructions containing integer, memory, and floating point operations. On subsequent cycles, a cluster, as demonstrated by cluster 0, may issue instructions from different V-Threads.

### 3.1. Single Cluster Execution

A *V-Thread* consists of at least one and up to three H-Threads, each running concurrently on a different cluster. The MAP has sufficient hardware resources to hold the state of six V-Threads (18 H-Threads), with each V-Thread occupying a *thread slot*. Three of the thread slots are *user* slots, two are for *events*, and one is the *exception* slot. User threads run in the user slots, handlers for asynchronous events run in the event slots, and handlers for synchronous exceptions detected and localized within a cluster, such as protection violations run in the exception slot. Message arrival is treated as an asynchronous event.

The H-Threads within the same V-Thread may communicate and synchronize via registers, while H-Threads of different V-Threads must synchronize and communicate through memory or messages.

On each cluster, up to six H-Threads (one from each V-Thread) are interleaved dynamically over the cluster's resources on a cycle-by-cycle basis. A synchronization pipeline stage holds the next instruction to be issued from each of the six H-Threads until all of its operands are present and all of the required resources are available, similar to the architecture described in Ref. 4.



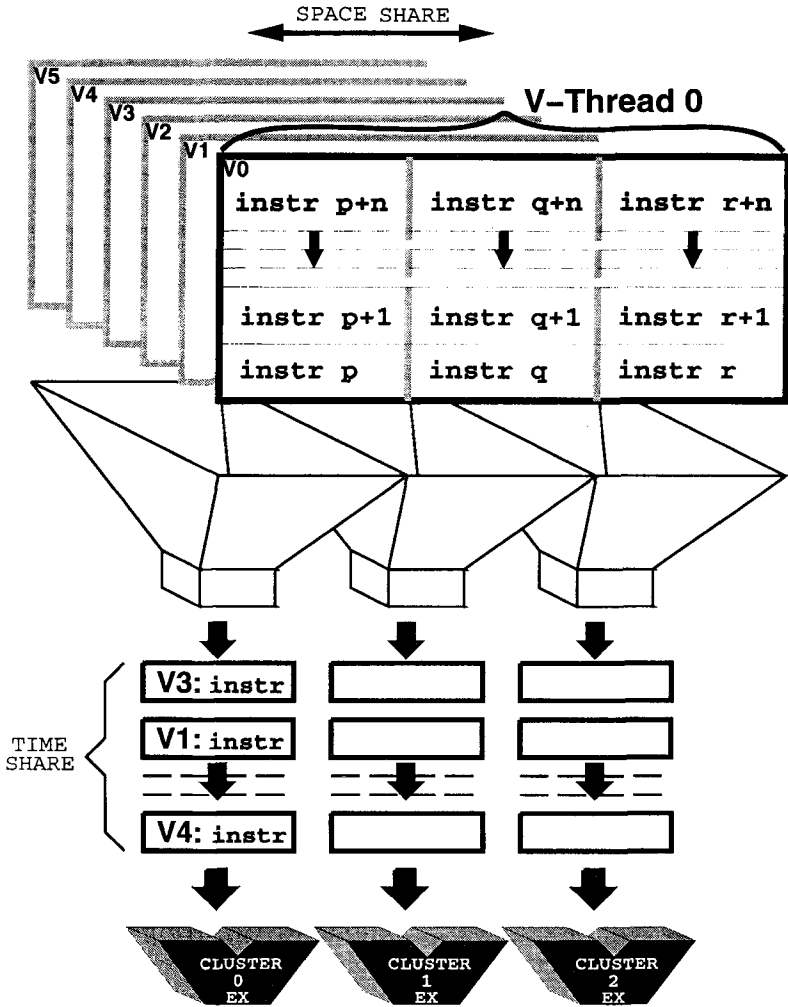


Fig. 3. Multiple V-Threads are interleaved dynamically over the cluster resources. Each V-Thread consists of 3 H-Threads which execute on different clusters.

On every cycle, the SZ stage decides which instruction to issue from those which are ready to execute. The SZ stage checks each instruction to determine if all of its operands are present, all of the resources it requires are available, and all of the barrier conditions are satisfied. The MAP chip implements both a memory barrier (MBAR) and a cluster barrier (CBAR) instruction. The MBAR instruction causes a thread to stall until all of its outstanding memory references complete, while CBAR stalls until the other

two clusters each reach their own CBAR instruction. After determining which instructions are ready to execute, the SZ stage selects one based on priority and deadlock avoidance criteria. A one bit per thread priority scheme is used to allow high priority threads to use more of the execution slots. A pre-emption counter ensures that threads of low priority get at least some portion of the execution slots.

An H-Thread that is stalled waiting for data or resource availability consumes no resources other than the thread slot that holds its state. Multiple H-Threads may be interleaved with zero delay, allowing task switching to mask even very short pipeline latencies as well as longer communication and synchronization latencies. As long as its data and resource dependencies are satisfied, a single thread may issue an instruction every cycle. Therefore, single thread performance is not penalized as a result of the M-Machine's support for multithreading.

### 3.2. Multicluster Execution

An *H-Thread* runs on a single cluster and executes a sequence of operation triplets (zero or one operation for each of the 3 ALUs in the cluster) that are issued simultaneously. Within an H-Thread, instructions are guaranteed to issue in order, but may complete out of order. An H-Thread may communicate and synchronize via registers with the two other H-Threads within the same V-Thread. An H-Thread may only read operands from its own register file, but can write directly into the register files of its collaborating H-Threads that are within the same V-Thread on other clusters.

The H-Thread mechanism can support multiple execution models. H-Threads can execute as independent threads with possibly different control flows to exploit loop-level or thread-level parallelism. Alternatively, the compiler can schedule the three H-Threads in a V-Thread as a unit to exploit instruction level parallelism, as in a VLIW machine. In this case the compiler may insert explicit register-based synchronization operations or employ the-CBAR instruction to enforce instruction ordering between H-Threads. Unlike the lock-step execution of traditional VLIW machines, H-Thread synchronization occurs only as required by data or resource dependencies. While explicit synchronization incurs some overhead, it allows H-Threads to slip relative to one other in order to accommodate variable-latency operations such as memory accesses.

Figure 4 shows a simple illustrative example of the instruction sequences of a program fragment on 1 and 2 H-Threads. The program is the body of the inner loop of a "smoothing" operation using a 7-point stencil on 3-D grid. On a particular grid point, the smoothed value is given by:

$$u_* = u_* + a \times r_* + b \times (r_u + r_d + r_n + r_s + r_e + r_w) \quad (3.1)$$

where  $r_*$  is the residual value at that point, and  $r_u, r_d, r_n, r_s, r_e$  and  $r_w$  are the residuals at the neighboring grid points in the six directions UP, DOWN, NORTH, SOUTH, EAST and WEST respectively. In order to better illustrate the use of H-Threads, advanced optimization (such as software pipelining) is not performed.

Figure 4a shows the single H-Thread program, with a 12 long instruction stream which includes all of the memory and floating point operations. The weighting constants  $a$  and  $b$  are kept in registers. Figure 4b shows the instruction streams for two H-Threads working cooperatively. Each H-Thread performs four memory operations and some of the arithmetic calculations. Instruction 7 in H-Thread 0 calculates a partial sum and transmits it directly to register  $t_2$  in H-Thread 1. The empty instruction on H-Thread 1 is used to prepare  $t_2$  for H-Thread synchronization; H-Thread 1 will not issue instruction 7 until the data arrives from H-Thread 0 as explained later.

The use of multiple H-Threads reduces the static depth of the instruction sequences from 12 to 8. On a larger 27-point stencil, the depth is

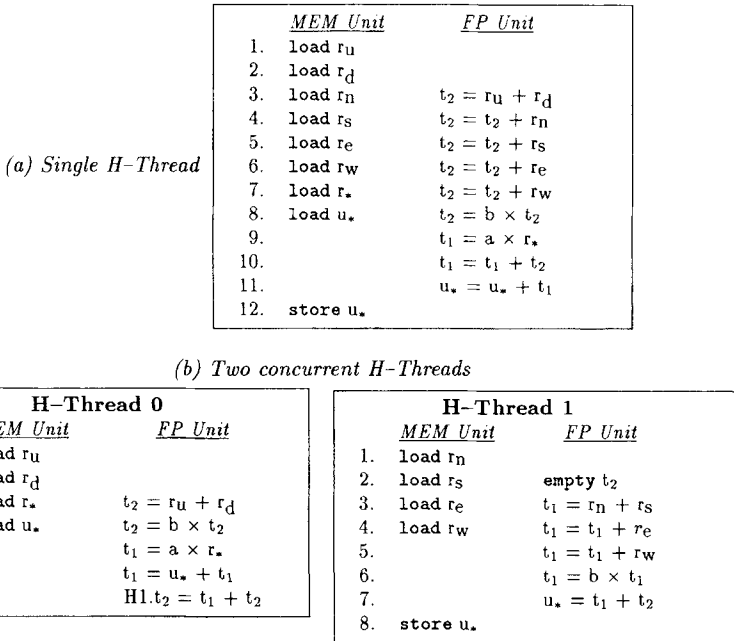


Fig. 4. Example of H-Threads used to exploit instruction level parallelism: (a) single H-Thread; (b) two H-Threads. The computation is a smoothing operator using a 7-point stencil on a 3-D grid:  $u_* = u_* + a \times r_* + b \times (r_u + r_d + r_n + r_s + r_e + r_w)$ .

reduced from 32 to 15 when run on 3 H-Threads. The actual execution time of the program fragments will depend on the pipeline and memory latencies.

### 3.2.1. Thread Synchronization

As shown in the previous example, H-Threads within the same V-Thread may synchronize with one another through registers. A scoreboard bit associated with the destination register is cleared (empty) when a multi-cycle operation, such as a load, issues and is automatically set (full) by MAP hardware when the result is available. An operation that uses the result will not be selected for issue until its register operands are present with the corresponding scoreboard bits set.

All inter-cluster data transfers require explicit register synchronization. To prepare for intercluster data transfers, the receiving H-Thread executes an EMPTY operation to mark empty a set of destination registers. As each datum arrives from the transmitting H-Thread over the C-Switch, the corresponding destination register is automatically set full by MAP hardware. An instruction in the receiving H-Thread that uses the arriving data will not be eligible for issue until its data is available. Therefore, explicit synchronization operations required by VLIW style execution across H-Threads may be overlapped with the inter-cluster data transfers inherent in the executing program. In order for the clusters to coordinate the emptying and writing of registers, three main synchronization mechanisms are provided: global condition registers, the CBAR instruction, and a tightly coupled mode.

### 3.2.2. Global Condition Registers

Each V-Thread has an independent bank of global condition code (CC) registers. Each bank is composed of three sets of four single-bit global CC registers and is used to broadcast binary values between H-Threads within a V-Thread. Similar to data registers, each global CC register has an accompanying scoreboard bit. The MAP global CC registers are physically replicated on each of the clusters, instead of being centrally located. An H-Thread may broadcast to other H-Threads of the same V-Thread using one of its writable global CC registers (only one of the three sets per H-Thread is writable), but may read and mark empty its local copy of *any* global CC register in its bank. Using these registers, all three H-Threads can execute conditional branches and assignment operations based on a comparison performed by a single H-Thread.

The scoreboard bits associated with the global CC registers may be used to rapidly synchronize among the H-Threads within a V-Thread. Figure 5 shows an example of two H-Threads synchronizing at loop

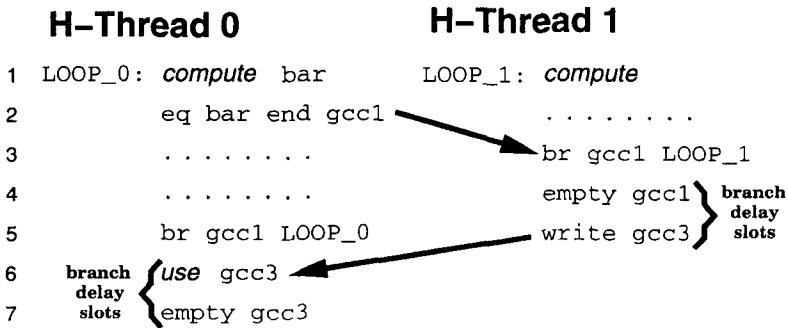


Fig. 5. Loop synchronization between two H-Threads using MAP global condition code (CC) registers.

boundaries. Two registers are involved in the synchronization, in order to provide an interlocking mechanism which ensures that neither H-Thread rolls over into the next loop iteration.

H-Thread 0 computes `bar`, compares it (using `eq`) to `end`, and broadcasts the result by targetting `gcc1`. H-Thread 1 uses `gcc1` to determine whether to branch, marks `gcc1` empty again, and writes to `gcc3` to notify H-Thread 0 that the current value of `gcc1` has been consumed. H-Thread 0 blocks until `gcc3` is full, and then empties it for the next iteration. Neither thread can proceed with the next iteration until both have completed the current one. Due to the multicopy structure of MAP global CC registers, this protocol can easily be extended to perform a fast barrier among 3 H-Threads executing on different clusters, without combining or distribution trees.

### 3.2.3. Cluster Barrier Instruction

Although the global condition registers can be used to implement pairwise barriers and quickly broadcast the results of comparison operations to be used by conditional branches on other clusters, global barriers do incur noticeable overheads. The cluster barrier (CBAR) instruction is intended to be used for fast barrier synchronization among the clusters. When executed on one H-Thread, the CBAR instruction stalls in the synchronization (SZ) pipeline stage until the H-Threads on the other 2 clusters execute CBAR instructions. Then all 3 clusters may proceed past the CBAR. Thus, a global barrier can be executed with only one instruction overhead and none of the communication overhead required using data or global condition registers. Since the CBAR operations stall until ready, no execution unit cycles are used waiting for the other H-Threads to

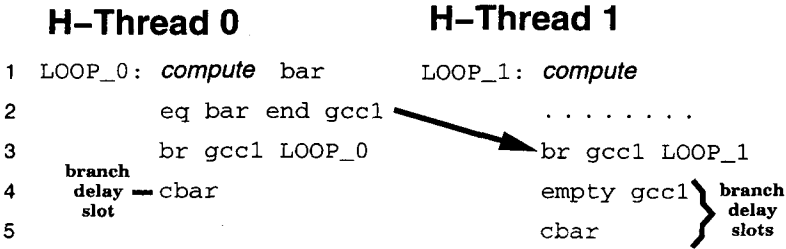


Fig. 6. Loop synchronization between two H-Threads using the CBAR instruction.

reach the barrier, and the arithmetic units on that cluster may be used by the other resident V-Threads.

Figure 6 shows the same loops as Fig. 5, but using CBAR instead of the GCC registers for synchronization. The register gcc1 is still used for broadcasting the result of the comparison operations, but the overhead of using gcc3 is removed, by using the CBAR instruction to enforce instruction ordering across the different H-Threads.

### 3.2.4. Tightly Coupled Mode

The overhead of executing the CBAR instruction may become significant for code that requires frequent synchronization across all of the clusters. Tightly coupled mode provides implicit synchronization between all of the H-Threads on the 3 clusters, essentially by making every instruction an implicit CBAR. The MAP chip then behaves like a 9-wide VLIW processor, in which no operation from instruction  $i + 1$  may issue (leave the SZ pipeline stage) until all of the operations from instruction  $i$  have issued. However, it is not identical to a VLIW as the cycle-by-cycle scheduling decisions on each cluster are independent. Since latencies are not completely known at compile time, the compiler must still coordinate register transfers between clusters by emptying the register on the destination cluster. However, traditional VLIW compilation techniques can be used to generate the code.

## 3.3. Asynchronous Exception Handling

Exceptions that occur outside the MAP cluster are termed *events* and are handled asynchronously by generating an *event record* and placing it in a hardware event queue. Local TLB misses, block status faults, memory synchronization faults, and message arrivals are events that are handled asynchronously. These events are precise in the sense that the faulting operation and its operands are specifically identified in the event record,

but they are handled asynchronously, without stopping the thread. Each H-Thread in the *event1* V-Thread slot handles one class of events. Local TLB misses are handled on cluster 0, and arriving messages are handled on clusters 1 and 2, depending on the priority of the message. Memory synchronization and status faults are handled in the *event0* slot and can use all 3 H-Threads in the slot to execute the event handler.

The dedicated handler located in each H-Thread of the *event* V-Thread slots processes event records to complete the faulting operations. The event handler loops, reading event records from the register-mapped queue and processing them in turn. A read from the queue will not issue if the queue is empty. For example, on a local TLB miss, the hardware formats and enqueues an event record containing the faulting address as well as the write data or read destination. The software TLB miss handler reads the record, places the requested page table entry in the TLB, and restarts the memory reference. The thread that issued the reference does not block until it needs the data from the reference that caused the miss. Similarly, inter-node message arrival is also treated as an event in which the contents of the message are written into the appropriate event queue (which serves as the message queue).

Handling events asynchronously obviates the need to cancel all of the issued operations which follow an operation that faults, a significant penalty in a 9-wide machine with deep pipelines. Dedicating H-Threads to this purpose accelerates event handling by eliminating the need to save and restore state, and allows concurrent (interleaved) execution of user threads and event handlers. Asynchronous event handling does require sufficient queue space to handle the case where every outstanding instruction generates an exception. To reduce queue size requirements, those exceptions that can be detected in the first execution cycle, such as protection violations and some arithmetic exceptions, stall all user H-Threads in the affected cluster, and are handled synchronously by the local H-Thread of the *exception* V-Thread. User H-Threads executing on neighboring clusters are unaffected.

### 3.4. Discussion

There are two major methods of exploiting instruction level parallelism. Superscalar processors execute multiple instructions simultaneously by relying upon runtime scheduling mechanisms to determine data dependencies.<sup>(6,7)</sup> However, they do not scale well with increasing number of function units because a greater number of register file ports and connections to the function units are required. In addition, superscalars attempt to schedule instructions at runtime (much of which could be done at compile

time), but can only examine a small subsequence of the instruction stream to do so.

Alternatively, Very Long Instruction Word (VLIW) processors such as the Multiflow Trace series<sup>(6)</sup> use only compile time scheduling to manage instruction-level parallelism, resource usage, and communication among a partitioned register file. However, the strict lock-step execution is unable to tolerate the dynamic latencies found in multiprocessors.

Processor Coupling, as originally introduced in Ref. 4, used implicit synchronization between the clusters on every wide instruction. Relaxing the lock-step synchronization, as described in this section, has several advantages. First, it is easier to implement, because control is localized completely within the clusters. Second, it allows more slip to occur between the instruction streams running on different clusters (H-Threads), which eliminates the automatic blocking of one thread on long latency operations of another, providing more opportunity for latency tolerance. Finally, the H-Threads can be used flexibly to exploit both instruction and loop level parallelism. When H-Threads must synchronize, they may do so explicitly through registers, at a higher cost than implicit synchronization. However, fewer synchronization operations are required, and many of them can be included in the data transfer between clusters which are inherent in the executing program. In addition, global condition code registers, the CBAR instruction, and tightly coupled mode provide lower cost synchronization mechanisms when synchronization latency is critical.

Using multiple threads to hide memory latencies and pipeline delays has been examined in several different studies and machines. Gupta and Weber explore the use of multiple hardware contexts in multiprocessors,<sup>(9)</sup> but the context switch overheads they used are too large to mask pipeline latencies. MASA<sup>(10)</sup> as well as HEP,<sup>(11)</sup> and TERA<sup>(12)</sup> use fine grain multithreading to issue an instruction from a different context on every cycle in order to mask pipeline latencies. However, with the required round-robin scheduling, single thread performance is degraded by the number of pipeline stages. The zero cost switching among V-Threads and the pipeline design of the MAP provide fast single thread execution as well as latency tolerance for better local memory bandwidth utilization. Furthermore, none of the previous multithreaded machines have multiple clusters for exploiting wide instruction level parallelism.

Various machines optimized for dataflow languages<sup>(13-15)</sup> provide hardware support for fine grained synchronization between threads (usually via memory synchronization bits), but they do not exploit instruction level parallelism, nor do they provide low cost register-based synchronization between threads. The XIMD architecture<sup>(16)</sup> uses multiple ALUs to exploit instruction level parallelism as well as thread level parallelism. However, it



uses a single global register file and does not interleave multiple threads over the same execution units. Two approaches that do exploit instruction level parallelism using multiple threads and multiple ALUs include Ref. 17 and 18.

#### 4. INTER-NODE CONCURRENCY MECHANICS

The M-Machine provides a fast, protected, user-level message passing substrate. A user program may communicate and synchronize by directly sending messages or by reading and writing remote memory using a coherent software shared memory system layered on the message-passing substrate. Direct messaging provides maximum performance data transfer and synchronization while shared memory support simplifies programming. Remote memory access is implemented using fast trap handlers that intercept load and store operations which reference remote data. These handlers send messages to other nodes to complete remote memory references transparently to user programs. Additional hardware and software mechanisms allow remote data to be cached locally in both the cache and external memory.

##### 4.1. Message Passings Support

The M-Machine provides hardware support for injecting a message into the network, determining the message destination, and dispatching a handler on message arrival. For example, Fig. 7 shows the M-Machine instruction sequences for both the sending and receiving components of a remote memory store. The message sending sequence (Fig. 7a) loads the

###### (a) Message Send

```
LOAD   A[0], MC1      ; load A[0] into register 1
SEND   Raddr, Rdip, #1 ; send a 3 word remote store
                               ; message to the processor
                               ; containing VA in Raddr
```

###### (b) Message Receive

```
loop:
JMP    Rnet           ; jump to DIP (remote write)
;start of remote write code
MOVE   Rnet, R1       ; move virtual address into R1
STORE  Rnet, R1       ; store word to memory
BRANCH loop          ; branch to message dispatch code
```

Fig. 7. Example of M-Machine code implementing a remote store: (a) Sending a 3 word remote store message; (b) Receiving and performing the store. On the receiving end Rnet is the register mapped to the head of the message queue.

data to be stored into general register MC1. The SEND instruction takes three arguments, the target address (Taddr) contained in Raddr, the dispatch instruction pointer (DIP) in Rdip, and the message body length (#1). When the SEND issues, the Global Translation Lookaside Buffer (GTLB) translates virtual address Raddr into a physical node identifier, and a 3 word message containing DIP, Taddr, and the contents of MC1 is sent to that node. When the message arrives at the destination (Fig. 7b) hardware enqueues it in the priority 0 message queue. An H-Thread dedicated to message handling jumps to the handler via the DIP contained in the first word of the message, executes a store operation, and branches back to the dispatch portion of the code. Two message priorities are provided: user messages are sent at priority zero, while priority 1 is reserved for system level message reply, thus avoiding deadlock.

#### 4.1.1. Message Address Translation

As described in Ref. 19, the explicit management of processor identifiers by application programs is cumbersome and slow. To eliminate this overhead, the MAP implements a Global Translation Lookaside Buffer (GTLB), backed by a software Global Destination Table (GDT), to hold mappings of virtual address regions to node numbers. These mappings may be changed by system software. The user specifies the destination of a message with a virtual address, which the network output interface hardware uses to access the GTLB and calculate the physical destination node.

A range of virtual addresses (called a page-group) is mapped across a region of processors with a single GTLB entry. In order to simplify encoding, the page-group must be a power of 2 pages in size, where each page is 512 words. The mapped processors must be in a contiguous 3-D rectangular region with a power of 2 number of nodes on a side. This information is encoded in a single GTLB entry as shown in Fig. 8. The *virtual page* field is used as the tag during the fully associative GTLB lookup. The *starting node* enumerates the coordinates of the origin of the region of mapped processors, while the *extent* specifies the base 2 logarithm of the X, Y, and Z

				Extent		
Virtual Page	Starting Node	Page-group Length	Pages/Node	Z	Y	X
42 bits	16 bits	6 bits	6 bits	3 bits each		

Fig. 8. Format of a Global Destination Table (and GTLB) entry, used to determine a physical node identifier from a virtual address.

dimensions of the region. The *page-group length* field specifies the number of local pages that are mapped into the page group. The *pages-per-node* field indicates the number of pages placed on each consecutive processor, and is used to implement a spectrum of block and cyclic interleaving

#### 4.1.2. Message Reception

At the destination node, an arriving message is automatically placed in a hardware message queue. The head of the message queue is mapped to a register accessible by an H-Thread (in either cluster 1 or 2, depending on message priority) in the *event* V-Thread. The message dispatch handler code running in that H-Thread stalls on the empty register until a message arrives, marking the register full; the handler then reads the dispatch instruction pointer (DIP) from the register and jumps to it. This starts execution of the specific handler code to perform the action requested in the message. Some of the actions include remote read, remote write, and remote procedure call. The message need not be copied to or from memory, as it is accessible via a general register. In order to avoid overflow of the fixed size message queue and back up of the network, only short, well-bounded tasks are executed by message handlers. Longer tasks are enqueued to be run as a user process on a user V-Thread.

#### 4.1.3. Protection

The M-Machine communication substrate provides fully protected user-level access to the network. The SEND instruction atomically launches a message into the network, preventing a user from occupying the network output indefinitely. The automatic translation provided by the GLTB ensures that a program may only send messages to virtual addresses within its own address space. Finally, restricting the set of user accessible DIPs prevents a user handler from monopolizing the network input. If an illegal DIP is used, a fault will occur on the sending thread before the message is sent.

#### 4.1.4. Throttling

In order to prevent a processor from injecting messages at a rate higher than they can be consumed, the M-Machine implements a *return-to-sender* throttling protocol. A portion of a local node's memory is used for returned message buffering. When a message is sent, a counter is automatically decremented, which reserves buffer space for that message, should it be returned. If the counter is zero, no buffer space is available and no additional messages may be sent; threads attempting to execute a SEND instruction will stall. When the message reaches the destination a reply is sent indicating whether the destination was able to handle the message.

If the message was consumed, the reply instructs the source processor to increment its counter, deallocating the buffer space. Otherwise, the reply contains the contents of the original message which are copied into the buffer and sent again later.

#### 4.1.5. Discussion

The M-Machine provides direct register-to-register communication, avoiding the overhead of memory copying at both the sender and the receiver, and eliminating the dedicated memory for message arrival, as is found on the J-Machine.<sup>(20)</sup> Register-mapped network interfaces have been used previously in the Mars Machine,<sup>(21)</sup> J-Machine, and iWarp,<sup>(22)</sup> and have been described by \*T<sup>(23)</sup> as well as Henry and Joerg.<sup>(24)</sup> However, none of these systems provide protection for user-level messages.

Systems, like the J-Machine, that provide user access to the network interface without atomicity must temporarily disable interrupts to allow the sending process to complete the message. The M-Machine's atomic SEND instruction eliminates this requirement at the cost of limiting message length to the number of cluster registers. Most messages fit easily in this size, and larger messages can be packetized and reassembled with very low overhead.

Automatic translation of virtual processor numbers to physical processor identifiers is used in the Cray T3D.<sup>(25)</sup> The use of virtual addresses as message destinations in the M-Machine has two advantages. When combined with translation hardware, it provides protection for user initiated messages, without incurring the overhead of operating system invocation, as messages may not be sent to processors mapped outside of the user's virtual address space. It also facilitates the implementation of global shared memory. The interleaving performed by the GTLB, although not as versatile as the CRAY T3D address centrifuge or the interleaving of the RP3,<sup>(26)</sup> provides a means of distributing ranges of the address space across a region of nodes.

In contrast to both \*T and FLASH<sup>(21)</sup> which use a separate communication coprocessor for receiving incoming messages, the M-Machine incorporates that function on its already existing execution resources, an H-Thread in the event V-Thread. This avoids idling a dedicated processor when it is not in use. During periods of few messages, user threads may make full use of the cluster's arithmetic and memory bandwidth.

## 4.2. Non-Cached Shared Memory

Fast access to remote memory is provided transparently to the user with a combination of hardware and software mechanisms. When a load or store

operation to a global virtual address causes a Local Translation Lookaside Buffer (LTLB) miss, a software trap is signalled. Like the hardware dedicated to message arrival, one H-Thread in the *event1* V-Thread is reserved for handling LTLB misses. The LTLB miss handler code probes the GTLB to determine where the requested data is located, and if necessary, sends a message to the destination node. If the data is in fact local, the LTLB miss handler fetches the required page table entry and places it in the LTLB. Using a small portion of the execution resources for fast trap handling reduces the latency of both local LTLB misses and remote data access.

The sequence of operations required to satisfy a remote memory load is shown below. The labels *HW* and *SW* indicate whether the action is performed by hardware or software.

1. *HW*: Memory operation accesses the cache and misses (2 cycles).
2. *HW*: LTLB miss occurs, enqueueing an event (2 cycles).
3. *SW*: Software accesses the local page table (LPT), probes the GTLB, and composes and sends a message containing the referenced and return addresses (48 cycles).
4. *HW*: Message delivered to remote node (5 cycles).
5. *SW*: Message handler fetches requested data from memory, formats a reply message, and sends it (29 cycles).
6. *HW*: Return message delivered (5 cycles).
7. *SW*: Message handler decodes the original load destination register and writes the data directly there (41 cycles).

Timelines for both remote read and write accesses are shown in Figs. 9 and 10. These measurements are based on prototype message and event handlers written in assembly code and running on the M-Machine simulator. A user level program running on node 0 makes read and write requests to memory on neighboring node 1. Except for the message handler that runs on demand, node 1 is idle. All references to memory system data structures in the software handlers are assumed to cache hit.

Table I shows a comparison of preliminary results of local and remote access latencies (in cycles), for single word accesses. A read is completed when the requested data has been written into the destination register. A write is completed when the line containing the data has been fully loaded into the cache. The remote read and write accesses are larger than their local counterparts due to the software intervention required to send the message to the remote node. However, the time to perform a remote read that hits in the cache is only a twice as large as a local read that

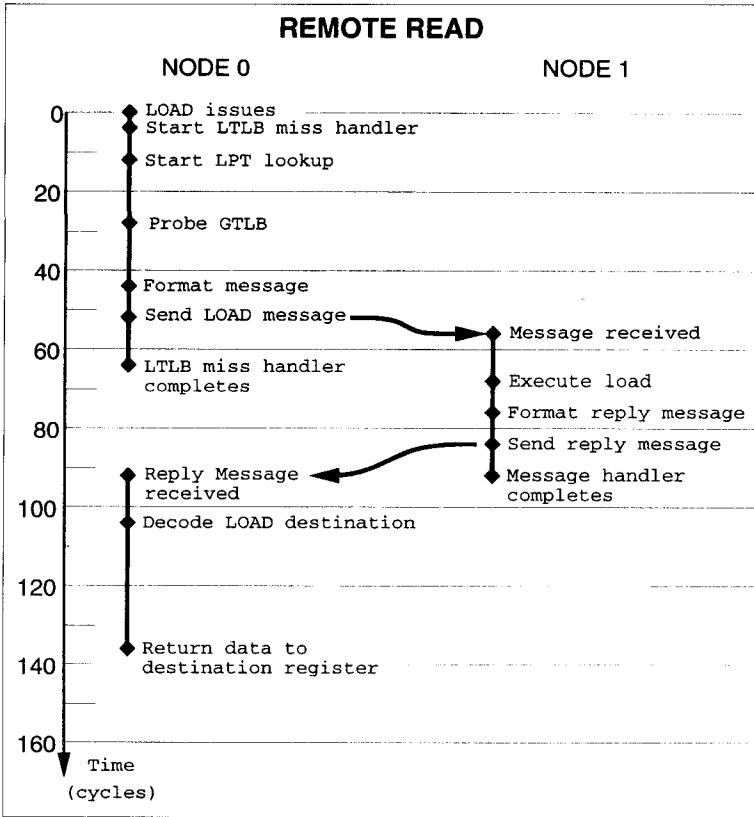


Fig. 9. Remote read access.

**Table I. Comparison of Local and Remote Access Times, Assuming No Resource Contention**

Access Type	Access Times (cycles)	
	READ	WRITE
Local Cache Hit	3	2
Local Cache Miss	13	19
Local LTLB Miss	61	67
Remote Cache Hit	138	74
Remote Cache Miss	154	90
Remote LTLB Miss	202	138

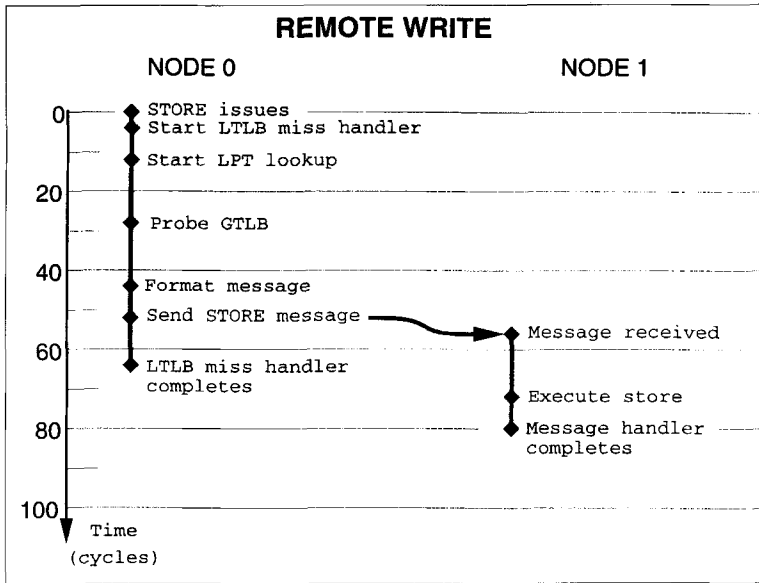


Fig. 10. Remote write access.

requires software intervention (LTLB miss). For the remote write, which does not require return data, the difference is only 10%.

The primary contributors to remote access latency in the M-Machine are searching for the faulting address in the local page table and decoding the reply message (about 40 cycles each). The page-table search is required only when accessing the first block of a page. Accesses to subsequent blocks cause block-status faults (rather than page faults) which skip the page-table access. The reply decode could be accelerated by prohibiting the faulting V-Thread from swapping out during the memory operation.

### 4.3. Caching and Coherence

Even though remote accesses are fast, their latency is still large compared to local memory references. This overhead reduces the ability of the map to use the network and remote memory bandwidth effectively. To reduce overall latency and improve bandwidth utilization, each M-Machine node may use its local memory to cache data from remote nodes.

In addition to the virtual to physical mapping, each LTLB (and LPT) entry contains 2 status bits for each cache block in the page. These *block status* bits are used to provide fine-grain control over 8 word blocks,

allowing different blocks within the same mapped page to be in different states. This fine-grain control over data is similar to that provided in hardware based cache coherent multiprocessors, and alleviates the false sharing that exists in other software data coherence systems.<sup>(28)</sup> The two block status bits are used to encode the following four states:

- **INVALID:** The block may not be read, written, or placed in the hardware cache.
- **READ-ONLY:** The block may be read, but not written.
- **READ/WRITE:** The block may be read or written.
- **DIRTY:** The block may be read or written, and it has been written since being copied to the local node.

One software policy that uses the block status bits fetches remote cache blocks on demand. When a memory reference occurs, the block status bits corresponding to the global virtual address are checked in hardware. If the attempted operation is not allowed by the state of the block, a software trap called a *block status* fault occurs. The trap code runs in the *event0* V-Thread. The block status handler sends a message to the home node, which can be determined using the GTLB, requesting the cache block containing the data. The home node logs the requesting node in a software managed directory and sends the block back. When the block is received, the data is written to memory and the block status bits are marked valid. If the virtual page containing the block is not mapped to a local physical page, a new page table entry is created and only the newly arrived block is marked valid. The remote data may be loaded into the on-chip cache, and modifications to the data will automatically mark the block state dirty. More complex coherence schemes can map blocks from different virtual pages into the same physical page, reducing the amount of unmapped physical memory.

The software handlers used to transmit data from node to node may implement a variety of coherence policies and protocols. This code is easily incorporated within the remote read and write handlers described in Section 4.2. Using local memory as a repository will allow more remote data to be cached locally than could fit in the on-chip cache alone.

#### 4.3.1. Discussion

Directory-based, cache coherent multiprocessors such as Alewife<sup>(29)</sup> and DASH<sup>(30)</sup> implement coherence policies in hardware. This improves performance at the cost of flexibility. Like the M-Machine, FLASH<sup>(27)</sup> implements remote memory access and cache coherence in software, but uses a coprocessor. However, this system does not provide block status bits



in the TLB to support caching remote data in local DRAM. The subpage status bits of the KSR-1<sup>(31)</sup> perform a function similar to that of the block status bits of the M-Machine.

Implementing remote memory access and coherence completely in software on a conventional processor would involve delays much greater than those shown in Table I, as evidenced by experience with the Ivy system.<sup>(28)</sup> The M-Machine's fast exception handling in a dedicated H-Thread avoids the delay associated with context switching and allows the user thread to execute in parallel with the exception handler. The GTLB avoids the overhead of manual translation and the cost of a system call to access the network. Finally, the M-Machine provides memory-mapped addressing of thread registers to allow the operation to be completed in software.

## 5. M-MACHINE SOFTWARE

The M-Machine addresses the problem of parallel software by supporting an incremental approach to parallelization. Unlike conventional parallel machines, the M-Machine is designed to efficiently run a sequential program that uses all the machine's memory, including that on remote nodes. A shared address space, high-performance messaging, and caching remote data in local DRAM provide fast access to remote data. The programmer can then incrementally improve program performance by adding parallelism. The cache coherence mechanisms enable efficient sharing of data across processors. The high-speed messaging network and runtime system support allow for low-overhead task parallelism. The ability to support fine-grain parallelism increases the number of suitable tasks and allows extraction of more parallelism from small problems. Support for synchronizing memory operations and global addressing simplifies user-level communication and synchronization between tasks and reduces overhead. Caching in DRAM automates much of the data placement and migration. For the cases where a programmer wants to extract the maximum performance fast, protected, user-level messages may be employed.

The M-Machine software is being designed and implemented jointly with the Scalable Concurrent Programming Laboratory at Caltech. The Multiflow compiler<sup>(32)</sup> has been ported to the M-Machine to generate long instructions spanning multiple clusters. The Multiflow compiler is designed to generate VLIW instructions from a sequential source program using Trace Scheduling. The modifications required to generate multicluster code for the M-Machine consist of partitioning the graph (DAG) of the trace into sub-DAGs that may be executed on different clusters with minimal communication. The sub-DAGs are then scheduled for each cluster using

the greedy instruction scheduler of the Multiflow compiler. Explicit synchronization is required to take the place of the implicit synchronization of a true VLIW. Communication is implemented by writing to remote registers, while the global condition registers and cluster barrier instruction are used to implement explicit barrier synchronization. An algorithm that might be used to discover the synchronization points is described in Ref. 33. The compiler currently generates code for a single cluster, and a prototype for generation multi-cluster code is being tested.

A prototype runtime system consisting of primitive message and event handlers has also been implemented. Approximately 90 percent of the runtime system code is implemented in C, compiled using the Multiflow compiler, and runs on the M-Machine simulator; the remaining 10 percent are assembly code routines which access hardware features not exposed to the compiler. The runtime system consists of independent modules which manage virtual memory allocation, physical memory allocation, memory coherence between nodes, and multiple threads on a single node and across nodes. The implementation of the runtime system is described more fully in Ref. 34.

## 6. CONCLUSION

In this paper we have described the architecture of the M-Machine with an emphasis on its novel features. The M-Machine, currently under development, is a 3-D mesh, each node of which contains a multi-ALU processor (MAP) and 8 Mbytes of synchronous DRAM. Each MAP chip consists of three 64-bit 3-issue clusters connected by a cluster switch, a 4-way interleaved on-chip cache, an external memory interface, and on-chip network interfaces and routers.

Instruction level parallelism is exploited both within a cluster and across clusters using H-Threads. An H-Thread may communicate and synchronize through registers with H-Threads on different clusters but within the same V-Thread. A 27 point stencil computation on 3 H-Threads (9-wide issue) has half the static instruction count of 1 H-Thread (3-wide issue).

To increase use of the local memory and execution bandwidth, multiple tasks, called V-Threads, are interleaved on a cycle-by-cycle basis independently on each of the clusters. Each cycle, a different thread may be selected for execution, or if only one V-Thread is resident, it may issue an instruction every cycle on each cluster.

The M-Machine has a user-level, protected, fast message passing substrate to reduce communication and remote memory latencies. Messages are composed in general registers and sent via a user level SEND instruction.

Arriving messages are extracted by system-level software message dispatch handlers, which are always resident in the *event1* V-Thread. The message contents are accessed via a register mapped queue. The message need not be copied to or from memory on either the sending or receiving side. Two level translation is used to independently relocate objects in the physical address space on a node, and in the processor namespace.

The fast message system is used to provide the user with transparent access to remote memory. When a user's load or store instruction traps to software on a LTLB miss, a message is sent to a remote node to perform the access. While slower than local accesses, a remote load can be satisfied in 138 cycles, while a remote store can be satisfied in 74 cycles. In order to facilitate local caching of remote data, 2 status bits for each block (8 words) in a page are added to the LTLB and page table entries. When an invalid block is accessed, a trap to software occurs which can retrieve the missing block from a remote node, copy it into local memory, and mark the status bits valid.

Both a C language based architectural simulator and a Verilog based, cycle accurate, RTL model have been completed and are being used for software development and hardware validation. A combination of manual, random, and compiler generated tests have been used to validate both the RTL model and the circuit schematics against the architectural simulator. The hardware design of the MAP is nearly complete; all the modules have been designed, 100% of the schematics are done, and 95% of the datapath layout is complete. The MAP will be fabricated on a single integrated circuit with a projected area of 18 mm × 18 mm in 0.5 μm CMOS with 5 metal layers. The target clockrate is 100 Mhz and tapeout is expected in the middle of 1997. After tapeout, the runtime system will be completed and the optimization and scheduling components of the compiler will be improved. When the manufactured MAP chips are returned, a single node system will be built to test the communication and synchronization mechanisms between H-Threads, and a 16-node system will be built to evaluate the inter-node communication and memory systems.

The M-Machine addresses the growing imbalance between memory system capacity and bandwidth, making all of memory close to some processor and increasing the aggregate bandwidth to memory. By employing multiple processors and multiple ALUs within a processor, the M-Machine enables parallelism to be used to both access more memory simultaneously, and keep the expensive communication channels between a given processor and its local memory busy.

The M-Machine increases the percentage of chip area devoted to the processor from 0.13% to 16% for a typical 1996 system. A 32-node (96 clusters) M-Machine with a total of 256 Mbytes of memory requires 50%

more chip area than a uniprocessor with the same amount of memory but provides 96 times as much peak performance, a 64:1 improvement in peak-performance/area. More importantly it provides 32 times the bandwidth between the processors and memory. The 64:1 improvement in peak-performance/area makes the increased parallelism of the M-Machine cost effective even in cases where only a small fraction of its peak performance is realized.

We expect that the architecture concepts demonstrated in the M-Machine will be useful in machines ranging from single-node personal computers, through workstations with tens of nodes, to servers with hundreds to thousands of nodes. Memory bandwidth and capacity are becoming the dominant factor in the cost and performance of systems of all scales. By changing the processor/memory ratio, providing methods for extracting parallelism at all levels, and supporting an incremental approach to parallelism, the M-Machine's mechanisms will lead to more cost effective and programmable machines across the price-performance spectrum.

## ACKNOWLEDGMENTS

We would like to thank the members of the Scalable Concurrent Programming Laboratory at Caltech for their feedback on the M-Machine architecture and for their work on the M-Machine software effort: Steve Taylor, Daniel Maskit, and Bryan Chow. In addition we would also like to thank Kathy Knobe and the anonymous referees for their comments on various versions of this paper. The research described in this paper was supported by the Advanced Research Projects Agency and monitored by the Air Force Electronic Systems Division under contract F19628-92-C0045.

## REFERENCES

1. J. L. Hennessy and N. P. Jouppi, Computer Technology and Architecture: An Evolving Interaction. *Computer*, pp. 18-29 (September 1991).
2. C. A. Mead, L. A. Conway, *Introduction to VLSI Systems*. Addison- Wesley, Reading, Massachusetts, (1980).
3. L. Gwennap, New MIPS Chip Targets Windows NT Boxes. *Microprocessor Report* (November 18, 1992).
4. S. W. Keckler, and W. J. Dally, Processor Coupling: Integrating Compile Time and Runtime Scheduling for Parallelism. *Proc. 19th Int'l. Symp. Computer Archit.*, Queensland, Australia, ACM, pp. 202-213 (May 1992).
5. N. P. Carter, S. W. Keckler, and W. J. Dally, Hardware Support for Fast Capability-Based Addressing. *Proc. Sixth Int'l. Conf. on Archit. Support Progr. Lang. Oper. Syst. (ASPLO VI)*, Association for Computing Machinery Press, pp. 319-327 (October 1994).

6. R. Tomasulo, An Efficient Algorithm for Exploiting Multiple Arithmetic Units. *IBM J.* **11**, 25–33 (January 1967).
7. W. M. Johnson, *Superscalar Microprocessor Design*. Prentice Hall, Englewood Cliffs, New Jersey (1991).
8. R. P. Colwell, W. E. Hall, C. S. Joshi, D. B. Papworth, P. K. Rodman, and J. E. Tornes, Architecture and Implementation of a VLIW Supercomputer. *Proc. Supercomputing*, IEEE Computer Society Press, pp. 910–919 (November 1990).
9. A. Gupta, and W.-D. Weber, Exploring the Benefits of Multiple Hardware Contexts in a Multiprocessor Architecture: Preliminary Results. *Proc. 16th Ann. Symp. Computer Archit.* IEEE, pp. 273–280 (May 1989).
10. R. H. Halstead, and T. Fujita, MASA: A Multithreaded Processor Architecture for Parallel Symbolic Computing. *15th Ann. Symp. Computer Archit.* IEEE Computer Society, pp. 443–451 (May 1988).
11. B. J. Smith, Architecture and Applications of the HEP Multiprocessor Computer System. *SPIE Vol. 298 Real-Time Signal Processing IV*, Denelcor, Inc., Aurora, Colorado, pp. 241–248 (1981).
12. R. Alverson *et al.*, The Tera Computer System. *Proc. Int'l. Conf. Supercomputing*, ACM SIGPLAN Computer Architecture News, pp. 1–6 (September 1990).
13. R. S. Nikhil, G. M. Papadopoulos, Arvind, \*T: A Multithreaded Massively Parallel Architecture. Computation Structures Group Memo 325-1, Laboratory for Computer Science, Massachusetts Institute of Technology (November 1991).
14. H. H. Hum *et al.*, A Design Study of the EARTH Multiprocessor. *Int'l. Conf. Parallel Archit. and Compilation Techn.*, pp. 59–68 (1995).
15. S. Sakai, Y. Kodoma, and Y. Yamaguchi, Prototype Implementation of a Highly Parallel Dataflow Machine em-4. *Proc. Fifth Int'l. Parallel Processing Symp.*, IEEE Computer Society, pp. 278–286 (May 1991).
16. A. Wolfe, and J. P. Shen, A Variable Instruction Stream Extension to the VLIW Architecture. *Proc. Fourth Int'l. Conf. Archit. Support for Progr. Lang. Oper. Syst.*, ACM Press, pp. 2–14 (April 1991).
17. G. S. Sohi, S. E. Breach, and T. Vijaykumar, Multiscalar Processors. *Proc. 22nd Int'l. Symp. Computer Archit.*, pp. 414–425 (May 1995).
18. D. M. Tullsen, S. J. Eggers, and H. M. Levy, Simultaneous Multithreading: Maximizing On-Chip Parallelism. *Proc. 22nd Int'l. Symp. Computer Archit.*, pp. 392–403 (May 1995).
19. M. D. Noakes, D. A. Wallach, and W. J. Dally, The J-Machine Multicomputer: An Architectural Evaluation. *Proc. 20th Int'l. Symp. Computer Archit.*, San Diego, California, IEEE, pp. 224–235 (May 1993).
20. W. J. Dally *et al.*, The J-Machine: A Fine-Grain Concurrent Computer. *Proc. the IFIP Congress* G. Ritter, (ed.), North-Holland, pp. 1147–1153 (August 1989).
21. P. Agrawal, W. Dally, W. Fischer, H. Jagadisich, A. Krishnakumar, and R. Tutundjian, A. Mars, A Multiprocessor-Based Programmable Accelerator. *IEEE Design Test* **4**:28–36 (October 1987).
22. S. Borkar *et al.*, Supporting Systolic and Memory Communication in Iwarp. *Proc. 17th Int'l. Symp. Computer Archit.*, pp. 70–81 (May 1990).
23. G. M. Papadopoulos, G. A. Boughton, R. Grainer, and M. J. Beckerle, \*T: Integrated Building Blocks for Parallel Computing. *Proc. Supercomputing*, IEEE, pp. 624–635 (1993).
24. D. S. Henry, and C. F. Joerg, A Tightly-Coupled Processor-Network Interface. *Fifth Int'l. Conf. Archit. Support for Progr. Lang. Oper. Systems (ASPLOS V)*, ACM, pp. 111–122 (October 1992).
25. Cray Research, Inc., *Cray T3D System Architecture Overview*. Chippewa Falls, Wisconsin (1993).

26. G. Pfister *et al.*, The IBM Research Parallel Processor Prototype (RP3): Introduction and Architecture. *Proc. Int'l. Conf. Parallel Processing*, pp. 764–771 (1985).
27. J. Kuskina, D. Ofelt, M. Heinrich, J. Heinlein, R. Simoni *et al.*, The Stanford FLASH Multiprocessor. *Proc. 21st Int'l. Symp. Computer Archit.*, IEEE, pp. 302–313 (April 1994).
28. L. K. Ivy, A Shared Virtual Memory System for Parallel Computing. *Int'l. Conf. Parallel Processing*, pp. 94–101 (1988).
29. A. Agarwal *et al.*, The MIT Alewife Machine: A Large-Scale Distributed-Memory Multiprocessor. *Scalable Shared Memory Multiprocessor*, Kluwer Academic Publishers, (1991).
30. D. Lenoski, J. Laudon, T. Joe, D. Nakahira, L. Stevens, A. Gupta, and J. Hennessy, The DASH prototype: Implementation and Performance. *Proc. 19th Ann. Int'l. Symp. Computer Archit.*, IEEE, pp. 92–103 (1992).
31. S. J. Frank *et al.*, Multiprocessor Digital Data Processing System. United States Patent No. 5,055,999 (October 8 1991).
32. P. G. Lowney, S. G. Freudenberger, T. J. Karzes, W. D. Lichtenstein, R. P. Nix, J. S. O'Donnell, and J. C. Ruttenberg, The Multiflow Trace Scheduling Compiler. *J. Supercomputing* 7(1/2):51–142 (May 1993).
33. A. Zaafrani, H. G. Dietz, and M. O'Keefe, Static Scheduling for Barrier MIMD Architectures. *Int'l. Conf. Parallel Processing* (1990).
34. Y. Gurevich, The M-Machine Operating System. Master of Engineering Thesis, Massachusetts Institute of Technology, Department of Electrical Engineering and Computer Science (September 1995).