# The M-Machine Operating System

by

Yevgeny Gurevich

Submitted to the Department of Electrical Engineering and
Computer Science
in partial fulfillment of the requirements for the degree of

Master of Engineering in Electrical Engineering and Computer
Science

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

September 1995

© 1995 Yevgeny Gurevich. All rights reserved.

Author . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Department of Electrical Engineering and Computer Science
August 11, 1995

Certified by . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
William J. Dally
Associate Professor
Thesis Supervisor

Accepted by . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Frederick R. Morgenthaler
Chairman, Department Committee on Graduate Theses

# The M-Machine Operating System

by

## Yevgeny Gurevich

## Abstract

This document details the design and implementation of an operating system written specifically for the M-Machine, a multicomputer currently being designed at MIT. The operating system is designed to be lightweight and flexible, able to support a UNIX-like operating system layer interface to higher-level code, while at the same time exposing machine primitives to user programs in a safe and efficient manner. The operating system's central features are its support for fast and efficient thread creation and built-in memory-coherence to present the view of global virtual memory to user-level programs as well as higher-level protected subsystems. Four core components are presented - the physical and virtual memory managers, the thread manager, and the memory-coherence manager.

# Acknowledgments

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

The M-Machine is a new multicomputer currently being designed at the MIT AI Lab. The machine's hardware features, some radically different from conventional architectures, require a custom operating system. The operating system is meant to define a small collection of powerful primitives which may be used to construct interface layers in order to emulate existing, familiar operating systems. At the same time, this low-level system code attempts to expose novel hardware features to user-level code in a safe manner. For this reason, the low-level OS needs to be efficient and flexible. Flexible, in terms of providing a framework of very general primitives, and efficient in order to allow other system personalities to reside as higher-level layers without significantly impacting performance.

Following the general trend in operating system design, the M-Machine's OS (M-Machine Runtime System or MARS) is a loose collection of managers which work in concert, instead of a single monolithic kernel such as in traditional implementations UNIX. These managers provide the minimum functionality necessary for an operating system - memory-management and process (thread) management. They form the basis for allowing user-level programs to execute on the machine in protected, stable manner. These low-level managers execute on each node of the M-Machine, similar to the microkernel of Amoeba, which performs process and memory management tasks. As in Mach, the collective managers are designed to enable the implementation of a UNIX-like API which sits above the low-level OS layer. Such an implementation

can be efficient and at the same time, provide a common and familiar programming environment.

As in Mach and Amoeba, the OS presented in this thesis supports lightweight thread creation which may be used a basis for the much heavier and often inefficient UNIX fork, although newer implementations of UNIX have moved towards this design as well, providing more lightweight process-creation functions.

A novel addition to this operating system is low-level memory-coherence management. Even operating systems such as Mach, which were designed in part to run on multicomputers, do not integrate a core global shared virtual memory system. A distributed shared memory server in Mach would operate at a higher-level. Furthermore, the global shared virtual memory, supported by hardware-based capabilities, allow the OS to employ a single machine-wide memory map which is identical for all processes. This differs from operating systems like UNIX, Amoeba, and Mach, where virtual memory maps depend on the currently-executing process. A single address space simplifies sharing and writing parallel programs. The use of capabilities free the OS from having to use complicated software-based capability schemes to enforce protections on shared memory. Shared virtual memory is as inexpensive to access, and as safe from errant and malicious threads, as a thread's private memory. Like Amoeba but unlike Mach and UNIX, the M-Machine OS does not include a pager. Such an addition requires additional complexity, and design time for a complete I/O system as well.

This thesis is divided into three general sections. In the first, chapter 2 provides a quick overview of particular aspects of the M-Machine architecture which will both shape the design of the operating system and enable it to perform its duties in an efficient manner. Chapter 3 then provides a high-level picture of the M-Machine operating system's structure.

The second section presents more detailed design and implementation of the four central subsystems within the operating system - the physical memory manager, the virtual memory manager, the thread manager, and the memory-coherence manager. These are covered in chapers 4, 5, 6, and 7.

In the last section, the interface for user programs to the system is described in chapter 8, some performance figures are given in chapter 9, and chapter 10 concludes with project status and future work which needs to be done.

# Chapter 2

# Target Hardware Overview

This section presents a brief description of the machine architecture targetted by MARS - the M-Machine. The M-Machine is a shared-memory superscalar multi-computer being designed at the MIT Artificial Intelligence Laboratory. A detailed architectural design is provided in [4]. At the high level, the M-Machine consists of a mesh of *nodes* serviced by a high-speed network substrate. Each node consists of four *clusters*. Clusters contain an integer, memory, and floating-point unit capable of issuing in parallel on each clock cycle. Multiple register files and other thread state support up to six thread contexts in hardware simultaneously. The clusters may communicate with each other through a dedicated cluster-switch, and to four cache-banks through a memory switch. With this design, the machine's peak issue rate is twelve operations per clock cycle, with up to 12 outstanding memory references being serviced by the individual cache banks at any one time. Several of the machine's distinctive features which greatly affect its operating system design are presented in this chapter. They include (1) hardware primitives to support global shared virtual memory, (2) operations for atomic memory access, (3) hardware-enforced capabilities for memory protection, (4) support for fast context-switching, and the ability to concurrently maintain several thread contexts in hardware, (5) support for message-send primitives at the instruction level, and (6) mechanisms for accessing hardware state through a memory-mapped configuration space.

## 2.1   A shared-address-space multicomputer

The M-Machine supports a global 54-bit virtual address space across all of its nodes.
Local on-node caches are virtually-addressed, a global translation lookaside buffer
(GTLB) maintains mappings of virtual addresses to their home nodes, and system
software is required to maintain memory coherence between nodes. Memory refer-
ences which miss in the cache are handled by an external memory interface (EMI)
which probes an on-node local translation lookaside buffer (LTLB) to determine which
physical page provides backing for the referenced virtual address. Due to the design
of the memory system, lines in the cache must be backed by a local physical page so
that they may be flushed to external memory by hardware in the event of a cache-line
conflict. References which pass to the LTLB but miss there as well result in an event
record being generated which allows software intervention.

### 2.1.1   Hardware primitives for implementing shared-memory

In order to support an implementation of software-based coherent shared memory,
the M-Machine architecture maintains two status bits for each 8-word block of vir-
tual memory. These *block-status bits*, signifying whether a line is invalid, read-only,
exclusive-clean, or exclusive-dirty, are maintained by hardware in the local translation
lookaside buffer and cache. The memory system prevents an access from completing
if it violates the block-status bits, instead generating an event record which allows
system software to intervene and satisfy the access. There are three fault types: write
to an invalid line, read to an invalid line, and write to a read-only line. All events
are handled by a dedicated system thread as explained in the next chapter. System
software is expected to replicate and manage these block status bits in node page
tables when altering entries in the LTLB.

### 2.1.2   Atomic Test-and-set Memory Operations

In order to support access to global shared data structures in the face of concur-
rency, each word of the machine's memory includes a lock bit which is referenced

in atomic test-and-set memory operations. These synchronizing memory operations allow programs to perform loads or stores conditional upon the status of the lock bit (the precondition), and set the bit to a known value if they succeed (postcondition). Conditional synchronizing memory operations return a condition which is true if the test-and-set succeeded and the memory operation completed, and false otherwise. Unconditional synchronizing memory operations generate events if the preconditions they require are not met. Details of instructions are in [5].

## 2.1.3 Hardware-supported capabilities

In order to maintain global shared virtual memory without the use of access lists, capabilities are used to enforce memory protection. Words may be tagged as pointers to memory segments of a power of two bytes in length by system software, and given out to user-level processes. User processes are only allowed to copy the pointers as is or to modify their address portion so as to change their offsets within the memory segment that the pointer represents. In this way, it is not necessary for the operating system to maintain separate page tables for each process in a multi-process environment since pointers may not be forged. As explained in [2], this presents a problem when a thread deallocates a segment of virtual memory since the operating system does not know a priori which threads may have been passed this pointer[1]. A global garbage-collection of allocated virtual memory must be performed to find and destroy any clones of pointers to virtual segments before such segments may be deemed clean and available once again for allocation. However, as [2] shows that for large address spaces, reclamation may be performed extremely infrequently. [1] presents a more detailed description of the M-Machine's capabilities. The use of different pointer types to allow efficient user access to system code will be revisited in chapter 8.

---

[1]Threads may pass pointers around in messages, by writing them into memory shared by other threads, or direct intercluster register-file writes

**User vthreads (slots 0 –3)**

**Exception vthread (slot 5)**

**Event vthread (slot 4)**

| cluster 0 | cluster 1 | cluster 2 | cluster 3 |

LTLB Miss
Event
Record

hardware
event
queue

priority 0
message
input
queue

priority 1
message
input
queue

Figure 2-1: Thread Slots and Clusters

## 2.1.4  V- and H-Threads

A thread which executes on the M-Machine is identified as a V-Thread. It occupies one of six hardware thread slots and may be composed of up to four decoupled H-Threads, each running on a different cluster on the node. H-Threads communicate with each other either through memory or intercluster register-file writes. More details of these mechanisms are given in [4]. At any instant, any combination of four H-Threads from the six different V-Thread slots may be issuing instructions down the piplines of the clusters. V-Threads are round-robin scheduled by the hardware to allow each fair access to machine resources. Four of the hardware thread slots are intended for user-level threads. The remaining two thread slots are meant for system-level handlers. Certain registers in the system-level thread slots are mapped to hardware resources such as the event and network input queues, as shown in figure 2-1. These system thread slots form the core of the M-Machine operating system as will be described in chapter 3. The hardware support for several thread slots allows for efficient context switching among available user threads for better latency tolerance. In addition, there is no expensive penalty for invoking system handlers to

18

respond to events and messages since suspension and eviction of a user thread to make room for a system handler is not required. The handlers are always active, sleeping until an event or incoming message requires their attention. Finally, the controlled manner in which system handlers are invoked - similar to a protection violation invoking kernel mode in traditional operating systems - protects handlers from errant threads since no direct function call is involved.

## 2.1.5 Support for Efficient Message-passing

Primitive hardware instructions to perform an atomic message send allow threads to inject messages into the M-Machine's internode network without needing to call system software. At the user-level, this allows threads to invoke handlers on a particular node if they obtain (1) an entry pointer into a message-handler routine and (2) a pointer to a virtual memory segment mapped to the destination node. The requirement for a pointer to a message-handler routine ensures that incoming messages are serviced by trusted code which will not lock up the network input queue. At the system level, threads are allowed to send messages directly to physical node numbers instead of using virtual addresses. This message-send primitive is employed by the memory-coherence management software as explained in chapter 7. A ten-word message size limit is adequate for the system software's requirement of shipping 8-word cache lines with some extra status information.

## 2.1.6 Memory-mapped Access to Hardware State

Threads on the M-Machine may access hardware state through load and store memory operations which target configuration space. Pointers tagged with the configspace type identify such accesses and requests are passed to configuration space controllers in each cluster. Machine state such as the LTLB, portions of the instruction-cache, hardware thread contexts, and status registers may be read and modified by system software. Since the configuration address space is 54 bits, hardware state is laid out sparcely so as to simplify hardware decoding of requested addresses. As will

become evident in the rest of this document, configuration-space access will be one of the central tools employed by the runtime system to manage the machine. Since configuration-space pointers allow such powerfull access, these pointers are never given out to user-level threads, and are constructed when needed by priviledged threads or generated directly by hardware state machines on event-record generation.

# Chapter 3

# Runtime System Overview

The M-Machine Runtime System (MARS) is split into two distinct pieces - system functions which are invoked by user-level threads and execute within the caller's thread slot, and low level handlers which perform physical memory management, memory-coherence, and thread management. System functions allow protected system code execution within a user thread's context with the help of the capabilities mentioned in the previous chapter. A detailed description of system entry and exit is provided in [1]. The operating system presented here is not truly complete, as it lacks a design for I/O, among other components.

The current runtime system implementation uses its own data pointer when invoked as system code, but for simplicity still borrows the caller's stack for intermediate values and spill space. This, however, is actually a potential security leak since a malicious user-level thread may pass a copy of its stack pointer to a confederate (perhaps another H-Thread within its own slot) which may then overwrite portions of the stack or snoop on the stack contents, hoping to encounter a pointer it is not normally allowed to access. MARS can be modified to counter this security problem. A more secure system would employ distinct system stacks inaccessible by the user-level caller. Such stacks may be maintained as linked lists of virtual or physical segments allocated by the OS at boot time, and popped for use by system code when a system function is first invoked. System handlers are invoked indirectly through fault mechanisms and are therefore more secure, using their own dedicated stacks and

```
┌─────────────────────────────────────────────────────────────────────────┐
│              System functions executed in the 4 user thread slots         │
│                                                                           │
│        Thread Manager                      Virtual–Memory Manager          │
│     tSpawn                                                                 │
│               tSignal                   vmem_alloc      vmem_dealloc       │
│       tFork                                                               │
│               tSleep                                                      │
│        tExit                                                              │
│                                                                           │
│               ▼         Stubs for Handler Entry    ▼                      │
│         add_eh_job                    unmap_page   map_page               │
│                                          lookup_page                       │
│               ▲                    ▲                                      │
├───────────────────────────────────────────────────────────────────────────┤
│            Low–level handlers executing in System thread slot             │
│  Event Handler   │ LTLB          │ P0 Message      │ P1 Message           │
│                  │ Miss Handler  │ Handler         │ Handler              │
│                  │               │                 │                      │
│  Thread          │ Physical      │ Thread          │ Thread               │
│  Management:      │ Page          │ Management:      │ Management:           │
│  scheduling      │ Management    │ Remote          │ Remote               │
│                  │               │ tSleep/tSignal  │ tSleep/tSignal       │
│                  │               │                 │                      │
│  Memory          │               │ Memory          │ Memory               │
│  Coherence       │               │ Coherence       │ Coherence            │
│  requests        │               │ requests        │ requests             │
│                  │               │                 │                      │
└───────────────────────────────────────────────────────────────────────────┘
```

Figure 3-1: OS Components Overview by Hardware

data segment pointers.

All components of the runtime system are designed with several key principles in mind - handlers and system functions are meant to be lightweight, tolerate concurrency, and be flexible and general enough to support a variety of high-level operating systems built from the primitives that they provide.

# 3.1   Differences from a Traditional Operating System

Unlike traditional operating systems such as the UNIX and Mach variants, MARS is designed from the ground up to support concurrent execution of lightweight threads in a single global virtual address space, and provide a view of coherent virtual memory even to higher-level operating system components. Unlike Mach, where message-

passing provides a secure interface for communication, communication among and between user-level processes and operating system components is accomplished through function calls and memory access. This is in great part due to the hardware capabilities which enable protected low-cost shared-memory access, and the underlying memory-coherence protocol. MARS does, however, share several design concepts with Mach, such as the support for synchronization primitives, inexpensive IPC, and lightweight threads/processes.

Instead of a kernel or microkernel of linear code which is invoked by user code through a page-fault mechanism, most system calls in MARS are handled within a user thread slot through protected system entry. This has several key advantages - while system calls are being performed by one user thread, other user threads are not prevented from executing their code. In addition several system calls may be active at any time, even the same system function (system calls must be designed to be re-entrant and use locks when accessing shared data structures.) Finally, more critical services such as the handling of memory protection violations and page faults are still invoked in the traditional fault-reponse manner, but with two key differences. First, while a system event handler is executing other threads may continue issuing until a conflict in a hardware resource arises, in which case the system threads are given higher priority. Second, even the faulting thread may make progress until it requires the use of data which is being serviced by the fault mechanisms. More importantly, even the invocation of critical OS services may be performed concurrently, with three and sometimes four different system-level fault handlers being able to service independent requests at the same time.

A more detailed view of the low level handlers which reside in the system thread slot is provided in the next section.

## 3.2 System Thread Components

The system V-Thread, running in thread slot four on each node of the M-Machine, is composed of four decoupled H-Threads effectively providing four independent han-

dlers which may concurrently satisfy system events. The four H-Threads are the Event Handler (EH), LTLB Miss Handler, Priority 0 Message Handler (P0MH), and Priority 1 Message Handler (P1MH). All thread components remove events from hardware-based event FIFO queues and process each event in turn. If no events are present, handlers simply block until an event arrives, thereby allowing other threads to issue and not stealing any execution cycles on the machine. An overview of these system components is shown in figure 3-1. Each thread blocks on a different hardware FIFO, allowing up to four events to be handled at a time. Events are usually of fixed-length and are inserted into queues by hardware state machines. Each of the H-Threads in the system V-Thread has integer register 14 mapped to the head of its respective event queue. When that register is used as a source for an operation, the head word of the next event in the hardware queue is used as the data source. Integer register 15 maps to the body of the event, used to access all remaining words in a hardware event. Once an event word is read out of the hardware queue, it is effectively popped from the queue and may not be recovered. Therefore, most handlers store away event words if they are intended to be used multiple times. The following subsections describe each of the four system H-Threads.

## 3.2.1 Event Handler

The Event Handler responds to block-status miss events, global translation lookaside buffer misses (GTLB Miss), and synchronization misses (SYNC Miss). GTLB misses occur when user-level message-send instructions target virtual addresses which contain no address to node-number mapping in the GTLB. Synchronization misses occur when a synchronizing memory operation fails to proceed because the referenced memory location does not have the requested precondition. The MARS system has not been designed to handle the two latter cases, although future work makes extending the event handler quite simple. In addition, software job queues are used by other components of the OS to request that the Event Handler perform certain tasks, such as evicting or installing threads. A SIGNAL event wakes the event handler to ensure that it gets a chance to examine these software job queues. The handling of block

24

status miss events is a part of the memory-coherence protocol described in chapter 7.

## 3.2.2   LTLB Miss Handler

The LTLB Miss Handler is the most critical component of the MARS system. Due to the nature of the M-Machine memory system, a miss in the LTLB locks up the external memory interface until that miss is serviced. Other threads may continue issuing operations until such time as they cause a cache miss, in which case memory requests stack up until the EMI is freed by the LTLB Miss Handler. The LTLB Miss Handler itself has a separate path (the *bypass* path) into the EMI which insures that it may always access physical memory. Therefore, when the Miss Handler is executing, there is absolutely no guarantee that any other thread is active and able to make progress on the machine. This makes it especially important that the handler not access any data structures which may be locked by other system threads. Such locked data structures may never be released if the owner of the lock is blocked waiting for the LTLB Miss Handler to free up the EMI for memory accesses. In addition, the handler may only access physically-addressed memory, since a virtual address reference may miss in the LTLB itself, causing deadlock. In its normal mode of operation, the LTLB handler maintains the local page table which contains mappings from virtual to physical pages, and refills LTLB entries on an address miss. The handler may also be called through a faked miss mechanism by system software to create, remove, or lookup the mappings that it creates. This mechanism is decribed in more detail in section 4.5.3. Finally, since the LTLB miss handler cannot guarantee that other portions of the event-handling system are able to make progress, it cannot cause Block-Status, Sync, or GTLB misses, or send messages.

## 3.2.3   Message Handlers

The P1 and P0 Message handlers receive messages from the network destined for their node and respond by executing message-handling functions. Such functions may implement a variety of mechanisms including remote memory transfer, remote

25

procedure call, and thread spawn. Others form the core of the software memory-coherence implementation. In order to guarantee deadlock-free execution, all request messages are sent via priority 0, and acknowledgements are returned on priority 1. Priority 1 messages are intended to be handled unconditionally, eventually allowing the network to drain of all P1 traffic and allowing all message traffic to make forward progress. For every P0 message received, at most one P1 message should be returned as an acknowledgement.

## 3.2.4   Availability and Reentrancy

As mentioned previously, since the system handlers reside in an active system V-Thread there is no need to swap in their context and save or restore user thread state before they may begin fulfilling a request. This makes for fast and efficient reponses to what are effectively interrupts without slowing down user code which may be executing concurrently. In addition, there is no time wasted restoring the register-file contents or setting up thread state for the system thread.

The limitations of the LTLB miss handler have already been discussed. In general all event handlers are not reentrant since they are the sole mechanism available for fulfilling their respective event requests. The event handler may not cause and events such as block status, SYNC, or GTLB misses, to occur. In order to maintain the progress guarantees of the network, the P1 message handler may not itself send out messages.

A hardware mechanism prevents user-level threads (those executing in thread slots 0-3) from issuing if the hardware event queue for the event handler rises above a watermark. This mechanism is in place to bound the number of outstanding events in the system and prevent hardware queue overflow. For this reason, it may be possible that protected subsystems which execute within user thread slots may not be able to issue instructions until the event handler has serviced events in the hardware queue. This introduces another constraint upon event handler operation - code executing in the event handler may not wait for locks held by code executing in user thread slots.

## 3.2.5 Signalling the Event Handler

In some instances other OS components, including message handlers, may request that the Event Handler perform a function, such as a message resend, in proxy for them. This is especially true if a message needs to be sent in response to an ACK arriving at the P1MH. In these cases, a request record is added to a *software job queue* for the Event Handler to fulfill at a future time. A signal event is then issued. In order to avoid overflowing the hardware event queue, only a single signal event may be in the hardware event queue at a time. This is accomplished by keeping a word in memory (the event lockword) on which the handlers may synchronize. System code adding a request for the Event Handler (the producer) first adds the event to a software queue and then attempts to set the lock bit of the event lockword to full. If the word was previously full, the set fails and the producer goes on. If the word was previously empty, the producer adds a new SIGNAL event to the hardware event queue. For its part, the Event Handler always resets the event lockword each time it dequeues the SIGNAL event. This guarantees for the producers that if the lockword is set, a previously-issued SIGNAL is still in the hardware event queue (or recently popped from it) and will be examined by the Event Handler as detailed below.

There are two software job queues because handlers within the system thread slot, and protected subsystems within user thread slots may be attempting to add software jobs for the event handler. All code executing within user thread slots synchronizes access to the job queue so that there is a single producer at a time. A job queue is a ring buffer, with two global pointers into it - the *cur* pointer and the *free* pointer. The cur pointer is read and modified only by the consumer (the event handler thread). It is advanced each time a new event is read out and identifies which events have been read out of the job queue. The free pointer is read by both producer and consumer, but advanced only by the producer. Each time a producer wishes to add a new event to the job queue, it reads the free pointer, and begins storing new event words starting at the free pointer and moving down (wrapping around the end of the event buffer if necessary). As its last action, the producer advances the free pointer with a single store operation. This is the atomic action which signifies that a new event

is available. For its part, the event handler checks the cur pointer against the free pointer each time it looks for an event to service. If the cur and free pointers match, no new software events are in the job queue. Otherwise, the event handler may start reading off the cur pointer and advancing it - servicing the next request in the queue.

This mechanism allows the event handler to safely dequeue events without needing to acquire a lock. A second queue is used for the two message handlers to enqueue jobs with the event handler. They also synchronize among themselves to guarantee that there is only one thread adding events to the software job queue at a time.

Given sufficient buffer space to hold all requests, this mechanism is deadlock-free and guarantees that all events in the software queues will eventually be handled. The reason for using the SIGNAL event is to guarantee that requests in the software queue will be examined by the Event Handler if no hardware events are being generated and the Event Handler is blocked waiting for one. The SIGNAL effectively wakes the thread so that it may look at the events stacked up in its software queues. A producer which is unable to set the lockword and therefore add the SIGNAL event is guaranteed that the SIGNAL word is either still in the hardware queue or is just being removed by the Event Handler. In either case, since it had enqueued the request in a software queue prior to attempting a SIGNAL, the producer is guaranteed that the Event Handler will wake up and take a look at the recently added request, as long as the Event Handler runs through the entire software queue before attempting to sleep again. Finally, the lockword also insures that at most one signal has been placed in the hardware queue at a time, preventing queue overflow.

### 3.2.6  System Call handling in User Thread Slots

Despite the great deal of infrastructure developed for handling events in the dedicated system thread slot, many higher-level system calls may be handled by trusted software running within a caller's user thread slot (V-Thread slots 0 to 3). Such system functions include virtual memory allocation, thread and process creation (but not scheduling), invocation of remote functions or spawning remote threads, bulk memory transfer, and others. In short, most routines made available by high-level operating

28

systems which do not require direct manipulation of low-level data structures such as page tables or memory-coherence directories may be safely executed within a user thread slot. In addition, system calls which work with protected data structures that reside in virtual memory may take full advantage of the memory-coherent global memory supplied by low-level OS components. This layered design provides a lot of flexibility.

### 3.2.7 Capabilities and Protection

Capabilities enable user threads to enter system functions in a protected manner. The runtime system "exports" a collection of system functions during the loading of user executables. User programs containing references to system functions are patched with entry pointers to runtime system functions by a trusted loader. Entry pointers may be loaded and used in jump instructions, but may not have their addresses changed. This provides a safe entry mechanism since the system functions which are exported are guaranteed to be entered at well-defined points. Since the setting of the pointer bit is a priviledged operation, user programs may not forge entry pointers of their own. This also means that as the OS evolves, the exact entry points and number of available system functions may change, but legacy programs will still execute correctly since patching is performed at load and not link time. In order for system functions to gain access to the runtime's data segment and associated system-level data structures, a system data segment pointer is stored within the system's code segment by the boot code. As a user-level thread enters a system function, the entry pointer is changed to an execute-system-mode instruction pointer which points into the system code segment. This allows the callee system function to load the system data segment pointer by offsetting from its IP (now allowed since the IP is no longer an entry pointer) and performing a load, overwritting the user data segment pointer. On return to the caller, the user's data segment pointer is restored and a jump to a return pointer switches the thread back to user mode. This process is explained in chapter 8.

29

## 3.2.8 Page Table Design

The global virtual memory supported by the machine allows the system software to use a single inverted page table to maintain virtual-to-physical page mappings for all allocated memory on each node. The M-Machine uses a 4-Kbyte page size for both physical and virtual pages. An open-hashing page table on a node with 16Mbytes of physical memory requires only 8192 entries to be twice as large as necessary for maximum capacity.[1] Assuming 4 64-bit words per entry, this works out to a 0.2% overhead for an inverted page table. Once again, the advantages of maintaining a single page table for all processes running on the node are clear - no switching of tables is necessary on context switches, speeding up multiprocessing performance. In addition, since capabilities prevent user threads from forging pointers, no additional mechanisms are required to prevent a process from accessing virtual memory allocated for other processes. Finally, no special support is required for shared virtual memory. Once a process gives out a virtual pointer to another thread, the virtual segment may be read and written by both. Several flavors of protected pointers allow processes to set up read-only or read-write shared segments.

Since virtual segments do not necessarily need to be backed by contiguous regions of physical memory, a chained list of physical pages is used by the physical memory manager to dole out backing pages to virtual segments. As physical pages become available for allocation, they are added to the free page chain in a FIFO manner. To speed up allocation, a background process may be used to clean physical pages before they become available for allocation. Physical memory management is detailed in the next chapter.

---

[1] The M-Machine currently being designed is expected to have 8MBytes of on-node physical memory.

Accessible by User Threads

**Virtual Segment Manager**          **Thread Manager**

```
vmem_alloc                    tFork              _getMyTC
vmem_dealloc                  tExit              _getParent
                              tSpawn             _getDP
                              hSpawn
                              tSignal
                              tSleep
```


Accessible by Protected Subsystems or User Thread faults

**Physical Page Manager**          **Cache–Coherence Manager**

```
PPM_init                      ccrequest
PPM_map                       ccinvalidate
PPM_unmap                     ccreturnStore
PPM_reclaim_local             ccreturnLoad
PPM_reclaim_remote            ccNackRO
                              ccNackRW
```

**Virtual Segment Manager**

```
vmem_prime
```

Figure 3-2: OS Components Overview by Function

# 3.3 Breakdown of MARS into Functional Components

The previous section approached the MARS design from the point of view of hardware resources used for runtime implementation. This section provides an overview of the runtime system as it is broken down into functional components. The runtime system can be viewed as a collection of managers running in a largely autonomous manner to satisfy requests. At times, managers may call upon each other to fulfill certain requests. This is most commonly the case when system threads require access to physical memory - they call upon the physical memory manager to allocate new physical pages or return information on existing virtual-to-physical mappings.

### 3.3.1  Physical Memory Management

Physical memory management - the maintenance of the local page table and the LTLB - is handled exclusively by the LTLB Miss Handler. Since the handler thread is not allowed access to data structures which may be locked by other threads (as explained in section 3.2.2), there is effectively no overlap in the information which it maintains with that of any other thread. Access to the LTLB Miss Handler is performed in a fault-response manner similar to traditional OS's as described above. In general, misses to reserved virtual addresses which are kept unmapped by the LTLB handler are used as triggers to invoke specific handler functions - such as removal of a particular virtual-physical mapping, creation of a new one, or return of information about an existing mapping.

### 3.3.2  Virtual Memory Management

Virtual memory is doled out in segments by a Virtual Segment Manager which is composed of a series of system functions accessible by user threads. The VSM does not allocate physical backing to the segments which it gives out, simplifying its design and allowing it to run independent of other pieces of the system software within user thread slots - it does not require access to hardware tables, registers, or other machine state. At boot time, the managers on each node of the M-Machine are primed with virtual segments which they may give out and effectively manage independently. The underlying data structure used for tracking allocated and available segments is the buddy list. Details of the VSM and Buddy list allocation are given in chapter 5.

### 3.3.3  Memory-coherence Management

The software memory coherence implementation of the M-Machine is centered around the actions performed by Event and Message handlers on each node. The Event Handler on a requesting node sends P0 requests to home nodes in response to local block-status miss events. The P0 message handler on a shared data item's home node receives requests for cache lines, updates a memory-coherence directory and ships out

blocks of memory as P1 acknowledgements. The P1 message handler on the original requesting node receives the remote cache line (an implicit acknowledgement to its request) and installs it locally. In the event of cache-line conflicts or flush requests, the Event Handler on a node sharing remote data may be required to invalidate and, in the case of dirty lines, return cache lines to their home nodes. The memory-coherence implementation is detailed in chapter 7.

### 3.3.4 Process Management

The management of user processes is broken down into two pieces. System calls invoked within user thread slots are used to fork user threads and add them to lists of ready-to-run threads. Other system calls allow threads to sleep, or signal other sleeping threads. The actual manipulation of hardware thread slots for evicting threads and/or installing new ones is performed by the event handler. This localizes access to the machine hardware so that it is performed by a single thread which is guaranteed to be always active. Although not strictly necessary, this localization simplifies aspects of the memory-coherence implementation as detailed in later chapters.

# Chapter 4

# Physical Memory Management

The M-Machine physical memory manager (PMM) is responsible for maintaining virtual-to-physical page mappings on each node and keeping track of available and allocated physical page frames. Physical memory (sometimes referred to as consisting of backing pages) is usually the ultimate target of memory operations issued on the M-Machine.[1] In the M-Machine memory hierarchy, each node requires a PMM to maintain mappings between virtual pages and their associated physical backing store within a page table. Without a page frame to back it, a virtual address reference cannot be completed. To increase memory-system performance, a 64-entry cache for these mappings is maintained in hardware (the LTLB). The PMM is responsible for keeping the LTLB in sync with the mappings found in the page table. Hardware events notify the PMM when a mapping was not found in the LTLB - an LTLB Miss Event. The PMM must find a mapping within the page table and place it in the LTLB, perhaps evicting a conflicting mapping for a different virtual page. This chapter first introduces a functional interface to the memory-management functions, describes the data structures employed, and details the implementation of the memory manager. As described briefly in the previous chapter, the LTLB Miss Handler is solely responsible for these functions. Section 4.3 explains the rationale behind this design decision.

---

[1]Exceptions are I/O addresses which are memory-mapped into virtual address space, and configuration-space which is a totally separate address space.

# 4.1  System Calls

The PMM performs three different functions as part of its management duties -
creating virtual-physical mappings, removing these mappings, and returning existing
mapping and status information. Interface definitions are shown in table 4.1. These
system calls are meant for protected subsystem use and are only exposed to other OS
components, not user-level threads.

| Function | Description |
|---|---|
| `void PPM_init(int initword)` | Initializes the physical memory manager. The low halfword of *initword* contains the physical page number of the start of unallocated physical memory (the runtime system resides in pages below this page). The high halfword contains the number of pages to add to the local physical page pool (the size of each node's external memory minus the number of page frames consumed by the runtime system). |
| `int PPM_map(int vpn)` | Creates a mapping between virtual page *vpn* and an available, unallocated physical page frame. There are two pools from which to draw page frames - one for backing local virtual addresses, and one for backing virtual addresses mapped to remote nodes. The page frame is taken from either the local or remote-memory pool, depending on the whether the virtual page is local or remote. Returns the page frame number assigned to the new mapping. |
| `void PPM_unmap(int vpn)` | Destroys the mapping of virtual page *vpn* with its physical page frame. |
| `int PPM_reclaim_local(int ppn)` | Returns the page frame *ppn* to the local frame pool. |
| `int PPM_reclaim_remote(int ppn)` | Returns the page frame *ppn* to the frame pool used for remote backing pages. |
| `int PPM_lookup(int vpn)` | Returns the number of the frame backing virtual page *vpn*. Returns -1 if no mapping is found. |

Table 4.1: Physical Page Manager exported functions

35

VPN

probe

VPN match –> use entry

Virtual Page #

Physical Page #

Status bits 0-31

Status bits 32-63

No Match –> reprobe

Figure 4-1: PPM Hash Table Structure

## 4.2 Data Structures

Two main data structures are employed by the PMM. First, the page table is an open hash, used to store virtual to physical mappings and block-status information. The hash table is initialized at machine boot time and sized so that it has room for twice as many mappings as there are page frames on a node. Since the page table is not hierarchical but a hash table, having a large, potentially sparse table is critical for performance reasons - too small a table will result in many conflicts and longer lookup times. Open hashing tables, as described in [3], tolerate entry conflicts without employing chains (linked lists of entries which map to the same location in the hash table) thereby increasing average-case performance. Each hash table entry consists of four words - the actual virtual page number used to define the mapping, its associated physical page number, and 128 block status bits[2] packed into two words (see figure 4-1.)

In order to actually allocate backing pages for virtual pages, the PMM needs to

---

[2]Each page contains 512 words divided into 64 8-word cache lines. 2 Block status bits for each of the 64 cache lines require a total of 128 bits.

maintain a list of all unallocated page frames. The most efficient data structure is a chain of page frame numbers which resides within the unallocated frames themselves. The PMM maintains a single 64-bit word which contains the page frame number of the next unallocated frame which may be used as a backing page. The first word within that frame itself contains the page frame number of the next frame to use. Thus, a chain of available frame numbers is maintained within the unallocated frames. A page frame chain is terminated by a -1 which is never expected to be a valid page frame number. In figure 4-2, the free page frame chain starts at page 15, and terminates at page 2. There are a total of 6 frames in the chain. Popping a new frame for use simply requires reading out the frame number from the frame about to be used and substituting it in the pointer to the next available frame. A list of allocated frames is not required since that information is implicitly stored within the hash table. Any frame popped from the free frames chain must be used in a hash table entry. Conversely, removing a frame from the hash table requires that it be added to the free frames chain. A second 64-bit word stores the frame number which is the last available frame in a chain. This makes returning pages to the page frame chain very simple - the tail frame's next frame entry is modified from *-1* to the frame being added to the chain. The tail frame number is then changed to reflect a new end-of-chain page frame number.

Since the memory-coherence system is so closely tied to the OS, special provision for frames which are used as backing for shared cache lines is made in the PMM. Instead of maintaining a single chain of available page frames, two chains are employed. The first is used to allocate normal backing pages for local data. The second is a limited collection of frames, perhaps some fraction of total on-node memory, which may be used as backing pages for shared cache lines. Once this pool is exhausted, shared cache lines must be evicted until an entire frame is freed up, at which time it will become available for allocation as a backing page of remote data again. The PPM_reclaim_remote call explicitly tells the PPM that a particular frame has been cleaned and should be added to the remote backing page pool. This particular aspect of page management is discussed further in chapter 7.

```
freePage      17
lastPage       2

 1
     16
 2
     -1
 3
   1
 4


15
     3
16
     2
17
     15
18
```

Figure 4-2: PPM Free Page List

# 4.3  Design Rationale

The reason for placing physical memory management in the hands of a system level handler instead of trusted code which may execute in user-level thread slots is tied closely to the M-Machine's memory system design. Since the LTLB miss handler needs access to the page table in order to insert and remove LTLB entries, no other software components may lock the page table data structure (as explained in section 3.2.2). Since locking the local page table is not allowed, there is only a single software component which remains able to access the page table - the LTLB miss handler itself. Because access to the miss handler is restricted to hardware-generated events which occur only on TLB misses, a few system-level routines act as wrappers around the special miss-response interface. These wrappers allow other system components to make standard function calls which in turn result in forced LTLB misses to reserved virtual page numbers. This implementation is detailed in section 4.5.3.

## 4.4 Page Allocation Policy

The PMM employs on-demand page frame allocation. That is, if the LTLB Miss Handler does not find a virtual-physical mapping in the on-node page table for a memory access which touched a particular page, it is assumed that a new mapping needs to be created. This allows efficient use of very large and sparse virtual segments - allocation of a virtual segment does not mean that physical backing needs to be created immediately. Instead, individual LTLB misses to virtual pages cause page frames to be allocated. In fact, the current runtime system only employs the PPM_map function when performing memory-coherence management, since the common case is for mappings to be created on-the-fly by this automatic allocation policy.

In standard operating systems, a memory reference to an unmapped virtual page is considered a disallowed memory access (a segmentation violation or bus error) which needs to be terminated. On the M-Machine, capabilities are used to control memory access. Since threads may not generate pointers on their own, they may not access arbitrary memory locations. All memory accesses which are issued by the memory functional unit on each node's cluster have had their capabilities verified. Therefore, a page fault is not considered a disallowed access on the machine, but rather an access to previously unmapped memory which is still a valid memory reference.

## 4.5 Implementation

The PMM is implemented as a low-level handler written in assembly which pops LTLB Miss events off the hardware event queue and passes them on to the handler body, which is written in C. Once a new event has come in, the assembly stub moves the four words of the event (referenced address, event header word, associated data, and configspace pointer to the faulting thread) into argument registers as defined by the M-Machine compiler and runtime system, and calls the body function. The handler body then determines whether the virtual address reference which caused the LTLB Miss is a fake virtual address and requires that special handling be employed,

or whether it is a standard reference. Upon return, the stub restores its stack and returns to waiting for the next event to arrive.

The body of the handler is written in C, as shown in appendix D (ltlb_body.c). It calls on functions which manipulate the data structures outlined at the beginning of this chapter. The data structure code is also written in C and shown in appendix D. Initialization code is assumed to set up these structures when the LTLB handler is first spawned.

### 4.5.1 Event Format

The hardware-composed LTLB Miss Event consists of four words, shown in table 4.2. The address word identifies the referenced address which caused the LTLB miss. The header word encodes information such as the opcode which was used in the address reference, the issuing V-Thread slot, cluster, and source and destination registers. If the operation was a store, the opdata word contains the data which was attempted to be stored. Finally, the faultcp is the configspace pointer to be used by the software to write thread registers when fulfulling memory requests in software. If the faulting operation was a load op, the pointer offsets directly into the configspace-mapped location of the destination register of the load operation. If it was a store operation that faulted, the configspace pointer identifies a location which updates the faulting thread's *membar* counter. Conditional synchronizing operations have the faultCP identify the their destination cc register.

| Word | Description |
|---|---|
| Header | Encodes information regarding the operation which caused the LTLB Miss and the issuing thread |
| Virtual Address | Virtual address which was not found in the LTLB |
| Opdata | Contains the 64-bit value which was attempted to be stored if the faulting operation was a store op. |
| faultCP | Configspace pointer to thread state for the faulting V-Thread. |

Table 4.2: LTLB Miss Event Format

The low-level handler simply moves the message header and body words into

40

argument registers and calls the manager's C-based handling function.

## 4.5.2   Initial Page Lookup

The handler body code extracts the virtual page number from the miss address that it is passed. This is a simple procedure of shifting off the 12 least significant address bits and masking off the 10 high protection/length bits to retain just the 42-bit page number. The virtual page number is then used to probe the page table by calculating the hash function and indexing into the table. A thorough study of good hash functions has not been performed. In the current implementation, the hash function is an XOR of a 16-bit constant and the rearranged bytes of the low 32 bits of a virtual address. C code is shown in figure 4-3.

```
result = (
        (((vpn >> 24) & 0xffL) << 16) |
        (((vpn >> 16) & 0xffL) << 24) |
        ((vpn >>  8) & 0xffL) |
        ((vpn & 0xffL) << 8)
        ) ^ 0x134aL;
```

Figure 4-3: Hash Function Calculation

Using the algorithm of open hashing, if the handler finds an entry marked *deleted*, or a valid entry whose vpn does not match the vpn being probed, the vpn is rehashed and probing continues. If a vpn match is found, the LTLB is accessed through configspace to determine which existing LTLB entry is to be evicted to make room. Since block-status bits for the virtual page whose mapping is to be evicted may have been modified, they have to be written back into the page table before the mapping can be evicted. Therefore, the existing entry is read from the LTLB through configspace load operations. The vpn of the evicted entry is used to probe into the page table and the block-status bits for that page are copied back into the hash table. Finally, the vpn-ppn mapping and associated block status bits for the page which is to be added to the LTLB are written into the LTLB through configspace stores, overwritting the evicted entry. The EMI is then unlocked through a configspace

41

store, allowing the instruction which caused the miss to be retried automatically by the EMI, this time presumably hitting in the LTLB. Throughout this entire procedure the actual faulting operation is not retried by the handler, and the virtual address never used as the target of a memory operation - the LTLB Miss Handler operates only on physical addresses or configspace addresses.

If continual probing does not find a virtual page match, the page is determined to have no physical backing, and a new backing page needs to be allocated. A GTLB probe is performed to determine whether the virtual address which was referenced is mapped to the node handling the LTLB Miss (a locally-mapped page). If it was a local page reference, local handling is invoked. Otherwise, remote handling is performed. These two distinct cases are described below.

**Local Handling**

The free page chain pointer of normal (not memory-coherence) backing page frames is read to determine the next available frame which may be allocated. The first word of that page is copied into the free page chain pointer, effectively popping off the backing page. This page number is then added to the page table, creating a new virtual-physical mapping. Block-status bits are set to `exclusive` (read/write) for all lines in the page and also written into the page table. Finally, this entry is added to the LTLB so that the next memory access to this page does not cause another LTLB miss.

**Remote Handling**

An initial reference to a remote virtual page requires that a physical page from the pool of memory-coherence pages be used for the mapping. In the simple case, a physical page is available and is popped off the memory-coherence backing chain, in a manner similar to that described in the section above. The only difference is that the block-status bits for the page are set to `invalid` since the node does not yet contain any remote data. When the memory operation is retried, the memory reference no longer causes an LTLB Miss (since the mapping was written added to

the LTLB and page table) but causes a Block-Status miss instead, which then results in the invocation of software leading to the local installation of a remote cache line.

If no physical page frames are available in the backing pool, a special physical frame number, -1, is used as a marker, identifying the fact that no backing frames are available. Since the block-status bits for the new mapping are still set to invalid, there are no problems involved in using the same mapping for all virtual pages which do not have available backing pages.[3] As section 7.3.2 details, this marker page is used as a trigger for the memory-coherence manager to perform a cleanup of existing shared pages and make room for new ones.

### 4.5.3 Fake Miss Interface

As explained in the beginning of this chapter, certain virtual pages which are always unmapped on the machine (trigger pages) are used to request direct manipulation of the PMM data structures by the LTLB Miss Handler. Since user threads are never given pointers to these special pages (and cannot create ones on their own), the miss handler is guaranteed that calls to it through misses are made by trusted subsystem code.

This mechanism, which involves threads generating memory faults to trigger actions by low-level components of the operating system, is similar to standard kernel-entry methods of other operating systems. As outlined in the previous chapter, performance is improved on the M-Machine because no actual thread swapping and context switching is performed.

The actual virtual page numbers used as triggers for the PPM_map, PPM_unmap, and PPM_lookup functions are compile-time constants in the kernel source code and may be picked rather arbitrarily, so long as they are not a subset of the virtual pages which may be allocated by the Virtual Segment Manager (see chapter 5). In the current runtime implementation, these virtual page numbers start at 0x80000. To invoke the

---

[3]The only occasion, in fact, when multiple virtual pages may be mapped to the same physical page within a single node is when a physical backing page is unavailable, in which case all block status bits are set to invalid.

LTLB Miss Handler, a thread issues a conditional synchronizing store instruction, targetting one of the three trigger addresses as shown below:

```
               /* cause a fault              */
instr memu stscnd <data-register>, <trigger-address-register>, <cc>;
               /* block until fault completes */
instr memu ct <cc> ...
```

The store instruction allows the thread to pass 64 bits of data to the LTLB Miss Handler. In most cases, this contains the virtual page number which is to be used as an argument to the PMM functions. By issuing an instruction conditioned on the value of the cc register in the trigger instruction, the requesting thread blocks until the LTLB Miss Handler has completed the request and fills the cc register. Since the functions all have full 64-bit integer return values, the LTLB handler needs to have a simple way to return data to the requesting thread. One mechanism is to overwrite the integer register conventionally used as the return argument register by the compiler - integer register 6 - with the return value. This may be done through a configuration space store operation. In this case, the functional wrapper around the PPM_unmap primitive may look like:

```
PPM_unmap::
    /* i6 contains the argument to this function */
    instr memu stscnd i6, <trigger>, cc1;        -- trigger the LTLB Handler
    instr ialu ct cc1 jmp RETIP;                 -- wait for completion
    instr ;                                      -- i6 (the return register)
    instr ;                                      --   is already set properly
    instr ;                                      --   at this point
```

The current implementation instead writes a single physical memory location (called _ltlb_data_for_mh) which is then loaded by the caller to retrieve the appropriate value. In order to prevent concurrent accesses to this location, a lock is used by callers to serialize access.

## 4.5.4  Reclaiming Pages

The local and remote reclamation functions simply take the supplied page number and add that page to tail of the proper physical page chain. Cleaning needs to be

performed before pages may be considered reclaimed and ready for reallocation.

## 4.5.5  Unmapping

When mappings need to be destroyed, the virtual page number argument to the PPM_unmap function is used to probe the page table until its entry is found. At that time, the entry is removed from the page table and replaced by a *deleted* marker, as necessary for an open-hashing table. In addition, cache lines may need to be flushed, and the LTLB modified to remove the virtual-physical mapping.

Note that no provisions are made for when virtual-physical mappings should be torn down. Higher level OS components make calls upon the PMM to create or eliminate mappings, but the PMM does not need to employ any policy for when mappings should be removed from the page table. Usually, this will be the work of the Virtual Segment Manager, which needs to deallocate physical backing once a virtual segment has been freed. See section 5.3.2 for more details.

# Chapter 5

# Virtual Memory Management

The Virtual Segment Manager (VSM) doles out segments of virtual address space for use by user threads and other portions of the operating system. Segments are a power of two bytes in length, with length protections enforced by the segment length field in pointers. A buddy list allocator is used for the implementation of the underlying allocation mechanism. Since a copy of the low-level operating system runs on each node of the M-Machine, each node's VSM executes independently of all others. This section describes the interface to the VSM, explains the data structures which are employed, and then details the rather simple implementation. Design issues conclude this chapter.

## 5.1  System Calls

The VSM exports a total of three functions. The first, **vmem_prime**, is accessible only to other protected subsystems and allows the bootstrap to initialize the allocator with segments of virtual memory which are available for allocation. The **vmem_alloc** call returns a segment of a requested size, while the **vmem_dealloc** call accepts a segment for deallocation. Table 5.1 provides a brief overview.

| Function | Description |
|---|---|
| `void vmem_prime(void *segment_ptr)` | Identifies the virtual segment pointed to by *segment_ptr* as available for allocation. The pointer length field in the pointer's capabilities explicitly identifies the size of the segment. |
| `void *vmem_alloc(int bytecount)` | Returns a pointer to a clean segment of virtual memory of at least *bytecount* bytes in size. Returns a NULL pointer if no such segment may be allocated. |
| `void vmem_dealloc(void *segment_ptr)` | Deallocates the segment of virtual memory identified by *segment_ptr*. |

Table 5.1: Virtual Segment Manager exported functions



Figure 5-1: Buddy List Structure

## 5.2  Data Structures

The VSM maintains two buddy lists for memory allocation and deallocation. A buddy list is essentially an array of sorted linked-lists of segments (see figure 5-1). The array has as many entries as there are segment sizes available on the machine. On the M-Machine, an array of 51 entries allows segments to range from 8 bytes to $2^{54}$ bytes in length.

The free-segment buddy list stores information on segments which are available for allocation. Initially empty, the list is first primed with segments by the `vmem_prime` call. Subsequent `vmem_alloc` calls remove entries from this free segment list, perhaps

modifying it in the process, and return newly-available segments.

The dirty-segment buddy list records segments which user programs or protected subsystems have asked to be deallocated - those that have been passed to the vmem_dealloc call. This allows deallocated segments to be collected and coallesed into larger segments for bookkeeping pending garbage-collection. A deallocated segment cannot be moved directly from the dirty to the clean list unless a garbage-collection phase has ensured that there are no clones of this virtual address remaining on the entire machine. Therefore, the dirty buddy list essentially stores segments which are candidates for a garbage-collection phase. The current implemention of MARS does not perform garbage-collection. Therefore, the dirty list is just a repository for segments which will never be given out. In fact, it is possible to make the vmem_dealloc function a null function.

Finally, a statically-defined linked-list of nodes for use in both buddy lists allows the allocator to function without needing dynamic memory-allocation itself, although this limits the number of memory segments which may reside in the buddy lists to a compile-time constant.

## 5.3   Implementation

The MARS bootstrap splits the M-Machine's global virtual memory into segments and assigns them to different nodes. Each node's boot code calls vmem_prime with its assigned segment, priming the free-segment buddy list data structure and allowing subsequent allocation calls to use that memory.

Once the priming is complete, the VSM accepts calls from user-level as well as system-level code. Since the VSM code may be running within many thread slots at once, a lock is used to serialize access to global data structures.

### 5.3.1   Allocation

When an allocation call is made, the requested number of bytes is used to calculate the smallest power-of-two-byte segment which contains at least that many bytes. As [7]

is the on-node page table. It may be inefficient for the VSM to run through all of the virtual pages within a segment and make PPM_unmap and PPM_reclaim_local calls for each, especially if the segment is large and the number of actual pages provided as backing for it is unknown. For small segments, the VSM cannot deallocate the mapping because other segments within the same virtual page (and hence mapped to a common physical page frame) may be still active. The VSM splits the deallocation problem into three cases.

The dirty-segment buddy list is used in the reverse manner of allocation - segments which are deallocated are coalesced with their buddies to try to form a single segment of as large a size as possible. In the case of deallocating segments smaller than a single virtual page, no unmapping is performed by the VSM initially. Instead, as the segment is coalesced with other dirty segments, the VSM waits (perhaps for many deallocations to follow) until a segment the size of a virtual page is finally formed. This means that many small segments which had resided within the same virtual page have all been finally deallocated. The VSM can make a single pair of calls (PPM_reclaim_local(PPM_unmap(vpn)) - that is, return the page previously used to back the deallocated virtual page to the local page pool) at this point, passing to the Physical Memory Manager the virtual page number of the segment.

For segments of moderate size (smaller than the number of physical pages on a node) which are deallocated, the above design will not work since the segment already spans many pages. Instead, unmap and reclaim calls are made for each page within the segment. Finally, for very large segments, the PMM must be called to unmap the entire segment, which requires that the PMM search for all entries in the page table which match a *range* of virtual pages (not just a single one) and remove all such mappings. This is especially efficient for very large segments, since the number of pages which need to be tested is limited by the amount of on-node physical memory.

and [6] explain, this may lead to wasting both virtual and physical address space since a little less than half of the entire segment may go to waste (e.g. a call to allocate a segment of 129 bytes will return a segment of at least 256 bytes in size). The waste of virtual address space is not a great problem, given the large size of the machine's address space. Physical address waste is limited a maximum of a single page, due to the policy of on-demand page allocation. That is, since only those virtual addresses which are targetted by a memory operation require physical backing, having a large number of allocated but unused virtual pages at the end of a segment does not cause wasted physical page frames to be allocated.

A search of available segments in the clean buddy list then begins for a segment of the appropriate size. This is a simple procedure given that the requested segment size is known - the free-segment array is indexed to find if any segments of the needed size exist. If the array entry is non-NULL, the linked list of segments of the requested size is modified as a segment is popped off the list. The pointer to the newly-allocated segment is returned to the caller. If no segments of the needed size exist, the allocator begins looking for larger segments, simply by moving up in the free-segment array, looking for linked-lists of larger and larger segments. Any larger segment that is found can be repeatedly split into two until the correct size segment is once again available. Leftover segments are added to the buddy list in the process, for later allocation.

The on-demand allocation of page frames simplifies VSM implementation since no mappings from virtual pages within the allocated segment to frames need occur - the LTLB Miss Handler will perform those tasks as each virtual page is touched.

## 5.3.2   Deallocation

As mentioned previously, virtual segments which are deallocated are added to the dirty-segment buddy list and need to pass a garbage-collection phase before being added to the clean list. However, an initial unmapping phase must occur to remove any virtual-physical mappings used by the segment, in order to free up physical memory. The need to deallocate physical backing from a segment actually poses a problem because the only data structure which lists all allocated physical page frames

49

## 5.4 Design Issues

For a machine with a large virtual address space, such as the M-Machine, buddy list allocators are quite efficient because they can quickly manage segments of memory which vary greatly in size. The fact that segment-size is encoded directly into all M-Machine pointers make this scheme even more efficient - a call to deallocate a segment uses the segment-size field to determine which low address bits of the pointer to ignore, and which high bits to use when searching for a segment's buddy.

The unmapping of backing page frames for segments seems the most inefficient aspect of the VSM design, and can be improved if a bitmap of which virtual segments actually have physical backing is maintained for each segment which is allocated. This increases overhead, however, and requires more storage space for such bitmaps. It is not clear, for example, how to maintain a mapping-bitmap for a segment of $2^{24}$ bytes. Such a segment spans 4096 page frames, requiring 64 words for a bitmap. The advantages of the current design lie in the fact that once a segment has been allocated and returned to a user thread, all information pertaining to its existence is no longer maintained by the VSM, until a call is made to deallocate it. At that point, the segment itself is provided to the VSM, which may add it to the dirty-segment list and start keeping track again. In fact, the reason for maintaining the dirty segment list (as opposed to a more simple linked-list of deallocated segments) is not for performance improvements, but rather for storage efficiency - fewer individual segment-information nodes need to be maintained if dirty segments are naturally coalesced. If two buddies are combined to form a larger segment in the dirty segment list, this frees up one more segment-information node which may be reused by the system.

# Chapter 6

# Thread Management

In traditional operating systems, a *process* represents a basic vehicle for executing code. Processes may be composed of threads which cooperate and share an address space and any special structures assigned to their collective process by the operating system. In MARS, there is no real concept of a heavyweight process. Since all privileges are granted through pointers given out by the system, all threads are protected from each other, yet any subset may cooperate on a task as well.

The MARS thread mananger is responsible for allocating and destroying user-level threads, scheduling threads to run in the available user thread slots, and managing interthread synchronization through the `tSignal` and `tSleep` interfaces. Threads which synchronize through explicit message-passing or shared-memory have no need for the thread manager to aid in their communication. The sleep and signal interface allows multiple threads to sleep on a single signal and be all awakened when it arrives, and even for a single thread to provide a signal mask, so that it is possible to group signals into categories and allow threads to pick which types of signals they wish to receive.

Instead of using some integer to identify each process (thread) which has been created by the thread manager, a *context pointer* is used instead. A context pointer is a pointer of type **key** whose address portion names a virtual memory segment which contains state information about the thread (the thread context.) Since this pointer cannot be used to read or write memory, it may be returned to user-level threads

as a magic cookie, identifying a particular thread. When an operation needs to be performed on the underlying thread state, a privileged system function may simply modify protections on the pointer from *key* to *read-write*, without the need to index into a process table. As will become evident, context pointers are used extensively within the thread management system to identify and track threads.

## 6.1  System Calls

This section describes the system calls available to user threads for accessing thread manager functionality. All of these system calls may safely execute as priviledged code in user thread slots since they do not modify any hardware state. Table 6.1 lists the common thread manager calls.

This set of system calls provides a great deal of functionality to user threads with a very simple interface. An example of how these calls are used is given in figure 6-1. In this example, the main parent thread spawns off a child to execute the function foo and then sleeps on a T_CHILD_EXIT signal, waiting for the child to complete. There is no explicit tSignal, because the signal is performed by the tExit function which the parent passed to the child. The _getDP function returns the parent thread's data pointer so that the child may share all of the parent's data structures.

More complex examples of system call use will be given later in this chapter.

In addition to the above system calls, several internal functions of the thread manager are invoked by the Event Handler and Message handlers. These include actual scheduling, and low-level signal reception.

## 6.2  Data Structures

The thread manager uses a structure called a thread context to store information about each live thread on its node. A signal table is used to manage the signal/sleep interface. Finally, pointers to chains of thread contexts maintain information on active threads. These structures are described in this section.

53

| Function | Description |
|---|---|
| `void *tFork(void *IP, void *DP, void *retIP, int numargs, void *parent, ...)` | Creates a new thread which will begin execution at address `IP`. The thread's data pointer is set to `DP`. When the thread exits, it will jump to `retIP`. The number of arguments passed to the function at `IP` is given in *numargs*, followed by the arguments themselves. *parent* is usually left NULL. This function returns a key pointer identifying the newly-created thread. |
| `void tExit(int retval)` | Standard exit procedure usually passed as the retIP to `tFork`. Signals its parent thread with a T_CHILD_EXIT signal and return-value *retval*. |
| `void *tSpawn(int numargs, void *IP, int node, ...)` | Forks a thread on remote node given by *node*. The data pointer is the same as the thread which called this function. The forked thread will start executing at `IP` and signal its parent when done. The number of arguments being passed to the function at `IP` and the argument list itself is also given. Returns a key pointer identifying the spawned thread. |
| `int tSleep(void *sigword, int mask)` | Puts the calling thread to sleep until a signal arrives which targets the *sigword*. The *mask* allows the calling thread to only be wakened by a subset of all signal arriving for the signal word. A mask of 0 will always match a signal. Returns the data which was send to the signal word (see `tSignal`.) The signal word must be a key pointer. |
| `int tSignal(void *sigword, int data)` | Attempts to wake all threads sleeping on *sigword*. The *data* is the data returned to all matching sleepers. If no sleepers are found, a dormant signal is recorded. The signal word must be a key pointer. |

Table 6.1: Thread Manager system calls

## 6.2.1  Thread Contexts

The thread manager defines a thread context data structure which is used to store information about each live thread. Several linked-lists of thread contexts group these threads into collections of running, pending, and kill threads. Running threads are the user-level threads actually occupying V-Thread slots on the manager's node. Pending threads are waiting to be scheduled to run on the hardware. Blocked threads are sleeping on a signal and should not be swapped into a thread slot until wakened

```c
int foo(int i, int j) {
    int x = 0;

    printf("This is function foo!\n");

    printf("let's calculate i + j : %d\n", i + j);
    printf("foo exiting");
    return i + j;
}

int main(int argc, char **argv) {
    char *mydp;
    void *child;
    int i;

    mydp = _getDP();       /* _getDP returns the thread's own data pointer */
    child = tFork(foo, mydp, tExit, NULL, 2, 1, 10);
    printf("main: forked foo (child pointer is %p)\n", child);

    i= tSleep(child, T_CHILD_EXIT);
    printf("main: woken with signal 0x%x\n", i);
    return 0;
}
```

Figure 6-1: Sample Thread Management system call usage

(they are stored implicitly in a signal table described later). Kill threads are waiting to be garbage-collected and removed from service. Together with running threads, kill threads may occupy hardware thread slots, but should be evicted by the thread scheduler.

Figure 6-2 shows a C structural definition of a thread context. The main sections of the context structure are the individual H-Thread contexts, (which define the entire register state of the H-Threads that compose the user thread), global thread state information, and linkages to other contexts.

The HContext structure simply contains space for all of the integer, floating-point, and condition registers of a particular H-Thread, the four restart instruction-pointers (used when installing a thread for execution), hardware and software memory-barrier counters (count how many memory references the thread still has outstanding in the system), and a scoreboard of which registers are vacant.

```
struct ThreadContext {
    struct ThreadContext *Next;
    struct ThreadContext *Parent;
    struct ThreadContext *Sibling;
    struct ThreadContext *Children;

    struct HContext hthreads[4];   /* register state for each H-Thread */

    int VSlot;
    int flags;                     /* hFull and hIssue bits IIIIFFFF */
    int SCC  ;                     /* stall-cycle counter            */
    int SCL  ;                     /* stall-cycle limit              */

    int signalData;               /* data passed when thread woken   */
    int need_to_block;            /* thread is blocked for a signal  */
    int need_to_wake;             /* signal has arrived              */
    int need_to_sleep;            /* thread has asked to sleep       */
};
```

Figure 6-2: Thread Context data structure

Global state information records which H-Threads of the user thread are active and may issue. When a thread is first forked, only the first H-Thread is active. If the thread spawns other H-Threads to neighboring clusters, this value will change. Thread flags are composed of eight bits in two 4-bit bitmaps - called hFull and hIssue. The hFull bitmap records which H-Threads are part of the V-Thread represented by the thread context. The hIssue bitmap is used as a mask to tell hardware which H-Threads may issue operations down their cluster pipelines. Special state information used in the signal/sleep implementation is also part of global thread state. The signalData field records the data word with which a thread was wakened. The three state bits of need_to_block, need_to_wake, and need_to_sleep are used by the scheduler to help decide which of the pending/running lists is to receive this thread. These state bits will be discussed in detail in the section on signalling. Finally, the thread Stall Cycle Limit (SCL) and Stall Cycle Counter (SCC) are used by the the M-Machine hardware to generate events if a particular user-level thread has been stalled and unable to issue for a certain number of cycles.

The linkages (Next, Parent, Sibling, and Children) allow thread contexts to

56

Figure 6-3: Context Linkages

be threaded onto several linked-lists at once. The main pointer is **Next**, which is used in the running, pending, and kill lists mentioned above and described in detail in a later section. The Parent pointer points to the thread's parent. Usually, the parent is the thread which tFork'ed the thread, although a different parent may be substituted (this is the **parent** argument to the tFork call). The Sibling pointer is a secondary linked list, which winds itself through all of the children of a particular parent thread. That is, even if the children of a particular parent are strewn around different pending/running/kill lists, this single list can identify all of the children of the parent regardless of where they are. This makes it easy to find and kill all children of a particular parent thread, without needing to look through all lists of threads (looking for contexts with a particular parent). Finally, the children pointer is the head of the Sibling list, which resides with the parent. Figure 6-3 makes this structural arrangement more explicit.

In this example, the first thread on the pending list is the parent of three threads - one also on the pending list, and two others that are running. One of its children is the parent of a thread which is on the kill list.

Thread contexts reside in virtual address space, and are dynamically allocated by the tFork call. Since all virtual addresses are unique across the entire machine, a thread context unambiguously identifies a thread to all operating system components across all nodes of the machine. All threads may access their own context pointer through a call to _getMyTC, and the context pointer of their parent with _getParent. The pointers that are returned are key-type pointers, to prevent user threads from actually modifying thread state.

## 6.2.2   Signal Table

In order to maintain information on which threads have performed signals and which threads have tried to sleep, the thread manager uses a chained hash table of signal entries.

A signal entry records information about a thread which has asked to be put to sleep, or a signal which has been made before any thread has slept on it (see figure 6-4.)

```
typedef struct se {
    struct se *next;
    int signal_word;
    int signal_data;
    struct ThreadContext *sleeper;
} signal_entry;
```

Figure 6-4: Signal Entry

If a thread has slept on a signal word, the two arguments to the sleep call (signal_word and mask) are recorded along with the thread context of the thread making the sleep call (*sleeper* in the signal entry. If the entry is recording a signal for which no thread has slept yet, sleeper is NULL and the signal_data is the actual data passed to the tSignal call.

58

```
signal         Signal Hash Table
word

                                  chain of entries – searched for singal_word match

probe

              next                  next
              signal_word           signal_word
              signal_data           signal_data
              sleeper                sleeper


              next
              signal_word
              signal_data
              sleeper
```

Figure 6-5: Signal Hash Table Structure

Signal entries are split into chains and referenced from the signal hash table, to improve lookup speed. The unique signal_word is hashed and identifies the chain, which may then be searched for matching entries. Signal entries may be dynamically allocated in a manner similar to thread contexts, or a fixed number may be statically allocated at compile-time into the runtime system (similar to what is done by the virtual segment manager.)

## 6.2.3 Thread Lists

The low-level scheduler employs thread lists, headed by pointers to *Pending*, *Running*, and *Kill* lists. All threads active on a node belong to one of these lists, or have *sleeper* entries in some signal table (effectively the collection of blocked threads). This guarantees that thread manager components have a way to find all active threads on the node by following these structures. The *Next* pointer in a context lets it be threaded in one of these lists. A thread may be in only one of these lists at a time.

# 6.3  Implementation

When the thread manager is initialized, it sets up a blank signal table and resets the running, pending, and kill thread lists to contain a single running thread - the bootstrap. Calls to the manager's system calls will begin modifying these structures. It was briefly noted that the thread manager is really composed of system calls executing in user thread slots and a low-level scheduler tied into the event-handler system. For this reason, the thread management implementation uses a producer-consumer model for servicing requests. User-accessible system calls invoke functions which set up and sometimes modify thread state. After certain global data structures are modified, the event handler is signalled through its software job queue to perform the low-level scheduling tasks. This two-phase design simplifies the implementation of individual thread manager components. It also allows thread manager subsystems to execute in conjunction with the scheduler without relying on locks to serialize access to common data structures - all data structures which are modified by the portion of the thread manager which runs in the event handler slot do not interfere with other thread manager functions.

In general, producers create or modify thread contexts which are then added to the running, pending, and kill lists by the scheduler (this list modification is performed when the event handler responds to certain signals). The scheduler examines these lists each time it is invoked and performs lowl-level functions such as thread eviction and installation. The following sections describe the producer's contribution to handling system calls. It is important to note here that in all critical sections of the portion of the Thread Manager that runs in user thread slots, a lock called the userthreadLock is used to serialize access to global data structures among user threads. This lock is not accessed by the low-level scheduler, and hence does not cause it to block in any of its activities.

### 6.3.1 tFork

The tFork function needs to allocate a new thread context by calling on the virtual segment manager, and fill an initial H-Thread with information passed to it. It allocates a new thread stack, again calling on the VSM, and pushes arguments on the stack exactly as the called thread expects to see them. The return pointer is set up as well, so that the exit function passed to the tFork is the last function executed. Parent/child/sibling linkages are updated to reflect the fact that a new thread has been created and that it belongs to some parent. If the parent pointer passed to the fork call is NULL, the thread executing the fork call is considered the parent (this is the common case). Remote parents are a special case, which are handled within the tSpawn implementation. Finally, the event handler is signalled to add the new thread context to the pending list. This signifies that the thread is ready to execute and is waiting to be scheduled into an available thread slot.

### 6.3.2 tExit

The tExit call must mark its own thread for termination since it is executed within the very user thread which is trying to exit. First, the thread calls tSignal on its own context pointer with a return value of T_CHILD_EXIT. Any thread waiting for this particular child to exit (most likely its parent) will be wakened.

The sibling list is modified to reflect the termination of this thread. The event handler is then signalled to add the thread to the kill list. The event handler removes the thread from the running or pending lists and adds it to the kill list. Finally, tExit blocks on an empty register to prevent stealing any more execution cycles. Eventually, the scheduler will be invoked and terminate the thread which had been added to the kill list.

### 6.3.3 tSignal

The tSignal system call is used by a thread to signal another thread, passing it a 64-bit data word. Signals are made upon *signal words*, which are key pointers given

out by the operating system. The most common signal words are the thread context pointers exchanged by the parent and child during a tFork call. Other signal words may be obtained simply by calling on the operating system to demote the protections of a virtual-memory read-only or read-write pointer to key.

The tSignal call takes a signal word and a 64-bit data word as arguments and determines which signal table to examine. If the address defined by the signal word is mapped to the thread manager's own node, the local signal table is examined. Otherwise, a message is sent to the node where the signal word is mapped, and a TM local to that signal table is invoked. The TM determines whether an address is remote or local by making a call to _sysGPRB, a function which performs a GTLB probe and returns the node number to which an address is mapped. This allows threads on different nodes to signal each other and for all thread managers to quickly decide which signal table needs to be referenced.

Once a local TM is invoked to examine the signal table, the signal word is used as the input to a hash function and an index into the signal hash table is calculated. This index identifies a chain of signal table entries which is to be searched to find a match or matches (for multiple sleepers) on the signal word. In order for a signal to match an entry, it must meet three criteria.

1. the *signal_word* field of the entry must match the signal word passed to tSignal

2. the signal_data [mask] field in the entry bitwise ANDed with the signal_data passed to tSignal must be nonzero (unless the mask is 0, in which case this criterion is always considered satisfied)

3. the *sleeper* field of the entry must be non-NULL.

For each match that is made, the thread identified by the sleeper context pointer is wakened (this process is described below.) Once all sleepers have been wakened, the signal operation has completed. If no sleepers were found, a *dormant* signal entry is added. This means that the signal is added to the signal hash table and waits for a sleeper to come along, at which point the thread which attempted to sleep on the signal is automatically wakened. Such dormant signals are added to the ends of

62

the signal chains, to handle cases where multiple dormant signals for the same signal word are added. In these cases, the signals are meant to be popped off in a FIFO manner, until they are all used up.

For each thread context which needs to be wakened, the tSignal system call must decide whether the wakening occurs locally or remotely. Once again the TM probes the GTLB, this time to determine whether or not the sleeper thread context is mapped to the local node. If the thread context is remote, a *Wake* message is sent to the appropriate home node of the thread. Otherwise, the event handler is signalled to set a thread's wake data. This causes the thread context's *signalData* field to be written with the signal data passed to tSignal, and the *need_to_wake* field set to true, signifying that if the thread happens to be blocked, the scheduler should move it to the pending list.

### 6.3.4 tSleep

A user thread calls tSleep when it wishes to block, stopping execution until a signal wakes it. This is especially useful when a thread has spawned off some children which are to perform long-latency operations and wishes to be informed when these operations have completed. Although it is possible for the parent thread to spin on global memory locations waiting for child thread to modify them, this is extremely inefficient if the child processes are expected to take a long time to complete their operations, and the parent has no other work to perform.

For this reason, the calling thread identifies itself as sleeping on a particular signal word, and also passes a mask as data. This mask is used to filter out certain signals to the signal word which the sleeping thread does not wish to see (as described above). As in the case of tSignal, the signal word is used to probe the GTLB to find the home node of the signal table. If the signal table is remote, the thread asks to be put to sleep locally and sends a message to be added to the remote signal table. It is important to note here that it is possible for the message to arrive and a dormant signal to be found which would cause a wake message to be returned, all before the local TM is able to put this thread to sleep. The need_to_sleep, need_to_wake, and need_to_block

63

Pending  tInstall  →  Running

tEvict

SYStSleep

eh(EVENT_SLEEP)
tPutToSleep

Pending
need_to_sleep     tInstall  →     Running
need_to_sleep    eh(EVENT_WAKE)  →

Running
need_to_sleep
need_to_wake

tEvict

eh(EVENT_SLEEP)
tPutToSleep

eh(EVENT_SLEEP)
tPutToSleep

eh(EVENT_SLEEP)
tPutToSleep

Pending
need_to_sleep
need_to_wake

need_to_block

Running
need_to_block

tEvict

eh(EVENT_WAKE)

eh(EVENT_WAKE)

eh(EVENT_SLEEP)
tPutToSleep

Pending
need_to_wake
need_to_block

Running
need_to_wake
need_to_block

tHandleSignals

tHandleSignals

Figure 6-6: State Transitions in Signal/Sleep Implementation

bits define state-transitions to handle such cases. Figure 6-6 shows a state-transition diagram where a thread state is a function if its *need_to_xxx* bits and whether it is running, pending, or neither. Transitions occur as a result of the low-level scheduler performing routine scheduling tasks, or being invoked as a result of signals to the event handler (EVENT_SLEEP and EVENT_WAKE). Certain functions are automatically invoked as a result of these signals (such as tPutToSleep and tHandleSignals).

Finally, whether as a result of a Sleep message from a remote TM or the fall-through case of a local tSleep call, the TM needs to add a sleeper for the signal word to the signal table. Again as in the tSignal case, the signal word is used as the hash input to find a chain of signals. The chain is examined for any *dormant* signals to this

word. If a dormant signal is found and the data within it filters through the mask provided by the calling thread, the thread is immediately wakened. If the thread was local to the signal table, the data is returned directly to the thread without the thread having ever been put to sleep. Otherwise, a wake message is sent to the home node of the sleeper thread.

If no dormant signal entries are found, a new sleeper entry is made. Finally, if the TM is still executing locally, it makes a call to sysSignalSleep, which asks the scheduler to move the thread off the running list (if possible) and consider it blocked until a signal arrives. At the same time, this action causes the thread to empty the return-value register and block on it. Whenever this register is written (as a result of the scheduler restarting a thread which is being wakened by a signal) the thread will resume execution and return a value to the caller of tSleep.

Figures 6-7 and 6-8 show an example of the use of signal and sleep calls for in-terthread synchronization. The parent thread forks a child called longprint, which in turn forks off longprint_child. Longprint then waits for its child to signal it. Mean-while, the main parent sleeps on a signal from longprint. longprint_child signals its parent and then goes to sleep, waiting for longprint to signal it. At this time, both main and longprint are sleeping on the same signal word. When longprint is wakened by its child's signal, it signals to its own threadcontext pointer, waking both its child and its parent. Finally, longprint waits for its child to exit before exiting itself. The main thread waits for longprint to exit.

## 6.3.5  tSpawn

The tSpawn system call is a good example of how lower-level thread manager prim-itives may be composed to form a more useful function. A tspawn is essentially a request by the user to fork a thread on a remote node and still have the child's thread context be returned to the parent. The tSpawn implementation first creates a nonce which will be used for a signal/sleep pair. In the current implementation, this nonce

65

```
#include <stdio.h>
#include "syscalls.h"
#include "tsignal.h"

int main(int argc, char **argv) {
    void *child;
    int i;

    printf("Sample signal/sleep program\n");
    child = tFork(longprint, _getDP(), tExit, 6, NULL, 1, 2, 3, 4, 5, 6);
    printf("main: forked off %p\n", child);
    i = tSleep(child, T_ALL_SIGNALS);
    printf("main: woken with 0x%x\n", i);
    /* wait for a while */
    for (i = 0; i < 900; i++) ;
    i = tSleep(child, 0x100);
    printf("main: woken with 0x%x from child %p exit\n", i, child);
    return 0;
}
```

Figure 6-7: Sample signal and sleep system call usage: main thread

is simply a newly-allocated segment of virtual memory used and then discarded.[1] A message is then generated, and the nonce and arguments to the spawn are sent to the destination node. Finally, the calling thread performs a tSleep on the nonce, waiting to be notified when the new thread has been created. It expects the return value of the tSleep (the data when it is signalled with tSignal) will be the thread context of the new child.

On the receiving node, a message-handler dispatch function processes the tSpawn request. The Spawn message is unpacked and arguments formatted for a tFork call. This time, instead of a parentTC of NULL being passed, the TC of the remote parent is substituted (this was passed in the message, along with argument list, IP, and so on), allowing linkages to be set up correctly. After the tFork completes and returns a thread context, the message-handler performs a tSignal on the nonce passed within the spawn request message, passing the child thread context as data. This eventually

---

[1]Since the VSM returns pointers as Read/Write, a *demote* call is made to change the protections to key pointer.

```
int longprint_child(int i, int j) {
    int sleepval;

    printf("longprint3_child: i * j = %d\n", i * j);
    tSignal(_getSelfTC(), 0x112);

    /* now wait until longprint signals me */
    printf("longprint_child: going to wait for longprint to signal me\n");

    sleepval = tSleep(_getParent(_getSelfTC()), T_ALL_SIGNALS);

    printf("longprint_child: woken with 0x%x and exiting\n, sleepval);
    return 4;
}

int longprint(int i, int j, int k, int l, int m, int n) {
    int x = 0;
    void *child;
    int sleepval;

    printf("longprint: %d, %d, %d, %d, %d, %d\n", i, j, k, l, m, n);

    child = tFork(longprint_child, _getDP(), tExit, 2, NULL, 5, 11);
    if (child) {
printf("longprint: forked off %p, and sleeping on it\n", child);
sleepval = tSleep(child, T_ALL_SIGNALS);
printf("longprint: woken with 0x%x from child %p\n", sleepval, child);

        for (x = 0; x < 200; x++)
  if (!(x % 20))
    printf("longprint: %d\n", x);

tSignal(_getSelfTC(), 0x223);

/* sleep on child exiting */
sleepval = tSleep(child, T_CHILD_EXIT);
printf("longprint: child %p exited\n", child, sleepval);
    }

    printf("longprint exiting");
    return 1;
}
```

Figure 6-8: Sample signal and sleep system call usage: child threads

wakens the calling parent who receives the child thread context just like the return value of a tFork.

## 6.3.6 Scheduler

The scheduler portion of the Thread Manager runs as part of the event handler - responding to requests placed in the software job queues. Requests are summarized in table 6.2. The generic EVENT_SCHEDULE is the most interesting to cover because it encompasses the important tasks of installing and evicting threads.

| Request | Arguments | Description |
|---------|-----------|-------------|
| EVENT_SCHEDULE | | Perform generic scheduling: wakes threads which have **need_to_wake** set. Terminates threads on the kill list. Attempts to install threads on the pending list, perhaps evicting running threads to make room. |
| EVENT_SLEEP | tc | Puts thread identified by thread context pointer **tc** into a blocked state. If the thread is already running, it is moved to the front of the running queue so it is the first to be swapped out if an eviction is necessary. If thread is on the pending list, it is removed from the list so as not to be mistakenly installed during scheduling. Sets thread's **need_to_block** state bit. |
| EVENT_FORK | tc | Adds thread identified by thread context pointer **tc** to the Pending list. |
| EVENT_WAKE | tc data | Sets the **need_to_wake** state bit of the thread identified by **tc**. Sets the thread's **signalData** field to **data**. If the thread is not currenly occupying a thread slot (running) it is added to the pending list. |
| EVENT_KILL | tc | Adds the thread identified by **tc** to the kill list. |

Table 6.2: Thread Manager system calls

The scheduler completes three tasks when asked to perform scheduling.

## Cleaning Killed Threads

First, all threads in the kill list are popped and terminated, if possible. Their thread context is freed, the hardware thread slot state that they occupy (if they are still

installed in a thread slot) is reset and the thread slot marked as unoccupied. If threads which are popped off the kill list still have outstanding memory events which are to be resolved in software or outstanding hardware events, the threads may not be terminated and are added back to the kill list. A check in the code which runs through the kill list makes sure that recirculating threads into the kill list does not cause an infinite loop of pushes and pops.

### Signal Handling

The thread scheduler then deals with outstanding signal-handling. A thread which is (1) in the pending or running lists, (2) has its need_to_wake state bit set, and (3) has its need_to_sleep bit unset, is set active by copying signalData into the appropriate return register. If it is occupying a thread slot, the thread's return-register (i10) is written with the contents of the context's signalData field directly (using a configuration-space write). Otherwise, the register is modified within the thread context and the empty bit for that register set to full so that the register can be read the next time that the thread is installed into a thread slot. In both cases, the need_to_wake bit is reset.

### Installing Threads

In its third task, the scheduler pops a thread off the pending list (the candidate) and attempts to install it into a free user thread slot. If no free thread slots exist, a thread is popped off the running list and evicted (if possible). Eviction involves halting all H-Threads which are issuing within the V-Thread - accomplished by writing to the thread flags region of configuration space mapped to the hardware thread slot which the thread occupies. The thread flags are modified to zero out the hIssue bits for the thread. Then, for each active H-Thread within the V-Thread, all of the register-file state is copied into the thread context. Four H-Thread IP's for use in the thread-restart process are read out from each cluster. Finally, state like software and hardware membar counters are updated. Once eviction succeeds, the thread context is pushed to the end of the pending list.

69

When a free thread slot has been found for the candidate, a reverse of the eviction process begins. First, the candidate's hFull thread flags are written into the configuration space mapped to the thread slot into which it is being installed. These flags set the hFull bits for all H-Threads which are to run within the candidate. This has the effect of resetting all thread state within individual clusters. This is a safe procedure since no hIssue bits are set, so the thread will not attempt to issue from a non-existing IP. Then, individual H-Thread state is updated by reading thread context data and writing into the thread slot through configspace. After all register-file and membar counter state has been written, a series of 4 IP writes are made for each H-Thread. These writes prime a hardware restart engine which fetches instructions and can restart a thread. Lastly, the candidate is pushed to the end of the running list.

# Chapter 7

# Memory-Coherence Management

This chapter details the M-Machine's software-based memory-coherence protocol. As mentioned in previous chapters, the software implementation is closely tied to other OS components, such as the Physical Memory Manager and Thread Manager. The memory-coherence system provides the view of a single globally-shared virtual address space which is accessible by user threads independent of the node on which they execute. That is, any thread which performs a memory-reference to a word of virtual memory will have that request satisfied even if the segment of virtual memory is not mapped to the thread's *home node*. Each word of virtual memory is mapped, through the GTLB and a software Global Page Table (not implemented in the current runtime system), to an M-Machine node - the home node of that data. For purposes of the memory-coherence protocol described in this paper, the granularity is on an 8-word block basis (words in each 8-word block of memory must have the same home node in common). The term "memory block" (or just "block") refers to an 8-word section of virtual memory, the size of an individual cache-line, which may be shared among several nodes. In the rest of this chapter, the *home node* means the node to which a particular block of memory is mapped, and a *requesting node* is used to identify a node which wishes to access data from the home node. In rare instances, the home and reqesting nodes may be the same.

In broad terms, the memory-coherence manager allows threads to transparently read and modify blocks of memory which are not mapped to their local nodes. Load

and store operations which attempt to access off-node data fault to software with *block-status* misses (BSM). A portion of the memory-coherence manager (MCM) which runs in the event-handler thread enqueues BSMs into a software event table, and sends out request messages for accessed blocks. Message-handing functions in the P0 and P1 Message Handler threads respond to request messages by modifying local coherence directories, local cache, and the LTLB, and send blocks to requesting nodes. Local message handlers on requesting nodes accept responses to the MCM requests sent out by the event handler and install blocks locally. The cache and LTLB of the requesting node is modified, and events pending to the block which were enqueued in the software event table are popped and satisfied at this time.

The following sections briefly describe the internal functions used by the MCM, present data structures employed by the home and requesting sides of the coherence protocol, and details the MCM implementation, including a state-machine model for tracking individual memory blocks.

# 7.1 Internal Functions

The MCM is split into three components which run as part of the event handler, and the two message handler threads. Table 7.1 lists the functions executed by the event handler thread. These functions may be grouped into three categories - functions which are executed as part of the requesting node's initial handling of blocks-status misses, functions which are executed in proxy for a requesting node's P1 Message Handler, and functions which are executed in proxy for a home node's P1 Message Handler. The proxy functions are actually wrapped up in the event handler's routine which services the software job queue, and are therefore shown in a stylized manner which does not actually appear in the source code.

The home node's MCM handles incoming requests for blocks, as well as acknowledgements for block invalidations which it sends out. These functions are outlined in table 7.2.

Lastly, the requesting node's MCM handles home node responses to the requests

that were sent out by its own event handler. It also responds to invalidation messages coming from the home node. These functions are outlined in table 7.3.

## 7.2 Data Structures

Each node's MCM uses two data structures - one for managing blocks for which the node is a home node, and the other for tracking requests for blocks which the node makes in its capacity as a requesting node. The home-node information is stored in a coherence directory, while requested blocks are stored in a software event table.

### 7.2.1 Coherence Directory

The coherence directory is simply a linked list of lists. Each toplevel entry in the list contains the address of a block of memory which is shared by at least one node, state information about the block, and a list of nodes which share that block (these are nodes to which this block has been sent). Blocks may be in one of three states. Read shared blocks may have multiple nodes which share them. Exclusive shared blocks may only be held by a single node. Transitioning blocks are in the process of being revoked from all sharers because a conflicting request for them has been made (a request for a readonly or exclusive copy for a block which was held exclusive by a different node, or a request for an exclusive copy if the block was held readonly by at least one node).

Functions are provided to add a new sharing node for a particular block (`CCDirectory_addSharing`) to the directory, and remove a sharing node from the list of nodes sharing a particular block (`CCDirectory_popSharing`). Other functions access and modify block state.

This current implementation is not efficient in terms of search time. Future implementations of the directory should use a chained hash table to access shared addresses with greater speed.

73

**Software Event Table**

```
                        (  .-- -- --- -  |   Event Table Entry
                        |                |   .........................
                        |                |   :  VPN                  :
         Number of      |                |   |                       |
         physical       {                |   : status                :
         pages used     |  .             |   |                       |
         for backing    |                |   : queue pointer         :
                        |  .  .          |   :.......................:
                        |                |
                        |.--  .          |
                        (
```

**▼ Event Queue Node**

```
┌─────────────────┐  ┌─────────────────┐  ┌─────────────────┐  ┐one event queue
│ next            │  │ next            │  │ next            │  │entry per shared
│ address         │  │ address         │  │ address         │  │8-word memory
│ state           │  │ state           │  │ state           │  }block
│ invalidate ptr  │  │ invalidate ptr  │  │ invalidate ptr  │  │
│ events          │  │ events          │  │ events          │  │
│ tail            │  │ tail            │  │ tail            │  ┘
└─────────────────┘  └─────────────────┘  └─────────────────┘
```

```
►┌─────────┐  ►┌─────────┐  ►┌─────────┐
 │ next    │   │ next    │   │ next    │
 │ header  │   │ header  │   │ header  │
 │ address │   │ address │   │ address │
 │ data    │   │ data    │   │ data    │
 │ CP      │   │ CP      │   │ CP      │
 └─────────┘   └─────────┘   └─────────┘
 ‿‿‿‿‿‿‿‿
 Individual
 Block-status
 Miss event
```

Figure 7-1: MCM Software Event Table

## 7.2.2 Software Event Table

The software event table is used by the cache-coherence manager to record block-status miss events which are being handled in software and maintain information about the status of blocks which have been requested from a home node. The table contains three-word entries and implictly maps physical page frame numbers to virtual page numbers and queues of requests. That is, the $ith$ entry in the table refers to the $ith$ page frame on the local node which is used as backing for remote virtual memory blocks. This table is statically-sized at link time, or at the time that the Physical Memory Manager is asked to reserve a range of frames for backing of remote memory with the PPM_local2remote function. Figure 7-1 shows event table layout.

The event table is probed with both a virtual address and a physical page frame number to access event queues for that block. The frame number is used to directly index into the table and locate a table entry. The table entry's virtual page number

74

field is compared against the page number portion of the virtual address. If the numbers match, the pointer to the entry's queue of requests is followed (the structure of the queue is described below). If no vpn match is made, the frame number is considered stale, and a page-table probe (PPM_lookup) must be performed. In this way, the software event table functions almost like a reverse page table, except that information that it holds may be stale and inconsistent with the local page table. State information is associated with each table entry as well. Currently the only state information is a bit which informs the caller that the physical frame associated with the entry is marked for eviction, and no new events should be added to its queue.

The last component of the event table entry is the software queue entry pointer. This identifies the head of a linked list of queue entries. Each queue entry represents an 8-word memory block for which event information is stored. There may be at most 64 such entries in any linked list since there are at most 64 different blocks within a virtual page. Each entry contains information on the state of the block (to be discussed later), a 64-bit invalidation pointer if the home node has requested that this block be invalidated and returned [1], an address field which is used to identify which of the 64 blocks this block represents [2], and pointers to the head and tail of an event list for this block. The event list is a collection of entries which represent block-status miss events which have been removed from the hardware event queue by the event handler. Each miss event entry contains all four words which compose a block status miss, and a *next* pointer for use in linked lists.

Use of these data structures will be explained when implementation is detailed. To obviate the need for dynamic memory allocation of these structures, a collection of software queue entries and miss event entries are statically allocated at compile-time and initialized into lists of available entries at runtime. Entries are popped from the lists of free entries when needed, and returned to these lists when no longer used in the event table. Since the event table is statically-sized at compile time, it also does not need any dynamic memory allocation.

---

[1]Invalidation pointers are pointers to a *yankbuffer* structure, described later in this chapter.

[2]Although the current implementation uses a full 64 bits, only 6 are necessary since the rest may be reconstructed from the virtual page number of the containing event table entry.

# 7.3  Implementation

A memory-coherence protocol needs to handle a variety of common-case memory-sharing requests, and deal properly with a number of more unusual cases which are a result of the asynchronous nature of multinode execution. This section first presents a simplified view of common-case operation of the coherence protocol, introducing how the different handlers interact and employ the data structures that were presented in the last section. The motivation for employing a state-machine model of block states is presented, along with the model. Further sections then explain handling of more subtle coherence cases.

## 7.3.1  Simplified Roundtrip Coherence Path

Figure 7-2 is helpful in clarifying the mechanisms introduced in this section.

All nodes initially start execution without sharing any remote data. Threads which reference off-node data begin the process of remote-block fetching and installation. The process begins when a thread causes an LTLB Miss, since while a page of virtual address space may have physical backing on its home node, a remote node will not have such backing. A thread (called the *faulting* thread in the rest of this section) which references off-node memory will cause an LTLB Miss with its memory reference which will invoke the Physical Memory Manager as described in chapter 4. The PMM will determine that the virtual address is a remote-address and create a new page-table entry mapping the virtual page to a new backing page frame take from the remote backing pool. Block status bits for all blocks within the page will be set to *invalid*. When the hardware retries the memory-reference, an LTLB entry will be found, but block-status bits for the block containing the referenced address will be invalid. The hardware will therefore generate a Block Status Miss event and add it to the hardware event queue. The event, similar to the LTLB Miss Event, will contain a header word, faulting address, source data if the operation was a store, and a configuration space pointer into thread state for the faulting thread. A 20-bit field within the header word contains the frame number retrieved from the LTLB at the

Figure 7-2: End-To-End Communication in Simple-Path Coherence Protocol

time that the block-status miss was generated. See table 7.4 for the event header format.

## Sending a Request

When the event handler pops the block-status miss event from the hardware queue, it determines the type of the event from the low four bits of the header word. Finding that it is a block-status miss, the event handler dispatches the event to the _BSM_xx functions which interface assembly-coded portions of the event handler with higher-level functions written in C. The assembly code then calls _EH_handle_bsm, passing it all four event words. This function uses the header's encoded physical page frame number to index into the event table and find an entry. Initially, all entries within the table will contain invalid mappings (virtual page numbers of -1). Therefore, the

event handler will not find a match between the faulting address' page and the page in the table entry. At this point, the handler decides that the page information is stale (it could have been changed between the time that the hardware determined the mapping from the LTLB and the time that the event handler had removed the event from the hardware queue) and performs a page table lookup (calls PPM_lookup). The resulting page is again used to probe into the table and again a match will not be found. At this point, the handler must deduce that the event table entry is not current, and creates a mapping, simply by writing the faulting address' virtual page number into the entry's vpn field.

Having found a valid table entry for the fault, the event handler examines the backing page's state information, to make sure that the page is not marked for eviction. Since it is not (the table is initialized so), the handler attempts to enqueue the block-status miss event. Since the page table's queue pointer is null, a new software queue entry is popped from the list of free entries and added as the head of a new list. Its address field is set to that of the faulting address with the low 6 bits masked off (indicating an entire 8-word block). A new miss event entry is also popped, initialized with the event words, and added to the event queue for the block in which the faulting address resides. The function returns certain flags which enable the caller to determine what actions to take. The *send_message* flag is set because a new software queue entry was added, and therefore this was the first reference to this block. The calling function (the event dispatch handler) then decides to send a message to the home node of the faulting address, requesting that the remote block be sent back. A MSG_ccrequest priority 0 message is sent, containing the header word and virtual address. At this point, the work of the requesting node's event handler is complete. The node must now wait for an acknowledgement to its request.

All further events targetting the block in the meantime are added to the event queue for that block so that spurious request messages are not sent. As long as there are events remaining in the software queue for a particular block, new events are added but no messages are sent.

78

## Fulfilling Requests

When it receives a MSG_ccrequest message, the home node's priority 0 message handler removes the message arguments from the message queue, packages them as function arguments, and calls the ccrequest function of the MCM. ccrequest examines the event header which was sent in the message and determines whether the request was for a readonly or an exclusive block based on the opcode of the operation that faulted on the requesting node. A ld operation results in a call to ccrequest_ld while a st or any of the synchronizing ld/st variants result in a ccrequest_st, ccrequest_stsu or ccrequest_ldsu being called.

In any case, the home node checks the coherence directory to determine what is the state the requested block. Assuming that this is the first coherence request to be serviced, the directory will return the fact that the block is unshared. In this case, the directory is modified to have the requesting node as a sharer for the block in question. If this was a store request, the store which was requested to be performed is performed locally (the opdata passed in the request message is used as the data source of the store operation). Block-status bits for the block are then changed to INVALID, and the block is read out and sent as an acknowledgement to the requesting node. In response to a load request, the block-status bits are changed to READONLY since the home node's thread can continue reading the block, and the block is read out and sent to the requesting node.

## Installing Remote Data

On the return path, the acknowledgement to the a block request returns to the requesting node as a ccreturnLoad or ccreturnStore, depending on the type of sharing which was granted (exclusive or readonly). In either case, the address and header which return in the acknowledgement are used by the MCM to index into the event table in the same manner as performed by the event handler. This time, there is a match between the entry's vpn and the vpn of the requested address (since this was correctly updated by the event handler prior to the request message being sent) and the entry's software queue pointer is followed and the queue entry for the appro-

79

priate block is found. The block contents are read out of the message queue by an assembly function and written into local memory (a backing page exists since there is a mapping in the event table from the ppn listed in the header, and the vpn in the faulting address). Block-status bits for the virtual address of the block are set properly (READONLY or READWRITE, depending on the type of sharing allowed). All events stacked up for the requested block are then handled in turn, by performing the faulted memory operations, this time on memory which has been installed locally. After all events have been processed, the event entries and software queue entry are returned to their free pools, and, if no other cache blocks have been requested for that particular virtual page, the pointer to the software queues in the table entry for the backing frame is reset to NULL.

## 7.3.2   Diverging from the Simple Case

This section begins to explore the more interesting cases which must be dealt with by the MCM. Each section will identify a case not covered in the above simplified example and ammend the actions taken by affected components. The cases will parallel the previous section in the order of the components that are introduced - starting with the event handler.

### Out of Backing Pages

In the previous section, the page frame number located by the M-Machine memory system was assumed to be a valid physical page frame. As mentioned in section 4.5.2, the ppm will create a mapping of a virtual page number to physical page frame -1 if no backing frames for remote data remain. This information may be returned in the event header of a block-status miss. **It is the policy of the MCM not to send requests for remote blocks unless physical backing is obtained first.** Therefore, the MCM first performs a PPM_lookup to make sure that a mapping hasn't been created since the block-status miss first occured. If the lookup returns a valid page, the event handler can perform the probe as before and continue processing.

On the other hand, if an invalid mapping is returned again, the event handler

makes note that cleaning of shared pages must be performed to free up a backing frame, and adds the entire event to a local software queue, effectively recirculating it so that it may continue taking a look at the event from time to time and being able to finally satisfy it when physical backing is obtained. Meanwhile, to prevent user threads from continuing to cause block-status misses and overfilling the recirculation queue, all user threads are prevented from issuing instructions (the event handler turns off their hIssue thread state bits).

In order to find pages suitable for reuse, the event handler may run through the event table, looking for entries which have no pointers to software queues of events. Such pages are ripe for eviction since no outstanding requests to their pages remain and therefore all of the shared blocks within these pages may be evicted (and sent back to their home nodes if dirty). In order to evict a shared page, the event handler performs the following actions:

1. Performs 64 *putcstat* operations, setting block-status bits for each block within the page to *invalid*. Putcstat's return value, the previous block-status bits, are used to check whether each block was dirty. Every dirty block is shipped back to the home node with a sysPushDirty call, which sends the address and the 8-words of the block to the home node in a MSG_ccreturnDirty message.

2. Calls PPM_unmap to remove old virtual-physical mapping for the virtual page being evicted.

3. Returns the backing page to the backing page chain with a call to PPM_reclaim_remote.

After a virtual page has been evicted and the backing frame is returned for reuse, the event handler makes a PPM_map call to give physical backing to a new virtual page, which was missing backing previously. Finally, the entry corresponding to the newly-acquired backing page is modified to reflect a new virtual page number, and the process of adding a new software queue entry may continue as before.

If no pages may be evicted right away (each entry in the event table has a valid software queue pointer, signifying that there is at least one outstanding event per page waiting for a block to be returned), some pages are chosed for eviction and their

state bits in the event table are set, indicating that no new events are to target these pages since they must be evicted.

In order to prevent running out of backing page frames, the event handler is designed to examine the number of page frames remaining after each event is handled. If the frame count is below a watermark, the handler must perform preemptive page eviction to free up backing frames. This may be accomplished by keeping a pointer into the event table which is advanced until a suitable candidate frame (one with a valid VPN mapping, but no queue pointer) is found. This frame undergoes the eviction process described in the steps above and may be added to the backing pool.

## Backing Page is Marked for Eviction

The case in the previous section presents another problem for the event handler. If it finds an event table entry for the faulting address and the virtual pages match, the physical page frame may be locked. If the state bit for that entry is set, the event handler is prevented from adding the new event to the software queue (although one optimization is to allow it to add the event if it targets an existing block, so that the event will be handled with all other events for that block as soon as the home node returns the necessary data) and must recirculate it. This case becomes analogous to the event handler not having an appropriate backing page, although in this particular instance no search for new backing pages is required.

A modification to the priority 1 message handler which deals with returning blocks must be made as well. When the last software queue entry for a particular virtual page has been freed and the event table entry's software queue pointer set to NULL, the message handler must check the status bit of that entry. If the status bit is set, the page is ready for eviction. Since the P1 message handler is not allowed to send out messages (this is to avoid deadlock in the machine's network) and message-sends of dirty blocks may be required when performing a page eviction, the P1MH enqueues an eviction job with the event handler in the handler's software job queue. Some time in the future, the event handler will respond to the eviction request and perform the same type of operations in evicting a page as mentioned in the previous subsection.

## Invalidations Required

Moving to the home node of requested data, the case of incompatible block sharing arises. As mentioned briefly at the beginning of this chapter, when the home node probes the coherence directory, it may discover that several nodes are sharing a block which has just been requested as an exclusive copy; or a node other than the requesting node may have an exclusive copy of the block. In both cases, all of the nodes currently sharing the block must have their shared copies revoked, before the latest request can be satisfied.

The home node performs the invalidation with the help of a new data structure - the yankbuffer. The yankbuffer records information about the request which caused the invalidation to be performed, and the number of invalidation messages outstanding. A circular buffer of pointers to free yankbuffers is accessed to acquire a new yankbuffer. This circular buffer is then used to return a yankbuffer for reuse once the invalidation process has completed. The invalidation protocol begins as follows: a new yankbuffer is acquired and the four words of request information written into it. The requesing node number is written as well, so that the MCM knows which node sent this request. Lastly, the number of nodes which currently share the block is written into the yankbuffer.

With the yankbuffer initialized, the message handler sets the state of the block in the coherence directory from shared exclusive or shared readonly, to TRANSITION-ING, signifying the fact that an invalidation of this block is in progress. The message handler begins popping nodes from the coherence directory list of sharers for the requested block and sends an `MSG_ccinvalidate` message to each. The block address and yankbuffer address are sent in each message. Once all messages have been sent, the message handler's immediate task is complete, and it is ready to handle the next incoming message. Other portions of the MCM will respond to the invalidations and cause the block to be sent to the requesting node which caused the invalidations.

As acknowledgements to the invalidation messages arrive at the P1 message handler (invoking the `ccreturnYank` and `ccreturnyankFull` functions), the yankbuffer pointer that is sent along is used to decrement the invalidation count within the

buffer. Dirty blocks which are returned in acknowledgements are copied into home node local memory.

All requests for blocks which come in while the blocks are in the transitioning state are NACKed back to their senders. This frees the home node from buffering requests for blocks locally, and instead places the burden of buffering on the network and requesting nodes, as NACKS are returned to home nodes, buffered, and new requests sent out.

Once the invalidation count reaches zero, the state of the requested block on the home node may be returned to the exclusive-copy state since (1) all node which had previously shared the block have acknowledged that they no longer hold copies, and (2) no new copies were given out since any new requests are met with a NACK. The state in the coherence directory remains transitioning, however. The original event is read out of the yankbuffer and added as a job to the event handler so that the full block-request code may be executed. The event cannot be handled directly by the P1MH since a reponse to a block request involves a message-send, which it not allowed for the P1MH. The state of the block in the coherence directory remains transitioning, to make the window of vulnerability when another request may come in an acquire rights to the block ahead of the original request as small as possible. The yankbuffer is returned to the circular buffer of free yankbuffers.

As the event handler performs the request procedure, it removes the block from the coherence directory (since no node is sharing the block) and calls the ccrequest function (normally called by message-dispatch code) directly, passing it the event information enqueued in its software job queue entry. At this point, the entire invalidation procedure is complete and the request which originally caused the invalidations gets another chance to acquire the block.

### Receiving NACKs

In the previous subsection, the home node was shown to be capable of sending NACKs in response to block requests. This section describes how the requesting node's P1MH must deal with NACKs. Since the events which caused the request messages to be sent

are still enqueued in software, the MCM does not need to perform another lookup in the event table when it receives a NACK. Intead, it needs to add a job for the event handler to resend the NACKed request. The actual NACK message which is sent by the home node contains the entire contents of the original request message. This makes it quite a simple task for the P1MH to add a resend request for the event handler - it passes all of the words of the NACK message to the EH. The event handler will dequeue the request some time in the future and retransmit the request. Once again, the reason that the P1MH cannot retransmit the request on its own is to avoid deadlock in the network - the P1MH is not allowed to send out any messages. Figure 7-3 summarizes the invalidation protocol.

**Performing Block-Invalidation**

Another task that the MCM must now perform is invalidating shared blocks in response to invalidation requests from the home node. When an invalidation message arrives, it bears only the virtual address, and does not contain any physical page frame information as events do. Therefore, the P0MH which handles the invalidation request must perform an explicit `PPM_lookup` to determine the local page frame which is used for backing the virtual page in question. The `putcstat` operation is performed on the virtual address to set block-status bits to *invalid* and return the previous state of the block. If the block was dirty, the page frame number is used along with the low 12 bits of the virtual address to determine the offset within the page frame where the block resides, and to read the block out into an acknowledgement message to be sent to the home node. In any case, the invalidation is acknowledged with either a simple ACK or an ACK bearing the contents of a dirty shared block. The invalidate ACK also contains the yankbuffer pointer which was passed in the invalidate message. As described above, this yankbuffer pointer is used on the return trip by the home node's P1MH to decrement the invalidate counter and decide when all nodes which shared the block have relinquished their copies.

Figure 7-3: Block Invalidation in Memory Coherence Protocol

## Dealing with Orderless Messages and Asynchrony

The protocol design presented so far seems to handle a variety of special cases, but the more interesting remain to be covered in this section. Particular problems arise when guarantees on message-ordering don't exist[3], and when asynchronous invalidation and NACK messages must be dealt with.

A requesting node may receive an invalidation message while it is still installing a newly-acquired block. Should the original ACK message to the block request be crossed with a later invalidation message, the requesting node may even receive the invalidation message before the actual data ACK arrives. To handle these cases, the

---

[3]At the time of the coherence protocol design, the M-Machine did not guarantee message ordering. The machine hardware has since been ammended to allow in-order messages to be used.

coherence protocol employs a state-machine model for memory blocks. That is, each block which has an entry in the event table has associated with it a state. This state helps MCM components decide what to do when messages or events concerning that block arrive. A block state is represented using five bits which encode the history of requests and responses targetting that block. These bits are:

1. PX : Pending Read/Write Request
2. PR : Pending Read Request
3. I  : Block Needs to be Invalidated
4. AX : ACK to R/W Request Received
5. NX : NACK to R/W Request Received

Initially, a software queue entry for a block gets its state set to PX or PR, depending on whether a readonly or readwrite copy of the block was requested from the home node. This records the fact that a request for the block has been sent to the home node and the requesting node is waiting for a NACK or ACK to return. In some instances, both PX and PR bits will be set - this occurs when first a read-invalid block-status miss is handled and the event handler sends out a request for a readonly copy of the block. Later, store to the same block will cause a write-invalid miss which will require that an exclusive copy of the block be requested. The EH will alter the state of the block from PR to (PR | PX) to note that two requests have been sent.

Should an invalidate message arrive before the actual data returns, the I state will be added to the block state. This will allow the MCM to keep track of the fact that after the request is ACKed or NACKed, an invalidation should be performed. The MCM cannot invalidate its block immediately after the invalidation message arrives because the invalidation and reponse to a block request could have gotten crossed, resulting in a block coming back later which should have been invalidated. If an invalidate message arrives when the block state is zero (meaning no software queue entry even exists for the block), it is safe to perform the invalidation immediately since no request messages for that block have been sent.

Figure 7-4: State Transition Diagram for Requested Blocks

Transitions which are performed when new messages arrive at the requesting node, or when events occur, can then be defined on these states. A state-transition diagram which is to be followed for each block is shown in figure 7-4. In this figure, transitions are triggered by the arrival of messages (NACK(X), ACK(X), NACK(R), ACK(R)) and new events (RI for read to an invalid block, WI for write to an invalid block, WR for write to a readonly block). This transition diagram clarifies the job of the MCM. When an invalidate message arrives for a block whose state is PR or PX, for instance,

88

the invalidate bit is added to the state and the yankbuffer pointer is added to the software queue entry for that block (hence the need for an invalidate pointer entry in the event queue data structure). If an ACK for the block is returned, the block will be installed, all of the events pending to it will be resolved, and then a job enqueued with the event handler will request that the handler perform an invalidation phase. The previously-stored yankbuffer pointer will be used when the block is invalidated and shipped, if necesary, to the home node. If a NACK arrives, a job for the event handler will be enqueued so that first the block is invalidated, and then a new request for the block is sent (since the MCM always resends requests if it receives a NACK).

This state-machine model tolerates out-of-order messages and asynchronous invalidations by imposing a rigid flow of control on the MCM and only allowing actions to be taken if the block is in a known and consistent state with the action being performed: for instance invalidated if no more messages are pending for that block.

With the state-machine model in place, the MCM design becomes more complete. When an event entry is enqueued for a particular block, the software queue entry's state is updated with proper PX and PR bits. An invalidation message handler (the code in ccInvalidate) first checks the state of a block to determine which state-transition to perform. Similarly, the P1MH checks block state when ACKs and NACKs arrive, to determine which actions to take. Usually in response to ACKs, this involves installing the block and then transitioning to a completed state or enqueuing jobs with the event handler to perform latent invalidations. In response to NACKs the actions are to enqueue jobs with the event handler, the nature of the jobs dependent on the block state - either invalidate-and-send-request, or send-request if no invalidation is required.

## Dealing With Concurrency

Since all of the threads which work in concert to provide memory-coherence need to access the MCM data structures (the event table on the requesting node, and the coherence directory on the home node) locks are used to enforce serialized access. The current implementation uses extremely coarse interlocks - a single lock is assigned

to each data structure (`sqlock` for the event table, and `ccdirlock` for the coherence directory). These locks must be acquired before functions which access and/or modify their associated data structures may be called. It is important to note here that regardless of lock granularity, the system must be implemented in such a way, that the event handler and P0 message handler may not hold locks which will prevent the P1MH from making progress at the time that they send out messages. As mentioned several times before, this is to prevent deadlocks from occuring - the P1MH must always be able to make progress and service its message queue even if other OS components such the the MCM running in an event handler slot are blocked, waiting to send a message into a saturated network.

When blocks are being installed, it is sometimes worthwhile for the P1MH to unlock the event table each time that it pops a new event from the event table which targets that block. This allows the event handler which is popping block status misses from the hardware event queue to add the event to the event table even while previous events are being popped off. This prevents spurious messages from being sent for blocks which are already installed locally. However, the system must be able to handle the case that a spurious message is sent. This may occur if the P1MH runs through all of the event table entries for a block that it received and then removes all traces that the block was installed, by deallocating the software queue entry for that block. Meanwhile, a latent block-status miss to the newly-installed block may be popped by the event handler, and a new event entry will be created. Since no software queue entry will have been found, the event handler decides to send a request message for the block. The home node will notice that the requesting node had already been listed as a sharer of the block, but will oblige with another copy. This allows requesting nodes to flush their shared blocks without having to inform the home node. The only side-effect of this flushing is that unnecessary invalidation messages may sometimes be sent by the home node.

Normally, installing a duplicate block is not a problem. However, if the original block was installed as exclusive, it may already be dirty by the time the second (and stale) home node's copy of the block arrives. This means that the requesting node,

90

when executing the code in `ccreturnStore` may not blindly install the block that it received in the ACK message. Instead, it must check the existing block-status bits of the block which was previously installed (local block) and determine whether the block is dirty or not. If the local block is dirty or readwrite, the stale copy is not installed. If, however, the local block is in the a readonly or invalid state, the block in the message is installed. In either case, all events pending to the block are satisfied as before.

| Function | Type | Description |
|---|---|---|
| INVALIDATE(int node, void *address, void *yankbuffer) | Request Proxy | Invalidates the memory block identified by *address* from the local cache and sets block-status bits to *invalid* in the LTLB and/or page table. Sends an acknowledgement to the home node *node*, sending along the *yankbuffer*. If the block was dirty, sends the dirty block within the acknowledgement. |
| RESENDSTORE(int header, void *address, int opdata, void *faultCP | Request Proxy | Sends a ccrequest message to the home node of *address*. |
| RESENDLOAD(int header, void *address, int opdata, void *faultCP | Request Proxy | Sends a ccrequest message to the home node of *address*. |
| INV_STORE(...) | Request Proxy | Combination of the INVALIDATE and RESENDSTORE/RESENDLOAD cases above. First invalidates a block and returns it to the home node. Then sends a request for it. |
| INV_LOAD(...) | Request Proxy | Same as above |
| REQUEST(int header, void *address, int opdata, void *faultCP) | Home Proxy | Executes the function ccrequest as if a request message for the block identified by *address* was received. |
| EH_handle_bsm(int header, void *address, int opdata, void *faultCP | BSM Handling | Responds to a local Block-Status Miss event. Enqueues the event (composed of the 4 argument words) into the software event table and returns status flags which tell the calling function what type of request message (if any) to send out. Returns 0 on error. A flag of 0x1 means no failure was detected. A flag of 0x2 means that a ccrequest message should be sent to the home node of *address*. A flag of 0x4 requests that the thread which caused the event be prevented from issuing any more instructions. A flag of 0x8 means that teh event request should be recirculated and tried again later. |

Table 7.1: Event Handler's MCM functions

| Function | Type | Description |
|---|---|---|
| ccrequest(void *address, int header, int opdata, void *faultCP, int node) | priority 0 | Processes a request for the block containing *address* from node *node*. Dispatches to helper functions ccrequest_st and ccrequest_ld depending on the type of operation encoded in *header*. May also call ccyankline if a shared block needs to be revoked from current sharing nodes. Sends a response to *node*, bearing the requested memory block or a NACK, or has the event handler do so in proxy at a later time. |
| ccreturnyankFull(int *yank_buffer) | priority 1 | Processes an acknowledgement to an invalidation message. The acknowledgement contains a dirty block which must be installed locally. Once installed, the original request which lead to the invalidation is processed in proxy by the event handler. |
| ccreturnYank(int *yank_buffer) | priority 1 | Processes an acknowledgement to an invalidation message. Decrements an invalidation counter for each such acknowledgement received. If the counter reaches zero, the block is considered unshared again and the request which lead to the invalidation is processed in proxy by the event handler. |

Table 7.2: Home Node MCM functions

93

| Function | Type | Description |
|---|---|---|
| `ccNackRO(void *address, int header, int opdata, void *faultCP, int node)` | priority 1 | Deals with a NACK returned by the home node in response to a readonly sharing request. Usually, the event handler is asked to resend the original request, to the home node, *home*. The first four arguments to this function are the arguments which were returned in the NACK, and used in the repeat request by the event handler. |
| `ccNackRW(...)` | priority 1 | Same as above, except that the NACK is in response to an exclusive block request. |
| `ccinvalidate(void *address, void *bufPtr, int node)` | priority 0 | Responds to an invalidation request from the home node, *node*, of the block identified by *address*. Takes steps to invalidate the block locally and, if dirty, to ship it back to the home node. |
| `ccreturnLoad(void *address, int header, int node)` | priority 1 | Installs the block which is returned in response to a readonly sharing request from node *node*. The 8 words of the block remain in the hardware message queue and are read out by an assembly-level helper function. |
| `ccreturnStore(...)` | priority 1 | Same as above, except that the block is installed for read/write as an exclusive copy. |

Table 7.3: Requesting Node MCM functions

| Description | bits |
|---|---|
| OP_ACTION | 56 - 63 |
| issuing thread slot | 48 - 55 |
| issuing functional unit | 42 - 47 |
| issuing cluster | 40 - 41 |
| target register file | 36 - 39 |
| target register | 32 - 35 |
| target cc | 28 - 31 |
| precondition | 26 - 27 |
| postcondition | 24 - 25 |
| physical page frame number | 4 - 23 |
| event type | 0 - 3 |

Table 7.4: Event Header Format

# Chapter 8

# Exposing System Calls to User Threads

The runtime system managers mentioned in previous chapters need to export certain system calls to user programs. This is accomplished through the use of jump tables and load-time program patching - mechanisms described in this chapter.

In order to allow user programs to safely access certain system function entry points, the programs need to be given *entry* pointers into runtime system code which they may then use to perform jmp instructions. The runtime system currently uses an object file called syscall.o which is linked with every user-level executable. This file contains stubs for all exported system calls which the program may wish to use. The stubs are simply functions which load system entry pointers from memory and jump on them. Entry pointers are loaded from locations in the data segment which are flagged to the loader as needing to be patched. This simplifies interfacing with the M-Machine compiler, since the compiler has no notion of which functions are system functions. Therefore, it expects to be able to place references to external system functions and have them resolved at link time. Again, this is already accomplished by having syscall.o contain stubs for all system functions, which means that from the point of view of the compiler and linker, a user-level executable has all of its symbols resolved before it is loaded. Figure 8-1 shows an example of a stub written in M-Machine assembly.

```
_tFork::
        GET_FRAME
        LOAD_FAR_LABEL(_tFork_ptr, itemp0, DStart)   /* load ldptr value */
        instr ialu jmp itemp0;                       /* jump to system code */
        instr ;
        instr ;
        CALC_RETIP

        RETURN                                       /* return to caller */
```

Figure 8-1: Sample syscall.m stub

Since stubs load system entry pointers from memory and the values of these entry pointers are known only at load time, the syscall.o object file contains magic numbers and relocation entries within it which signify that certain locations of its data segments need to be patched with pointers at load time. These pointers are called ldptr in the assembly language, and have their own relocation type. The trusted loader reads the object file, looking for ldptr relocations and replacing the contents of the data segment where the ldptrs are stored with entry pointers into system code. The magic numbers stored where the ldptr's are defined are used to determine which system function entry pointer needs to be stored there. The trusted loader is passed a table of associations between magic numbers and system entry pointers. This allows the syscall.o to create a table of ldptr values in its data segment, and use the stubs to load these values and jump on them. This patching is safe, since the user cannot trick the loader into giving out privileged information - any entry pointer which can be given out defines a protected entry point, and only entry pointers which the OS is willing to give out are passed to the loader. Examples of ldptr usage from syscall.m are shown in figure 8-2.

The entry pointer table passed to the loader is constructed at boot time, with values which are taken from system call function stubs, offset from the runtime IP. These are usually physical addresses. System call function stubs exist for each actual system function and act as an interface to the system function. Once called, the stubs perform two tasks. First, they issue an *mbar* instruction, which insures that

```
data;
align 0 mod 8;
_tFork_ptr::        ldptr 0x0000ffffaaaaaab0;
_tExit_ptr::        ldptr 0x0000ffffaaaaaaab;
```

Figure 8-2: Sample syscall.m ldptr usage

```
_tForkX::
        instr memu mbar;        -- issue mbar right away to keep registers safe
        GET_FRAME
        PUSH(DStart)            -- save caller's data segment pointer
        instr ialu imm __SYSTEM_UDAT_PTR, itemp0;   -- offset where system's
                                                    --   data pointer is stored
        instr ialu leab IP, itemp0, DStart;         -- create a pointer
                                                    --   to this offset
        instr memu ld DStart, DStart;               -- load system's data
                                                    --   pointer off the IP
        FCALL(_SYStFork)                            -- call the actual runtime
                                                    --   system function
        SPOP(DStart)                                -- restore user's data ptr
        RETURN                                      -- return to caller
```

Figure 8-3: Sample runtime stub

any memory operations which the caller performed will complete and overwrite any registers before the stub continues execution. This prevents a malicious user from issuing memory operations which may overwrite the register set of structed code as it begins execution. Secondly, while the IP of the executing system code points into runtime system space (as opposed to the user-level caller's space), the data segment pointer still points to the caller's data segment. The system function stub saves away the existing data segment pointer, and then loads the runtime system data pointer off its IP. The runtime data segment pointer is stored there at system boot time, for the express purpose of making it available to system-function callees. A runtime stub for the tFork system call is shown in figure 8-3.

# Chapter 9

# Performance Measurements

This chapter presents performance measurements of some runtime system components. It should be noted that although cycle-counts are included, these numbers are the result of executing a runtime system which was compiled with a compiler still under development and with absolutely no optimizations being performed. The more interesting numbers to examine are the breakdown of cycle-counts within long-latency operations to determine where most of the time is being spent.

## 9.1 The LTLB Miss Handler and Physical Memory Management

Tables 9.1 and 9.2 list the cycle counts of performing physical memory management tasks by the LTLB Miss Handler. Note that table lookups are quite fast, but the time to create a new mapping, which involves acquiring a new page frame from the free page list, is the largest component of an LTLB Miss.

## 9.2 Virtual Memory Allocation

The virtual segment manager takes an average of 950 cycles to allocate and return a virtual segment. A selected run is shown in table 9.3.

| Subcomponent | Cycles | Notes |
| --- | --- | --- |
| Initial LPT lookup | 283 | Lookup fails |
| Create new mapping | 1398 | Creates new virtual-physical mapping |
| Second Lookup | 236 | Added entry now found |
| Find conflicting LTLB Entry | 266 | For evicting existing LTLB entry |
| Writing new LTLB Entry | 231 | Evict old entry and write new one |
| Other | 1423 | |
| Total | 3837 | Total time to handle a miss to an unallocated page |

Table 9.1: Cycle count breakdown of LTLB Miss Handling

| Function | Cycles | Notes |
| --- | --- | --- |
| PPM_lookup | 1281 | Lookup a mapping in the page table |
| PPM_unmap | 1789 | Remove a mapping from both LTLB and the page table |

Table 9.2: Cycle counts for selected PPM functions

# 9.3  Thread Management

Table 9.4 shows that aside from thread context allocation and initialization, forking off a thread is quite inexpensive. This suggests that keeping available thread contexts around after they are destroyed may help improve performance.

# 9.4  Memory-Coherence

Table 9.9 shows a cycle-breakdown for handling a block-status-miss by the event handler. Note that while the event table is being updated, the update is not being directly simulated. It is expected that this time will be quite substantial. Cycle counts

| Subcomponent | Cycles | Notes |
| --- | --- | --- |
| Jump to protected subsystem | 95 | Including an mbar and restoring system data ptr |
| Allocate new segment | 801 | Actual buddy list allocation |
| Return from subsystem | 31 | Includes restoring user's data ptr |
| Total | 927 | Total time to allocate a virtual segment |

Table 9.3: Cycle count breakdown of Virtual Memory Allocation

| Subcomponent | Cycles | Notes |
|---|---|---|
| Subsystem entry | 89 | |
| Allocate new thread context | 935 | See VSM times in previous section |
| Initialize thread context | 7279 | |
| Allocate thread stack | 1248 | |
| Add job to EH job queue | 1033 | Tells EH to add thread to pending list |
| Add job to EH job queue | 1025 | Tells EH to perform scheduling |
| Return from subsystem | 126 | |
| Other | 1837 | |
| Total | 13572 | Total time to fork off a thread |

Table 9.4: Cycle count breakdown of tFork

| Subcomponent | Cycles | Notes |
|---|---|---|
| Pop from pending list | 162 | Get a new candidate |
| Install candidate | 2725 | Includes copying entire register state |
| Other | 348 | |
| Total | 3235 | Total time to install a thread into empty slot |

Table 9.5: Cycle count breakdown of tInstall

| Subcomponent | Cycles | Notes |
|---|---|---|
| Subsystem entry | 104 | |
| Signal T_CHILD_EXIT | 4685 | Includes allocating signal entry |
| Add EH job | 788 | Add EXIT signal |
| Other | 1316 | |
| Total | 6893 | Total time for a thread to call tExit and block |

Table 9.6: Cycle count breakdown of tExit

| Subcomponent | Cycles | Notes |
|---|---|---|
| Send spawn message | 2018 | Includes nonce allocation (1718 cycles) |
| Perform tSleep on nonce | 2800 | |
| Return Signal Message Processing | 4453 | Time to wake from when signal arrives |
| Total | 9217 | Does not include time that thread was sleeping |

Table 9.7: Cycle count breakdown of sender tSpawn

| Subcomponent | Cycles | Notes |
|---|---|---|
| Perform local fork | 9267 | |
| Perform signal on nonce | 564 | Sends message to spawner's node |
| Other | 326 | |
| Total | 10157 | Doesn't include time that remote caller was sleeping |

Table 9.8: Cycle count breakdown of receiving tSpawn request

for handling BSM's which don't require message-sends average about 410. This means that there is about a 700-cycle premium to sending out a request message, putting a thread to sleep, and performing other bookkeeping.

| Subcomponent | Cycles | Notes |
|---|---|---|
| Assembly prologue | 37 | Time to call C handler function |
| Add to event table | 174 | This is not simulated |
| Stop thread from issuing (icache miss) | 261 | |
| Request Message send | 126 | Read out data and send request message |
| Other | 495 | |
| Total | 1093 | Time to handle a block-status miss |

Table 9.9: Cycle count breakdown of handling a BSM

Table 9.10 shows cycle breakdowns for handling a coherence request by the home node. Note that as above, the coherence directory code is not being simulated and is expected to be a substantial portion of the total execution time. The total roundtrip time from block-status-miss to completion of line installation is about 8400 cycles, or about 1050 cycles per event to that line (the cycles of adding events after the initial request has been sent overlap the response times).

| Subcomponent | Cycles | Notes |
|---|---|---|
| Page-table lookup | 1471 | |
| Reading and sending cache line | 106 | |
| Other | 1182 | Includes coherence directory modification |
| Total | 2759 | Time to handle a cc request |

Table 9.10: Cycle count breakdown of home node's handling a ccrequest

| Subcomponent | Cycles | Notes |
|---|---|---|
| Read line from message and install | 75 | |
| Pop and satisfy 8 events to line | 3326 | (415 cycles/event) |
| Other | 1116 | |
| Total | 4517 | Time to handle a cc ACK |

Table 9.11: Cycle count breakdown of requesting node's handling an ACK

# Chapter 10

# Status and Future Directions

In this chapter, I present a broad overview of the currently-implemented MARS components and chart a course for what work remains to be done to develop MARS into a truly robust system.

## 10.1   Key OS Features and Contributions

The operating system presented in this thesis is quite novel. This is in great part due to the unique hardware platform to which MARS is tailored. The M-Machine's support for multiple thread contents, hardware-based capabilities, and configuration-space access to hardware state has been presented. What sets MARS apart most strongly from existing operating systems is its reliance on a collection of concurrently-executing managers to perform OS functions, instead of a single monolithic kernel or even microkernel. Most systems, regardless of light or heavyweight nature of the kernel, still require user-level programs to fault into a single-threaded kernel. With up to four system-level handlers able to execute at the same time, (and several additional protected subsystems in user slots) MARS is a truly decentralized operating system. The highest priority thread - the PMM - is still just a single thread performing only physical page management.

The use of capabilities by the OS can dramatically enhance performance. By turning thread context pointers used by MARS into Key pointers and giving them away

to user-level threads, the OS is able to obviate the use of more levels of indirection in order to protect threads. At the same time, once the thread context pointer is passed to a trusted OS component, the conversion of pointer type allows the system to access thread state very quickly, without requiring a lookup table. Capabilities are also used by the loader and runtime system to export system calls to user threads. Again, because no fault is required to enter a trusted subsystem, and because system-level code may execute in a user-level thread slot, performance of other threads is not affected.

By coupling a single virtual address space (in itself not a novel idea) with capabilities, MARS is able to provide efficient shared memory for all user and higher-level system threads. No special provisions are required to map virtual address spaces independently for each thread. A single virtual address map simplifies page and thread management. Context switches need only deal with register contents and other localized thread state.

Finally, the low-level support for coherent memory across the nodes of a multiprocessor makes MARS quite unique. Although operating systems like Mach may rely on hardware-based coherence, or allow a software coherence layer to be built independently using add-on memory managers, MARS takes a middle-ground. This results in memory-coherence more flexible than if built into hardware, at a performance cost. Because the coherence system is built on such a low level - within the message and event handlers - higher-level components are free to execute in such an environment. For example, the system-level loader can easily distribute a data segment of a newly-loaded executable over several nodes without requiring explicit message-passing. Simply storing a large array into virtual memory striped across several nodes will transparently distribute it. This makes the task of writing not only user-level programs, but also other system routines much simpler. This is certainly demonstrated by the ease with which multithreaded shared-memory code may be written under this OS (as shown by the example programs in appendix E).

## 10.2 Existing Components

The MARS system is composed of a collection of assembly and C source code files which compiled, assembled, and linked into a single executable. This executable is loaded into the M-Machine Simulator for testing and development work, and runs completely in physical memory.

The bootstrap - the boot.m assembly file - is the first to execute. It spawns off remaining system threads and performs initialization of the four managers presented in previous chapters. Each of the four system-level handlers contains an assembly-level portion which sets up arguments by popping events from hardware-mapped registers and calls on higher level functions written in C. The handlers are the event handler (event.m), P0 and P1 message handlers (both in message_event.m), and LTLB miss handler (ltlb_event.m). Sthe syscall.m assembly file contains stubs which allow user-level programs to call on exported system calls. This file is assembled and linked with user programs and is not linked into the runtime system.

The components written in C are divided on a roughly functional basis.

The physical memory manager is composed of the ltlb_body.c, ppm.c, lpt.c, and pplist.c files.

The virtual-memory manager is composed of stubs in vmem.m and actual routines in buddy.c.

The thread manager is divided into tmanager.c, tmanager2.c, and tsignal.c, with certain stubs written in boot.m.

Cache-coherence code is in cc_home.c and cc_request.c, with stubs in boot.m to handle message-sends and line-installation.

The cache-coherence data structures are actually compiled into the M-Machine simulator instead of being part of the runtime system. Both the cache-coherence directory and the event table are some of the largest components which remain to be fully implemented within the runtime system. Table 10.1 shows the breakdown of OS components by source file.

105

| File | Description |
|------|-------------|
| boot.m | Main system bootstrap. Also includes several assembly stubs for special instructions and message sends. |
| event.m | Event handler H-Thread source code. Marshalls arguments before calling code in eh.c |
| ltlb_event.m | LTLB Miss Handler H-Thread/PMM source code. Marshalls arguments before calling code in ltlb_body.c |
| message_event.m | P0 and P1 Message Handler source code. Interfaces to routines in memory-coherence and thread management functions. |
| sysloader.m | Loads user programs into memory and executes them. |
| vmem.m | Assembly stubs for VMM. Calls VSM functions in buddy.c |
| buddy.c | Virtual Segment Manager source code. |
| cc_home.c | Home node end of memory-coherence functions. |
| cc_request.c | Requesting node end of memory-coherence functions. |
| eh.c | Event Handler source code - for dealing with the software job queue, as well as responding to block-status miss events. |
| lpt.c | Local Page Table management tasks of the physical memory manager. |
| ltlb_body.c | Core LTLB Miss Handler code written in C. |
| pplist.c | Code to manage free page chains. Written by Andy Shultz from design by the author. |
| ppm.c | Physical Page Manager code for dealing with individual map/unmap/reclaim calls. Written by Andy Shultz from design by the author. |
| sq.c | Event Table code for use in memory-coherence. Currently incorporated directly into the M-Machine simulator and not linked into the runtime executable. |
| tmanager.c | Code for the thread manager dealing mostly with forking, evicting, and installing threads. |
| tmanager2.c | Additional code for the thread manager, dealing mostly with exiting a thread and maintaining parent/child linkages. |
| tsignal.c | Thread manager code dealing with signal/sleep. |

Table 10.1: MARS Sources Files

106

## 10.3　Future Work

Additional debugging and testing still needs to be performed on the runtime system to iron out bugs, although several test programs which have excercised all aspects of the runtime system, from memory-management to thread creation and communication to memory-coherence, have been successfully executed. These programs include the tfork suites (tfork2.c, tfork3.c, and tfork4.c), the matmul parallel matrix multiply programs (matmul1.c and matmul2.c), and iterative Jacobian relaxation programs (jacoby.c, jacoby2.c, jacoby3.c, jacoby4.c, jacoby5.c, and jacoby6.c).

### 10.3.1　Loader

The system's loader is an assembly stub which calls into the M-Machine simulator to perform actual program-loading. This component should be implemented as a protected subsystem able to run completely in virtual memory and load other processes without requiring low-level interaction with the runtime system - aside from the I/O aspect of accessing an executable's raw contents, calls to vmem_alloc and tFork are all that are required.

### 10.3.2　Memory-Coherence

The memory-coherence data structures and code for manipulating them should be moved out of the simulator and into the runtime system directly. This includes porting the implementations of the SSQEnqueue, SSQDeqeuue, SSQGetState, SSQSetState, and other such functions. The work should be relatively simple because the existing implementation is already written in C. The more involved development work must deal with the implementation of the backing-page invalidation and eviction strategy which was presented in the memory-coherence chapter. This will also require that the event handler call upon the physical page manager to determine the number of available backing pages. A low watermark will require preemptive evictions of shared lines to make more pages available should they become necessary. Speed optimizations to improve average-case performance for directory lookups will require modifying the

existing memory-coherence directory code to use a chained hash table instead of a simple linked-list of memory-block addresses.

### 10.3.3  Virtual Memory Management

The deallocation of virtual segments and underlying garbage-collection phase needs to be designed. This involves collecting dirty virtual segments in the dirty buddy list on each node and then performing a garbage-collection phase at very infrequent intervals. The actual garbage-collection will involve several phases. First, an initial round of communication needs to be performed so that all nodes enter into a garbage-collection phase, and prevent user threads from issuing any operations. In addition, all event and message queues need to be drained to remove any latent events and messages which may contain pointers to dirty segments. In a second phase, all local register files and physical memory needs to be examined to look for references to dirty segments. Any pointers which are found need to be replaced (perhaps with *errval* pointers) or NULL pointers). The system must be careful to avoid physical memory used by the OS itself. After local cleanup is completed, references to dirty segments must also be removed from all other nodes on the machine, so the garbage-collector needs to contact all other nodes and ask them to perform a local cleaning. Upon completion of the cleaning phase, another round of communication needs to inform nodes that garbage-collection is complete, and user threads may issue.

### 10.3.4  UNIX Personality

An entire UNIX system-call layer may be written using the low-level system primitives. This will present a familiar system-level interface for programmers to target without sacrificing general system performance. Thread and process-creation calls would be most interesting to implement in terms of MARS calls. Process creation calls like `fork` and `exec` would require little additional work and may be written in terms of primitives like `tFork`. The `signal` and `waitpid` would perhaps be the most challenging. The UNIX idea of letting programs install system handlers to dispatch

on signal events can be extended in the MARS system to allow dispatch threads to run, which absolves the runtime system of needing to save away current program state when handling a signal. Synchronization between the main thread and its signal handlers will need to be designed, however.

In terms of memory-allocation, it is quite likely that the UNIX sbrk call may be a NULL call if user threads are given enough virtual address space for code, data, and stack at the outset. Giving threads very large address spaces does not introduce a tremendous inefficiency problem since on-demand backing of virtual pages with physical page frames allows threads to have access to large address spaces without wasting physical memory.

# Appendix A

# MARS Messages

This appendix chapter lists the messages employed by MARS.

| Message IP | Message Words | Description |
|---|---|---|
| MSG_ccreturnDirty | address word1 ... word8 | Returns the a dirty block named by *address* to the home node. Words 1-8 are the contents of the block. |
| MSG_ccreturnyankFull | yankbuf word1 ... word8 | Returns a dirty block as a response to an invalidation message. The block is named by the address stored at the home node in the *yankbuf*. Words 1-8 are the contents of the block. |
| MSG_ccreturnyank | yankbuf | Acknowledges an invalidation request with the information that a shared line is no longer at the requesting node. The *yankbuf* sent in the original invalidation message is returned to the home. |
| MSG_ccinvalidate | address yankbuf | Sends an invalidation for a block identified by *address* to a node which shares that block. The *yankbuf* pointer to a local yankbuffer structure is passed as well. This pointer is returned in the ACK to the invalidation. |
| MSG_ccNackRO | address header data fcp | Sends a NACK message to a requesting node in response to a request for a readonly copy of a line. The contents of the request message are bounced back to the sender. |
| MSG_ccNackRW | address header data fcp | Similar to above, except message is in response to a request for an exclusive copy of a line. |
| MSG_ccreturnLoad | address header word1 ... word8 | Sends a readonly copy of a block from a home node to a requesting node. The block starts at *address* and consists of the 8 data words. The *header* sent in the original request is returned as well. |
| MSG_ccreturnStore | address header word1 ... word8 | Same as above, only an exclusive line is returned. |

Table A.1: Memory Coherence Messages

111

| Message IP | Message Words | Description |
|---|---|---|
| MSG_tWake | tc signal_data | Sends a message to invoke the SYStWake function on the home node of the context *tc*. The thread identified by *tc* is to be wakened with the *signal_data*. |
| MSG_tSleep | signal_word tc data_mask | Invokes a SYStSleep function on the home node if *signal_word*, adding a sleeper entry for thread context *tc* with a mask of *data_mask*. |
| MSG_tsignal | signal_word signal_data | Invokes a SYStSignal function at the home node of *signal_word*. |
| MSG_tspawn | nargs dp ip arg1 ... arg5 | Spawns a thread executing the function at *ip* with up to *nargs* number of arguments. The thread's data pointer is *dp*. |

Table A.2: Thread Management Messages

# Appendix B

# MARS Header Files

This chapter contains the header files used by the assembly routines and C functions in the M-Machine runtime system.

```
typedef struct {
    int header;
    void *address;
    int opdata;
    int *faultCP;
} eventBuffer;
```

```
#define CC_READONLY     0
#define CC_EXCLUSIVE    1
#define CC_UNSHARED     2
#define CC_INVALID      3
#define CC_TRANSITION   4

#define BSB_INVALID     0
#define BSB_READONLY    1
#define BSB_EXCLUSIVE   2
#define BSB_DIRTY       3
```

```
#define PRINTF(format_string, data_ptr)                          \
CONSTRUCT_LONG_PTR(format_string, Intarg0, data_ptr)             \
PUSH(AP)                                                          \
PUSH(Intarg0)                                                     \
Instr lalu mov SP, AP;                                           \
LIBCALL(_printf)                                                 \
POP(Intarg0)                                                     \
POP(AP)                                                          \
Instr lalu mov Intarg0, Intarg0;

#define SYSGETLOCK(temp_label, lock_addr)                        \
  CONSTRUCT_LONG_PTR(lock_addr, Itemp0, DStart)                  \
temp_label:                                                      \
  Instr memu stscnd cf, 1, IP, Itemp0, LCC1;                     \
  Instr lalu cf LCC1 br temp_label;                              \
  Instr ; Instr ; Instr ;

#define SYSPUTLOCK(lock_addr)                                    \
  CONSTRUCT_LONG_PTR(lock_addr, Itemp0, DStart)                  \
  Instr memu stscnd ct, 0, IP, Itemp0, LCC2;                     \
  Instr lalu ct LCC2 mov Intarg0, Intarg0;

/* uses Intarg0 and Itemp0.   Fills FMCIP */
#define MAKE_XM_PTR(address)                                     \
  CONSTRUCT_LONG_LABEL(address, Intarg0)                         \
  Instr lalu leab IP, Intarg0, Intarg0;                         \
  Instr lalu lsh Intarg0, #4, Intarg0                            \
    memu mov #P_EXMSG, Itemp0;                                   \
  Instr lalu lsh Intarg0, #-4, Intarg0;                         \
  Instr lalu lsh Itemp0, #60, Itemp0;                           \
  Instr lalu or Intarg0, Itemp0, Intarg0;                       \
  Instr lalu setptr Intarg0, Intarg0                            \
    lalu empty ##0x80000;                                        \
  Instr lalu mov Intarg0, FMCIP;
```

```
/* BASE USER */
#define IP         11
#define SP         12
#define Itemp0     13
#define RETIP      14
#define DStart     15
#define IRetVal    16
#define IntArg0    16
#define IntArg1    17
#define IntArg2    18
#define IntArg3    19
#define IntArg4    110
#define IntArg5    111
#define AP         112
#define RP         113
#define Itemp1     114
#define Itemp2     115

#define ITEMP0_EMPTY_MASK   0x0008
#define IRETVAL_EMPTY_MASK  0x0040
#define INTARG0_EMPTY_MASK  0x0040
#define INTARG1_EMPTY_MASK  0x0080
#define INTARG2_EMPTY_MASK  0x0100
#define INTARG3_EMPTY_MASK  0x0200
#define INTARG4_EMPTY_MASK  0x0400
#define ITEMP1_EMPTY_MASK   0x4000
#define ITEMP2_EMPTY_MASK   0x8000

/* regular system */
/* #define sstart  111 */

/* Event System */
/* regular event */
#define evHead    114
#define evBody    115

/* LTLB */
#define evtemp1   110
#define LTLBHDR   112
#define VADDR     113
#define OPDATA    114
#define FAULTCP   115

#if 0
#define EVTEMP1_EMPTY_MASK  0x0400
#define LTLBHDR_EMPTY_MASK  0x1000
#define VADDR_EMPTY_MASK    0x2000
#define OPDATA_EMPTY_MASK   0x4000
#define FAULTCP_EMPTY_MASK  0x8000
#endif

/* System Message */
#define mstemp1   110
#define MsgHead   114
#define MsgBody   115

#define MSTEMP1_EMPTY_MASK  0x0400

/* FP Registers */
#define FpArg0    f1
#define FRetVal   FpArg0
#define FpArg1    f2
#define FpArg2    f3
#define FpArg3    f4
#define FpArg4    f5
```

```
#define FpArg5    f6
#define FpArg6    f7
#define FpArg7    f8
#define FpTemp0   f9
#define FpTemp1   f10
#define FpTemp2   f11
#define FpTemp3   f12
#define FpTemp4   f13
#define FpTemp5   f14
#define FpTemp6   f15

/* message composition registers */
#define FMC0      f4
#define FMC1      f5
#define FMC2      f6
#define FMC3      f7
#define FMC4      f8
#define FMC5      f9
#define FMC6      f10
#define FMC7      f11
#define FMC8      f12
#define FMC9      f13
#define FMCDest   f14
#define FMCIP     f15

#define MC0       14
#define MC1       15
#define MC2       16
#define MC3       17
#define MC4       18
#define MC5       19
#define MC6       110
#define MC7       111
#define MC8       112
#define MC9       113
#define MCDest    114
#define MCIP      115

/* condition registers */
#define LCC0      cc0
#define LCC1      cc1
#define LCC2      cc2
#define LCC3      cc3

/* remember. FCALL consumes ITemp0 */
#ifndef FCALL
#define FCALL(function)                                         \
CONSTRUCT_LONG_LABEL(function, Itemp0)                          \
Instr lalu leab IP, Itemp0, Itemp0;                            \
Instr lalu jmp Itemp0;                                          \
Instr ; Instr ; CALC_RETIP                                      \
#endif

#define CALC_RETIP                                              \
Instr lalu lea IP, #4, RETIP;

#define PUSH(x)                                                 \
Instr lalu lea SP, #-8, SP;                                    \
Instr memu st x, SP;

#define POP(x)                                                  \
Instr memu ld    SP, #8, x;

#define SPOP(x)                                                 \
Instr memu ld    SP, #8, x;
```

```
Instr lalu mov  x, x;

#define CPUSH(x,TYPE,CC)                              \
Instr memu TYPE CC lea SP, #-8, SP;                   \
Instr memu TYPE CC st x, SP;

#define CPOP(x,TYPE,CC)                               \
Instr memu TYPE CC ld    SP, #8, x;

#define CPUSHX(x,TYPE,CC)                             \
memu TYPE CC st x, #-8, SP

#define GET_FRAME                                     \
PUSH(RETIP)

#define GET_RETIP                                     \
POP(RETIP)

#define FREE_FRAME \
Instr ;  \
Instr ;

#ifndef LINCALL
#define LINCALL(function)                             \
CONSTRUCT_LONG_LABEL(function,ltemp0)                 \
Instr lalu lea DStart, ltemp0, ltemp0;                \
Instr memu ld  ltemp0, ltemp0;                        \
Instr lalu jmp ltemp0;                                \
Instr ; Instr ; CALC_RETIP
#endif

#ifndef RETURN
#define RETURN                                        \
GET_RETIP                                             \
Instr lalu jmp RETIP;                                 \
FREE_FRAME                                            \
Instr ;
#endif
```

```
#define OPCODE_LD           47
#define OPCODE_ST           48
#define OPCODE_FST          49
#define OPCODE_LDS          50
#define OPCODE_LDSU         51
#define OPCODE_STS          52
#define OPCODE_STSU         53
#define OPCODE_STSCHD       54
#define OPCODE_FSTS         55
#define OPCODE_FSTSU        56
#define OPCODE_FSTSCHD      57
```

```
#define P_READ       0x0
#define P_RW         0x1
#define P_EXUSER     0x2
#define P_EXSYS      0x3
#define P_ENTERUSER  0x4
#define P_ENTERSYS   0x5
#define P_EXMSG      0x6
#define P_CONFIG     0x7
#define P_KEY        0x8
#define P_PHYSICAL   0x9
#define P_ERROR      0xA
```

```
#define EVENT_SIGNAL_INVALIDATE      0x1
#define EVENT_SIGNAL_RESENDSTORE     0x2
#define EVENT_SIGNAL_RESENDLOAD      0x3
#define EVENT_SIGNAL_INV_STORE       0x4
#define EVENT_SIGNAL_INV_LOAD        0x5
#define EVENT_SIGNAL_REQUEST         0x8
#define EVENT_SIGNAL_EVICT           0x9

#define EVENT_SCHEDULE               0x10
#define EVENT_SLEEP                  0x11
#define EVENT_FORK                   0x20
#define EVENT_KILL                   0x21
#define EVENT_WAKE                   0x22
```

eh.h          Fri Jul 21 15:34:39 1995          1

void add_eh_job(...);

```
#ifndef LPT_H
#define LPT_H


#ifndef LONG64

#define LTLBHashN1l        0x0
#define LTLBHashDeleted    0x1
#ifndef ULong64
#define ULong64 unsigned long
#endif
#ifdef Long64
#define Long64 long int
#endif

#else
#define LTLBHashN1l        0x0LL
#define LTLBHashDeleted    0x1LL
#include "stddefs.h"
#endif

#ifndef NUMPHYSPAGES
#define NUMPHYSPAGES    128   /*was 128*/
#endif

#ifndef LPTSIZE
#define LPTSIZE         128   /*was 256*/
#endif

#define IS_ANSI 0

/* LPT management */

typedef struct lpt_entry {
        int vpn;
        int ppn;
        int status1;
        int status2;
} LPTEntry;

typedef struct {
        LPTEntry htab[LPTSIZE];
        long table_hits[LPTSIZE];
        long numnlls;
        long numdels;
        long numcleans;
        long numlookups;
        long numhits;
} LPTable;

typedef LPTEntry *pLPTEntry;
typedef LPTable  *pLPT;

#if IS_ANSI
LPTEntry        LPTEntry_init(void);
void            LPT_init(pLPT);
int             LPT_calchash(int vpn, long l);
int             LPT_insert(pLPT table, int vpn, int ppn, int, int);
int             LPT_remove(pLPT table, int vpn);
int             LPT_lookup(pLPT table, int vpn);
pLPTEntry       LPT_findEntry(pLPT table, int vpn);
void            LPT_unparse(pLPT table);
void            LPT_stats(pLPT table);
#else
LPTEntry LPTEntry_init();
LPT_init();
int LPT_calchash();
int LPT_insert();
int LPT_remove();
int LPT_lookup();
pLPTEntry LPT_findEntry();
LPT_unparse();
LPT_stats();
#endif
#endif
```

```
/* this code was written by Andy Schultz */

#define PPNULL -1
#define PAGE_SIZE 4096

#define IS_ANSI 0

#if IS_ANSI
void PPList_init(pplist thelist, int start, int end);
int PPList_getpage(pplist thelist);
void PPList_putpage(pplist thelist, int thepage);
void PPList_unparse(pplist thelist);
#else
PPList_init(); /*which list, start, end*/
int PPList_getpage();        /*which list*/
PPList_putpage();   /*which list, thePage*/
PPList_unparse();        /*which list*/
#endif

/*these are in pointer.m*/
#if IS_ANSI
void * createPointer(int, int, int);
#else
void * createPointer();
#endif

typedef struct{
    int first;
    int last;
} PPList;
typedef PPList *pplist;
```

```c
#define IS_ANSI 0

#if IS_ANSI
int PPM_lookup(int vpn);    /*returns PPN for that VPN*/
int PPM_unmap(int vpn);     /*removes virtual-phsyical*/
                            /*and returns ppn*/
int PPM_map(int vpn);       /*maps in either local or remote list*/
int PPM_reclaim_local(int ppn); /*returns PPN to local pool*/
int PPM_reclaim_remote(int ppn); /*ditto for remote*/
int PPM_init(int start); /*initializes structures*/
int PPM_local2remote(int num); /*moves (up to) num pages.
                                 /*returns number moved*/
int PPM_remote2local(int num); /*ditto in the other direction*/
int PPM_remote_left(); /*number of remote pages left*/
int PPM_local_left();  /*number of local pages left*/
#else
int PPM_lookup();  /*returns PPN for that VPN*/
int PPM_unmap();   /*removes virtual-phsyical*/
                   /*and returns ppn*/
int PPM_map();     /*maps in either local or remote list*/
int PPM_reclaim_local(); /*returns PPN to local pool*/
int PPM_reclaim_remote(); /*ditto for remote*/
int PPM_init(); /*initializes structures*/
int PPM_local2remote(); /*moves (up to) num pages.
                          /*returns number moved*/
int PPM_remote2local(); /*ditto in the other direction*/
int PPM_remote_left(); /*number of remote pages left*/
int PPM_local_left();  /*number of local pages left*/
#endif
```

```c
/* an ENode is an event node, recording information about an individual
   memory event */
typedef struct line_event_entry {
    struct line_event_entry *next;
    ULong64 header;
    MWord address;
    ULong64 data;
    MWord CP;
} ENode;

/* line-state information.  Just a bit vector */
typedef struct {
    unsigned int px : 1;
    unsigned int pr : 1;
    unsigned int inv : 1;
    unsigned int ax : 1;
    unsigned int nx : 1;
} SQ_STATE;

/* a software queue node maintains event chains for a particular cache
   line.  It contains the cache-line state, invalidation information,
   etc. */
typedef struct sqn {
    struct sqn *next;
    ULong64 address;              /* cache-line address (low 6 bits are zero) */
    SQ_STATE state;
    char emptied;
    ENode *events;
    ENode *tail;
    ULong64 invalidate_ptr; /* stores a single yankbuffer ptr per cache line */
} SQNode;

/* an entry in the event table - contains the virtual page number,
   head of the an SQNode list, and status information about the
   physical page which is used for backing */
typedef struct request_entry {
    SQNode *head;
    ULong64 vpn;
    int status;
} REntry;

#ifndef NUM_NODES
#define NUM_NODES       10
#endif

#ifndef BACKING_SIZE
#define BACKING_SIZE 128
#endif

void ENodeList_init();
ENode* ENodeList_pop();
void ENodeList_push(ENode *node);
void ENode_unparse(ENode *node);
void SQNodeList_init();
SQNode* SQNodeList_pop();
void SQNodeList_push(SQNode *node);
SQ_STATE SQNode_getState(SQNode *node, ULong64 address);
int SQNode_setState(SQNode *node, ULong64 address, int newState);
int backing_setState(REntry *entry, ULong64 address, int newState);
void REntry_unparse(REntry *entry);
void backing_unparse(REntry *BackingHashTable);
void backing_init(REntry *BackingHashTable);
int backing_addInvalidate(REntry *entry, ULong64 address, ULong64 ptr);
ULong64 backing_getInvalidate(REntry *entry, ULong64 address);
int backing_addEvent(REntry *entry, ENode *node);
ENode *backing_popEvent(REntry *entry, ULong64 address);
ENode *backing_firstWriteEvent(REntry *entry, ULong64 address);
void BackingTable_unparse(int n);

int header_isWriteEvent(ULong64 header);
```

```
/* return line-state information for address */
/* if address is NULL, returns whether page ppn has any events to it */
int SqGetState(void *address, int ppn);

/* sets the line state for address */
int SqSetState(void *address, int new_state, int ppn);

/* pop off the next event targetting address.  If event buffer is
   NULL, returns the invalidate pointer. */
void * SqDequeue(void *address, eventBuffer *eb, int ppn);

/* returns the first store targetting the address
   (writes is into the eb struct, actually) */
int SqGetFirstW(void *address, eventBuffer *eb, int ppn);

/* pushes a new event targetting the address. */
/* returns flags which help caller decide whether to send
   a current message and stuff like that */
int SqEnqueue(int header, void *address, void *opdat, void *faultcp, int ppn);

extern int sqlock;
```

```
/* return own data-segment pointer */
char *_getDP();

/* given a thread context, return the parent of that thread context */
int * _getParent(int *mytc);

/* return own thread context */
int * _getSelfTC();

void tRnparse();

void *tFork(void *threadIP, void *DataPtr, void *returnIP, void *parent,
            int numargs, ... );
int tExit(int);
int tSleep(int *signal_word, int mask);
void *tspawn(int numargs, void *function, int node_num, ...);
int hspawn(int numargs, void *threadIP, int dest_cluster, ...);
```

```c
/* Information for a particular hthread context */
struct HContext {
    int int_reg_file[16];         /* 16 integer registers      */
    float fp_reg_file[16];        /* 16 floating-point registers */
    int local_cc[4];              /* 4 local cc registers      */
    int global_cc[8];             /* the eight global cc registers */
    int hardware_mbar_counter;    /* hardware memory-barrier counter */
    int software_mbar_counter;    /* software memory-barrier counter */
    char *restartIPvector[4];     /* IP's to restart the thread */
    int empty_scoreboard;         /* register empty-state scoreboard */
};

/* A live thread context */
struct ThreadContext {
    struct ThreadContext *Next;
    struct ThreadContext *Parent;
    struct ThreadContext *Sibling;
    struct ThreadContext *Children;

    struct HContext hcthreads[4];

    /* generate global thread information */
    int vslot;                    /* the slot that the thread occuples */
    int flags;                    /* hFull and hIssue bits IIIIFFFF    */
    int SCC;
    int SCL;

    /* Information used for signal/sleep */
    int signalData;               /* data passed to you when you wake  */
    int need_to_block;
    int need_to_wake;
    int need_to_sleep;
};

struct GlobalThreadState {
    int occupied;
    struct ThreadContext *Pending;
    struct ThreadContext *PendingEnd;
    struct ThreadContext *Kill;
    struct ThreadContext *KillEnd;
    struct ThreadContext *Running;
    struct ThreadContext *RunningEnd;
};

void HContext_init(struct HContext *context);
void ThreadContext_init(struct ThreadContext *context);
void ThreadContext_unparse(struct ThreadContext *context);
int ThreadContext_free(struct ThreadContext *context);

struct ThreadContext *tAlloc();
void                  tUnparse();

void                  tAddPending(struct ThreadContext *tc);
struct ThreadContext *tPopPending();

void                  tAddRunning(struct ThreadContext *tc);
struct ThreadContext *tPopRunning();

void                  tAddKill(struct ThreadContext *tc);
struct ThreadContext *tPopKill();

struct ThreadContext *tSelfTC();

void *SYSfork(void *IP, void *DP, void *retIP, int numargs, ... );
int   SYSExit(int returnvalue);
```

```c
void SYSUnparse();
int  tSchedule();
void tKill(struct ThreadContext *);
void tPutToSleep(struct ThreadContext *);
```

```c
/* make a configuration space pointer into thread slot */
int *sysMakeCP(int slot);

/* create a pointer from the raw bits */
int *sysSetPtr(int bits);

/* return the home node of the address */
int sysGetHN(void *address);

/* set blocks-status bits of address */
int sysPUTCSTAT(void *address, int new_bits);

/* acquire a lock */
int sysGetLock(void *lock_addr);

/* release a lock */
int sysPutLock(void *lock_addr);

void mbarLoadUpdate(int, void *, int, int);

void sysSMSMessage(int *);
void sysSendLine(char *, void *, int);
void sysSendInvalidateAck(void *, int);
void sysSendInvalidate(int, void *, int *);

void update_ombc();
int *get_offsetptr_into_ppn(int, void *);

#ifndef FALSE
#define FALSE 0
#endif

#ifndef TRUE
#define TRUE 1
#endif
```

```
#define T_CHILD_EXIT    0x100
#define T_KILL          0x200
#define M_FAULT         0x400
#define T_ALL_SIGNALS   0x000

int SYStSignal(struct ThreadContext *target,
               struct ThreadContext *waker,
               int data);
int SYStSleep(int mask, int *buffer);
void SigTab_Init();
```

# Appendix C

# MARS Assembly Code

This chapter contains the assembly routines for the M-Machine runtime system.

```
/*********************************************************
 *
 * M-Machine runtime system boot code.  Contains code to initialize
 * system variables, and start up system handlers.
 *
 * Written by        Yevgeny Gurevich
 * Version           0.95
 * Modification Date 9/26/94
 * Modification Date 8/7/95
 *
 * --------------------------------------------------
 * Bugs:  None
 * --------------------------------------------------
 *********************************************************/

#include "newreg.h"
#include "helpmacros.h"
#include "ccdefs.h"
#include "pointers.h"
#include <libcall.h>

#define PAGE_SEG_LENGTH     16
#define POINTER_LENGTH      6
#define POINTER_ADDRESS     54
#define PAGE_TABLE_LENGTH   2   /* In K. Multiply by 1024 to get bytes */

#define CALLFAIL       Instr lalu br Fail; Instr ; Instr ; Instr ;

/* the following are offsets from the global thread state area for a vthread */
#define CONFIG_CP       #0x00000
#define CONFIG_FLAGS    #0x00008
#define CONFIG_VACANT   #0x00010
#define CONFIG_SCC      #0x00018
#define CONFIG_SLC      #0x00020

/* base of the GTLB in configspace */
#define CONFIG_IN_GTLB_BASE     0x204000

/* use virtual memory */
#define VMEM            1

data;
    _SysStackEnd:    ptr    r.0.0;    -- address of last valid system
                                      -- stack address
    _Sys_LPT_Start:: ptr    r.0.0;    -- Starting address of
                                      -- local page table
    _Sys_SFH_Stack:  ptr    r.0.0;    -- Starting address of system fault
                                      -- handler stack
    _Sys_LMH_Stack:  ptr    r.0.0;    -- Starting address of ltlb miss
                                      -- handler stack
    _Sys_DH_Stack:   ptr    r.0.0;    -- Starting address of system
                                      -- dispatch handler stack

    _Sys_Memory_End:: ptr    r.0.0;    --
    _Sys_VMemory_End:: ptr   r.0.0;    --

    _format_string1:  asciz  "M-Machine System Runtime v0.95\n";
    _format_acquire_failed: asciz  "handler setup has failed\n";

/* place holders for system load pointers - this array is passed to the loader */
_SysCallStart::
    ptr  r.0.1;
    ptr  r.0.2;
    ptr  r.0.3;
    ptr  r.0.4;
    ptr  r.0.5;
    ptr  r.0.6;
    ptr  r.0.7;
    ptr  r.0.8;
    ptr  r.0.9;
    ptr  r.0.10;
    ptr  r.0.11;
    ptr  r.0.12;
    ptr  r.0.13;

_SystemStackSize := 0x1000;        -- size, in bytes, of the system
                                   -- stack segment

Config_Global_Base := CONFIG_IN_GTLB_BASE;

/* the following are used because the runtime system links in its own
   copy of clib.o, where references are made to these two.  */
_rpcx:  p64 MSG_InvokeRPC;
_vmemx::p64 vmem_alloc;

text;

/*      First, set up a location in system code space where we store
        away the system's data pointer for future use when returning
        from user code */

align 0 mod 8;
_SYSTEM_UDAT_PTR__:: ptr k.0.11;       -- system's data pointer

/***********************************************
 *                                             *
 *         Beginning of System Boot code       *
 *                                             *
 ***********************************************/

Main::    /* the first step is to set up memory and registers correctly */

          /* store away pointer to system's USER data segment */
          /* and set DStart (15) to the beginning of USER */

          /* ON ENTRY -
             12 --> ptr to beginning of data segment
             13 --> ptr to end of entire area, including BSS segment */
          Instr lalu empty #0x0002;
          Instr lalu mov 13, f1;            -- save away end-of-data area

          LOAD_FAR_LABEL(USER, ltemp0, 12)
          Instr lalu lea 12, ltemp0, DStart;
          Instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
          Instr lalu leab IP, ltemp0, ltemp0;
          Instr memu st DStart, ltemp0;

          /* NOTE: Initially, 12 is set to the beginning of the data segment,
             and this conflicts with the logical use of 12 as the SP.
             Therefore, make sure that once SP is written, 12 is no longer
             needed, by copying all of its usefulness to other registers
             or within a data word accessible by astart or DStart. */

          /* find the system end point and make that the beginning of stack. */
          /* this value is now origially given to us in 13, we saved it away */
          /*    into f1 and now we restore and use it                        */
          Instr lalu empty #(ITEMP0_EMPTY_MASK);
```

```
Instr (alu mov fi, ltemp0;

/* set the size properly */
Instr lalu lsh ltemp0, #10, ltemp0;
Instr lalu lsh ltemp0, #-10, ltemp0;
Instr lalu lmm #0x97c0, SP;
Instr lalu shoru 0x0000, SP;
Instr lalu shoru 0x0000, SP;
Instr lalu shoru 0x0000, SP;
Instr lalu or ltemp0, SP, SP;
Instr lalu setptr SP,SP;

Instr lalu lmm __SystemStackSize, ltemp0;          -- Size of system's
                                                   -- stack
Instr lalu lea SP, ltemp0, ltemp0;                 -- ptr to end of stack
Instr lalu lea ltemp0, #-8, ltemp0;                -- last word of stack
Instr lalu lmm __SysStackEnd, ltemp1;
Instr lalu lea DStart, ltemp1, ltemp1;
Instr memu st  ltemp0, ltemp1;                     -- store address of stack end

/* flag the end of the system stack by writing
                0x99998888 into last word */
LOAD_FAR_LABEL(_SysStackEnd, ltemp0, DStart)
Instr lalu lmm #0x9999, ltemp1;
Instr lalu shoru #0x8888, ltemp1;
Instr memu st  ltemp1, ltemp0;

/* now SINCE STACK GROWS UPWARD we change the SP */
Instr lalu mov  ltemp0, SP;
Instr lalu lea  SP, #-8, SP;
Instr lalu mov  SP, AP;

/* now actually ready for a frame */
Instr lalu mov 10, RETIP;
GET_FRAME

/* perform a C library call, printing some information about the
        runtime system */
Instr lalu mov SP, intarg1;
Instr lalu mov AP, intarg2;
Instr lalu mov DStart, intarg3;
LOAD_FAR_LABEL(_SysStackEnd, intarg4, DStart)
PRINTF(_format_string1,DStart)

PRINTF(BuildDate, DStart)

FCALL(_syscall_setup)                 -- set up system call library
FCALL(_physical_memory_setup0)        -- just enough for handlers to use
FCALL(_handler_setup)
FCALL(_physical_memory_setup1)        -- set up physical memory on nodes

/* check for return value */
Instr lalu leq IRetVal, 10, cc1;
Instr lalu ct cc1 br _sync_acquire_failed;
Instr ; Instr ; Instr ;

FCALL(_user_thread_setup)
Instr memu mbar;

#if VMEM
FCALL(_virtual_memory_setup)          -- set up virtual memory on nodes
#endif

FCALL(_tinit)    /* needs to be done after virtual memory setup */
FCALL(_ccinit)   /* print out debugging information on location of
                    the yankbuf data structure */
```

```
FCALL(INIT_Lib)

/* case out on own node number to execute node-specific code. */
/* this might be obviated when the loader can place code and
        data on differing nodes.  For now, since all code and data
        is placed on each node in the same manner, we have to case
        out in the code itself. */
LIBCALL(NodeId)
Instr lalu leq IRetVal, 10, cc0;
Instr lalu ct cc0 br _node0code;
Instr ;
Instr ;
Instr ;

Instr lalu leq IRetVal, #1, cc0;
Instr lalu ct cc0 br _node1code;
Instr ;
Instr ;
Instr ;

Instr lalu leq IRetVal, #2, cc0;
Instr lalu ct cc0 br _node2code;
Instr ;
Instr ;
Instr ;

Instr lalu leq IRetVal, #3, cc0;
Instr lalu ct cc0 br _node3code;
Instr ;
Instr ;
Instr ;

Instr lalu br fail;
Instr ;
Instr ;
Instr ;

_node1code:
_node2code:
_node3code:
        /* terminate bootstrap on all nodes other than node0 */
        Instr lalu lmm #0x4444, intarg0;
        FCALL(_SYSExit)
        Instr ;
        Instr ;
        Instr lalu empty #((TEMP0_EMPTY_MASK);
        Instr lalu mov ltemp0, ltemp0;

_node0code:
_SYSTEM_LOADER_LOOP:
        FCALL(_system_loader)
        Instr lalu leq IRetVal, #0, cc1;
        Instr lalu ct cc1 br _loader_failed;
        Instr ;
        Instr ;
        Instr ;

        Instr lalu br _SYSTEM_LOADER_LOOP;
        Instr ;
        Instr ;
        Instr ;

data;
_loader_failed_string: asclz "Loader Failed.\n";
```

```
text:
_loader_failed:
        PRINTF(_loader_failed_string, DStart)

        Instr lalu br _boot_fail;
        Instr ;
        Instr ;
        Instr ;

_sync_acquire_failed:
        PRINTF(_format_acquire_failed, DStart)

        Instr lalu br _boot_fail;
        Instr ;
        Instr ;
        Instr ;

_boot_fail:
_cleanup:
        FREE_FRAME
        Instr lalu br Fail;
        Instr ;
        Instr ;
        Instr ;

/***************************************
/*                                     */
/*      Setup of System Call jump table */
/*                                     */
/***************************************

data:
_scall_string1: asclz 'Setting up syscall jump table\n';
text:
_syscall_setup:
        GET_FRAME
#if VERBOSE_SYSCALLSETUP
        PRINTF(_scall_string1, DStart)
#endif

        CONSTRUCT_LONG_PTR(_SysCallStart, ltemp0, DStart)

        CONSTRUCT_LONG_LABEL(_tSleepX, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(_tExitX, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(sysmalloc, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(vmem_alloc, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(MSG_invokeRPC, ltemp1)
        Instr    lalu leab IP, ltemp1, ltemp1;
        Instr    lalu lsh ltemp1, #4, ltemp1;
        Instr    lalu lsh ltemp1, #-4, ltemp1;
        Instr    menu mov #P_EXMSG, lntarg0;
        Instr    lalu lsh lntarg0, #60, lntarg0;      -- creating XM ptr
        Instr    lalu or ltemp1, lntarg0, ltemp1;     -- creating XM ptr
```

```
        Instr   lalu setpcr ltemp1, ltemp1, ltemp1;    -- creating XM ptr

        Instr   menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(_tUnparseX, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(_tForkX, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(_tSignalX, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(_getParentX, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        CONSTRUCT_LONG_LABEL(_getSelfTCX, ltemp1)
        Instr lalu leab IP, ltemp1, ltemp1;
        Instr menu st ltemp1, #8, ltemp0;

        Instr lalu mov #1, lRetVal;
        GET_RETIP
        Instr lalu jmp RETIP;
        FREE_FRAME
        Instr ;

/***************************************
*                                     *
* Initiating System Fault and Message handlers *
*                                     *
***************************************/

data:
_shs_string1: asclz 'Setting up System Handlers...\n';
text:

_handler_setup:
        GET_FRAME
#if VERBOSE_HANDLERSETUP
        PRINTF(_shs_string1, DStart)
#endif

        --
        -- set up system THREAD4 - EH, LTLB, PONH, PIMH --
        --

        Instr lalu imm #(0x7440), lntarg0;
        Instr lalu shoru #0x0000, lntarg0;
        Instr lalu shoru #0x0008, lntarg0;
        Instr lalu shoru #0x0000, lntarg0;
        Instr lalu setpcr lntarg0, lntarg0;

        PUSH(lntarg0)                                  -- mark active threads manually

        /* one thing to do is construct a cspace ptr to global thread state */
        Instr lalu imm #1, ltemp0;
        Instr lalu lsh ltemp0, #16, ltemp0;
        Instr lalu or lntarg0, ltemp0, ltemp0;
        Instr lalu setpcr ltemp0, ltemp0;

        /* ltemp0 is ptr to global thread state */
```

```
Instr lalu lea ltemp0, #CONFIG_CP, ltemp0;

/* store the system data pointer in the CP location */
/* this will eventually change, of course          */
Instr menu st 12, #0x18, ltemp0;               -- Write thread CP

Instr menu st 10, #0x8, ltemp0;                -- Write 0 into SCC
Instr lalu lmm #0x400, ltemp1;
Instr menu st ltemp1, #0x10, ltemp0;           -- Write 1024 into SCL

/* now modify local thread state */
POP(ltarg0)
PUSH(ltarg0)

/* now need to write threads flag -- that is in global thread state */
Instr lalu lmm #1, ltemp0;
Instr lalu lsh ltemp0, #16, ltemp0;
Instr lalu or  ltarg0, ltemp0, ltemp0;
Instr lalu setptr ltemp0, ltemp0;

/* modifies THREADS flag in CONFIG SPACE */
Instr menu mov #0x0f, ltarg1;
Instr lalu lea ltemp0, #CONFIG_FLAGS, ltemp0;  -- Mark all hthreads
Instr menu st ltarg1, #0, ltemp0;

/* set ip's for each of the threads in the 4 clusters */
LOAD_FAR_LABEL(EVENT_IP, ltemp0, DStart)
Instr lalu lmm #0x0700, ltarg1;
Instr lalu lea ltarg0, ltarg1, ltarg1;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr menu st ltemp0, ltarg1;

LOAD_FAR_LABEL(LTLB_IP, ltemp0, DStart)
Instr lalu lmm #0x1700, ltarg1;
Instr lalu lea ltarg0, ltarg1, ltarg1;         -- write C0 IP
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr menu st ltemp0, ltarg1;

LOAD_FAR_LABEL(DISPATCH_IP, ltemp0, DStart)
Instr lalu lmm #0x2700, ltarg1;
Instr lalu lea ltarg0, ltarg1, ltarg1;         -- write C1 IP
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr menu st ltemp0, ltarg1;

LOAD_FAR_LABEL(DISPATCH_IP, ltemp0, DStart)    -- write C2 IP
Instr lalu lmm #0x3700, ltarg1;
Instr lalu lea ltarg0, ltarg1, ltarg1;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
```

```
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, #8, ltarg1;
Instr lalu lea ltemp0, #4, ltemp0;
Instr menu st ltemp0, ltarg1;                   -- write C3 IP

/* so     C0 of slot 4 is event handler
 *        C1 of slot 4 is ltlb fault handler
 *        C2 of slot 4 is message handler
 *        C3 of slot 4 is message handler
 *        Generate some stacks and give them pointers  */

/* stack for Event Handler */
Instr lalu lmm #0x2000, ltarg0;
FCALL(mem_alloc)

Instr lalu lmm #(0x2000 - 0x08), ltarg1;        -- set to last word
Instr lalu lea IRetVal, ltarg1, IRetVal;        -- in the segment

/* IRetVal is now pointer to beginning of sync fault handler stack */
CONSTRUCT_LONG_PTR(_Sys_SFH_Stack, ltarg1, DStart)
Instr menu st IRetVal, ltarg1;

/* store in configspace, in register SP (12) */
POP(ltemp0)              -- general thread configspace pointer that we saved up
PUSH(ltemp0)

Instr lalu lmm #0x0010, ltarg1;                 -- C0.12 offset
Instr lalu lea ltemp0, ltarg1, ltarg1;
Instr menu st IRetVal, #24, ltarg1;             -- C0.12 <- stack ptr
                                                -- advance to C0.15
Instr menu st DStart, #48, ltarg1;              -- C0.15 <- data ptr
                                                -- advance to C0.111
Instr menu st DStart, ltarg1;                   -- C0.111 <- DSTART

/* now generate pointer to ltlb handler stack */
Instr lalu lmm #0x1000, ltarg0;
FCALL(mem_alloc)

Instr lalu lmm #(0x1000 - 0x08), ltarg1;        -- set to last word
Instr lalu lea IRetVal, ltarg1, IRetVal;        -- in the segment

/* IRetVal is now pointer to beginning of ltlb handler stack */
CONSTRUCT_LONG_PTR(_Sys_LMH_Stack, ltarg1, DStart)
Instr menu st IRetVal, ltarg1;

/* store in configspace, in register SP (12) */
POP(ltemp0)
PUSH(ltemp0)

Instr lalu lmm #0x1010, ltarg1;                 -- C1.12 offset
Instr lalu lea ltemp0, ltarg1, ltarg1;
Instr menu st IRetVal, #24, ltarg1;             -- C1.12 <- stack ptr
                                                -- advance to C1.15
Instr menu st DStart, #48, ltarg1;              -- C1.15 <- data ptr
                                                -- advance to C1.111
Instr menu st DStart, ltarg1;                   -- C1.111 <- DSTART

-- Dispatch Handler --
Instr lalu lmm #0x2000, ltarg0;
FCALL(mem_alloc)

Instr lalu lmm #(0x2000 - 0x08), ltarg1;        -- set to last word
Instr lalu lea IRetVal, ltarg1, IRetVal;        -- in the segment

/* IRetVal is now pointer to beginning of dispatch handler stack */
```

```
CONSTRUCT_LONG_PTR(_Sys_DM_Stack, Intarg1, DStart)
Instr menu st IRetVal, Intarg1;

/* store in configspace, in register SP (12) */
POP(Itemp0)
PUSH(Itemp0)

Instr lalu lmm #0x2010, Intarg1;                          -- C3.12 offset
Instr lalu lea Itemp0, Intarg1, Intarg1;
Instr menu st IRetVal, #24, Intarg1;

Instr menu st DStart, #48, Intarg1;
Instr menu st DStart, Intarg1;

-- priority 1 message handler --
Instr lalu lmm #0x2000, Intarg0;
FCALL(mem_alloc)

Instr lalu lmm #(0x2000 - 0x08), Intarg1;                 -- set to last word
Instr lalu lea IRetVal, Intarg1, IRetVal;                 -- in the segment

/* store in configspace, in register SP (12) */
POP(Itemp0)
Instr lalu lmm #0x2010, Intarg1;                          -- C4.12 offset
Instr lalu lea Itemp0, Intarg1, Intarg1;
Instr menu st IRetVal, #24, Intarg1;

Instr menu st DStart, #48, Intarg1;
Instr menu st DStart, Intarg1;

Instr lalu mov #1, IRetVal;
RETURN


_user_thread_setup:
/* set the threadflags of all user thread slots to known and
   useful values                                                 */
GET_FRAME

Instr lalu lmm #(0x7440), Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu setptr Intarg0, Intarg0;

Instr lalu lmm #1, Itemp0;
Instr lalu lsh Itemp0, #16, Itemp0;
Instr lalu add Itemp0, #CONFIG_FLAGS, Itemp0;
Instr lalu lea Intarg0, Itemp0, Itemp0;
Instr lalu setptr Itemp0, Itemp0;

Instr lalu mov #0x11, Intarg1;
Instr menu st Intarg1, Itemp0;                            -- set self to issuing and full in c0

Instr lalu lmm #(0x7440), Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu shoru #0x0002, Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu setptr Intarg0, Intarg0;

Instr lalu lmm #1, Itemp0;
Instr lalu lsh Itemp0, #16, Itemp0;
Instr lalu add Itemp0, #CONFIG_FLAGS, Itemp0;
Instr lalu lea Intarg0, Itemp0, Itemp0;
```

```
Instr lalu setptr Itemp0, Itemp0;

Instr lalu mov #0x0, Intarg1;
Instr menu st Intarg1, Itemp0;                            -- set all to inactive and empty

Instr lalu lmm #(0x7440), Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu shoru #0x0004, Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu setptr Intarg0, Intarg0;

Instr lalu lmm #1, Itemp0;
Instr lalu lsh Itemp0, #16, Itemp0;
Instr lalu add Itemp0, #CONFIG_FLAGS, Itemp0;
Instr lalu lea Intarg0, Itemp0, Itemp0;
Instr lalu setptr Itemp0, Itemp0;

Instr lalu mov #0, Intarg1;
Instr menu st Intarg1, Itemp0;                            -- set others to inactive

Instr lalu lmm #(0x7440), Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu shoru #0x0006, Intarg0;
Instr lalu shoru #0x0000, Intarg0;
Instr lalu setptr Intarg0, Intarg0;

Instr lalu lmm #1, Itemp0;
Instr lalu lsh Itemp0, #16, Itemp0;
Instr lalu add Itemp0, #CONFIG_FLAGS, Itemp0;
Instr lalu lea Intarg0, Itemp0, Itemp0;
Instr lalu setptr Itemp0, Itemp0;

Instr lalu mov #0, Intarg1;
Instr menu st Intarg1, Itemp0;                            -- set others to inactive

Instr lalu mov #1, IRetVal;
RETURN

/*****************************************************/

_physical_memory_setup:
text:

GET_FRAME

/* load end of stack and align to page boundary */
LOAD_FAR_LABEL(_SysStackEnd, Itemp0, DStart)
Instr lalu lea Itemp0, #8, Itemp0;

/* need to align to page boundary */
Instr lalu ash Itemp0, #52, Itemp1;
Instr lalu ash Itemp1, #-52, Itemp1;
Instr lalu leq Itemp1, #0, cc1;
Instr lalu cf cc1 lmm #0x1000, Itemp1;
Instr lalu cf cc1 lea Itemp0, Itemp1, Itemp0;

/* zero out low 12 bits */
Instr lalu lmm #Oxffff, Itemp1;
Instr lalu shoru #0xffff, Itemp1;
Instr lalu shoru #0xffff, Itemp1;
Instr lalu shoru #0xf000, Itemp1;
Instr lalu and Itemp1, Itemp0, Itemp0;
Instr lalu setptr Itemp0, Itemp0;

CONSTRUCT_LONG_PTR(_Sys_Memory_End, Intarg1, DStart)
Instr menu sts  ua, 1, Itemp0, Intarg1;
```

```
    /* this _Sys_Memory_End is used as the base for internal
       physical-memory allocation until enough stuff is set up
       for the physical memory manager to operate. That is,
       this is enough to start allocating stacks for the system
       threads, and so on. */

    RETURN

/**********************************************************
 *
 * Set up On-Node Virtual Memory
 *
 **********************************************************/

data:
_vmem_allocated: asciz "Returned: $p\n";
_vmem_print2:    asciz "VPN: 0x%x\n";
_vmem_node:      asciz "Node #: 0x%x\n";
text:
_virtual_memory_setup:
    /* Create start of address space for virtual memory */
    /* dependent on which node we reside! */

    GET_FRAME

    LIBCALL(NodeId)
    Instr lalu leq lRetVal, 10, cc0;
    Instr lalu ct cc0 br _node0codeA;
    Instr ; Instr ; Instr ;

    Instr lalu leq lRetVal, #1, cc0;
    Instr lalu ct cc0 br _node1codeA;
    Instr ; Instr ; Instr ;

    Instr lalu leq lRetVal, #2, cc0;
    Instr lalu ct cc0 br _node2codeA;
    Instr ; Instr ; Instr ;

    Instr lalu leq lRetVal, #3, cc0;
    Instr lalu ct cc0 br _node3codeA;
    Instr ; Instr ; Instr ;

    Instr lalu br fall;
    Instr ; Instr ; Instr ;

_node3codeA:
    /* on node 3, address starts at 0x20000000 = 256 MB ??? */

    FCALL(_buddyInit)

    Instr lalu lmm #0x1540, ltemp0;
    Instr lalu shoru ##0x0000, ltemp0;
    Instr lalu shoru ##0x0160, ltemp0;
    Instr lalu shoru ##0x0000, ltemp0;
    Instr lalu setptr ltemp0, ltemp0;

    Instr lalu mov ltemp0, lntarg0;

    FCALL(_buddyPrime)

    Instr lalu br _nodeAdone;
    Instr ; Instr ; Instr ;

_node2codeA:
    /* on node 2, address starts at 0x20000000 = 256 MB ??? */
```

```
    FCALL(_buddyInit)

    Instr lalu lmm #0x1540, ltemp0;
    Instr lalu shoru ##0x0000, ltemp0;
    Instr lalu shoru ##0x0140, ltemp0;
    Instr lalu shoru ##0x0000, ltemp0;
    Instr lalu setptr ltemp0, ltemp0;

    Instr lalu mov ltemp0, lntarg0;

    FCALL(_buddyPrime)

    Instr lalu br _nodeAdone;
    Instr ; Instr ; Instr ;

_node1codeA:
    /* on node 1, address starts at 0x20000000 = 256 MB ??? */

    FCALL(_buddyInit)

    Instr lalu lmm #0x1540, ltemp0;
    Instr lalu shoru ##0x0000, ltemp0;
    Instr lalu shoru ##0x0120, ltemp0;
    Instr lalu shoru ##0x0000, ltemp0;
    Instr lalu setptr ltemp0, ltemp0;
    Instr lalu mov ltemp0, lntarg0;

    FCALL(_buddyPrime)

    Instr lalu br _nodeAdone;
    Instr ; Instr ; Instr ;

_node0codeA:
    /* on node 0, address starts at 0x10000000 = 256 MB ??? */

    FCALL(_buddyInit)

    Instr lalu lmm #0x1540, ltemp0;
    Instr lalu shoru ##0x0000, ltemp0;
    Instr lalu shoru ##0x0100, ltemp0;
    Instr lalu shoru ##0x0000, ltemp0;
    Instr lalu setptr ltemp0, ltemp0;
    Instr lalu mov ltemp0, lntarg0;

    FCALL(_buddyPrime)
    LIBCALL(CCInit)
    LIBCALL(SQInit)

_nodeAdone:
    /* ok, write gtlb into config space */
    Instr lalu mov P_CONFIG, ltemp0;           -- Set up prot field
    Instr lalu ash ltemp0, POINTER_LENGTH, ltemp0;   -- add length field
    Instr lalu add ltemp0, PAGE_SEG_LENGTH, ltemp0;  -- Get ready for
    Instr lalu ash ltemp0, POINTER_ADDRESS, ltemp0;  -- starting address

    CONSTRUCT_LONG_LABEL(Config_Global_Base, ltemp1)
    Instr lalu add ltemp1, ltemp0, ltemp0;
    Instr lalu setptr ltemp0, ltemp0;

    CONSTRUCT_LONG_PTR(GLOBAL_PAGE_TABLE, ltemp1, DStart)
    Instr menu ld ltemp1, #8, lntarg1;       -- first entry
    Instr menu st lntarg1, #8, ltemp0;
    Instr menu ld ltemp1, #8, lntarg1;
```

```
        instr memu st intarg1, #8, ltemp0;

        instr memu ld ltemp1, #8, intarg1;        -- second entry
        instr memu st intarg1, #8, ltemp0;
        instr memu ld ltemp1, #8, intarg1;
        instr memu st intarg1, #8, ltemp0;

_vmem_setup_done:
        RETURN

_physical_memory_setup1:
        GET_FRAME

        /* load end of stack and align to page boundary */
        LOAD_FAR_LABEL(_Sys_Memory_End, ltemp0, DStart)
        instr lalu lea ltemp0, #8, ltemp0;

        /* need to align to page boundary */
        instr lalu ash ltemp0, #52, ltemp1;
        instr lalu ash ltemp1, #-52, ltemp1;
        instr lalu leq ltemp1, #0, cc1;
        instr lalu cf cc1 lmm #0x1000, ltemp1, ltemp1;
        instr lalu cf cc1 lea ltemp0, ltemp1, ltemp0;

        /* zero out low 12 bits */
        instr lalu lmm ##0xffff, ltemp1;
        instr lalu shoru ##0xffff, ltemp1;
        instr lalu shoru ##0xffff, ltemp1;
        instr lalu shoru ##0xf000, ltemp1;
        instr lalu and ltemp1, ltemp0, ltemp0;
        instr lalu setptr ltemp0, ltemp0;

        instr lalu lsh ltemp0, #10, intarg0;       -- calculate base page to prime
        instr lalu lsh intarg0, #-22, intarg0;     -- the physical memory manager
        PUSH(intarg0)
        instr lalu mov SP, AP;                     -- Physical Page Manager
        FCALL(_PPM_Init)                           -- Initialization
        SPOP(intarg0)

        RETURN

/*******************************************************************
 *
 * These are stubs and helpful system routines
 *
 *******************************************************************/

_tSleepX::
        instr memu mbar;
        GET_FRAME
        PUSH(DStart)
        instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
        instr lalu leab IP, ltemp0, DStart;
        instr memu ld DStart, DStart;

        FCALL(_SYStSleep)

        SPOP(DStart)

        RETURN

_tUnparseX::
        instr memu mbar;
        GET_FRAME
        PUSH(DStart)
        instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
```

```
        instr lalu leab IP, ltemp0, DStart;
        instr memu ld DStart, DStart;

        FCALL(_SYStUnparse)

        SPOP(DStart)
        RETURN

_tForkX::
        instr memu mbar;
        GET_FRAME
        PUSH(DStart)
        instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
        instr lalu leab IP, ltemp0, DStart;
        instr memu ld DStart, DStart;

        FCALL(_SYStFork)

        SPOP(DStart)
        RETURN

_tExitX::
        instr memu mbar;
        GET_FRAME
        PUSH(DStart)
        instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
        instr lalu leab IP, ltemp0, DStart;
        instr memu ld DStart, DStart;

        FCALL(_SYStExit)

        instr ;
        instr ;
        instr lalu empty #(ITEMP0_EMPTY_MASK);
        instr lalu mov ltemp0, ltemp0;

        SPOP(DStart)
        RETURN

_tSignalX::
        instr memu mbar;
        GET_FRAME
        PUSH(DStart)
        instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
        instr lalu leab IP, ltemp0, DStart;
        instr memu ld DStart, DStart;

        FCALL(_SYStSignal)

        SPOP(DStart)
        RETURN

_getSelfTCX::
        instr memu mbar;
        GET_FRAME
        PUSH(DStart)
        instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
        instr lalu leab IP, ltemp0, DStart;
        instr memu ld DStart, DStart;

        FCALL(_SYSgetSelfTC)

        SPOP(DStart)
        RETURN

_getParentX::
```

```
        instr menu mbar;
        GET_FRAME
        PUSH(DStart)
        instr lalu imm __SYSTEM_UDAT_PTR__, itemp0;
        instr lalu leab IP, itemp0, DStart;
        instr menu ld DStart, DStart;

        FCALL(_SYSgetParent)

        SPOP(DStart)
        RETURN

_sysMakeCP::
        instr lalu loh intarg0, #17, itemp0;
        instr lalu imm #(0x74440), itemp0;
        instr lalu shoru #0x0000, itemp0;
        instr lalu shoru #0x0000, itemp0;
        instr lalu jmp RETIP;
        instr lalu shoru #0x0000, itemp0;
        instr lalu or intarg0, itemp0, intarg0;
        instr lalu setptr intarg0, intarg0;

_sysSetPtr::
        instr lalu jmp RETIP;
        instr lalu setptr intarg0, intarg0;
        instr ;
        instr ;

_sysGPRD::
        instr lalu jmp RETIP;
        instr menu gprb intarg0, itemp0, LCC1;
        instr lalu ct LCC1 mov itemp0, intarg0;
        instr lalu cf LCC1 mov #-1, intarg0;

_sysFLNE::
        instr lalu jmp RETIP;
        instr menu fine intarg0;
        instr ;
        instr ;

_sysPUTCSTAT::
        instr     lalu leq intarg1, #BSB_INVALID, LCC1;
        instr     menu ct LCC1 putcstat intarg0, #BSB_INVALID, intarg0
        instr     lalu leq intarg1, #BSB_READONLY, LCC2;
        instr     menu ct LCC2 putcstat intarg0, #BSB_READONLY, intarg0
        instr     lalu leq intarg1, #BSB_EXCLUSIVE, LCC3;
        instr     menu ct LCC3 putcstat intarg0, #BSB_EXCLUSIVE, intarg0
        instr     lalu leq intarg1, #BSB_DIRTY, LCC1;
        instr     menu ct LCC1 putcstat intarg0, #BSB_DIRTY, intarg0;
        instr     menu mbar
        instr     lalu mov #1, intarg2;
        instr     lalu llt intarg2, intarg1, LCC2;
        instr     lalu cf LCC2 br Fail;
        instr     lalu cf LCC2 jmp RETIP;
        instr     lalu cf LCC2 mov intarg0, intarg0;
        instr     ;
        instr     ;

_sysSMSMessage::
        PUSH(intarg0)
        instr menu sms MsgBody, #8, intarg0;
        instr menu sms MsgBody, #8, intarg0;
        instr menu sms MsgBody, #8, intarg0;
        instr menu sms MsgBody, #8, intarg0;
        instr menu sms MsgBody, #8, intarg0;
        instr menu sms MsgBody, #8, intarg0;
```

```
        instr menu sms MsgBody, #8, intarg0;
        instr menu sms MsgBody, #8, intarg0;
        instr menu sms MsgBody, #8, intarg0;
        POP(intarg0)
        instr menu fine intarg0;
        instr menu mbar;
        instr lalu jmp RETIP;
        instr ;
        instr ;

_sysPushDirty::
        /* the next words are filled with the cache line */
        PUSH(intarg0)
        instr     falu empty #0x0014;
        instr     lalu mov intarg1, FMC0;          -- destination node
        instr     lalu mov intarg2, #2;            -- destination node
        instr menu ld intarg0, #8, FMC1;
        instr menu ld intarg0, #8, FMC2;
        instr menu ld intarg0, #8, FMC3;
        instr menu ld intarg0, #8, FMC4;
        instr menu ld intarg0, #8, FMC5;
        instr menu ld intarg0, #8, FMC6;
        instr menu ld intarg0, #8, FMC7;
        instr menu ld intarg0, #8, FMC8;
        SPOP(intarg0)
        instr menu fine intarg0;
        instr menu mbar;
        /* message is cons'd up. Now just send it */
        MAKE_XM_PTR(MSG_ccreturnDirty)
        instr     falu fsndipt #9, #2, FMCIP, LCC3;    -- this is an ACK to P1
        instr     lalu ct LCC3 jmp RETIP;
        instr     ;
        instr     ;
        instr     ;

        CALLFAIL

_sysSendLine::
        /* the next words are filled with the cache line */
        PUSH(intarg0)
        instr     falu empty #0x0014;
        instr     lalu mov intarg1, FMC0;          -- yank buffer ptr
        instr     lalu mov intarg2, #2;            -- destination node
        instr menu ld intarg0, #8, FMC1;
        instr menu ld intarg0, #8, FMC2;
        instr menu ld intarg0, #8, FMC3;
        instr menu ld intarg0, #8, FMC4;
        instr menu ld intarg0, #8, FMC5;
        instr menu ld intarg0, #8, FMC6;
        instr menu ld intarg0, #8, FMC7;
        instr menu ld intarg0, #8, FMC8;
        SPOP(intarg0)
        instr menu fine intarg0;
        instr menu mbar;
        /* message is cons'd up. Now just send it */
        SYSPUTLOCK(_sqlock)
        MAKE_XM_PTR(MSG_ccreturnyankFull)
        instr     falu fsndipt #9, #2, FMCIP, LCC3;    -- this is an ACK to P1
        instr     lalu ct LCC3 jmp RETIP;
        instr     ;
        instr     ;
        instr     ;

        CALLFAIL
```

```
_sysSendInvalidateAck::
        SYSPUTLOCK(_sqlock)
        Instr   falu empty #0x0014;
        Instr   lalu mov intarg0, FMC0;        -- yank buffer ptr
        Instr   lalu mov intarg1, f2;          -- destination node
        MAKE_XM_PTR(MSG_ccreturnyank)
        Instr   falu fsndipt #1, f2, FMCIP, LCC3;    -- this is an ACK to P1
        Instr   lalu ct LCC3 jmp RETIP;
        Instr   ;
        Instr   ;
        Instr   ;
        CALLFAIL

_sysSendInvalidate::
        Instr   falu empty #0x0034;
        Instr   lalu mov intarg1, FMC0;
        Instr   lalu mov intarg2, FMC1;
        Instr   lalu mov intarg0, f2;
        MAKE_XM_PTR(MSG_ccinvalidate)
        Instr   falu fsndipt #2, f2, FMCIP, LCC3;
        Instr   lalu ct LCC3 jmp RETIP;
        Instr   ;
        Instr   ;
        Instr   ;
        CALLFAIL

_sysNackRO::
        SYSPUTLOCK(_ccdirlock)
        Instr   falu empty #0x00f4;
        Instr   lalu mov intarg0, FMC0;         -- address
        Instr   lalu mov intarg1, FMC1;         -- header
        Instr   lalu mov intarg2, FMC2;         -- opdata
        Instr   lalu mov intarg3, FMC3;         -- faultcp
        Instr   lalu mov intarg4, f2;           -- faultcp
        MAKE_XM_PTR(MSG_ccNackRO)
        Instr   falu fsndipt #4, f2, FMCIP, LCC3;    -- this is an ACK to P1
        Instr   lalu ct LCC3 jmp RETIP;
        Instr   ;
        Instr   ;
        Instr   ;
        CALLFAIL

_sysNackRW::
        SYSPUTLOCK(_ccdirlock)
        Instr   falu empty #0x00f4;
        Instr   lalu mov intarg0, FMC0;         -- address
        Instr   lalu mov intarg1, FMC1;         -- header
        Instr   lalu mov intarg2, FMC2;         -- opdata
        Instr   lalu mov intarg3, FMC3;         -- faultcp
        Instr   lalu mov intarg4, f2;           -- faultcp
        MAKE_XM_PTR(MSG_ccNackRW)
        Instr   falu fsndipt #4, f2, FMCIP, LCC3;    -- this is an ACK to P1
        Instr   lalu ct LCC3 jmp RETIP;
        Instr   ;
        Instr   ;
        Instr   ;
        CALLFAIL

_sysReadAndSendRO::
        PUSH(intarg0)
        Instr   falu empty #0x34;
        Instr   lalu mov intarg1, FMC0;         /* address */
        Instr   lalu mov intarg2, f2;
        Instr   lalu mov intarg3, FMC1;         /* header */
```

```
        Instr   memu ld intarg0, #8, FMC2;
        Instr   memu ld intarg0, #8, FMC3;
        Instr   memu ld intarg0, #8, FMC4;
        Instr   memu ld intarg0, #8, FMC5;
        Instr   memu ld intarg0, #8, FMC6;
        Instr   memu ld intarg0, #8, FMC7;
        Instr   memu ld intarg0, #8, FMC8;
        Instr   memu ld intarg0, #8, FMC9;
        POP(intarg0)
        Instr   memu flne intarg0;
        SYSPUTLOCK(_ccdirlock)

        MAKE_XM_PTR(MSG_ccreturnLoad)
        Instr   falu fsndipt #10, f2, FMCIP, LCC3;
        Instr   lalu ct LCC3 jmp RETIP;
        Instr   ;
        Instr   ;
        Instr   ;
        CALLFAIL

_sysReadAndSendX::
        /* perform the store locally */
        Instr   lalu and intarg1, #0x1f, ltemp0;
        Instr   lalu or ltemp0, intarg0, ltemp0;
        Instr   lalu setptr ltemp0, ltemp0;
        Instr   memu st intarg2, ltemp0;

PUSH(intarg0)
        Instr   falu empty #0x34;
        Instr   lalu mov intarg1, FMC0;
        Instr   lalu mov intarg3, f2;
        Instr   lalu mov intarg4, FMC1;
        Instr   memu ld intarg0, #8, FMC2;
        Instr   memu ld intarg0, #8, FMC3;
        Instr   memu ld intarg0, #8, FMC4;
        Instr   memu ld intarg0, #8, FMC5;
        Instr   memu ld intarg0, #8, FMC6;
        Instr   memu ld intarg0, #8, FMC7;
        Instr   memu ld intarg0, #8, FMC8;
        Instr   memu ld intarg0, #8, FMC9;
        POP(intarg0)
        Instr   memu flne intarg0;
        SYSPUTLOCK(_ccdirlock)

        MAKE_XM_PTR(MSG_ccreturnStore)
        Instr   falu fsndipt #10, f2, FMCIP, LCC3;
        Instr   lalu ct LCC3 jmp RETIP;
        Instr   ;
        Instr   ;
        Instr   ;
        CALLFAIL

data;
_threadlock::
        u64 0;

text;

_sysGetLock::
_sysgetlock_point1:
        Instr   memu stscnd cf, 1, RETIP, intarg0, LCC1;
        Instr   lalu ct LCC1 jmp RETIP;
        Instr   lalu cf LCC1 br _sysgetlock_point1;
        Instr   ; Instr ; Instr ;
```

```
_sysPutLock::
    Instr menu mbar;
    Instr menu stscnd ct, 0, RETIP, Intarg0, LCC1;
    Instr ;
    Instr lalu ct LCC1 jmp RETIP;
    Instr ; Instr ; Instr ;
    CALLFAIL

_sysSignalSleep::
    GET_FRAME

    /* EMPTY out l10 so that we are guaranteed to block as soon
       as EH starts running in response to a software job request */
    /* for this reason we MUST make sure that after l10 is emptied,
       no one writes it in any further calls we make! */
    Instr lalu empty #(INTARG4_EMPTY_MASK);

    SYSGETLOCK(generic_signal_eh2,_event_lockword)
    CONSTRUCT_LONG_PTR(_event_buf_free, ltemp0, DStart)  -- free buffer ptr
    Instr menu ld ltemp0, ltemp0;                         -- type of event
    Instr menu st Intarg0, #8, ltemp0;

    CONSTRUCT_LONG_PTR(_event_buf_end, Intarg2, DStart)
    CONSTRUCT_LONG_PTR(_event_buf_start, Intarg3, DStart)

    Instr    lalu leq ltemp2, Intarg2, LCC1;
    Instr    lalu ct LCC1 mov Intarg3, ltemp0;

    Instr    menu st Intarg1, #8, ltemp0;                 -- TC
    Instr    lalu leq ltemp0, Intarg2, LCC1;
    Instr    lalu ct LCC1 mov Intarg3, ltemp0;

    CONSTRUCT_LONG_PTR(_event_buf_free, Intarg2, DStart)
    Instr menu st ltemp0, Intarg2;

    SYSPUTLOCK(_event_lockword)

    /* now here we are to try to enqueue the event into the hardware
       queue, only if a global state variable is unset */
    CONSTRUCT_LONG_PTR(_event_signalword, Intarg0, DStart)
    Instr menu stscnd cf, 1, Intarg0, Intarg0, LCC1;
    Instr lalu cf LCC1 br _sysSignalSleepDone;
    Instr ; Instr ; Instr ;

    /* now we can signal the event */
    Instr lalu lmm #0x7fff, Intarg0;
    Instr lalu lmm #0x7ff0, Intarg1;
    LIBCALL(EventcqAdd)

_sysSignalSleepDone:
    Instr lalu mov Intarg4, Intarg0;
    RETURN

_sysHardwareSignalEH::
    GET_FRAME

    /* now here we are to try to enqueue the event into the hardware
       queue, only if a global state variable is unset */
    CONSTRUCT_LONG_PTR(_event_signalword, Intarg0, DStart)
    Instr menu stscnd cf, 1, Intarg0, Intarg0, LCC1;
    Instr lalu cf LCC1 br _sysHardwareSignalEHDone;
    Instr ; Instr ; Instr ;

    /* now we can signal the event */
    Instr lalu lmm #0x7fff, Intarg0;
```

```
    Instr lalu lmm #0x7ff0, Intarg1;
    LIBCALL(EventcqAdd)

_sysHardwareSignalEHDone:
    RETURN

_sysGetNodeId::
    GET_FRAME
    LIBCALL(NodeId)
    RETURN

_sysSendWake::
    /* send a wake message to the home node of the context */
    Instr lalu empty #0xc030;
    Instr lalu mov Intarg1, FMC0;
    Instr lalu mov Intarg2, FMC1;
    Instr lalu mov Intarg0, FMCDest;

    MAKE_XM_PTR(MSG_tWake)                /* set FMCIP appropriately */

    Instr lalu fsndOpt #2, FMCDest, FMCIP, LCC1;
    Instr lalu jmp RETIP;
    Instr lalu cf LCC1 mov #0, Intarg0;
    Instr lalu ct LCC1 mov #1, Intarg0;
    Instr ;

_sysSendSleep::
    /* send a sleep request to the home node of a signal word */
    Instr lalu empty #0xc070;
    Instr lalu mov Intarg0, FMC0;
    Instr menu gprb Intarg0, ltemp0, LCC1;
    Instr lalu cf LCC1 br fall;
    Instr lalu ct LCC1 mov ltemp0, FMCDest;
    Instr lalu mov Intarg1, FMC1;
    Instr lalu mov Intarg2, FMC2;

    MAKE_XM_PTR(MSG_tSleep)              /* set FMCIP appropriately */

    Instr lalu fsndOpt #3, FMCDest, FMCIP, LCC1;
    Instr lalu jmp RETIP;
    Instr lalu cf LCC1 mov #0, Intarg0;
    Instr lalu ct LCC1 mov #1, Intarg0;
    Instr ;

_sysSendSignal::
    /* send a signal to the home node of a signal word */
    Instr lalu empty #0xc030;
    Instr lalu mov Intarg1, FMC0;
    Instr lalu mov Intarg2, FMC1;
    Instr lalu mov Intarg0, FMCDest;

    MAKE_XM_PTR(MSG_tSignal)            /* set FMCIP appropriately */

    Instr lalu fsndOpt #2, FMCDest, FMCIP, LCC1;
    Instr lalu jmp RETIP;
    Instr lalu cf LCC1 mov #0, Intarg0;
    Instr lalu ct LCC1 mov #1, Intarg0;
    Instr ;

/*************************************************/

data:
LTLB_IP:        p64 LTLB_HANDLER_START;
EVENT_IP:       p64 EVENT_HANDLER_START;
```

```
    Instr lalu lmm #0x7ff0, Intarg1;
    LIBCALL(EventcqAdd)

_sysHardwareSignalEHDone:
    RETURN
```

```
DISPATCH_IP:    p64 DISPATCH_HANDLER_START;
GLOBAL_PAGE_TABLE::
u64 0x8000000000000000;  /* virtual page 0 */
u64 0x00000000000b9002;  /* mapped across 4x1x1 */
                         /* 256 Pages/node, total page length = 2^23 bytes */
u64 0x8000000000010000;  /* virtual page 0x1000 = 1024 */
u64 0x00000000000c1202;  /* mapped across 4x1x1 */
                         /* 512 Pages/node, total page length = 2^24 bytes */

#if 0
u64 0x8000000000000000;
u64 0x00000000000b9600;
u64 0x8000000000010001;
u64 0x00000000000b9600;
u64 0x8000000000020002;
u64 0x00000000000b9600;
u64 0x8000000000030003;
u64 0x00000000000b9600;
u64 0x8000000000040020;
u64 0x00000000000b9600;
u64 0x8000000000050021;
u64 0x00000000000b9600;
u64 0x8000000000060022;
u64 0x00000000000b9600;
u64 0x8000000000070023;
u64 0x00000000000b9600;
u64 0x8000000000080040;
u64 0x00000000000b9600;
u64 0x8000000000090041;
u64 0x00000000000b9600;
u64 0x80000000000a0042;
u64 0x00000000000b9600;
u64 0x80000000000b0043;
u64 0x00000000000b9600;
u64 0x80000000000c0060;
u64 0x00000000000b9600;
u64 0x80000000000d0061;
u64 0x00000000000b9600;
u64 0x80000000000e0062;
u64 0x00000000000b9600;
u64 0x80000000000f0063;
u64 0x00000000000b9600;
#endif
u64 0x1111111122222222; /* end of boot system data */
end;
```

```
/* **************************************************
 *
 * M Machine runtime system Physical-Page Manager helper functions
 *
 * Written by       Yevgeny Gurevich
 * Version          0.80
 * Modification Date    06/21/95
 * Modification Date    08/08/95
 *
 * **************************************************/

#include "newreg.h"
#include <libcall.h>
#include "helpmacros.h"
#include "opcodes.h"
#include "pointers.h"

#define CALLFAIL        Instr lalu br Fail; Instr ; Instr ; Instr ;

data;
/* ensures that only one process is calling on the PPM at any one time. */
lpt_lookup_lock:
        u64 0x0;

text;
_PPM_reclaim_remote::
        SYSGETLOCK(lpt_lookup_point4, lpt_lookup_lock)

        /* page 80005 is reclaim_remote request directed at LTLB Thread */
        Instr lalu lmm #0x1700, ltemp0;
        Instr lalu shoru ##0x0000, ltemp0;
        Instr lalu shoru ##0x8000, ltemp0;
        Instr lalu shoru ##0x5000, ltemp0;
        Instr lalu setptr ltemp0, ltemp0;
        Instr memu stsu ct, 1, lntarg0, ltemp0, LCC1;
        Instr lalu cf LCC1 br Fail;                   -- block for LTLB to
                                                      -- complete

        LOAD_FAR_LABEL(_ltlb_data_for_MH, IRetVal, DStart)   -- return
        SYSPUTLOCK(lpt_lookup_lock)
        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

_PPM_unmap::
        SYSGETLOCK(lptlookup_point3, lpt_lookup_lock)

        /* page 80000 is unmap request directed at LTLB Thread */
        Instr lalu lmm #0x1700, ltemp0;
        Instr lalu shoru ##0x0000, ltemp0;
        Instr lalu shoru ##0x8000, ltemp0;
        Instr lalu shoru ##0x0000, ltemp0;
        Instr lalu setptr ltemp0, ltemp0;
        Instr memu stsu ct, 1, lntarg0, ltemp0, LCC1;
        Instr lalu cf LCC1 br Fail;                   -- block for LTLB to
                                                      -- complete

        LOAD_FAR_LABEL(_ltlb_data_for_MH, IRetVal, DStart)   -- return
        SYSPUTLOCK(lpt_lookup_lock)
        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

_PPM_local2remote::
        SYSGETLOCK(lptlookup_point2, lpt_lookup_lock)
```

```
        /* page 80006 is local2remote request directed at LTLB Thread */
        Instr lalu lmm #0x1700, ltemp0;
        Instr lalu shoru ##0x0000, ltemp0;
        Instr lalu shoru ##0x8000, ltemp0;
        Instr lalu shoru ##0x6000, ltemp0;
        Instr lalu setptr ltemp0, ltemp0;
        Instr memu stsu ct, 1, lntarg0, ltemp0, LCC1;
        Instr lalu cf LCC1 br Fail;                   -- block for LTLB to
                                                      -- complete

        LOAD_FAR_LABEL(_ltlb_data_for_MH, IRetVal, DStart)   -- return
        SYSPUTLOCK(lpt_lookup_lock)
        Instr lalu jmp RETIP;
        Instr ; Instr ;

_PPM_lookup::
        /* lntarg0: virtual address                              */
        /* Lookup the physical page backing address in arg0      */
        /* returns ppn in location _ltlb_data_for_MH             */
        /* may change this later                                 */

        SYSGETLOCK(lptlookup_point1, lpt_lookup_lock)

        /* page 80002 is lookup request directed at LTLB Thread */
        Instr lalu lmm #0x1700, ltemp0;
        Instr lalu shoru ##0x0000, ltemp0;
        Instr lalu shoru ##0x8000, ltemp0;
        Instr lalu shoru ##0x2000, ltemp0;
        Instr lalu setptr ltemp0, ltemp0;
        Instr memu stsu ct, 1, lntarg0, ltemp0, LCC1;
        Instr lalu cf LCC1 br Fail;                   -- block for LTLB to
                                                      -- complete

        LOAD_FAR_LABEL(_ltlb_data_for_MH, IRetVal, DStart)   -- return
        SYSPUTLOCK(lpt_lookup_lock)
        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

_get_offsetptr_into_ppn::
        /* lntarg0: physical page number                          */
        /* lntarg1: virtual address with low 6 bits having the cache
                    line offset                                   */
        /* returns physical pointer into the backing page with the proper
           offset (takes the middle 6 bits of the VPN to calculate the
           proper offset. NOT the low 12 bits!)                    */

        Instr lalu lsh lntarg0, #12, lntarg0;        -- ppn
        Instr lalu lmm ##0x9700, ltemp0;             -- protection bits
        Instr lalu lsh ltemp0, #48, ltemp0;          -- protection bits
        Instr lalu or lntarg0, ltemp0, lntarg0;

        Instr lalu lsh lntarg1, #52, ltemp0;
        Instr lalu lsh ltemp0, #-58, ltemp0;
        Instr lalu lsh ltemp0, #6, ltemp0;           -- just cache line addr
        Instr lalu or lntarg0, ltemp0, lntarg0;

        Instr lalu setptr lntarg0, lntarg0;
        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

end;
```

```
#include </home/mm/tools/mslm/v1.5/yev/mslm_lib.h>
#include "/home/mm/tools/mslm/v1.5/yev/ccmagic.h"
#include "newreg.h"
#include <libcall.h>
#include "pointers.h"
#include "helpmacros.h"

#define MALLOCTEST 1
#define LOADER_READY 1

/* sets up the jump pointers in the jump table       */
/* call this once, before any uses of the C library  */
/* pstart is assumed to be the pointer to the USER    */
/* named data segment.  */

/* generates the magic jump pointer to hook in with simulator */
#define MAKE_MAGIC_PTR(x)                               \
Instr lalu lmm #*(P_PHYSICAL << 12), ltemp0;            \
Instr lalu shoru #0, ltemp0;                            \
Instr lalu shoru #*(x & 0xffff0000) >> 16), ltemp0;     \
Instr lalu shoru #*(x & 0x0000ffff), ltemp0;            \
Instr lalu sotptr ltemp0, ltemp0;

/************************************************/

.text:
INIT_Lib::
        /* Initialize the C library.  So far, nothing to do */
        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

/************************************************/

_sin: /* double sin(double x) */
        Instr lalu mov fpArg0, fpArg0;
#define SIN_DELAY 100
_sin_delay_loop:
        Instr      lalu mov #(SIN_DELAY / 5 - 1), ltemp0;
        Instr      lalu leq ltemp0, #0, cc0;
        Instr      lalu cf cc0 br _sin_delay_loop;
        Instr      lalu sub ltemp0, #1, ltemp0;
        Instr      lalu nop;
        Instr      lalu nop;

MAKE_MAGIC_PTR(SIN_MAGIC)
        Instr lalu jmp ltemp0; /* hook to simulator */
        Instr;
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

_cos: /* double cos(double x) */
        Instr lalu mov fpArg0, fpArg0;
#define COSINE_DELAY 100
_cos_delay_loop:
        Instr      lalu mov #(COSINE_DELAY / 5 - 1), ltemp0;
        Instr      lalu leq ltemp0, #0, cc0;
        Instr      lalu cf cc0 br _cos_delay_loop;
        Instr      lalu sub ltemp0, #1, ltemp0;
        Instr      lalu nop;
        Instr      lalu nop;

MAKE_MAGIC_PTR(COS_MAGIC)
        Instr lalu jmp ltemp0; /* hook to simulator */
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_exit: /* void exit(int status) */
        Instr menu mbar;
        LOAD_FAR_LABEL(_cleanup, ltemp0, DStart)
        Instr lalu jmp ltemp0;
        Instr ;
        Instr ;

_printf: /* int printf(char *format, arg ...) */
        /* assume Arg0 contains a pointer to the format string */
        Instr menu mbar;
_kprintf::
        PUSH(lntarg0)
        Instr menu ld lntarg0, ltemp0;
_printf_loadloop:
        Instr lalu extb ltemp0, #7, ltemp0;
        Instr lalu leq ltemp0, 10, LCC1;
        Instr lalu cf LCC1 br _printf_continue;
        Instr lalu cf LCC1 br _printf_loadloop;
        Instr lalu cf LCC1 lea lntarg0, #8, lntarg0;
        Instr menu cf LCC1 ld lntarg0, ltemp0;
        Instr ;
_printf_continue:
        SPOP(lntarg0)

----    Instr menu ld lntarg0, ltemp0;    -- fetch it in to make sure...
----    Instr menu mbar;
        Instr menu ld AP, ltemp0;
        Instr lalu mov ltemp0, ltemp0;
MAKE_MAGIC_PTR(PRINTF_MAGIC)
        Instr lalu jmp ltemp0;    /* hook to simulator */
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_scanf: /* int scanf(char *format, arg ...) */
        /* assume Arg0 contains a pointer to the format string */
        args 1-5 contain pointers to memory addresses
        Instr menu mbar;
        Instr lalu mov lntarg0, lntarg0;    -- validate format argument
MAKE_MAGIC_PTR(SCANF_MAGIC)
        Instr lalu jmp ltemp0;    -- hook to simulator
```

```
        Instr;
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_strchr:: /* char *strchr(const char *, int c)  */
        Instr memu mbar;
        Instr lalu mov Intarg0, Intarg0;        -- validate argument
        Instr lalu mov Intarg1, Intarg0;        -- validate argument
        MAKE_MAGIC_PTR(STRCHR_MAGIC)
        Instr lalu jmp Itemp0;                  -- hook to simulator
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_fopen:: /* FILE* fopen(const char *filename, const char *mode)   */
        /* assume Arg0 contains a pointer to the filename string */
        /* assume Arg1 contains a pointer to the mode string */
        Instr lalu mov Intarg0, Intarg0;        -- validate filename argument
        Instr lalu mov Intarg1, Intarg1;        -- validate mode argument
        MAKE_MAGIC_PTR(FOPEN_MAGIC)
        Instr lalu jmp Itemp0;                  -- hook to simulator
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_fclose:: /* int fclose(FILE *file)  */
        /* assume Arg0 contains a FILE pointer */
        Instr lalu mov Intarg0, Intarg0;        -- validate file argument
        MAKE_MAGIC_PTR(FCLOSE_MAGIC)
        Instr lalu jmp Itemp0;                  -- hook to simulator
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_fprintf:: /* int fprintf(FILE *file, char *format, arg . . .)   */
        /* assume Arg0 contains a pointer to the file */
        /* assume Arg1 contains a pointer to the format string */
        Instr memu mbar;
        Instr lalu mov Intarg0, Intarg0;        -- validate file argument
        Instr lalu mov Intarg0, Intarg1;        -- validate format argument
        MAKE_MAGIC_PTR(FPRINTF_MAGIC)
        Instr lalu jmp Itemp0;                  -- hook to simulator
        Instr;
        Instr;
        Instr;
```

```
        Instr;
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_gets:: /* char *gets(char *)   */
        /* assume Arg0 contains a pointer to the dest string */
        Instr memu mbar;
        Instr lalu mov Intarg0, Intarg0;        -- validate format argument
        MAKE_MAGIC_PTR(GETS_MAGIC)
        Instr lalu jmp Itemp0;                  -- hook to simulator
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_getchar:: /* int getchar()   */
        Instr memu mbar;
        MAKE_MAGIC_PTR(GETCHAR_MAGIC)
        Instr lalu jmp Itemp0;                  -- hook to simulator
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_putchar:: /* int putchar(int)   */
        Instr memu mbar;
        MAKE_MAGIC_PTR(PUTCHAR_MAGIC)
        Instr lalu jmp Itemp0;                  -- hook to simulator
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

/* string functions */
_strcpy:: /* char *strcpy(char *, const char *)   */
        Instr memu mbar;
        Instr lalu mov Intarg0, Intarg0;        -- validate argument
        Instr lalu mov Intarg1, Intarg1;        -- validate argument
        MAKE_MAGIC_PTR(STRCPY_MAGIC)
        Instr lalu jmp Itemp0;                  -- hook to simulator
        Instr;
        Instr;
        Instr;
```

```
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_lseek:  Instr lalu mov Intarg0, Intarg0;
        Instr lalu mov Intarg1, Intarg1;
        Instr lalu mov Intarg2, Intarg2;
        MAKE_MAGIC_PTR(FSEEK_MAGIC)
        Instr lalu jmp Itemp0;              /* hook to simulator */
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_fwrite:: /* size_t fwrite(void *ptr,
                          size_t element_size,
                          size_t count,
                          FILE *stream) */
        /* arg0 contains pointer to the data */
        /* arg1 contains size of pointed-to data in bytes */
        /* arg2 contains number of such elements to write */
        /* arg3 contains file pointer */
        Instr lalu mov Intarg0, Intarg0;   /* validate argument */
        Instr lalu mov Intarg1, Intarg1;   /* validate argument */
        Instr lalu mov Intarg2, Intarg2;   /* validate argument */
        Instr lalu mov Intarg3, Intarg3;   /* validate argument */
        MAKE_MAGIC_PTR(FWRITE_MAGIC)
        Instr lalu jmp Itemp0;              /* hook to simulator */
        Instr;
        Instr;
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_fread:: /* size_t fread(void *ptr,
                          size_t element_size,
                          size_t count,
                          FILE *stream) */
        /* arg0 contains pointer to the data */
        /* arg1 contains size of pointed-to data in bytes */
        /* arg2 contains number of such elements to write */
        /* arg3 contains file pointer */
        Instr lalu mov Intarg0, Intarg0;   /* validate argument */
        Instr lalu mov Intarg1, Intarg1;   /* validate argument */
        Instr lalu mov Intarg2, Intarg2;   /* validate argument */
        Instr lalu mov Intarg3, Intarg3;   /* validate argument */
        MAKE_MAGIC_PTR(FREAD_MAGIC)
        Instr lalu jmp Itemp0;              /* hook to simulator */
        Instr;
        Instr;
        Instr;
```

```
        Instr;

        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_bcopy:: /* char* bcopy(char *src, char *dest, int len) */
        /* copy from the src to the destination length number of bytes */
        Instr memu lea Intarg0, Intarg2, Intarg2;     -- limit
_bcopyloop:
        Instr memu ld Intarg0, Itemp1;                -- Itemp1 <- memword
        Instr lalu extb Itemp1, Intarg0, Itemp1;      -- extract single byte

        Instr lalu and Intarg1, #0x7, Itemp0;         -- calc Insb offset
        Instr lalu lsh Itemp0, #32, Itemp0;           -- calc Insb offset
        Instr lalu or Itemp0, Itemp0, Itemp1;         -- set Insb bits

        Instr memu ld Intarg1, #0, Itemp0;            -- Itemp0 <- destword
        Instr lalu Insb Itemp0, Itemp1, Itemp0;       -- Itemp0 <- newdest
        Instr memu st Itemp0, #1, Intarg1;            -- store destword

        Instr lalu lea Intarg0, #1, Intarg0;          -- check bounds
        Instr lalu lle Intarg2, Intarg0, cc0;         -- check bounds
        Instr lalu ct cc0 jmp RETIP;
        Instr lalu cf cc0 br _bcopyloop;
        Instr ;
        Instr ;
        Instr ;

?I._bcopy:: /* char* bcopy(char *src, char *dest, int len) */
        /* copy from the src to the destination length number of bytes */
        Instr memu lea Intarg0, Intarg2, Intarg2;     -- limit
?I._bcopyloop:
        Instr memu ld Intarg0, Itemp1;                -- Itemp1 <- memword
        Instr lalu extb Itemp1, Intarg0, Itemp1;      -- extract single byte

        Instr lalu and Intarg1, #0x7, Itemp0;         -- calc Insb offset
        Instr lalu lsh Itemp0, #32, Itemp0;           -- calc Insb offset
        Instr lalu or Itemp0, Itemp0, Itemp1;         -- set Insb bits

        Instr memu ld Intarg1, #0, Itemp0;            -- Itemp0 <- destword
        Instr lalu Insb Itemp0, Itemp1, Itemp0;       -- Itemp0 <- newdest
        Instr memu st Itemp0, #1, Intarg1;            -- store destword

        Instr lalu lea Intarg0, #1, Intarg0;          -- check bounds
        Instr lalu lle Intarg2, Intarg0, cc0;         -- check bounds
        Instr lalu ct cc0 jmp RETIP;
        Instr lalu cf cc0 br ?I._bcopyloop;
        Instr ;
        Instr ;
        Instr ;

_bcopydone:
        Instr lalu jmp RETIP;
        Instr;
        Instr;
        Instr;

_Sloader: /* void* loader(FILE *file, int node, void *text, void *data) */
        Instr lalu mov Intarg0, Intarg0;   /* validate argument */
        Instr lalu mov Intarg1, Intarg1;   /* validate argument */
        Instr lalu mov Intarg2, Intarg2;   /* validate argument */
        Instr lalu mov Intarg3, Intarg3;   /* validate argument */
        MAKE_MAGIC_PTR(LOADER1_MAGIC)
```

```
    Instr lalu jmp Itemp0;                  /* hook to simulator */
    Instr;
    Instr;

    Instr lalu jmp RETIP;
    Instr;
    Instr;
    Instr;

_Sysloader:  /* ptr* loader(FILE *file, void* startIP ) */
    Instr lalu mov Intarg0, Intarg0;  /* validate argument */
    Instr lalu mov Intarg1, Intarg1;  /* validate argument */
    Instr lalu mov Intarg2, Intarg2;  /* validate argument */
    Instr lalu mov Intarg3, Intarg3;  /* validate argument */
    MAKE_MAGIC_PTR(LOADER2_MAGIC)
    Instr lalu jmp Itemp0;                  /* hook to simulator */
    Instr;
    Instr;
    Instr;

    Instr lalu jmp RETIP;
    Instr;
    Instr;
    Instr;

_SHashCreate:  /* void hashcreate(Int NumPhysicalPages, Int NumMappings) */
    Instr lalu mov Intarg0, Intarg0;  /* validate argument */
    Instr lalu mov Intarg1, Intarg1;  /* validate argument */
    MAKE_MAGIC_PTR(HASHSTART_MAGIC)
    Instr lalu jmp Itemp0;                  /* hook to simulator */
    Instr;
    Instr;
    Instr;

    Instr lalu jmp RETIP;
    Instr;
    Instr;
    Instr;

_SNodeId:  /* Int nodeid() */
    MAKE_MAGIC_PTR(NODEID_MAGIC)
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
    Instr; Instr; Instr;

_SCCInit:  /* void ccinit(void) */
    MAKE_MAGIC_PTR(CCINIT_MAGIC)
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
    Instr; Instr; Instr;

_SCCShare:  /* Int ccshare(Int address, Int node, Int sharing_type) */
    MAKE_MAGIC_PTR(CCSHARE_MAGIC)
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
```

```
    Instr; Instr; Instr;

_SCCShareInfo:  /* Int ccshareinfo(Int address) */
    MAKE_MAGIC_PTR(CCSHAREINFO_MAGIC)
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
    Instr; Instr; Instr;

_SCCPopShare:  /* Int ccpopshare(Int address) */
    MAKE_MAGIC_PTR(CCPOPSHARE_MAGIC)
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
    Instr; Instr; Instr;

_SSQInit:  /* void sqinit(void) */
    MAKE_MAGIC_PTR(SQINIT_MAGIC)
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
    Instr; Instr; Instr;

_SSQEnqueue:  /* void sqenqueue(Int, void*, Int, void*) */
    MAKE_MAGIC_PTR(SQENQUEUE_MAGIC)
    Instr lalu mov Intarg0, Intarg0;
    Instr lalu mov Intarg1, Intarg1;
    Instr lalu mov Intarg2, Intarg2;
    Instr lalu mov Intarg3, Intarg3;
    Instr lalu mov Intarg4, Intarg4;
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
    Instr; Instr; Instr;

_SSQDequeue:  /* void sqdequeue(void) */
    MAKE_MAGIC_PTR(SQDEQUEUE_MAGIC)
    Instr lalu mov Intarg0, Intarg0;
    Instr lalu mov Intarg1, Intarg1;
    Instr lalu mov Intarg2, Intarg2;
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
    Instr; Instr; Instr;

_SSQGetState:  /* void sqgetstate(void *address) */
    MAKE_MAGIC_PTR(SQGETSTATE_MAGIC)
    Instr lalu mov Intarg0, Intarg0;
    Instr lalu mov Intarg1, Intarg1;
    Instr lalu jmp Itemp0;             /* hook to simulator */
    Instr; Instr; Instr;

    Instr lalu jmp RETIP;
    Instr; Instr; Instr;

_SSQSetState:  /* void sqsetstate(void) */
    MAKE_MAGIC_PTR(SQSETSTATE_MAGIC)
    Instr lalu mov Intarg0, Intarg0;
    Instr lalu mov Intarg1, Intarg1;
    Instr lalu mov Intarg2, Intarg2;
```

```
          Instr lalu jmp ltemp0;                    /* hook to simulator */
          Instr; Instr; Instr; Instr;

          Instr lalu jmp RETIP;
          Instr; Instr; Instr;

_SSQget[irstW::  /* void sqget(irstwrite(void) */
          MAKE_MAGIC_PTR(SQGETFIRSTWRITE_MAGIC)
          Instr lalu mov Intarg0, Intarg0;
          Instr lalu mov Intarg1, Intarg1;
          Instr lalu mov Intarg2, Intarg2;
          Instr lalu jmp ltemp0;                    /* hook to simulator */
          Instr; Instr; Instr; Instr;

          Instr lalu jmp RETIP;
          Instr; Instr; Instr;

_SSQDsequeue:  /* void sqdequeueh(void) */
          MAKE_MAGIC_PTR(SQDEQUEUEH_MAGIC)
          Instr lalu mov Intarg0, Intarg0;
          Instr lalu mov Intarg1, Intarg1;
          Instr lalu jmp ltemp0;                    /* hook to simulator */
          Instr; Instr; Instr; Instr;

          Instr lalu jmp RETIP;
          Instr; Instr;

_SSQUpdate:  /* void squpdate(void *cache_buffer, int address) */
          MAKE_MAGIC_PTR(SQUPDATE_MAGIC)
          Instr lalu mov Intarg0, Intarg0;
          Instr lalu mov Intarg1, Intarg1;
          Instr lalu jmp ltemp0;                    /* hook to simulator */
          Instr; Instr; Instr; Instr;

          Instr lalu jmp RETIP;
          Instr; Instr;

_SEventqAdd:  /* void eventqadd(int, int) */
          MAKE_MAGIC_PTR(EVENTQ_ADD_MAGIC)
          Instr lalu mov Intarg0, Intarg0;
          Instr lalu mov Intarg1, Intarg1;
          Instr lalu jmp ltemp0;                    /* hook to simulator */
          Instr; Instr; Instr; Instr;

          Instr lalu jmp RETIP;
          Instr; Instr;

_hlock:
          Instr lalu empty #INTARG0_EMPTY_MASK;
          Instr lalu mov Intarg0, Intarg0;
          Instr ;
          Instr ;
          Instr ;

_malloc::
          GET_FRAME
          LIBCALL(_vmsmx)
          RETURN

_hspawn:   /*     Intarg0 - numargs
                  Intarg1 - thread lp
                  Intarg2 - dest cluster #
           */
          Instr memu mbar;
```

```
GET_FRAME

          Instr memu hfork Intarg1, Intarg2, LCC1;
          Instr lalu cf LCC1 br _hspawn_done;
          Instr lalu cf LCC1 mov #0, IRetVal;
          Instr ;
          Instr ;

LOAD_FAR_LABEL(_hexit, ltemp0, DStart)
          Instr lalu leq Intarg2, #1, LCC1;
          Instr lalu ct LCC1 mov ltemp0, h1.RETIP;
          Instr lalu leq Intarg2, #2, LCC1;
          Instr lalu ct LCC1 mov ltemp0, h2.RETIP;
          Instr lalu leq Intarg2, #3, LCC1;
          Instr lalu ct LCC1 mov ltemp0, h3.RETIP;

          Instr lalu leq Intarg2, #1, LCC1;
          Instr lalu ct LCC1 mov DStart, h1.DStart;
          Instr lalu leq Intarg2, #2, LCC1;
          Instr lalu ct LCC1 mov DStart, h2.DStart;
          Instr lalu leq Intarg2, #3, LCC1;
          Instr lalu ct LCC1 mov DStart, h3.DStart;

          Instr lalu leq Intarg0, #0, LCC1;
          Instr lalu ct LCC1 br _hspawn_0_args;
          Instr ; Instr ; Instr ;

          Instr lalu leq Intarg0, #1, LCC1;
          Instr lalu ct LCC1 br _hspawn_1_args;
          Instr ; Instr ; Instr ;

          Instr lalu leq Intarg0, #2, LCC1;
          Instr lalu ct LCC1 br _hspawn_2_args;
          Instr ; Instr ; Instr ;

          Instr lalu leq Intarg0, #3, LCC1;
          Instr lalu ct LCC1 br _hspawn_3_args;
          Instr ; Instr ; Instr ;

          Instr lalu leq Intarg0, #4, LCC1;
          Instr lalu ct LCC1 br _hspawn_4_args;
          Instr ; Instr ; Instr ;

_hspawn_0_args:
          PUSH(AP)
          PUSH(Intarg0)
          PUSH(Intarg1)
          PUSH(Intarg2)
          Instr lalu imm #0x4000, Intarg0;
          FCALL(_malloc)
          Instr lalu imm #(0x4000 - 0x08), Intarg2;
          Instr lalu lea Intarg0, Intarg2, Intarg0;
          POP(Intarg2)
          POP(Intarg1)

          Instr lalu leq Intarg2, #1, LCC1;
          Instr lalu ct LCC1 mov Intarg0, h1.SP;
          Instr lalu leq Intarg2, #2, LCC1;
          Instr lalu ct LCC1 mov Intarg0, h2.SP;
          Instr lalu leq Intarg2, #3, LCC1;
          Instr lalu ct LCC1 mov Intarg0, h3.SP;

          POP(Intarg0)
          POP(AP)
```

```
Instr lalu br _hspawn_done;
Instr lalu mov #1, IRetVal;
Instr ;
Instr ;
_hspawn_1_args:
PUSH(AP)
PUSH(Intarg0)
PUSH(Intarg1)
PUSH(Intarg2)
Instr lalu lmm #0x4000, Intarg0;
Instr lalu lmm #(0x4000 - 0x08], Intarg2;
Instr lalu lea Intarg0, Intarg2, Intarg0;
FCALL(_malloc)
POP(Intarg2)
POP(Intarg1)

Instr lalu leq Intarg2, #1, LCC1;
Instr lalu ct LCC1 mov Intarg0, h1.SP;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.SP;
Instr lalu leq Intarg2, #3, LCC1;
Instr lalu ct LCC1 mov Intarg0, h3.SP;

POP(Intarg0)
POP(AP)

Instr lalu lea AP, #24, AP;
Instr memu ld AP, Intarg0;
Instr lalu leq Intarg2, #1, LCC1;
Instr lalu ct LCC1 mov Intarg0, h1.Intarg0;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu leq Intarg2, #3, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.Intarg0;
Instr lalu ct LCC1 mov Intarg0, h3.Intarg0;

Instr lalu br _hspawn_done;
Instr lalu mov #1, IRetVal;
Instr ;
Instr ;
_hspawn_2_args:
PUSH(AP)
PUSH(Intarg0)
PUSH(Intarg1)
PUSH(Intarg2)
Instr lalu lmm #0x4000, Intarg0;
Instr lalu lmm #(0x4000 - 0x08], Intarg2;
Instr lalu lea Intarg0, Intarg2, Intarg0;
FCALL(_malloc)
POP(Intarg2)
POP(Intarg1)

Instr lalu leq Intarg2, #1, LCC1;
Instr lalu ct LCC1 mov Intarg0, h1.SP;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.SP;
Instr lalu leq Intarg2, #3, LCC1;
Instr lalu ct LCC1 mov Intarg0, h3.SP;

POP(Intarg0)
POP(AP)

Instr lalu lea AP, #24, AP;
Instr memu ld AP, #8, Intarg0;
Instr lalu leq Intarg2, #1, LCC1;
```

```
Instr lalu ct LCC1 mov Intarg0, h1.Intarg0;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.Intarg0;
Instr lalu leq Intarg2, #3, LCC1;
Instr lalu ct LCC1 mov Intarg0, h3.Intarg0;

Instr memu ld AP, #8, Intarg0;
Instr lalu leq Intarg2, #1, LCC1;
Instr lalu ct LCC1 mov Intarg0, h1.Intarg1;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.Intarg1;
Instr lalu leq Intarg2, #3, LCC1;
Instr lalu ct LCC1 mov Intarg0, h3.Intarg1;

Instr lalu br _hspawn_done;
Instr lalu mov #1, IRetVal;
Instr ;
Instr ;

_hspawn_3_args:
PUSH(AP)
PUSH(Intarg0)
PUSH(Intarg1)
PUSH(Intarg2)
Instr lalu lmm #0x4000, Intarg0;
FCALL(_malloc)
Instr lalu lmm #(0x4000 - 0x08], Intarg2;
Instr lalu lea Intarg0, Intarg2, Intarg0;
POP(Intarg2)
POP(Intarg1)

Instr lalu leq Intarg2, #1, LCC1;
Instr lalu ct LCC1 mov Intarg0, h1.SP;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu leq Intarg2, #3, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.SP;
Instr lalu ct LCC1 mov Intarg0, h3.SP;

POP(Intarg0)
POP(AP)

Instr lalu lea AP, #24, AP;
Instr memu ld AP, #8, Intarg0;
Instr lalu leq Intarg2, #1, LCC1;
Instr lalu ct LCC1 mov Intarg0, h1.Intarg0;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.Intarg0;
Instr lalu ct LCC1 mov Intarg0, h3.Intarg0;

Instr memu ld AP, #8, Intarg0;
Instr lalu leq Intarg2, #1, LCC1;
Instr lalu ct LCC1 mov Intarg0, h1.Intarg1;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.Intarg1;
Instr lalu leq Intarg2, #3, LCC1;
Instr lalu ct LCC1 mov Intarg0, h3.Intarg1;

Instr memu ld AP, #8, Intarg0;
Instr lalu leq Intarg2, #1, LCC1;
Instr lalu ct LCC1 mov Intarg0, h1.Intarg2;
Instr lalu leq Intarg2, #2, LCC1;
Instr lalu ct LCC1 mov Intarg0, h2.Intarg2;
Instr lalu leq Intarg2, #3, LCC1;
Instr lalu ct LCC1 mov Intarg0, h3.Intarg2;
```

```
        Instr lalu hr _hspawn_done;
        Instr lalu mov #1, IRetVal;
        Instr ;
        Instr ;

_hspawn_4_args:
        PUSH(AP)
        PUSH(Intarg0)
        PUSH(Intarg1)
        PUSH(Intarg2)
        PUSH(Intarg3)
        Instr lalu lmm #0x4000, Intarg0;
        FCALL(_malloc)
        Instr lalu lmm #(0x4000 - 0x08), Intarg2;
        Instr lalu lea Intarg0, Intarg2, Intarg0;
        POP(Intarg3)
        POP(Intarg2)
        POP(Intarg1)

        Instr lalu leq Intarg2, #1, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h1.SP;
        Instr lalu leq Intarg2, #2, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h2.SP;
        Instr lalu leq Intarg2, #3, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h3.SP;

        POP(Intarg0)
        POP(AP)

        Instr lalu lea AP, #24, AP;
        Instr memu ld AP, #8, Intarg0;
        Instr lalu leq Intarg2, #1, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h1.Intarg0;
        Instr lalu leq Intarg2, #2, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h2.Intarg0;
        Instr lalu leq Intarg2, #3, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h3.Intarg0;

        Instr memu ld AP, #8, Intarg0;
        Instr lalu leq Intarg2, #1, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h1.Intarg1;
        Instr lalu leq Intarg2, #2, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h2.Intarg1;
        Instr lalu leq Intarg2, #3, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h3.Intarg1;

        Instr memu ld AP, #8, Intarg0;
        Instr lalu leq Intarg2, #1, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h1.Intarg2;
        Instr lalu leq Intarg2, #2, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h2.Intarg2;
        Instr lalu leq Intarg2, #3, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h3.Intarg2;

        Instr memu ld AP, #8, Intarg0;
        Instr lalu leq Intarg2, #1, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h1.Intarg3;
        Instr lalu leq Intarg2, #2, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h2.Intarg3;
        Instr lalu leq Intarg2, #3, LCC1;
        Instr lalu ct LCC1 mov Intarg0, h3.Intarg3;

        Instr lalu hr _hspawn_done;
        Instr lalu mov #1, IRetVal;

        Instr ;
        Instr ;

_hspawn_done:
        RETURN

_hexit::
        Instr memu hexit;
        Instr ;
        Instr ;
        Instr ;

#if 0
_getSelfTC::
        /* eventually, this needs to get threadlock and move to system level */
        Instr lalu lll;
        Instr lalu lmm #0x77c0, Intarg0;
        Instr lalu shoru #0x0000, Intarg0;
        Instr lalu shoru #0x0051, Intarg0;
        Instr lalu shoru #0x0000, Intarg0;
        Instr lalu setptr Intarg0, Intarg0;
        Instr memu ld Intarg0, Intarg0;

        Instr lalu lmm #0x2000, ltemp0;
        Instr lalu lsh ltemp0, #3, ltemp0;
        Instr lalu lea Intarg0, ltemp0, Intarg0;
        Instr memu ld Intarg0, Intarg0;
        Instr lalu lsh Intarg0, #4, Intarg0;
        Instr lalu lsh Intarg0, #-4, Intarg0;
        Instr lalu mov #P_KEY, ltemp0;
        Instr lalu jmp RETIP;
        Instr lalu lsh ltemp0, #60, ltemp0;
        Instr lalu or Intarg0, ltemp0, Intarg0;
        Instr lalu setptr Intarg0, Intarg0;

_getParent::
        Instr lalu lll;
        Instr lalu lsh Intarg0, #4, Intarg0;
        Instr lalu lsh Intarg0, #-4, Intarg0;
        Instr lalu mov #P_RW, ltemp0;
        Instr lalu lsh ltemp0, #60, ltemp0;
        Instr lalu or Intarg0, ltemp0, Intarg0;
        Instr lalu setptr Intarg0, Intarg0;
        Instr lalu lea Intarg0, #16, Intarg0;
        Instr memu ld Intarg0, Intarg0;
        Instr lalu lsh Intarg0, #4, Intarg0;
        Instr lalu lsh Intarg0, #-4, Intarg0;
        Instr lalu mov #P_KEY, ltemp0;
        Instr lalu jmp RETIP;
        Instr lalu lsh ltemp0, #60, ltemp0;
        Instr lalu or Intarg0, ltemp0, Intarg0;
        Instr lalu setptr Intarg0, Intarg0;
#endif

_getDP::
        /* returns own data ptr */
        Instr lalu jmp RETIP;
        Instr lalu mov DStart, Intarg0;
        Instr ;
        Instr ;

#if 0
_vexit::
        /* try to get own CP ptr and exit stuff */
        Instr ;
        Instr lalu empty #((TEMP0_EMPTY_MASK);
        Instr ;
```

```
        Instr [a]); mov ltemp0, ltemp0;
        Instr ;
        Instr ;
        Instr ;
        Instr ;
        Instr ;
        Instr ;
#endif


data USER;
    __sin::        p64  _sin;
    __cos::        p64  _cos;
    __printf::     p64  _printf;
    __scanf::      p64  _scanf;
    __gets::       p64  _gets;
    __strchr::     p64  _strchr;
    __strcpy::     p64  _sstrcpy;
    __fopen::      p64  _fopen;
    __fclose::     p64  _fclose;
    __fprintf::    p64  _fprintf;
    __fwrite::     p64  _fwrite;
    __fread::      p64  _fread;
    __fseek::      p64  _fseek;
    __bcopy::      p64  _bcopy;
    __cleanup::    p64  _cleanup;
    __malloc::     p64  _malloc;
    __hspawn::     p64  _hspawn;
    __hexlt::      p64  _hexlt;
    Loader::       p64  _Sloader;
    SysLoader::    p64  _Sysloader;
    HashCreate::   p64  _SHashCreate;
    NodeId::       p64  _SNodeId;
    CCInit::       p64  _SCCInit;
    CCShare::      p64  _SCCShare;
    CCShareInfo::  p64  _SCCShareInfo;
    CCPopShare::   p64  _SCCPopShare;
    SQInit::       p64  _SSQInit;
    SQEnqueue::    p64  _SSQEnqueue;
    SQDequeue::    p64  _SSQDequeue;
    SQGetState::   p64  _SSQGetState;
    SQSetState::   p64  _SSQSetState;
    SQDequeueH::   p64  _SSQDequeueH;
    SQGetFirstW::  p64  _SSQGetFirstW;
    SQUpdate::     p64  _SSQUpdate;
    EventqAdd::    p64  _SEventqAdd;
/* add entries for more c functions here */
/* remember, data labels are of form
    _c-function-name::
    they are global! */

end;
```

```
#include "newarg.h"
#include <libcall.h>
#include "helpmacros.h"

data USER;
return_string_text:
        asciz "Program returns: %08x%08x\n";
_user_printargs:
        asciz "Arg %d: (%s)\n";

text:
Main::

        Instr lalu mov #0, ltemp0;
        Instr lalu mov #0, ltemp1;
        GET_FRAME

        PUSH(lntarg0)
        PUSH(lntarg1)

        FCALL(INIT_Lib)
        FCALL(INIT_SysLib)

#if 0
        POP(lntarg1)
        POP(lntarg0)

        PUSH(lntarg0)
        PUSH(lntarg1)
        Instr lalu mov lntarg0, ltemp1;
        Instr lalu mov lntarg1, lntarg3;
        Instr lalu mov 1, lntarg1;
_argvprint_loop:
        Instr memu ld lntarg3, #0, lntarg2;           -- target string
        Instr lalu mov lntarg2, lntarg2;
        PRINTF(_user_printargs, DStart)

        Instr lalu leq lntarg1, ltemp1, cc1;
        Instr lalu cf cc1 br _argvprint_loop;
        Instr lalu cf cc1 add lntarg1, #1, lntarg1;
        Instr ;
        Instr ;
#endif

        POP(lntarg1)
        POP(lntarg0)

        PUSH(AP)
        PUSH(lntarg0)
        PUSH(lntarg1)
        Instr lalu mov SP, AP;
        FCALL(_main)
        POP(lntarg1)
        POP(lntarg0)
        POP(AP)

_cleanup::
        Instr lalu mov lRetVal, lntarg2;
        Instr lalu lsh lRetVal, #-32, lntarg1;
        PRINTF(return_string_text, DStart)

        Instr lalu lll;
        GET_RETIP
        Instr lalu jmp RETIP;
        FREE_FRAME
        Instr ;

end;
```

```
/*********************************************************\
 *                                                       *
 *  General Event handler (running in v4-h0 on each node *
 *  Yevgeny Gurevich - 1/29/95 v1.0                      *
 *                   - 3/13/95 v1.5                      *
 *                                                       *
\*********************************************************/

#include "newreg.h"
#include <libcall.h>
#include "helpmacros.h"
#include "opcodes.h"
#include "pointers.h"
#include "ccdefs.h"
#include "signaldefs.h"

#define FAULT_BLOCK_RI        9
#define FAULT_BLOCK_WI        10
#define FAULT_BLOCK_WR        11

#define FAULT_SYNC_MISS_01    12
#define FAULT_SYNC_MISS_10    13
#define FAULT_SYNC_MISS_N0    14
#define FAULT_SYNC_MISS_N1    15

#define CALLFAIL        Instr lalu br Fail; Instr ; Instr ; Instr ;

#define VERBOSE_EH 0

data:
_event_intro:    asciz "General Event Handler Installed\n";
_event_error1:   asciz "Unknown Event Type\n";
_event_inform2:  asciz "\tptr = %p\n\tdata = %lx\n\tcp = %p\n";
_event_signalword::
    u64 0;

text:
EVENT_HANDLER_START:
    Instr lalu mov #0, RETIP;
    Instr lalu mov #0, AP;
    GET_FRAME

    /* zero out temp regs */
    Instr lalu mov #0, itemp0;
    Instr lalu mov #0, evtemp1;
    Instr lalu mov #0, intarg0;
    Instr lalu mov #0, intarg1;
    Instr lalu mov #0, intarg2;
    Instr lalu mov #0, intarg3;
    PRINTF(_event_intro, DStart)

EVENT_WAIT:
    /* if there are any jobs to perform from the software job queue, take
       care of them (first) */
    FCALL(_EH_SoftwareDequeueLoop)

    /* once done with the softwareq queue, can block on a hardware event */
    Instr lalu br HARDWARE_EVENT_WAIT;
    Instr ;
    Instr ;
    Instr ;

SOFTWARE_SIGNAL:
    /* signal arrived via software. */
```

```
    /* unset the signal flag and let's look at the software job queue */
    SYSPUTLOCK(_event_signalword)
    Instr lalu br EVENT_WAIT;
    Instr ; Instr ; Instr ;

HARDWARE_EVENT_WAIT:
    /* wait for an event by blocking on the queue head */
    Instr lalu mov intarg0, intarg0;         -- move out event header
    Instr lalu mov evHead, intarg0;

    Instr   lalu imm #0x7fff, intarg2         -- hardware event types
            memu and intarg0, #0x0f, intarg1;    --    are 4-bits wide
                                                  -- signal to look at
                                                  --    software job queue
    Instr   lalu leq intarg0, intarg2, LCC1;
    Instr   lalu ct LCC1 br SOFTWARE_SIGNAL;
    Instr   lalu leq intarg1, #FAULT_BLOCK_RI, LCC2;
    Instr   lalu leq intarg1, #FAULT_BLOCK_WI, LCC3;
    Instr   lalu ct LCC2 br _BSM_RI;
    Instr   lalu ct LCC3 br _BSM_WI;
    Instr   lalu leq intarg1, #FAULT_BLOCK_WR, LCC1;
    Instr   lalu ct LCC1 br _BSM_WR;
    Instr ; Instr ; Instr ;

    Instr lalu leq intarg1, #FAULT_SYNC_MISS_01, cc1;
    Instr lalu ct cc1 br _SM_01;
    Instr lalu leq intarg1, #FAULT_SYNC_MISS_10, cc2;
    Instr lalu ct cc2 br _SM_10;
    Instr ;
    Instr lalu leq intarg1, #FAULT_SYNC_MISS_N0, cc1;
    Instr lalu ct cc1 br _SM_N0;
    Instr lalu leq intarg1, #FAULT_SYNC_MISS_N1, cc2;
    Instr lalu ct cc2 br _SM_N1;
    Instr ;
    Instr ;
    Instr ;

    Instr lalu mov evBody, intarg1;          -- address ptr
    Instr lalu mov evBody, intarg2;          -- data
    Instr lalu mov evBody, intarg3;          -- cp ptr

    PRINTF(_event_error1, DStart)

    Instr lalu br EVENT_WAIT;
    Instr ;
    Instr ;
    Instr ;

text:

_BSM_RI:
_BSM_WI:
_BSM_WR:

    /* package up arguments and call the C function 'EH_handle_bsm' */
    Instr lalu mov evBody, intarg1;          -- address ptr
    Instr lalu mov evBody, intarg2;          -- data
    Instr lalu mov evBody, intarg3;          -- cp ptr
    PUSH(intarg0)
    Instr lalu mov SP, AP;
    PUSH(intarg1)
    PUSH(intarg2)
    PUSH(intarg3)
    FCALL(_EH_handle_bsm)
    STOP(intarg3)
    STOP(intarg2)
    STOP(intarg1)
```

```
    Instr lalu leq [RetVal, #-1, LCC1;
    Instr lalu ct LCC1 br EVENT_WAIT;
    CTOP(intarg0, ct, LCC1)
    Instr ; Instr ;

    /* If ((RetVal & 0x4], that is a flag telling us that the issuing
       thread must be halted - because the bsm was an icache-miss.  */
    Instr lalu and intarg0, #4, ltemp0;
    Instr lalu leq ltemp0, #4, LCC1;
    Instr lalu cf LCC1 br _continue_BSM_handling;
    CTOP(ltemp0, cf, LCC1)
    Instr ;
    Instr ;

_ichalt_stop_thread:
    PUSH(ltemp0)                -- save header
    PUSH(intarg0)               -- sq returnval
    PUSH(intarg1)
    Instr lalu lsh ltemp0, #-48, intarg0;     -- address
    Instr lalu and intarg0, #0xf, intarg0;    -- extract threadslot
                                              -- extract threadslot

    Instr lalu lsh ltemp0, #-40, intarg1;     -- extract cluster
    Instr lalu and intarg1, #0x3, intarg1;    -- extract cluster

    /* we must change the thread running state to NOT issue
       and then back to issue when data comes back.  This means
       that we must lock some data structure before trying to
       get the CP made for us
    SYSGETLOCK(ichalt_vac_loop, _threadLock)                  */

    /* now we cons up the configspace ptr to the thread */
    Instr lalu lsh intarg0, #17, intarg0;
    Instr lalu imm #(0x7440), ltemp0;
    Instr lalu shoru #0x0000, ltemp0;
    Instr lalu shoru #0x0000, ltemp0;
    Instr lalu shoru #0x0000, ltemp0;
    Instr lalu or intarg0, ltemp0, intarg0;
    Instr lalu setptr intarg0, intarg0;

    /* now offset to global thread state */
    Instr lalu imm #1, ltemp0;
    Instr lalu lsh ltemp0, #16, ltemp0;
    Instr lalu add ltemp0, #8, ltemp0;
    Instr lalu lea intarg0, ltemp0, intarg0;    -- current running/full bits
    Instr memu ld intarg0, ltemp0;

    PUSH(intarg0)
    Instr lalu imm #1, intarg0;
    Instr lalu lsh intarg0, intarg1, intarg0;
    Instr lalu lsh intarg0, #4, intarg0;                    --
    Instr lalu not intarg0, intarg0;
    Instr lalu and intarg0, ltemp0, ltemp0;
    Instr lalu leq intarg0, ltemp0, LCC1;
    SPOP(intarg0)

    /* kill the hissue bit for the faulting h.hread */
    Instr memu st ltemp0, intarg0;
    SYSPUTLOCK(_threadLock)
    STOP(intarg1)
    STOP(intarg0)

_continue_BSM_handling:
    /* ((RetVal & 0x2) == 2) => send CC request on our own */
    Instr lalu and intarg0, #2, ltemp0;
    Instr lalu leq ltemp0, #2, LCC1;
```

```
    STOP(intarg0)

    SYSPUTLOCK(_sqlock)          -- Important to unlock the data
                                 -- structure before sending a
                                 -- message, in case you block
                                 -- and the PIMH needs to get sqlock

    Instr lalu ct LCC1 br _SendCCMessage;
    Instr lalu ct LCC1 imm #EVENT_WAIT, RETIP;
    Instr lalu ct LCC1 leab IP, RETIP, RETIP;
    Instr lalu cf LCC1 br EVENT_WAIT;
    Instr ;
    Instr ;
    Instr ;

data;
    _sync_inform1:     asciz 'Sync 01 handler\n';
    _sync_inform2:     asciz 'Sync 10 handler\n';
    _sync_inform3:     asciz 'Sync N0 handler\n';
    _sync_inform4:     asciz 'Sync N1 handler\n';
text;

_SM_01:
    CONSTRUCT_LONG_PTR(_sync_inform1, intarg0, DStart)
    PUSH(AP)
    PUSH(intarg0)
    Instr lalu lea SP, #8, AP;
    LIBCALL(_printf)
    POP(intarg0)
    POP(AP)

    Instr lalu mov evBody, intarg1;      -- address ptr
    Instr lalu mov evBody, intarg2;      -- data
    Instr lalu mov evBody, intarg3;      -- cp ptr

    CONSTRUCT_LONG_PTR(_event_inform2, intarg0, DStart)
    PUSH(AP)
    PUSH(intarg0)
    Instr lalu lea SP, #8, AP;
    LIBCALL(_printf)
    POP(intarg0)
    POP(AP)

    Instr lalu br EVENT_WAIT;
    Instr ;
    Instr ;
    Instr ;

_SM_10:
    CONSTRUCT_LONG_PTR(_sync_inform2, intarg0, DStart)
    PUSH(AP)
    PUSH(intarg0)
    Instr lalu lea SP, #8, AP;
    LIBCALL(_printf)
    POP(intarg0)
    POP(AP)

    Instr lalu mov evBody, intarg1;      -- address ptr
    Instr lalu mov evBody, intarg2;      -- data
    Instr lalu mov evBody, intarg3;      -- cp ptr

    CONSTRUCT_LONG_PTR(_event_inform2, intarg0, DStart)
    PUSH(AP)
    PUSH(intarg0)
    Instr lalu lea SP, #8, AP;
    LIBCALL(_printf)
```

```
        POP(Intarg0)
        POP(AP)

_SM_I10:
        Instr lalu br EVENT_WAIT;
        Instr ;
        Instr ;
        Instr ;

        CONSTRUCT_LONG_PTR(_sync_Inform3, Intarg0, DStart)
        PUSH(AP)
        PUSH(Intarg0)
        Instr lalu lea SP, #8, AP;
        LIBCALL(__print)
        POP(Intarg0)
        POP(AP)

        Instr lalu mov evBody, Intarg1;       -- address ptr
        Instr lalu mov evBody, Intarg2;       -- data
        Instr lalu mov evBody, Intarg3;       -- cp ptr

        CONSTRUCT_LONG_PTR(_event_Inform2, Intarg0, DStart)
        PUSH(AP)
        PUSH(Intarg0)
        Instr lalu lea SP, #8, AP;
        LIBCALL(__print)
        POP(Intarg0)
        POP(AP)

        Instr lalu br EVENT_WAIT;
        Instr ;
        Instr ;
        Instr ;

_SM_I11:
        CONSTRUCT_LONG_PTR(_sync_Inform4, Intarg0, DStart)
        PUSH(AP)
        PUSH(Intarg0)
        Instr lalu lea SP, #8, AP;
        LIBCALL(__print)
        POP(Intarg0)
        POP(AP)

        Instr lalu mov evBody, Intarg1;       -- address ptr
        Instr lalu mov evBody, Intarg2;       -- data
        Instr lalu mov evBody, Intarg3;       -- cp ptr

        CONSTRUCT_LONG_PTR(_event_Inform2, Intarg0, DStart)
        PUSH(AP)
        PUSH(Intarg0)
        Instr lalu lea SP, #8, AP;
        LIBCALL(__print)
        POP(Intarg0)
        POP(AP)

        Instr lalu br EVENT_WAIT;
        Instr ;
        Instr ;
        Instr ;

_SendCCMessage::
        /*      Intarg0 - header
                Intarg1 - target address
                Intarg2 - data
                Intarg3 - fault cp
```

```
*/
        Instr    lalu empty #0x4010;
        Instr    lalu mov Intarg0, FMC0;          -- event header
        Instr    menu gprb Intarg1, ltemp0, LCC1  -- node ld
        Instr    lalu mov Intarg1, FMC1;          -- virtual address
        Instr    lalu mov Intarg2, FMC2;          -- op data
        Instr    lalu mov Intarg1, FMCDest;       -- virtual address
        Instr    lalu mov Intarg3, FMC3;          -- fault CP
        /* eventually have to check whether nodenum from gprb matches
           own node.  For now, ignore this... */

        Instr    lalu mov ltemp0, Intarg1;

        /* Intarg1 contains node number! */
        /* try to create a dispatch IP and send a message */
        /* Generate XM ptr */
        MAKE_XM_PTR(MSC_ccrequest)

        /* fp send */
        Instr    lalu [srnd0 #4, FMCDest, FMCIP, cc1;
        Instr    lalu ct cc1 jmp RETIP;
        Instr    ; Instr ; Instr ;

        CALLFAIL

_mbarStoreUpdate::
        /* Intarg0 - configspace ptr */
        /* update the hardware threadstate by decrementing the mbar counter */
        Instr    lalu lsh Intarg0, #-8, Intarg0;
        Instr    lalu jmp RETIP;
        Instr    lalu lsh Intarg0, #8, Intarg0;
        Instr    lalu setptr Intarg0, Intarg0;
        Instr    menu st Intarg0, Intarg0;

_mbarCCLoadUpdate::
        /* Intarg0 - configspace ptr    */
        /* Intarg1 - header word         */
        /* Intarg2 - data                */
        /* Intarg3 - condition           */
        Instr    lalu leq Intarg3, #1, LCC1;
        Instr    menu cf LCC1 jmp RETIP;
        Instr    menu cf LCC1 st Intarg3, Intarg0 -- If load failed
                 lalu ct LCC1 extb Intarg1, #4, Intarg1; -- set condition
        Instr    lalu ct LCC1 lmm #0x6ff, ltemp0;        -- dest register
        Instr    lalu ct LCC1 not ltemp0, ltemp0;        -- cleared mbar-bit

        Instr    lalu ct LCC1 and Intarg0, ltemp0, ltemp0; -- mask

        /* now calculate register offset */
        Instr    lalu and Intarg1, #0xf, Intarg1;
        Instr    lalu lsh Intarg1, #3, Intarg1;

        Instr    lalu or ltemp0, Intarg1, ltemp0;
        Instr    lalu jmp RETIP;
        Instr    lalu setptr ltemp0, ltemp0;
        Instr    menu st Intarg2, ltemp0;
        Instr    menu st Intarg3, Intarg0;

_mbarLoadUpdate::
        /* Intarg0 - header word -- tells us which register and stuff  */
        /* Intarg1 - configspace ptr into thread                       */
        /* Intarg3 - actual data to be written in thread's target register */
        Instr    menu mbat;
        Instr    menu st Intarg3, Intarg1;      -- store data into cp ptr
```

```
/* we now need to activate the thread which was deactivated
   as a result of the original (l1 icache) miss... */
Instr lalu lsh lntarg1, #-20, ltemp0;
Instr lalu and ltemp0, #0xf, ltemp0;
Instr lalu lsq ltemp0, #0x5, LCC1;
Instr lalu cf LCC1 jmp RETIP;
Instr ; Instr ; Instr ;
SYSRETBLOCK(_phar_vac_loop, _threadlock)   -- get threadlock because
                                           --  updating thead state

Instr lalu lsh lntarg0, #-48, ltemp0;      -- threadslot
Instr lalu and ltemp0, #0xf, ltemp0;       -- threadslot

Instr lalu lsh lntarg0, #-40, lntarg1;     -- cluster
Instr lalu and lntarg1, #0x1, lntarg1;     -- cluster

/* now we cons up the configspace ptr to the thread */
Instr lalu lsh ltemp0, #17, lntarg0;
Instr lalu lmm #(0x74401), ltemp0;
Instr lalu shoru #0x0000, ltemp0;
Instr lalu shoru #0x0000, ltemp0;
Instr lalu shoru #0x0000, ltemp0;
Instr lalu or lntarg0, ltemp0, lntarg0;
Instr lalu setptr lntarg0, lntarg0;

/* now offset to global thread state */
Instr lalu lmm #1, ltemp0;
Instr lalu lsh ltemp0, #16, ltemp0;
Instr lalu add ltemp0, #8, ltemp0;
Instr lalu lea lntarg0, ltemp0, lntarg0;
Instr memu ld lntarg0, ltemp0;            -- current running/full bits

PUSH(lntarg0)
Instr lalu lmm #1, lntarg0;
Instr lalu lsh lntarg0, lntarg1, lntarg0;
Instr lalu lsh lntarg0, #4, lntarg0;
Instr lalu or lntarg0, ltemp0, ltemp0;
STOP(lntarg0)

/* make hisuue true for the hthread so that it can continue */
Instr memu st ltemp0, lntarg0;
SYSPUTLOCK(_threadlock)
Instr lalu jmp RETIP;
Instr ; Instr ; Instr ;

.end;
```

```
/*lock procedures for use with buddy.c*/

#include "newreg.h"
#include <libcall.h>

data:

_buddy_data_lock:
        u64 0x666;

text:

_buddysetlock::
        /*sets up the lock into the proper locked state*/
        /*no matter what it was before*/

        instr memu mbar;

        instr memu lds ua, 1, ltemp0, 10;
        instr lalu jmp RETIP;
        instr ; instr ; instr ;


_buddylock::
        /*spins until the lock is free, then locks and returns*/

        instr memu mbar;

        CONSTRUCT_LONG_PTR(_buddy_data_lock, ltemp0, DStart)
_buddyLock1:
        instr memu ldscnd cf, 1, ltemp0, 10, cc3;
        instr lalu cf cc3 br _buddyLock1;
        instr ;
        instr ;
        instr ;

        instr lalu jmp RETIP;
        instr ; instr ; instr ;


_buddyUnlock::
        /*unlocks the lock and returns*/

        instr memu mbar;

        CONSTRUCT_LONG_PTR(_buddy_data_lock, ltemp0, DStart)
        instr memu ldscnd cf, 0, ltemp0, 10, cc3;
        instr lalu jmp RETIP;
        instr ; instr ; instr ;

end;
```

```
#include "newreg.h"
#include "helpmacros.h"
#include <libcall.h>

/*--------------------------------------------------------------------
LTLB fault handler

----------------------------------------------------------------------*/

.text;

LTLB_HANDLER_START::
        Instr   lalu mov #0, RETIP;
        Instr   lalu mov #0, Itemp0;
        Instr   lalu mov #0, evtemp1;
        Instr   lalu mov #0, Intarg0;
        Instr   lalu mov #0, Intarg1;
        Instr   lalu mov #0, Intarg2;

        Instr   lalu empty ##0x4000;                    -- Empty 113, 114, 115

        /* Wait for an LTLB miss, then call ltlb_body.c code */
wait_ltlb_miss:
        Instr   lalu lsh VADDR, #10, evtemp1;           -- stalls until ltlb miss
_ltlb_miss_found:
        Instr   lalu mov 112, Intarg0;
        Instr   lalu mov 113, Intarg1;
        Instr   lalu mov 114, Intarg2;
        Instr   lalu mov 115, Intarg3;
        PUSH(Intarg0)
        Instr   lalu lea SP, #8, AP;
        PUSH(Intarg1)
        PUSH(Intarg2)
        PUSH(Intarg3)
        FCALL(_PFM_handle_miss)
        SPOP(Intarg3)
        SPOP(Intarg3)
        SPOP(Intarg3)
        SPOP(Intarg3)
_ready_to_continue:
        /* once we return, we must store the returnvalue of the miss-handling
           (1 = EMI should retry the faulting operation, 0 = don't)
           and unlock the LTLB */
        Instr   lalu lnm   ##0x7000, Itemp0;
        Instr   lalu shoru ##0x0000, Itemp0;
        Instr   lalu shoru ##0x002f, Itemp0;
        Instr   lalu shoru ##0xfff8, Itemp0;

        Instr   lalu br wait_ltlb_miss;
        Instr   lalu setptr Itemp0, Itemp0;
        Instr   lalu empty ##0x400;                     -- Empty 113, 114, 115
        Instr   memu st RetVal, Itemp0;

        Instr   lalu br fail;
        Instr   ;
        Instr   ;
        Instr   ;

end;
```

```
/* ***************************************************
 *
 * M-Machine runtime system Message Handler code.
 * Handles p0 and p1 messages arriving at node.
 * Now mostly stubs which call functions written in C
 *
 * Written by        Yevgeny Gurevich
 * Version           0.05
 * Modification Date 06/11/95
 *
 * Version           0.10
 * Modification Date 07/18/95
 *
 * Version           0.20
 * Modification Date 08/08/95
 *
 * -----------------------------------------------
 * System handling doesn't deal well with large argument lists.
 *
 * ***************************************************/

#include "newreg.h"
#include <libcall.h>
#include "helpmacros.h"
#include "opcodes.h"
#include "pointers.h"
#include "ccdefs.h"
#include "signaldefs.h"

#define CALLFAIL          Instr Ialu br Fail; Instr ; Instr ; Instr ;

data:
_ccdirlock::
        u64 0;              -- lock word on cache-coherence directory
                            -- synchronizes access to the ccdir data structure

_sqlock::
        u64 0;
_ltlb_data_for_MH::
        u64 0;
_MH_Intro:
        asciz 'Message Handler Installed\n';

text:
/* ***************************************************
 *
 * Initialize all necessary registers and data structures and
 * start waiting for a message to arrive
 *
 * ***************************************************/

DISPATCH_HANDLER_START::
        Instr Ialu mov #0, RETIP;
        Instr Ialu mov #0, ltemp0;
        Instr Ialu mov #0, mstemp1;
        Instr Ialu mov #0, Intarg0;
        Instr Ialu mov #0, Intarg1;
        Instr Ialu mov #0, Intarg2;
        Instr Ialu mov #0, AP;

        PRINTF(_MH_Intro, DStart)

/* ***************************************************
 *
 * Wait for incoming message
 *
 * ***************************************************/

MSG_NEXTMESSAGE:
        Instr Ialu mov MsgHead, ltemp0;
        Instr empty #0x0016;            -- XM IP for message
                                        -- empty f1,f2,f4

_message_arrived:
        Instr Ialu jmp ltemp0;          -- jump to dispatch IP
        Instr Ialu mov MsgBody, f1;     -- arg count (ignored now)
        Instr Ialu mov MsgBody, f2;     -- return address (sending node)
        Instr Ialu mov MsgBody, f4;     -- referenced addr (VADDR_REG)

        CALLFAIL

/* ***************************************************
 *                                                   \
 *                                                   .
 * Cache-Coherence Code                              .
 *                                                   .
 *                                                   /
 * ***************************************************/

text;

/* ***************************************************
 *
 * NACK message arriving (this tells us to resend the request by enqueing a
 * job with the event handler
 *
 * ***************************************************/

MSG_ccNackRO::
        Instr Ialu empty #(INTARG4_EMPTY_MASK);
        Instr       Ialu mov f2, Intarg4
                    Ialu mov MsgBody, Intarg0;
        Instr Ialu mov MsgBody, Intarg1;
        Instr Ialu mov MsgBody, Intarg2;
        Instr Ialu mov MsgBody, Intarg3;
        PUSH(Intarg0)
        Instr Ialu mov SP, AP;
        PUSH(Intarg1)
        PUSH(Intarg2)
        PUSH(Intarg3)
        FCALL(_ccNackRO)
        SPOP(Intarg0)
        SPOP(Intarg0)
        SPOP(Intarg0)
        SPOP(Intarg0)
        Instr Ialu br MSG_NEXTMESSAGE;
        Instr ; Instr ; Instr ;

        CALLFAIL

MSG_ccNackRW::
        Instr Ialu empty #(INTARG4_EMPTY_MASK);
        Instr       Ialu mov f2, Intarg4
                    Ialu mov MsgBody, Intarg0;
        Instr Ialu mov MsgBody, Intarg1;
        Instr Ialu mov MsgBody, Intarg2;
        Instr Ialu mov MsgBody, Intarg3;
        PUSH(Intarg0)
        Instr Ialu mov SP, AP;
        PUSH(Intarg1)
        PUSH(Intarg2)
        PUSH(Intarg3)
        FCALL(_ccNackRW)
        SPOP(Intarg0)
        SPOP(Intarg0)
        SPOP(Intarg0)
        SPOP(Intarg0)
```

```
        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ; Instr ;
        CALLFAIL
;.....................................................
;.
;.  INVALIDATE message arriving at P0
;.
;.....................................................

MSG_ccInvalidate::
        Instr lalu empty #(INTARG2_EMPTY_MASK);
        Instr     lalu mov f2, Intarg2
                  lalu mov MsgBody, Intarg0;
        Instr lalu mov MsgBody, Intarg1;
        PUSH(Intarg0)
        Instr lalu mov SP, AP;
        PUSH(Intarg1)
        FCALL(_ccInvalidate)
        SPOP(Intarg0)
        SPOP(Intarg1)
        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ; Instr ;
        CALLFAIL
;.........................................................
;.
;.  Dirty line being pushed (evicted) back to home node - execution on home node
;.
;.........................................................

MSG_ccreturnDirty::
        Instr lalu mov MsgBody, Intarg0;
        PUSH(Intarg0)
        Instr lalu mov SP, AP;
        FCALL(_ccreturnEvicted)
        SPOP(Intarg0)
        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ; Instr ;
        CALLFAIL
;.........................................................
;.
;.  YANK ACK coming back to home node (this is a "long" ack - with dirty line)
;.
;.........................................................

MSG_ccreturnyankFull::
        Instr lalu mov MsgBody, Intarg0;
        PUSH(Intarg0)
        Instr lalu mov SP, AP;
        FCALL(_ccreturnyankFull)
        SPOP(Intarg0)
        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ;
        CALLFAIL
;.........................................................
;.
;.  YANK ACK coming back to home node (this is a "short" ack - no dirty line)
;.
;.........................................................

MSG_ccreturnyank::
        Instr lalu mov MsgBody, Intarg0;
```

```
        PUSH(Intarg0)
        Instr lalu mov SP, AP;
        FCALL(_ccreturnyank)
        SPOP(Intarg0)
        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ; Instr ;
        CALLFAIL
;.........................................................
;.
;.  REQUEST for a cache line coming in to a p0 MH
;.
;.........................................................

MSG_ccrequest::
        Instr lalu empty #(INTARG4_EMPTY_MASK);
        Instr     lalu mov f2, Intarg4
                  lalu mov MsgBody, Intarg1;
        Instr lalu mov MsgBody, Intarg0;
        Instr lalu mov MsgBody, Intarg2;
        Instr lalu mov MsgBody, Intarg3;
        PUSH(Intarg0)
        Instr lalu mov SP, AP;
        PUSH(Intarg1)
        PUSH(Intarg2)
        PUSH(Intarg3)
        FCALL(_ccrequest)
        SPOP(Intarg0)
        SPOP(Intarg0)
        SPOP(Intarg0)
        SPOP(Intarg0)
        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ; Instr ;
        CALLFAIL
;.........................................................
;.
;.  ACK(X)       -- received acknowledgement with the entire cache line
;.
;.........................................................

MSG_ccreturnStore::
        Instr lalu empty #(INTARG1_EMPTY_MASK);
        Instr lalu mov f2, Intarg2;
        Instr lalu mov MsgBody, Intarg0;
        Instr lalu mov MsgBody, Intarg1;    /* header */
        PUSH(Intarg0)
        Instr lalu mov SP, AP;
        FCALL(_ccreturnStore)
        SPOP(Intarg0)
        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ; Instr ;
        CALLFAIL
;.........................................................
;.
;.  ACK(R)       -- received acknowledgement with the entire cache line
;.
;.........................................................

MSG_ccreturnLoad::
        Instr lalu empty #(INTARG1_EMPTY_MASK);
        Instr lalu mov f2, Intarg2;
```

```
        Instr lalu mov MsgBody, Intarg0;          /* address */
        Instr lalu mov MsgBody, Intarg1;          /* header */
        PUSH(Intarg0)
        Instr lalu mov SP, AF;
        FCALL(_sstructunload)
        SPOP(Intarg0)

        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ;
        CALLFAIL

/*****************************************************
 *
 *  Update the outgoing message buffer counter
 *
 *****************************************************/

_update_ombc::
        Instr    lalu lmm #8, ltemp0;              /* Increment OMBC constant */
        Instr    lalu lmm #((P_CONFIG << 12) | (0xf << 8)), lntarg1;
        Instr    lalu shoru #0x0000, lntarg1;
        Instr    lalu shoru #0x0030, lntarg1;
        Instr    lalu shoru #0x0000, lntarg1;
        Instr    lalu jmp RETIP;
        Instr    lalu or lntarg1, ltemp0, lntarg1;
        Instr    lalu setptr lntarg1, lntarg1;
        Instr    memu st ltemp0, lntarg1;

/*****************************************************
 *
 *  Thread Spawn message arriving on target node
 *
 *****************************************************/

MSG_InvokeRPC::
MSG_tspawn::
        /* making a call to fork */
        /* fork(ip, data, return, numargs, arg1, arg2, etc) */
        Instr lalu mov MsgBody, lntarg3;              -- numargs;

        Instr lalu mov MsgBody, lntarg1;              -- data ptr
        Instr lalu mov MsgBody, lntarg0;              -- threadip
        CONSTRUCT_LCMG_LABEL(_tExitEx, lntarg2)
        Instr lalu leab IF, lntarg2, lntarg2;

        /* set up stack */
        /* save control block for retrieval later */
        Instr lalu mov SP, lntarg4;

        Instr lalu lsh lntarg3, #3, ltemp0;
        Instr lalu sub 10, ltemp0, ltemp0;
        Instr lalu lea SP, ltemp0, SP;
        Instr lalu lea SP, #-48, SP;

        Instr lalu mov SP, AF;
        Instr memu st lntarg0, #8,  SP;
        Instr memu st lntarg1, #8,  SP;
        Instr memu st lntarg2, #8,  SP;
        Instr memu st lntarg3, #8,  SP;
        /* this is where we are to store the parentTC */
        Instr lalu lea SP, #8, SP;

        Instr lalu mov lntarg3, ltemp0;
        Instr lalu leq ltemp0, #0, LCC1;
        Instr lalu cf LCC1 br _done_tspawn_loop;
```

```
        Instr lalu cf LCC1 sub ltemp0, #1, ltemp0;
        Instr ;
        Instr ;
_tspawn_loop1:
        Instr lalu leq ltemp0, #0, LCC1;
        Instr lalu cf LCC1 br _tspawn_loop1;
        Instr memu st MsgBody, #8, SP;
        Instr lalu sub ltemp0, #1, ltemp0;
        Instr ;
        Instr ;
_done_tspawn_loop:
        Instr lalu lea AF, #32, SP;
        Instr memu st MsgBody, #-32, SP;

        PUSH(lntarg4)
        FCALL(_tForkX)
        SPOP(lntarg1)
        Instr lalu mov lntarg1, SP;

        /* now the last message word is the target of a signal */
        Instr lalu mov lntarg0, lntarg1;       -- the dataword for
                                               -- the signal is the
                                               -- actual childTC
                                               -- the word to signal

        Instr lalu mov MsgBody, lntarg0;
        FCALL(_tSignalX)

        Instr lalu br MSG_NEXTMESSAGE;
        Instr; Instr; Instr;

/*****************************************************
 *
 *  Message for a thread to be wakened arriving at the
 *  thread's home node
 *
 *****************************************************/

MSG_tWake::
        /* a wake message arriving for a node */
        Instr lalu mov MsgBody, lntarg0;
        Instr lalu mov MsgBody, lntarg1;
        PUSH(lntarg0)
        Instr lalu mov SP, AF;
        PUSH(lntarg1)
        FCALL(_SYStWake)
        SPOP(lntarg0)
        SPOP(lntarg0)

        Instr lalu br MSG_NEXTMESSAGE;
        Instr ; Instr ;

data:
_rval_loc:
        u64 0;

text:

/*****************************************************
 *
 *  Message for a sleep entry to be made arriving at the
 *  signal_word's home node
 *
 *****************************************************/

MSG_tSleep::
        /* a TC asks to be put to sleep */
        Instr lalu mov MsgBody, lntarg0;       -- signal_word
```

```
Instr lalu mov MsgBody, Intarg1;          -- tc
Instr lalu mov MsgBody, Intarg2;          -- data mask
PUSH(Intarg1)

PUSH(Intarg0)
Instr lalu mov SP, AP;
PUSH(Intarg1)
PUSH(Intarg2)
CONSTRUCT_LONG_PTR(_rval_loc, Intarg3, DStart)
PUSH(Intarg3)
FCALL(_SYStSleepRemote)
SPOP(Intarg3)
SPOP(Intarg2)
SPOP(Intarg2)
SPOP(Intarg2)
SPOP(Intarg1)


CONSTRUCT_LONG_PTR(_rval_loc, Intarg3, DStart)

Instr lalu leq Intarg0, 10, LCC1;
Instr lalu ct LCC1 br MSG_NEXTMESSAGE;
Instr falu cf LCC1 empty #0x0010;
Instr memo cf LCC1 ld Intarg3, FMC1;          -- result of tsleepremote
Instr lalu cf LCC1 mov Intarg1, FMC0;          -- tc involved

MAKE_XM_PTR(MSG_tSleepContinuation)
Instr falu fsndlpt #2, {2, FMCIP, LCC1;

Instr lalu ct LCC1 br MSG_NEXTMESSAGE;
Instr ; Instr ; Instr ;
CALLFAIL
```

```
/* **********************************************
 *
 * Response to a MSG_tSleep message wherein
 * a thread is told that it may waken immediately
 *
 * **********************************************/
```

```
MSG_tSleepContinuation::
    /* If we got a continuation message, this means that we need to
       rewaken a thread. But since we are a P1 MII, we cannot
       wait for any locks. Need to add a job to the eh queue */
    Instr lalu mov MsgBody, Intarg0;          -- tc
    Instr lalu mov MsgBody, Intarg1;          -- data
    PUSH(Intarg0)
    Instr lalu mov SP, AP;
    PUSH(Intarg1)
    FCALL(_SYStWake)
    SPOP(Intarg1)
    SPOP(Intarg1)

    Instr lalu br MSG_NEXTMESSAGE;
    Instr ; Instr ;
```

```
/* **********************************************
 *
 * A signal message arriving at signal_word's
 * home node. Asking for a new signal entry to be made
 *
 * **********************************************/
```

```
MSG_tSignal::
    /* a signal for this word arrives - must handle it */
    Instr lalu mov MsgBody, Intarg0;
```

```
Instr lalu mov MsgBody, Intarg1;

PUSH(Intarg0)
Instr lalu mov SP, AP;
PUSH(Intarg1)
FCALL(_SYStSignal)
SPOP(Intarg1)
SPOP(Intarg1)

Instr lalu br MSG_NEXTMESSAGE;
Instr ; Instr ; Instr ;
```

end;

```
/*****************************************
 *
 * passfail.m
 *
 * User pass, fail exit routines
 * Author: Steve Keckler
 * Modified to global label by Marco Fillo April 1st 1994
 *****************************************/

#include    "/home/num/apps/stdhead/regdef.h"
#define PASS_CODE 0xaa
#define FAIL_CODE 0xff
#define SYSFAIL_CODE 0xdd

.text;

/*------------------------------------------

User Pass routine: halts simulator when called.

--------------------------------------------*/

pass::
Pass::

    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu mov #PASS_CODE, IntZero;
    instr lalu br Halt;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;

/*------------------------------------------

User Fail routine: halts simulator when called.

--------------------------------------------*/

fail::
Fail::

    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu mov #FAIL_CODE, IntZero;
    instr lalu br Halt;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;

/*------------------------------------------

Sysfail routine:  spin loop

--------------------------------------------*/

sysfail::
Sysfail::

    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
    instr lalu nop;
```

```
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu mov *SYSFAIL_CODE, Int?ero;
Instr lalu br Halt;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;

/*------------------------------------------------------

User Halt routine: spin loop

-----------------------------------------------------*/

Halt::
halt::

Instr lalu br Halt;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;
Instr lalu nop;

end;
```

```
/*pointer creation procedures for use with buddy.c*/

#include "newreg.h"
#include <libcall.h>

#define PFLISTTEST 0

#if PFLISTTEST

data:

/*_fake_phys:
        u64 0x333 (30);*/

_fake_phys:
        u64 0x333;

text:

_initfake::
        /*expects a pointer to the fake*/
        /*physical memory in Intarg0*/

        Instr memu mbar;
        GET_FRAME

        CONSTRUCT_LONG_PTR(_fake_phys, Intarg1, DStart)
        Instr memu st Intarg0, #0, Intarg1;

        RETURN

_createpointer::
        /*expects: address in Intarg0 (int)*/
        /*         size in Intarg1          */
        /*         protection in Intarg2    */
        /*fake version for use in testing*/
        /*rpflist code*/

        Instr memu mbar;

        GET_FRAME

        CONSTRUCT_LONG_PTR(_fake_phys, Intarg1, DStart)
        Instr memu ld Intarg1, #0, Intarg2;
        Instr lalu lea Intarg2, Intarg0, Intarg0;
        Instr lalu setptr Intarg0, IRetVal;

        RETURN

#else
text:

_createpointer::
        /*expects: address in Intarg0 (int)*/
        /*         size in Intarg1          */
        /*         protection in Intarg2    */

        Instr memu mbar;

        Instr lalu lsh Intarg2, #6, Intarg2;
        Instr lalu or Intarg1, Intarg2, Intarg1;
        Instr lalu lsh Intarg1, #54, Intarg1;
        Instr lalu or Intarg1, Intarg0, Intarg0;
        Instr lalu setptr Intarg0, IRetVal;
```

```
        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

#endif

_setpointer::
        /*expects: pointer in Intarg0 (int)*/
        /*without the pointer bit set*/
        /*sets it and returns in IRetVal*/

        Instr memu mbar;

        Instr lalu setptr Intarg0, IRetVal;

        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

_resetpointer::
        /*expects a pointer in Intarg0*/
        /*resets it to its starting address*/
        /*and returns in IRetVal*/

        Instr lalu leab Intarg0, #0, IRetVal;

        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

_breakpoint::

        Instr memu mbar;
        Instr lalu ill;

        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

#if 0
_alloc::
        /*expects a size in Intarg0*/
        /*just a wrapper for _vmemx*/

        GET_FRAME

        PUSH(AP)

        /*condom the SP*/
        Instr lalu lea SP, -8, SP;

        Instr lalu lea SP, #8, AP;
        LIBCALL(_vmemx)

        /*undo*/
        Instr lalu lea SP, #8, SP;

        SPOP(AP)

        GET_RETIP

        Instr lalu jmp RETIP;
        Instr ;
        Instr ;
        Instr ;
#endif
end;
```

```
#include "newreg.h"
#include "pointers.h"
#include "helpmacros.h"

#include <libcall.h>

#define SETUP_POINTER(x)
CONSTRUCT_LONG_LABEL(x,IntArg0)
Instr lalu lea DStart, IntArg0, IntArg0;    /* data address */
Instr memu ld IntArg0, IntArg1;             /* data value */
Instr lalu leab IP, IntArg1, IntArg1;       /* generate pointer */
Instr memu st IntArg1, IntArg0;

data USER
align 0 mod 8;
_tSleep_ptr::      ldptr 0x0000[[[aaaaaaaa;
_tExit_ptr::       ldptr 0x0000[[[aaaaaaab;
_vmmax::           ldptr 0x0000[[[aaaaaaad;
_ipcx::            ldptr 0x0000[[[aaaaaaae;
_tUnparse_ptr::    ldptr 0x0000[[[aaaaaaaf;
_tFork_ptr::       ldptr 0x0000[[[aaaaaab0;
_tSignal_ptr::     ldptr 0x0000[[[aaaaaab1;
_getParent_ptr::   ldptr 0x0000[[[aaaaaab2;
_getSelfTC_ptr::   ldptr 0x0000[[[aaaaaab3;

/*******************************************/

text;
INIT_SysLib::
        /* Initialize the library of pointers. */
        Instr lalu jmp RETIP;
        Instr ; Instr ; Instr ;

/*******************************************/

text;

_getParent::
        GET_FRAME
        LOAD_FAR_LABEL(_getParent_ptr, ltemp0, DStart)
        Instr lalu jmp ltemp0;
        Instr ;
        Instr ;
        CALC_RETIP

        RETURN

_getSelfTC::
        GET_FRAME
        LOAD_FAR_LABEL(_getSelfTC_ptr, ltemp0, DStart)
        Instr lalu jmp ltemp0;
        Instr ;
        Instr ;
        CALC_RETIP

        RETURN

_tSleep::
        GET_FRAME
        LOAD_FAR_LABEL(_tSleep_ptr, ltemp0, DStart)
        Instr lalu jmp ltemp0;
        Instr ;
        Instr ;
        CALC_RETIP

        RETURN

_tSignal::
        GET_FRAME
        LOAD_FAR_LABEL(_tSignal_ptr, ltemp0, DStart)
        Instr lalu jmp ltemp0;
        Instr ;
        Instr ;
        CALC_RETIP

        RETURN

_tExit::
        GET_FRAME
        LOAD_FAR_LABEL(_tExit_ptr, ltemp0, DStart)
        Instr lalu jmp ltemp0;
        Instr ;
        Instr ;
        CALC_RETIP

        RETURN

_tUnparse::
        GET_FRAME
        LOAD_FAR_LABEL(_tUnparse_ptr, ltemp0, DStart)
        Instr lalu jmp ltemp0;
        Instr ;
        Instr ;
        CALC_RETIP

        RETURN

_tFork::
        GET_FRAME
        LOAD_FAR_LABEL(_tFork_ptr, ltemp0, DStart)
        Instr lalu jmp ltemp0;
        Instr ;
        Instr ;
        CALC_RETIP

        RETURN

_sysPtrRemote::
        Instr lalu lsh IntArg0, #4, IntArg0;
        Instr lalu lsh IntArg0, #-4, IntArg0;
        Instr lalu mov #P_KEY, ltemp0;
        Instr lalu jmp RETIP;
        Instr lalu lsh ltemp0, #60, ltemp0;
        Instr lalu or ltemp0, IntArg0, IntArg0;
        Instr lalu setptr IntArg0, IntArg0;

_tspawn::
/*      IntArg0 - numargs
        IntArg1 - thread ip
        IntArg2 - dest node */

        Instr memu mbar;

        GET_FRAME

        Instr lalu mov #1, FMC3;
        Instr lalu mov #1, FMC4;
```

```
Instr   lalu mov #1, FMC5;
Instr   lalu mov #1, FMC6;
Instr   lalu mov #1, FMC7;
Instr   lalu mov #1, FMC8;
Instr   lalu mov #1, FMC9;

Instr   lalu empty #0xc070;

CONSTRUCT_LONG_PTR(_rpcx, ltemp0, DStart)
Instr   memu ld ltemp0, FMCIP;

/* perform an RPC */
Instr   lalu mov lntarg0, FMC0;           /* num args   */
Instr   lalu mov DStart, FMC1;            /* data ptr   */
Instr   lalu mov lntarg1, FMC2;           /* thread lp  */
Instr   lalu mov lntarg2, FMCDest;        /* dest node # */

/* try to allocate memory for a signalword */
PUSH(AP)
PUSH(lntarg0)
Instr   lalu mov #0x8, lntarg0;
FCALL(_malloc)

/* when malloc returns, denote the ptr into a KEY-only ptr */
FCALL(_sysPtrDemote)

/* now move it for later use */
Instr   lalu mov lntarg0, lntarg1;        -- sleep ptr

SPOP(lntarg0)
SPOP(AP)

PUSH(lntarg1)

/* load up arguments */
Instr   lalu lea AP, #24, ltemp0;
Instr   lalu mov #1, ltemp1;
Instr   lalu lle ltemp1, lntarg0, LCC1;
Instr   memu ct LCC1 ld ltemp0, #8, FMC3;
Instr   lalu mov #2, ltemp1;
Instr   lalu lle ltemp1, lntarg0, LCC1;
Instr   memu ct LCC1 ld ltemp0, #8, FMC4;
Instr   lalu mov #3, ltemp1;
Instr   lalu lle ltemp1, lntarg0, LCC1;
Instr   memu ct LCC1 ld ltemp0, #8, FMC5;
Instr   lalu mov #4, ltemp1;
Instr   lalu lle ltemp1, lntarg0, LCC1;
Instr   memu ct LCC1 ld ltemp0, #8, FMC6;
Instr   lalu mov #5, ltemp1;
Instr   lalu lle ltemp1, lntarg0, LCC1;
Instr   memu ct LCC1 ld ltemp0, #8, FMC7;

/* problem is that tspawn should really be a system mode function */
SYSRTLOCK(_spawn_point1, _threadLock)

/* ltemp0 is free at this point, so are lntarg1, lntarg2, etc */
Instr   lalu lmm #0x76c0, ltemp0;
Instr   lalu shoru 0x0000, ltemp0;
Instr   lalu shoru 0x0051, ltemp0;
Instr   lalu shoru 0x0000, ltemp0;
Instr   lalu setptr ltemp0, ltemp0;

Instr   memu ld ltemp0, ltemp0;            /* get own CP */
Instr   lalu lmm #0x2000, lntarg1;
Instr   lalu lsh lntarg1, #3, lntarg1;
```

```
Instr   lalu lea ltemp0, lntarg1, ltemp0;     /* get TC */
Instr   memu ld ltemp0, ltemp0;
Instr   lalu lsh ltemp0, #4, ltemp0;
Instr   lalu lsh ltemp0, #-4, ltemp0;
Instr   lalu lmm #P_KEY, lntarg1;
Instr   lalu lsh lntarg1, #60, lntarg1;
Instr   lalu or lntarg1, ltemp0, ltemp0;
Instr   lalu setptr ltemp0, lntarg1;          /* lntarg1 is TC */

SYSPUTLOCK(_threadLock)

---

SPOP(lntarg2)                                  -- sleepword
PUSH(lntarg2)

Instr   lalu lle lntarg0, #0, LCC1;
Instr   lalu ct LCC1 empty #0x0100;
Instr   lalu cf LCC1 mov lntarg1, FMC3;
Instr   lalu ct LCC1 mov lntarg2, FMC4;
Instr   lalu ct LCC1 {snd0pt #5, FMCDest, FMCIP, LCC2
Instr   lalu ct LCC1 lne 10, 10, LCC2;
Instr   lalu ct LCC1 br _tspawn_continue;
Instr   lalu ccand LCC2, LCC1, LCC3;
Instr   lalu ct LCC1 mov #0, lRetVal;
Instr   lalu ct LCC3 mov #1, lRetVal;

Instr   lalu lle lntarg0, #1, LCC1;
Instr   lalu ct LCC1 empty #0x0300;
Instr   lalu cf LCC1 mov lntarg1, FMC4;
Instr   lalu ct LCC1 mov lntarg2, FMC5;
Instr   lalu ct LCC1 {snd0pt #6, FMCDest, FMCIP, LCC2
Instr   lalu cf LCC1 lne 10, 10, LCC2;
Instr   lalu ct LCC1 br _tspawn_continue;
Instr   lalu ct LCC1 ccand LCC2, LCC1, LCC3;
Instr   lalu ct LCC1 mov #0, lRetVal;
Instr   lalu ct LCC3 mov #1, lRetVal;

Instr   lalu lle lntarg0, #2, LCC1;
Instr   lalu ct LCC1 empty #0x0600;
Instr   lalu cf LCC1 mov lntarg1, FMC5;
Instr   lalu ct LCC1 mov lntarg2, FMC6;
Instr   lalu ct LCC1 {snd0pt #7, FMCDest, FMCIP, LCC2
Instr   lalu cf LCC1 lne 10, 10, LCC2;
Instr   lalu ct LCC1 br _tspawn_continue;
Instr   lalu ct LCC1 ccand LCC2, LCC1, LCC3;
Instr   lalu ct LCC1 mov #0, lRetVal;
Instr   lalu ct LCC3 mov #1, lRetVal;

Instr   lalu lle lntarg0, #3, LCC1;
Instr   lalu ct LCC1 empty #0x0c00;
Instr   lalu ct LCC1 mov lntarg1, FMC6;
Instr   lalu ct LCC1 mov lntarg2, FMC7;
Instr   lalu ct LCC1 {snd0pt #8, FMCDest, FMCIP, LCC2
Instr   lalu cf LCC1 lne 10, 10, LCC2;
Instr   lalu ct LCC1 br _tspawn_continue;
Instr   lalu ct LCC1 ccand LCC2, LCC1, LCC3;
Instr   lalu ct LCC1 mov #0, lRetVal;
Instr   lalu ct LCC3 mov #1, lRetVal;

Instr   lalu lle lntarg0, #4, LCC1;
Instr   lalu ct LCC1 empty #0x1800;
Instr   lalu ct LCC1 mov lntarg1, FMC7;
Instr   lalu ct LCC1 mov lntarg2, FMC8;
Instr   lalu ct LCC1 {snd0pt #9, FMCDest, FMCIP, LCC2
Instr   lalu cf LCC1 lne 10, 10, LCC2;
Instr   lalu ct LCC1 br _tspawn_continue;
```

```
    Instr    lalu ct LCC1 ccand LCC2, LCC1, LCC1, LCC3,
    Instr    lalu ct LCC1 mov r0, IRetVal;
    Instr    lalu ct LCC3 mov r1, IRetVal;

    Instr    lalu mov r0, IRetVal;
    SPOP(Intarg1)
    RETURN

_tSpawn_continue::
    /* now that we have spawned off a thread, call tSleep to
       wait for the result to come back... */
    SPOP(Intarg0)                              -- signal word
    Instr lalu mov r0, Intarg1;                -- mask
    FCALL( tSleep)
    RETURN


    end;
```

```
#include "mmreg.h"
#include "helpmacros.h"
#include <libcall.h>
#include "pointers.h"

/* system loader */
data:
_LD_hello_string:
        asciz 'M-Machine Loader v0.20\n';
_LD_prompt:
        asciz 'Enter filename > ';
_LD_input_string:
        asciz '%s';
_LD_commandline_string:
        asciz '';

        :
_LD_filename:
        asciz                               '';
_LD_load_file:
        ptr r.0.0;
_LD_argc:
        u64 0;
_LD_argv:
        p64 _LD_commandline_string;
        p64 _LD_commandline_string;
        p64 _LD_commandline_string;
        p64 _LD_commandline_string;
        p64 _LD_commandline_string;
        p64 _LD_commandline_string;
_LD_printargs:
        asciz 'Arg %d: (%s)\n';

_LD_filemode:
        asciz 'r';
_LD_load_file:
        ptr r.0.0;
_LD_NULL_PTR:
        ptr k.0.0;

_Sys_User_IP:       ptr k.0.0;
_Sys_User_CP:       ptr k.0.0;
_Sys_User_SP:       ptr k.0.0;

_LD_ip_string:  asciz 'IP = 0x%x\n';
_LD_cp_string:  asciz 'CP = 0x%x\n';

text:
_system_loader::
        GET_FRAME

        PRINTF(_LD_hello_string, DStart)
        PRINTF(_LD_prompt, DStart)

        CONSTRUCT_LONG_PTR(_LD_commandline_string, Intarg0, DStart)
        LIBCALL(__gets)

        /* Intarg1 contains the entire 'command-line'. It must be
           parsed, pointers into the different arguments created, and
           an argc calculated */

        /* don't make copies, just store pointers to different parts of
           the string; Intarg0 is the ptr to the entire input string */
        LOAD_FAR_LABEL(_LD_argc, Intarg1, DStart)       -- Intarg1 contains argc
        CONSTRUCT_LONG_PTR(_LD_argv, Intarg1, DStart)   -- Intarg1 contains
                                                           -- ptr to argv list

        Instr memu st Intarg0, #0, Intarg1;
        Instr lalu add Intarg1, #1, Itemp1;
```

```
        /* now loop THROUGH THE INPUT STRING, reading out characters and
           storing new string starting points into argv! */
_argc_loop:
        Instr memu ld Intarg0, Intarg2;                     -- Intarg2 <- word
        Instr lalu and Intarg0, #0x7, Itemp0;               -- 3 lsbits
        Instr lalu exb Intarg2, Itemp0, Intarg1;            -- Intarg3 contains
                                                            -- string[index]

        Instr    lalu leq Intarg3, #32, cc0;
        Instr    lalu leq Intarg3, #0, cc1;
        Instr    lalu ct cc1 br _done_argcloop;
        Instr    lalu ct cc0 lsh Itemp0, #32, Itemp0;
        Instr    lalu ct cc0 or Itemp0, #0, Itemp0;
        Instr    lalu ct cc0 lnsb Intarg2, Itemp0, Intarg2;
        Instr    lalu cf cc0 br _argc_loop;
        Instr    memu ct cc0 st Intarg2, Intarg0;

        Instr    memu cf cc1 lea Intarg0, #1, Intarg0;      -- store BACK into
                                                            -- source string
        Instr    memu ct cc0 st Intarg0, #8, Intarg1        -- next argv string
                 lalu ct cc0 add Itemp1, #1, Itemp1;

_done_argcloop:
        /* have completed the copying */

        /* nice thing to do now is to store the argc back in, and
           also store the first argv into _LD_filename */
        CONSTRUCT_LONG_PTR(_LD_argc, Intarg0, DStart)
        Instr memu st Itemp1, Intarg0;
        CONSTRUCT_LONG_PTR(_LD_filename, Intarg0, DStart)
        LOAD_FAR_LABEL(_LD_argv, Intarg1, DStart)
        LIBCALL(__strcpy)

        CONSTRUCT_LONG_PTR(_LD_argv, Intarg3, DStart)
        Instr lalu mov 1, Intarg1;

_argvprint_loop:
        /* name of file is in filename buffer */
        /* open the file */
        CONSTRUCT_LONG_PTR(_LD_filename, Intarg0, DStart)  -- name address
        CONSTRUCT_LONG_PTR(_LD_filemode, Intarg1, DStart)  -- mode address
        LIBCALL(__fopen)                                   -- open file

        /* IRetVal contains the file pointer */
        /* check file pointer to make sure its valid */
        CONSTRUCT_LONG_PTR(_LD_load_file, Itemp1, DStart)  -- store handle
        Instr memu st IRetVal, Itemp1;

        LOAD_FAR_LABEL(_LD_NULL_PTR, Itemp0, DStart)
        LOAD_FAR_LABEL(_LD_load_file, Itemp1, DStart)
        Instr lalu leq Itemp0, Itemp1, cc1;

        CONSTRUCT_LONG_LABEL(_LD_bad_opening, Itemp1)
        Instr lalu leab IP, Itemp1, Itemp1;
        Instr lalu ct cc1 jmp Itemp1;
        Instr ; Instr ; Instr ;

        /* first, decide in linkage -
           generate IP
           generate CP
           generate Stack */

        Instr lalu lmm ##0x4000, Intarg0;          -- allocate 16K for user code
        Instr lalu lsh Intarg0, #3, Intarg0;       -- allocate 128K for user code
        FCALL(vmem_alloc)

        /* set the ip type to user execute ptr */
```

```
Instr lalu lsh [RetVal, #4, [RetVal;
Instr lalu lsh [RetVal, #-4, [RetVal;          -- kill RW protection
Instr lalu mov r_EXUSER, [templ;
Instr lalu lsh [templ, #60, [templ;
Instr lalu or [templ, [RetVal, [RetVal;        -- set EXUSER protection
Instr lalu setptr [RetVal, [RetVal;

CONSTRUCT_LONG_PTR(_Sys_User_IP, [templ, DStart)
Instr memu st [RetVal, [templ;                  -- store in sys data segment

/* now need to decide where to place the text and data segments */
/* [templ0 is an argument to the loader! */
Instr lalu lmm ##0X4000, [ntarg0;              -- 16K bytes for user data
Instr lalu lsh [ntarg0, #3, [ntarg0;           -- 128K bytes for user data
FCALL(vmem_alloc)

Instr lalu mov [RetVal, [ntarg1;

CONSTRUCT_LONG_PTR(_Sys_User_CP, [templ, DStart)
Instr memu st [ntarg1, [templ;                  -- store in sys data segment

/* arguments to _LD_load_file are the FILE*, and the IP and CP */
LOAD_FAR_LABEL(_LD_load_file, [ntarg0, DStart)
LOAD_FAR_LABEL(_Sys_User_IP, [ntarg1, DStart)
LOAD_FAR_LABEL(_Sys_User_CP, [ntarg2, DStart)
CONSTRUCT_LONG_PTR(_SysCallStart, [ntarg3, DStart)
LIBCALL(SysLoader)

/* return value is either NULL ptr or entry IP */
LOAD_FAR_LABEL(_LD_NULL_PTR, [temp0, DStart)
Instr lalu leq [temp0, [RetVal, ccl;
Instr lalu ct ccl br _LD_bad_objectfile;
Instr ; Instr ; Instr ;

/* open complete. */
PUSH([ntarg2)                                   -- save CP
PUSH([RetVal)                                   -- save IP

LOAD_FAR_LABEL(_LD_load_file, [ntarg0, DStart)
LIBCALL(_fclose)                                -- close user executable

POP([temp0)                                     -- restore IP
PUSH([temp0)                                    -- save IP
CONSTRUCT_LONG_PTR(_LD_ip_string, [ntarg0, DStart)
Instr lalu mov [temp0, [ntarg1;

/* set up a user stack */
Instr lalu lmm #0x4000, [ntarg0;               -- request 16384-byte stack
FCALL(vmem_alloc)

Instr lalu lmm ##(0x4000 - 0x08), [temp0;      -- set to last word
Instr lalu lea [RetVal, [temp0, [RetVal;       -- in the segment

/* [RetVal now contains a ptr to beginning of user stack */
POP([temp0)                                     -- restore USER IP
POP([templ)                                     -- restore USER CP

/* save system values on system stack */
PUSH([RETIP)
PUSH(AP)
PUSH(DStart)

Instr lalu mov [RetVal, AP;                     -- initially, AP points to
                                                -- first word on user stack
```

```
/* construct key-only Stack pointer */
/* kill protections and make it key-only */
Instr lalu lsh      SP,      #4,       SP;
Instr lalu lsh      SP,      #-4,      SP;
Instr lalu lmm      ##0x8000, [ntarg0;
Instr lalu shoru    #0x0000, [ntarg0;
Instr lalu shoru    #0x0000, [ntarg0;
Instr lalu or       SP,      [ntarg0, SP;

/* need to construct [ntarg0 (argc) and [ntarg1 (argv) */
LOAD_FAR_LABEL(_LD_argc, [ntarg0, DStart)
CONSTRUCT_LONG_PTR(_LD_argv, [ntarg1, DStart)

Instr lalu setptr SP,    [temp2;               -- key ptr to go on USER stack
Instr lalu mov AP,       SP;                    -- set user stack ptr
Instr memu st [temp2,    SP;                    -- push system SP on user SP

CONSTRUCT_LONG_LABEL(_RETURNCODE, RETIP)
Instr lalu leab IP, RETIP, RETIP;               -- set user RETIP back
                                                -- to system code

Instr lalu lea [templ, #0x10, DStart;           -- Get rid of user TOC and
                                                -- set DStart to USER data

/* now just need to jump to user code */
Instr lalu jmp [temp0;
Instr ;
Instr ;
Instr ;

/* if execution reaches this point, something is really wrong!*/
Instr lalu br sysfall;
Instr ;
Instr ;
Instr ;

_RETURNCODE:
data:
_LD_return_string:
    asciz "Return from user code complete.\nSystem Stack Pointer= %p\n";

text:

Instr menu mbar;
Instr lalu lll;
/* SP is still user stack ptr */
POP([ntarg0)
Instr lalu mov [ntarg0, SP;

/* now this is keyonly, make physical again */
Instr lalu lsh      SP,      #4,       SP;
Instr lalu lsh      SP,      #-4,      SP;
Instr lalu lmm      ##0x9000, [ntarg0;
Instr lalu shoru    #0x0000, [ntarg0;
Instr lalu shoru    #0x0000, [ntarg0;
Instr lalu or       SP,      [ntarg0, SP;
Instr lalu setptr SP,    SP;

/* now SP is the actual saved stack ptr. */
POP(DStart)
POP(AP)                                          -- restore system AP
POP(RETIP)

/* now can print something nice */
```

```
        Instr lalu mov SP, intarg1;
        PRINTF(_LD_return_string, DStart)
        Instr lalu mov #1, lRetVal;
        RETURN

_LD_bad_opening:        /* error opening the data file */
data;
        LD_bad_open_stringa: asclz 'Error opening ';
        _LD_bad_open_stringb: asclz '.\n';

text;

        PRINTF(_LD_bad_open_stringa, DStart)
        PRINTF(_LD_filename, DStart)
        PRINTF(_LD_bad_open_stringb, DStart)
        Instr lalu mov #0, lRetVal;
        RETURN

_LD_bad_objectfile:     /* error opening the data file */
data;
        _LD_bad_open_stringc: asclz 'Error reading ';
        _LD_bad_open_stringd: asclz '.\n';

text;

        PRINTF(_LD_bad_open_stringc, DStart)
        PRINTF(_LD_filename, DStart)
        PRINTF(_LD_bad_open_stringd, DStart)
        Instr lalu mov #0, lRetVal;
        RETURN

end;
```

```
/*....................................................
 *....................................................
 *....................................................
 *         systoff.m                                  .
 *         System Table of Offsets                    .
 *....................................................
 * Author: Marco Fillo                                .
 * Date: 5/4/94                                       .
 * Last Mod:    5/6/94                                .
 *....................................................
 *....................................................
 *..................................................*/

text;
/* table of system pointers */
align 0 mod 8;
TOCptr:: u64 0;  /* to be patched by init.m code */
GLOBAL_PAGE_OFFSET:: u64 GLOBAL_PAGE_TABLE;
instr lalu nop;

data;
         /*****************************************/
         /*      SYSTEM Table of Offsets        */
         /*****************************************/

/* because move allows immediate field max 11 bits, data set may be
   too large to mov label directly, use a table of offsets         */
SYSTOFF:          u64 0x7fffffffffffffff;
ESYSTOFF:         u64 0x7fffffffffffffff;

         /****************************/
         /*      END of TOFF         */
         /****************************/

end;
```

```
/* ....................................................
 *
 * vmem.m
 * Yevgeny Gurevich
 * August 26, 1994
 *
 * Modified: 08/08/95
 *
 * This file contains the following virtual memory management functions:
 * vmem_alloc:   allocates a range of bytes and returns a new segment ptr
 * vmem_dealloc: deallocates a virtual memory segment
 *
 * These functions are merely stubs for the buddy list virtual memory
 * allocator
 *
 * ....................................................
 */

#include "nuwreg.h"
#include <libcall.h>

#define VERBOSE_ALLOC 0

data:
_mem_alloc_string:
    asciz "mem_alloc:\tAllocating %d bytes of memory...\n";
_mem_alloc_error_string:
    asciz "Error allocating memory - stsu failed.\n";
_mem_alloc_segment_string:
    asciz "Requested Segment Size: %d\n";
_m_m_alloc_m3_string:
    asciz "Bitmask to check 0x%x%x is 0x%x%x\n";
_mem_notok_string:
    asciz "Segment not ok\n";

#define VERBOSE_ALLOC 0

text:
/* this routine is used for internal physical memory allocation before the
   PPM and VSM are up and running. Use sparingly to just help jumpstart
   the handlers and such */
mem_alloc::

GET_FRAME
/* Uses DStart pointer to access _Sys_Memory_End and update it.
   Needs to synchronize on _Sys_Memory_End because more than one
   system/user process may be modifying it at a time
   Args :
     Intarg0 = length in bytes of requested segment
   Returns :
     returns Physical (change to R/W later) pointer to allocated space
   Modifies :
     DStart[_Sys_Memory_End]      */

    PUSH(DStart)

    Instr lalu imm __SYSTEM_HEAT_PTR__, ltemp0;
    Instr lalu lssb IP, ltemp0, DStart;
    Instr memu ld DStart, DStart;

    Instr lalu mov Intarg0, Intarg1;

    /* Intarg1 now contains number of bytes to allocate. */
    /* (first read in pointer from _Sys_Memory_End.
       Will need to sync on it. */


CONSTRUCT_LONG_PTR(_Sys_Memory_End, Intarg0, DStart)

_loop_load:
    Instr memu ldsu ct, 0, Intarg0, Intarg2, ccl;
    Instr lalu cf ccl br _loop_load;
    Instr ;
    Instr ;
    Instr ;

    /* need to calculate a size for the pointer, and also align it to
       the correct boundary */
    /* arg2 contains pointer to all memory. Need to find closest aligned
       chunk of size Intarg1 */

    /* how do you find requested segment size ? Assume is zero, and
       continue shifting out the requested size until is zero. Use a
       counter */
    Instr lalu mov Intarg1, Intarg3;    -- number of bytes to allocate
    Instr lalu sub Intarg3, #1, Intarg3;    -- segment size
    Instr lalu mov #0, ltemp0;              -- segment size

shift_loop:
    Instr lalu leq Intarg3, l0, ccl;
    Instr lalu cf ccl br shift_loop;
    Instr lalu cf ccl lsh Intarg3, #-1, Intarg3;
    Instr lalu cf ccl add ltemp0, #1, ltemp0;
    Instr ;

    /* MUST check arg3 to see if alignment is right. To this end, must
       ensure that for an m-word-segment request, the m+3 low-order
       address bits of arg3 are zero. */
    Instr lalu mov #-1, mstemp1;
    Instr lalu lsh mstemp1, ltemp0, mstemp1;
    Instr lalu lsh mstemp1, ltemp0, mstemp1;
    Instr lalu not mstemp1, mstemp1;

#if VERBOSE_ALLOC
    PUSH(Intarg0)
    PUSH(Intarg1)
    PUSH(Intarg2)
    PUSH(Intarg3)
    PUSH(Intarg4)
    PUSH(ltemp0)
    PUSH(mstemp1)
    Instr lalu mov Intarg2, Intarg1;
    Instr lalu lsh Intarg1, #-32, Intarg1;
    Instr lalu mov Intarg2, Intarg2;
    Instr lalu mov mstemp1, Intarg3;
    Instr lalu lsh Intarg3, #-32, Intarg3;
    Instr lalu mov mstemp1, Intarg4;
    CONSTRUCT_LONG_PTR(_mem_alloc_m3_string, Intarg0, DStart)
    PUSH(AP)
    PUSH(Intarg0)
    Instr lalu lea SP, #8, AP;
    LIBCALL(__printf)
    SPOP(Intarg0)
    SPOP(AP)

    SPOP(mstemp1)
    SPOP(ltemp0)
    SPOP(Intarg4)
    SPOP(Intarg3)
    SPOP(Intarg2)
    SPOP(Intarg1)
    SPOP(Intarg0)
#endif
```

```
[PUSH(mstempl)]
        Instr lalu and mstempl, Intarg2, mstempl;
        Instr lalu leq mstempl, 10, ccl;
        Instr lalu ct ccl br _segment_ok;
        Instr ;
        Instr ;
        Instr ;

/* segment not aligned                                    */
/* zero out Intarg2 with not of bitmask and then lea twice! */
        Instr lalu not mstempl, mstempl;
        Instr lalu and Intarg2, mstempl, Intarg2;
        Instr lalu setptr Intarg2, Intarg2;

/* new beginning of segment is bitmask */
        SPOP(mstempl)
[PUSH(mstempl)]
        Instr lalu add mstempl, #1, mstempl;
        Instr lalu lea Intarg2, mstempl, Intarg2;

_segment_ok:

        SPOP(mstempl)

/* create correct seglength field and protections */
[PUSH(Intarg1)]
        Instr lalu lea Intarg2, Intarg1, Intarg1;

/* Intarg1 now contains pointer to end of allocated memory */
        Instr memu stsu cf, 1, Intarg1, Intarg0, ccl;

SPOP(Intarg1)

        Instr lalu lsh Intarg2, #10, Intarg2;
        Instr lalu lsh Intarg2, #-10, Intarg2;        -- destroy length and
                                                      -- prot fields.

        Instr lalu lmm #0x9, mstempl;                 -- create Phys bitmap
        Instr lalu lsh mstempl, #60, mstempl;         -- create Phys bitmap

        Instr lalu lsh ltemp0, #54, ltemp0;           -- create correct seg length
        Instr lalu or ltemp0, mstempl, mstempl;
        Instr lalu or Intarg2, mstempl, Intarg2;      -- new ptr bitmask
        Instr lalu setptr Intarg2, Intarg2;

        Instr lalu ct ccl br _return_mem_alloc;
        Instr; Instr;

CONSTRUCT_LONG_PTR(_mem_alloc_error_string, Intarg0, DStart)
        PUSH(AP)
        PUSH(Intarg0)
        Instr lalu lea SP, #8, AP;
        LIBCALL(_printf)
        SPOP(Intarg0)
        SPOP(AP)

        Instr lalu mov #-1, Intarg2;

_return_mem_alloc:
        /* need to return pointer to alloc'd memory.  for now, no
           protections and just return Intarg2 */
        /* need to change its size to not overrun???*/
        Instr lalu mov Intarg2, lRetVal;
```

```
        SPOP(DStart)
        RETURN

data:
_vmem_alloc_string:
        asciz "vmem_alloc:\tAllocating %d bytes of memory...\n";
_vmem_alloc_error_string:
        asciz "vmem_alloc:\t\terror allocating memory - stsu failed.\n";
_vmem_alloc_segment_string:
        asciz "vmem_alloc:\trequested Segment Size: %d\n";
_vmem_alloc_m1_string:
        asciz "vmem_alloc:\t\tbitmask to check 0x%x%x is 0x%x%x\n";
_vmem_notok_string:
        asciz "vmem_alloc:\tsegment not ok\n";

text:

/*******************************************************
 *
 * The segment deallocation routine
 * actually removes backing store for the pages within the
 * segment.  Then frees the actual segment itself.
 *
 *******************************************************/

vmem_dealloc::
        /* deallocate a segment of virtual memory */
        Instr memu mbar;

        GET_FRAME

        /*
            Args:       Intarg0 = segment ptr

            Returns :   None

            Modifies :  Internal BuddyList
                        LPT
        */

        PUSH(DStart)

        Instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
        Instr lalu leab IP, ltemp0, DStart;
        Instr memu ld DStart, DStart;

        Instr lalu lsh Intarg0, #10, Intarg2;
        Instr lalu lsh Intarg2, #-22, Intarg2;

        Instr lalu lsh Intarg0, #4, Intarg3;          -- [first vpn #
        Instr lalu lsh Intarg3, #-58, Intarg3;
        Instr lalu mov #1, ltemp0;                    -- seg length
        Instr lalu lsh ltemp0, Intarg3, Intarg3;      -- bytes in segment
        Instr lalu lmm #*4096, ltemp0;
        Instr lalu lle Intarg3, ltemp0, ccl;
        Instr lalu ct ccl mov #1, Intarg1;            -- # pages to dealloc
        Instr lalu ct ccl lsh Intarg3, #-12, Intarg1; -- # pages to dealloc

_dealloc_loop1:
        /* this removes the VPN-PPN mapping in the TLB */
        /* eventually, it will also result in the mapping
           being removed from the local page table. */
        Instr lalu lmm #0x1700, ltemp0;
        Instr lalu shoru #*0x0000, ltemp0;
        Instr lalu shoru #*0x8000, ltemp0;
        Instr lalu shoru #*0x0000, ltemp0;
        Instr lalu setptr ltemp0, ltemp0;
```

```
        Instr memu st  Intarg2, ltemp0;
        Instr lalu sub Intarg1, #1, Intarg1;
        Instr lalu leq Intarg1, l0, cc1;
        Instr lalu cf cc1 br _dealloc_loop1;
        Instr lalu cf cc1 add Intarg2, #1, Intarg2;
        Instr ;
        Instr ;

        FCALL(_buddyFree)

        SPOP(DStart)

        RETURN

/***************************************************
*
* The segment allocation routine
* does not provide backing store.
*
****************************************************/

vmem_alloc::
        /* allocates a segment of virtual memory using buddylists */
        Instr memu mbar;

        GET_FRAME

        /*  Args:   Intarg0 = length in bytes of requested segment
         *
         *  Returns : Virtual N/W pointer to allocated space
         *  Modifies :
         *          Internal BuddyList
         *          LPT
         */

        PUSH(DStart)

        Instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
        Instr lalu leab IP, ltemp0, DStart;
        Instr memu ld DStart, DStart;

        FCALL( buddyAlloc)

#if INITIALIZE_SEGMENT
        /* the following initializes the data in the segment to
                deadbeef values */
        Instr lalu lsh lRetVal, #4, Intarg2;
        Instr lalu lsh Intarg2, #-58, Intarg2;   -- seg length
        Instr lalu sub Intarg2, #3, Intarg2;     -- bytes/word
        Instr lalu mov #1, ltemp0;
        Instr lalu lsh ltemp0, Intarg2, Intarg1; -- number of words
        Instr lalu sub Intarg1, #1, Intarg1;
        Instr lalu mov Intarg0, ltemp0;
_buddy_loop1:
        Instr lle Intarg1, l0, cc1;
        Instr lalu cf cc1 lmm Intarg1, ltemp1;
        Instr lalu cf cc1 br _buddy_loop1;
        Instr lalu cf cc1 shoru #0xdead, ltemp1;   -- delay
        Instr lalu cf cc1 shoru #0xbeef, ltemp1;   -- delay
        Instr memu cf cc1 st ltemp1, #8, ltemp0     -- delay
        Instr lalu cf cc1 sub Intarg1, #1, Intarg1;
#endif

_return_vmem_alloc:
        SPOP(DStart)
```

```
        GET_RETIP
        Instr lalu jmp RETIP;
        FREE_FRAME
        Instr ;

/***************************************************
*
* Debugging routine to print out the buddylist internal
* data structures.
*
****************************************************/

vmem_buddyPP::
        GET_FRAME

        /*    Args:   NONE */

        PUSH(DStart)

        Instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
        Instr lalu leab IP, ltemp0, DStart;
        Instr memu ld DStart, DStart;

        /* compiler linkage */

        FCALL(_buddyPP)

        SPOP(DStart)

        GET_RETIP
        Instr lalu jmp RETIP;
        FREE_FRAME
        Instr ;

/***************************************************
*
* AS FAR AS I CAN SEE, THIS IS DEFUNCT - Yev
*
* The malloc call (sysmalloc)
*
* Not clear whether we can substitute calls to vmem_alloc all of the time
*
****************************************************/

data:
_sysmalloc_string: asciz 'malloc_base: 0x%08x%08x\n';
text:

align 0 mod 8;
_malloc_base::
        ptr k.0.0;
align 0 mod 8;
sysmalloc::
        /* as it is written, this code is not reentrant! */

        Instr memu mbar;

        GET_FRAME
        /* system malloc library */
        /* Intarg0 contains number of bytes to allocate */

        PUSH(DStart)

        Instr lalu lmm __SYSTEM_UDAT_PTR__, ltemp0;
```

```
Instr lalu leah It, itemp0, DStart;
Instr memu ld DStart, DStart;

CONSTRUCT_LONG_LABEL(_malloc_base, itemp0)
Instr lalu leah It, itemp0, itemp0;
Instr memu ld itemp0, itemp1;
PUSH(itemp0)
PUSH(itemp1)
PUSH(Intarg0)

CONSTRUCT_LONG_PTR(_sysmalloc_string, Intarg0, DStart)
Instr lalu lsh itemp1, #-32, Intarg1;
Instr lalu mov itemp1, Intarg2;
PUSH(AP)
PUSH(Intarg0)
Instr lalu lea SP, #8, AP;
LIBCALL(_print)
STOP(Intarg0)
STOP(AP)

STOP(Intarg0)
STOP(Itemp1)
STOP(Itemp0)

PUSH(Itemp1)
Instr lalu lea itemp1, Intarg0, itemp1;  -- advance top of mem
Instr memu st itemp1, itemp0;            -- save back in base

Instr lalu mov itemp1, Intarg2;
Instr lalu lsh itemp1, #-32, Intarg1;
CONSTRUCT_LONG_PTR(_sysmalloc_string, Intarg0, DStart)
PUSH(AP)
PUSH(Intarg0)
Instr lalu lea SP, #8, AP;
LIBCALL(_print)
STOP(Intarg0)
STOP(AP)

STOP(RetVal)
STOP(DStart)                             -- return original base

GET_RETIP
Instr lalu jmp RETIP;
FREE_FRAME
Instr;

end;
```

# Appendix D

# MARS C Code

This chapter contains the C source files for the M-Machine runtime system.

```c
/* ........................................................
  .........................................................

  buddy.c

        An attempt to write a buddylist manager
        in C for use in the runtime system
        author: Andrew Shultz
        date: 6/14/95

  ........................................................
  .........................................................  /

#include <stdio.h>
/* #define NULL createPointer(0, 0, P_KEY) */
#include '/home/mmm/apps/runtime/v2.0/yev/pointers.h'

#define POOL_SIZE 50
/* #define ARRAY_SIZE 52 */
#define ARRAY_SIZE 30

#define addrmask 0x003fffffffffffff
#define sizemask 0x0fc000000000000

buddyPP();
buddyPrintlist();
buddyInit();
int buddyPrime();
void * buddyAlloc();
int buddyFree();
int buddyInsert();
void * buddyAllocate();

/* these are in pointer.m */
void * createPointer();
void * setPointer();
void * resetPointer();

struct Frog{
        struct Frog *next;
        int address;
};
typedef struct Frog *frog;

/* data structures have to be external to the procedures? */
/* should these be static? */

/* lock data structure? lock commands? */
/* separate in an assembly language file */

/* pointer to pool */
frog poolPtr;

/* pointer to used */
frog usedPtr;

/* list 'o' frogs */
struct Frog pool[POOL_SIZE];

/* array of free */
frog freeArray[ARRAY_SIZE];

/* array of dirty */
frog dirtyArray[ARRAY_SIZE];

main()
{
        void *ook, *eep, *pbht;
```

```c
        printf("Initializing Buddylist\n");
        buddyInit();
        buddyPrime(createPointer(48, 4, P_RW));
        buddyPrime(createPointer(32, 4, P_RW));
        buddyPrime(createPointer(0, 5, P_RW));
        ook = buddyAlloc(7);
        eep = buddyAlloc(7);
        pbht = buddyAlloc(7);
        buddyFree(eep);
        buddyPP();
        buddyFree(ook);
        buddyFree(pbht);
        printf("Done with Testing\n");
}

buddyPP()
{
        int i;
        /* Pretty print the data structures */

        printf("\nPretty Printing Buddylist Arrays\n");

        buddyLock();

        for(i=0; i<ARRAY_SIZE; i++)
        {
            printf("Free chunks of size %d:\n", i+3);
            buddyPrintlist(freeArray[i], 0);
        }
        for(i=0; i<ARRAY_SIZE; i++)
        {
            printf("Dirty chunks of size %d:\n", i+3);
            buddyPrintlist(dirtyArray[i], 0);
        }

        printf("Used frogs:\n");
        buddyPrintlist(usedPtr, 1);

/*      printf("Pool of frogs:\n");
        buddyPrintlist(poolPtr, 0);
*/
        buddyUnLock();
}

buddyPrintlist(theList, size)
frog theList;
{
        /* print out the list starting at theList */
        if(theList == NULL)
        {
            printf("End\n");
        }
        else
        {
            if(size)
                printf("Pointer: %p\n", theList->address);
            else
                printf("Address: %x\n", theList->address);
            buddyPrintlist(theList->next, size);
        }
}

buddyInit()
```

```c
{
int i;

/*lock somehow.  not sure how*/
buddySetLock();

/*set up the pointers*/
poolPtr = pool;
usedPtr = NULL;

/*set up the pool as a list*/
for(i=0; i<POOL_SIZE-1; i++)
{
    pool[i].next = &pool[i+1];
    pool[i].address = 666;
#if VERBOSE
    printf("*");
#endif
}
#if VERBOSE
printf("\n");
#endif
pool[POOL_SIZE-1].next = NULL;
pool[POOL_SIZE-1].address = 666;
#if VERBOSE
printf("freeArray: %p, dirtyArray: %p\n", freeArray, dirtyArray);
#endif
for(i=0; i<ARRAY_SIZE; i++)
{
    freeArray[i] = NULL;
    dirtyArray[i] = NULL;
#if VERBOSE
    printf(".");
#endif
}
#if VERBOSE
printf("\n");
#endif
buddyUnlock();
}

int buddyPrime(ptr)
void * ptr;
{
frog temp;
int size;

/*again, perform the lock somehow*/
buddyLock();

/*pull a frog out of the pool*/
if (poolPtr == NULL)
{
    printf("No frogs for priming! YAHHHHHH!\n");
    return(0);
}

temp = poolPtr;
poolPtr = poolPtr->next;

temp->address = (int)ptr & addrmask;
size = ((int)ptr & sizemask) >> 54;

/*buddyInsert should handle unlocking*/
return(buddyInsert(temp, size, freeArray));
}

void * buddyAlloc(size)
int size;
{
int i = 0, mask, pos;
frog temp, new;

if (size < 8)
{
    size == 8;
}

size--;

while (size != 0)
{
    size = size >> 1;
    i++;
}

size = i;
pos = i - 3; /*to fix it in the zero-indexed array*/

buddyLock();

if(freeArray[pos] != NULL)
{
    /*buddyAllocate should unlock*/
    temp = freeArray[pos];
    freeArray[pos] = temp->next;
    return(buddyAllocate(temp, size));
}

i = pos;
while((freeArray[i] == NULL) && (i<ARRAY_SIZE))
{
    i++;
}

if(i==ARRAY_SIZE)
{
    buddyUnlock();
    printf("Just no memory at all!\n");
    return(NULL);
}

while(i>pos)
{
    temp = freeArray[i];
    freeArray[i] = temp->next;
    if(poolPtr == NULL)
    {
        return(buddyAllocate(temp, i-3));
    }
    new = poolPtr;
    poolPtr = new->next;
    new->next = NULL;
    temp->next=new;
    mask = 1;
    mask = mask << i + 2; /*to fix zero start, but one less*/
    new->address = (temp->address + mask);
    freeArray[i-1] = temp;
    i--;
```

```c
    temp = treeArray[1];
    treeArray[1] = temp->next;
    return(buddyAllocate(temp, 1-3));
}


int buddyFree(ptr)
    void * ptr;
{
    frog temp, old;
    int size;

    ptr = resetPointer(ptr); /*reset it using leab*/
    buddyLock();
    temp = usedPtr;

    if(temp == NULL) /*nothing to free. shouldn't happen*/
    {
        printf("Huh? I don't have nothin' to free!\n");
        buddyUnLock();
        return(0);
    }

    if(temp->address == (int)ptr)
    {
        size = ((temp->address & sizemask) >> 54);
        temp->address = temp->address & addrmask;
        usedPtr = temp->next;
        return(buddyInsert(temp, size, dirtyArray));
    }

    old = temp;
    temp = temp->next;
    while(temp != NULL)
    {
        if(temp->address == (int)ptr)
        {
            size = ((temp->address & sizemask) >> 54);
            temp->address = temp->address & addrmask;
            old->next = temp->next;
            return(buddyInsert(temp, size, dirtyArray));
        }

        old = temp;
        temp = temp->next;
    }

    printf("Ack, couldn't find it!\n");
    buddyUnLock();
    return(0);
}


int buddyInsert(putIn, size, array)
    frog putIn;
    int size;
    frog* array;
{
    frog left, right, oldleft;
    int mask, temp, pos = size - 3;

    left = array[pos];
    oldleft = array[pos];
    right = array[pos];
    if(left == NULL)
    {
        putIn->next = NULL;
        array[pos] = putIn;
        buddyUnLock();
        return(1);
    }

    mask = 1;
    mask = mask << size;
    temp = putIn->address ^ mask; /*that's XOR*/
    /*
    printf("array: %p\tpos = %d\t, array[pos] = %d\n",
    array, pos, array[pos]); */

    if(right->address > putIn->address)
    {
        array[pos] = putIn;
        putIn->next = right;
        if((int)right->address == temp)
        {
            array[pos] = right->next;
            right->next = poolPtr;
            poolPtr = right;
            return(buddyInsert(putIn, size + 1, array));
        }
        buddyUnLock();
        return(1);
    }

    oldleft = left;
    left = right;
    right = right->next;

    while(right != NULL)
    {
        if(right->address > putIn->address)
        {
            left->next = putIn;
            putIn->next = right;
            if((int)right->address == temp)
            {
                left->next = right->next;
                right->next = poolPtr;
                poolPtr = right;
                return(buddyInsert(putIn, size + 1, array));
            }
            if((int)left->address == temp)
            {
                if(oldleft == array[pos])
                    array[pos] = right;
                else
                    oldleft->next = right;
                putIn->next = poolPtr;
                poolPtr = putIn;
                return(buddyInsert(left, size + 1, array));
            }
            buddyUnLock();
            return(1);
        }
```

```
        oldleft = left;
        left = right;
        right = right->next;
    }
    left->next = putIn;
    putIn->next = NULL;
    if((int)left->address == temp)
    {
        if(oldleft == array[pos])
            array[pos] = right;
        else
            oldleft->next = right;
        putIn->next = poolPtr;
        poolPtr = putIn;
        return(buddyInsert(left, size + 1, array));
    }
    buddyUnlock();
    return(1);
}


void * buddyAllocate(theFrog, size)
frog theFrog;
int size;
{
    size = size; /*leave this out, and DEATH!*/
    /*some sort of write-after-write error, or what?*/
    /*hnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnnm*/
    theFrog->address = (int)createPointer(theFrog->address,
                                          size, P_RW);
    theFrog->next = usedPtr;
    usedPtr = theFrog;
    buddyUnlock();
    return(setPointer(theFrog->address));
}
```

```c
/*
 *
 * Home-node end of the memory-coherence protocol
 *
 ************************************************/

#include <stdio.h>
#include "sysfunc.h"
#include "ccdefs.h"
#include "cc_funcs.h"
#include "sqdefs.h"
#include "signaldefs.h"
#include "opcodes.h"

#define VERBOSE 0
#if VERBOSE
#define Vprintf printf
#else
#define Vprintf
#endif

/* lock on the coherence-directory */
extern int ccdirlock;

int yanklock;

int yankbuf[70];

/* circular buffer of pointers to yankbuffers.  This allows up to
   10 invalidations to be active per node.  May have to increase this
   number, although running out of yankbuffers is not a critical error,
   since you can keep NACKing requests until ones free up */
int *yankBufList[10] = { &yankBuf[0],
                         &yankBuf[7],
                         &yankBuf[14],
                         &yankBuf[21],
                         &yankBuf[28],
                         &yankBuf[35],
                         &yankBuf[42],
                         &yankBuf[49],
                         &yankBuf[56],
                         &yankBuf[63] };

int **yankBufEnd = &(yankBufList[9]);
int **yankBufFree = &(yankBufList[0]);
int **yankBufCur = NULL;

/* assembly procedures */
void MSG_send_NackRO(void *address, int header, int opdat,
                     void *faultCP, int node);
void MSG_send_NackRW(void *address, int header, int opdat,
                     void *faultCP, int node);

/* helper functions */
void ccrequest_ld(void *, int, int, int *, int);
void ccrequest_st(void *, int, int, int *, int);
void ccyankline(void *, int, int, void *, int, int);
void ccyankbody(int *);

void add_eh_job(...);

#define VERBOSE_YANK 0

/* helpful for debugging */

void ccInit() {
    printf("&yankLock = %p\n", &yankLock);
    printf("&yankBufCur = %p\n", &yankBufCur);
    printf("&yankBufFree = %p\n", &yankBufFree);
    printf("yankBufList = %p\n", yankBufList);
}

/**************************************************
 *
 * Home node messages
 *
 **************************************************/

/**************************************************
 *
 * Request(line)
 *
 **************************************************/

void ccrequest(void *address, int header, int opdata, int *faultCP, int node) {
    /* a request for a line is coming in */
    int opcode;

    opcode = (header >> 56) & 0xff;

    Vprintf("ccrequest(%p, ..., %d) called.\n", address, node);

    /* from the opcode, determine which kind of sharing is required */
    switch (opcode) {
    case OPCODE_LD:
        ccrequest_ld(address, header, opdata, faultCP, node);
        break;
    case OPCODE_ST:
    case OPCODE_FST:
        ccrequest_st(address, header, opdata, faultCP, node);
        break;

    default:
        printf("0x%x: unknown opcode in ccrequest\n", opcode);
        break;
    }

    /* need to implement the remaining cases for synchronizing
       memory operations */
}

void ccrequest_ld(void *address, int header, int opdata,
                  int *faultCP, int node) {

    int shareInfo;
    int ppn;
    int *phys_ptr;

    sysGetLock(&ccdirlock);

    shareInfo = SCCShareInfo(address);

    /* TODO: cover case when event caused by own node AND
       condition is no longer valid (i.e. event was read to invalid line
       but the line is no longer shared X and is OK)  That is, when we
       have the line locally now, whereas we didn't at the time the event
       was generated. */

    /* shareInfo returns number of sharers of the line in high halfword
       and type of sharing in low halfword */
    if ((shareInfo & 0xf) == CC_EXCLUSIVE) {
```

```c
            /* sharing is already exclusive, so try to yank lines back from
               those already sharing */

            /* the following function must remember to unlock ccdirlock */
            ccyankline(address, header, opdata, faultCP, node, shareInfo);
            /* line is transitioning, which means we NACK the request */
        } else if ((shareInfo & 0xf) == CC_TRANSITION) {

            Vprintf("ccrequest_id: sending NackRO to %d because %p is transitioning\n",
                    node, address);

            /* the following function must remember to unlock ccdirlock */
            sysNackRO(address, header, opdata, faultCP, node);
        } else {
            /* modify directory for readonly sharing of the cache line */
            /* can we make ccshare also automatically perform a putcstat
               of the line to readonly ?? */
            if (SCCShare(address, node, CC_READONLY) == 3) {
                /* this should never happen! */
                printf("****** ccrequest_id: sharing is wrong!\n");
                sysNackRO(address, header, opdata, faultCP, node);
                return;
            }

            /* demote local line to readonly since another node now will share it
               with us */
            sysPUTCSTAT(address, BSB_READONLY);

            /* find the physical page backing this line so that we can safely
               read it out now.  Can't first read virtually and then set
               to readonly because line may go dirty in the meantime if it was
               readwrite */
            ppn = PPM_lookup(address);
            if (ppn == -1) {
                printf("****** ccrequest_id: ppn mapping for %p not found!\n",
                       address);
                sysPutLock(&ccdirlock);
                return;
            }

            /* calculate the correct cache-line address */
            phys_ptr = get_offsetptr_into_ppn(ppn, address);

            /* the following function must rememebr to unlock the ccdir */
            sysReadAndSendRO(phys_ptr, address, node, header);
        }
}

void ccrequest_st(void *address, int header, int opdata,
                  int *faultCP, int node) {

    int shareInfo;
    int ppn;
    int *phys_ptr;

    sysGetLock(&ccdirlock);

    shareInfo = SCCShareInfo(address);

    /* shareInfo returns number of sharers of the line in high halfword
       and type of sharing in low halfword */
    if ((shareInfo & 0xf) <= 1) {
        /* sharing is already exclusive or readonly, so try to yank lines
           back from those already sharing */

        /* following function must unlock ccdirlock */
        ccyankline(address, header, opdata, faultCP, node, shareInfo);
    } else if ((shareInfo & 0xf) == CC_TRANSITION) {
        /* line is transitioning, which means we NACK the request */

        Vprintf("ccrequest_st: sending NackRW to %d because %p is transitioning\n",
                node, address);

        /* following function must unlock ccdir */
        sysNackRW(address, header, opdata, faultCP, node);
    } else {
        /* modify directory for exclusive sharing of the cache line */
        /* can we make ccshare also automatically perform a putcstat
           of the line to readonly ?? */
        if (SCCShare(address, node, CC_EXCLUSIVE) == 3) {
            printf("****** ccrequest_st: sharing is wrong!\n");
            sysNackRW(address, header, opdata, faultCP, node);
            return;
        }
        sysPUTCSTAT(address, BSB_INVALID);

        ppn = PPM_lookup(address);
        if (ppn == -1) {
            printf(">>>******<<< ccrequest_st: ppn mapping for %p not found!\n",
                   address);
            sysPutLock(&ccdirlock);
            return;
        }

        phys_ptr = get_offsetptr_into_ppn(ppn, address);

        /* following function must unlock ccdir */
        sysReadAndSendX(phys_ptr, address, opdata, node, header);
    }
}

/********************
    ccdirlock >> yanklock
 ********************/

/* perform the necessary sends of invalidation messages for address */
void ccyankline(void *address, int header, int opdata, void *faultCP,
                int node, int shareInfo) {
    /* ccdirlock is held upon entry */

    int *yank_buffer;
    int sharing_nodeid;

    Vprintf("yankline: need to yank line %p with shareInfo %lx\n",
            address, shareInfo);

    sysGetLock(&yankLock);
    if (yankBufFree == yankBufCur) {
        printf("*********** yankline: no yank buffers available.  Must nack requestor.\n"

        /* TODO: send nack */

        sysPutLock(&yankLock);
        sysPutLock(&ccdirlock);
        return;
    }

    /* set yank_buffer to point to the next available yank_buffer we can use */
    yank_buffer = *yankBufFree;
```

```
    if (yankBufOut == NULL)
        yankBufOut = yankBufFree;

    /* just to overwrite the pointer in the circular buffer with
       and int to make sure we never try to follow it until the
       yankbuffer is freed for reuse */
    *(int*)yankBufFree = -2;

    if (yankBufFree == yankBufEnd)
        yankBufFree = &(yankBufList[0]);
    else yankBufFree++;

    /* now that we have a yankbuffer we can use, fill it */
    /* number of invalidation ACK's expected */
    yank_buffer[0] = (shareInfo >> 16) & 0xff;

    /* node number of original requesting node */
    yank_buffer[1] = node;

    /* the four event words */
    yank_buffer[2] = (int)address;
    yank_buffer[3] = header;
    yank_buffer[4] = opdata;
    yank_buffer[5] = (int)faultCP;

    sysPutLock(&yankLock);

SCCShare(address, -1, CC_TRANSITION);
    /* pop the sharing information and send out invalidate messages for
       each node currently sharing the line */
    while(1) {
        sharing_nodeId = SCCPopShare(address);
        if (sharing_nodeId == -1) {
            /* if no more nodes left, we are done! */
            sysGetLock(&ccdirlock);
            sysPutLock(&ccdirlock);
            return;
        } else {
            if (sharing_nodeId == node) {
                /* if the line is shared by the same node whose request we
                   are currently fulfilling, don't send an invalidate message
                   to it, but instead decrement the yank counter as if it
                   sent us an ACK */
                sysGetLock(&yankLock);
                yank_buffer[0] -= 1;
                if (yank_buffer[0] == 0) {
                    /* if yank count is zero, this means that no more
                       invalidations are necessary and no more ACKs are
                       forthcoming. We call yankBody to continue the
                       return-path of the invalidation protocol. */
                    ccyankBody(yank_buffer);
                    sysPutLock(&ccdirlock);
                    sysPutLock(&yankLock);
                    return;
                }
                sysPutLock(&yankLock);
            } else {
                /* send an invalidate message to a sharing node, passing the
                   yank buffer address and the address which is to be
                   invalidated over */

                /* TODO: problem here is that we send out messages while
                   having ccdir locked! */
                Vprintf("yankline: sending invalidate to %d for %p\n",
                        sharing_nodeId, address);
```

```
            if (sharing_nodeId == sysGetNodeId()) {
                /* if you need to send an invalidate message AND you are
                   the node you have to send that message to, then invoke the
                   invalidation routine directly instead of performing
                   a message-send. This is NOT just a performance optimization!
                   The problem is that if your input queue is blocked, you
                   can deadlock by not being able to send yourself a
                   message, while other invalidations still remain to
                   be sent! */
                sysPutLock(&ccdirlock);
                ccInvalidate(address, yank_buffer, sharing_nodeId);
                sysGetLock(&ccdirlock);
            } else
                sysSendInvalidate(sharing_nodeId, address, yank_buffer);
        }
    }
}


void ccreturnYankFull(int *yank_buffer) {
    /* an ACK in response to an invalidate message has come back to the
       home node, bearing a dirty line to be reinstalled */

    Vprintf("ccreturnYankFull: got return for address %p\n",
            yank_buffer[2]);

    sysGetLock(&yankLock);

    /* we don't need to look at the counter because there may only be
       one sharer for a dirty (exclusive) line at a time. */
    int *phys_ptr;
    int ppn;

    /* find the local backing page frame for the address */
    ppn = PPM_lookup(yank_buffer[2]);

    /* calculate cache-line address */
    phys_ptr = get_offset_ptr_into_ppn(ppn, yank_buffer[2]);

    /* store the message locally */
    sysSMSMessage(phys_ptr);
    ccyankBody(yank_buffer);
    sysPutLock(&yankLock);
}


void ccreturnYank(int *yank_buffer) {
    /* an ACK in response to an invalidate message has come back to the
       home node */

    /* we need to look at the counter because there may be multiple sharers
       for a readonly line */
    int *phys_ptr;
    int ppn;

    Vprintf("ccreturnYank: got return for address %p\n", yank_buffer[2]);

    sysGetLock(&yankLock);
    yank_buffer[0] -= 1;
    if (yank_buffer[0] == 0) {
        /* all done */
        ccyankBody(yank_buffer);
        sysPutLock(&yankLock);
    } else {
        /* more invalidations must come back, so we're done for now */
        sysPutLock(&yankLock);
```

```c
}

void ccyankBody(int *yank_buffer) {
    /* yanklock is held upon entry */

    /* all remote holders of the line have relinquished it.
       We can now set the state of the line here back to read/write */
    sysPUTSTAT(yank_buffer[2], BSB_EXCLUSIVE);

    /* now we need to enqueue a 'request' message with the event handler
       to be handled later. */

    Remember that event REQUEST needs to set the ccdir status
    of the line from transitioning to Invalid */

    /* can't do the following because of a compiler-bug. Workaround is
       to enqueue just the yankbuffer pointer and have the routine
       read words out individually. */
    /* add_eh_job(EVENT_SIGNAL_REQUEST, 5, yank_buffer[2],
       yank_buffer[3], yank_buffer[4], yank_buffer[5],
       yank_buffer[1]); */

    Vprintf("ccyankBody: executing [or %p <%d>\n",
       yank_buffer[2], yank_buffer[1]);
    add_eh_job(EVENT_SIGNAL_REQUEST, yank_buffer);

    /* free up the yankbuffer and return to circular buffer */
    yank_buffer[0] = -1;
    yank_buffer[1] = -1;
    yank_buffer[2] = -1;
    yank_buffer[3] = -1;
    yank_buffer[4] = -1;
    yank_buffer[5] = -1;

    if (yankBufCur)
        *yankBufCur = yank_buffer;
    else
        printf("....... ccyankBody: ERROR.  yankBufCur expected to be non-null\n");

    /* now update the yankbuffer stuff */
    if (yankBufCur == yankBufEnd)
        yankBufCur = &(yankBufList[0]);
    else yankBufCur++;

    if (yankBufCur == yankBufFree)
        yankBufCur = NULL;

    /* TODO: must create a new state called CC_TRANSITION2 to be used for
       purposes of making the ccreturnEvicted work better (see below).
       The reason is that we need to have a state similar to
       CC_TRANSITION in that external requests to the line will
       result in NACKs, but also so that a line which was returned
       as dirty to us will not wait for an ACK to a yank to come
       back (so CC_TRANSITION means that some ACKs are still expected,
       but CC_TRANSITION2 means that no ACKs are expected any more, but
       keep NACKing requests to the line anyway) */
}

extern int *event_buf_cur;
void event_buf_advance();

void ccEventRequest() {
    /* event lockword is locked upon entry */
    void *address;
    int header;
    int opdata;
    int *faultCP;
    int node;

    address = (void *)(event_buf_cur[0]);
    event_buf_cur[0] = 0xf11d;
    event_buf_advance();

    header = (int)(event_buf_cur[0]);
    event_buf_cur[0] = 0xf11d;
    event_buf_advance();

    opdata = (int)(event_buf_cur[0]);
    event_buf_cur[0] = 0xf11d;
    event_buf_advance();

    faultCP = (int *)(event_buf_cur[0]);
    event_buf_cur[0] = 0xf11d;
    event_buf_advance();

    node = (int)(event_buf_cur[0]);
    event_buf_cur[0] = 0xf11d;
    event_buf_advance();

    /* get ccdirlock and change line from transitioning to something else */
    sysGetLock(&ccdirlock);
    SCCShare(address, -1, CC_INVALID);
    sysPutLock(&ccdirlock);

    Vprintf("ccEventRequest: calling ccrequest(%p, ..., %d)\n",
        address, node);

    /* now we can try the request which caused the whole invalidation
       business over again because no more nodes are sharing the
       line any more */
    ccrequest(address, header, opdata, faultCP, node);
}

void ccreturnEvicted(void *address) {
    /* a node has pushed on us a dirty line of theirs.
       This line should be installed locally again.  If the line is
       transitioning, then we already sent out an invalidation and should
       wait for the invalidation to come back.  Otherwise, we simply remove
       the node from sharing list for the line */

    int *phys_ptr;
    int ppn, shareInfo;

    Vprintf("ccreturnEvicted: got return for address %p\n", address);

    sysGetLock(&ccdirlock);

    ppn = PPM_lookup(address);
    phys_ptr = get_offsetptr_into_ppn(ppn, address);

    /* store this line locally because it was dirty from whoever
       sent it to us */
    sysSMMessage(phys_ptr);

    shareInfo = SCCShareInfo(address);
    if ((shareInfo & 0xf) == CC_TRANSITION) {
        /* if the line was transitioning, this means that we are still waiting
```

```
         for some reply to come back.  In this case, just keep waiting for
         an ACK */
         kprintf("recv-infinity: %p has been yanked.  Waiting for reply\n");
      } else {
         /* otherwise, remove the fact that this node is sharing our line */
         /* this is where the CC_TRANSITION2 (explained above) will come into play,
            because we won't wait forever for a nonexistent ACK */
         SCCPopShare(address);
         SCCShare(address, -1, CC_INVALID);
      }

      sysPutLock(&ccdirlock);
   }
}
```

```c
/* ***********************************************************
 *
 * Requesting-node end of the memory-coherence protocol
 *
 * ***********************************************************/

#include <stdio.h>
#include <varargs.h>
#include "cc_funcs.h"
#include "sysfunc.h"
#include "pointers.h"
#include "sqdefs.h"
#include "ccdefs.h"
#include "signaldefs.h"
#include "opcodes.h"
#include "eh.h"

#define VERBOSE 0

#if VERBOSE
#define Vprintf kprintf
#else
#define Vprintf
#endif

/* ************************************************************
 *
 * note: need to prove that process can't hold eventLockword
 * and later try for sqlock!!
 *
 *   order is sqlock >> eventLockword
 * ************************************************************/

/* -----------------------------------------------------------
 *
 * requesting node message-handlers
 *
 * -----------------------------------------------------------*/

/* ************************************************************
 *
 * NACK(R) - got a nack to a readonly line request
 *
 * ************************************************************/

void ccNackRO(void *address, int header, int opdata,
                    int *faultCP, int node) {

    int line_state;
    int *bufPtr;
    int ppn;
    eventbuffer eb;

    update_ombc();

    ppn = (header >> 4) & 0xffff;

    sysGetLock(&sqlock);

    line_state = SSQGetState(address, ppn);

    switch(line_state & 0xff) {
    case 0x8:
            /* resend load request */
            SSQGetState(address, 0x8, ppn);
```

```c
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_RESENDLOAD,  4, header, address, opdata,
                        faultCP);

        break;
    case 0x9:
    case 0x10:
        /* resend store request */
        SSQSetState(address, 0x10, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_RESENDSTORE,  4, header, address, opdata,
                        faultCP);

        break;

    case 0xa:
    case 0xe:
        /* signal invalidate so that the line will eventually be invalidated */
        SSQSetState(address, 0x0, ppn);
        bufPtr = (int*)SSQDequeue(address, NULL, ppn);
        if (bufPtr == NULL) {
            kprintf("ccNackRO: signal invalidate bufPtr is NULL!\n");
        }
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_INVALIDATE, 4, node, address, bufPtr, ppn);
        break;
    case 0xc:
        /* signal invalidate and resend load */
        SSQSetState(address, 0x8, ppn);
        bufPtr = (int *)SSQDequeue(address, NULL, ppn);
        sysPutLock(&sqlock);
        eb.header = header;
        eb.opdata = opdata;
        eb.faultCP = faultCP;
        add_eh_job(EVENT_SIGNAL_INV_LOAD, node, address, bufPtr, ppn, &eb);
        break;
    case 0xd:
    case 0x14:
        /* signal invalidate and resend store */
        SSQSetState(address, 0x10, ppn);

        /* if several of the first events were readonly, must find the first
           store attempt for this line and use that one as the request */
        SSQGetFirstW(address, &eb, ppn);
        bufPtr = (int*)SSQDequeue(address, NULL, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_INV_STORE, node, address, bufPtr, ppn, &eb);
        break;
    case 0x18:
        /* set state to 0x10 */
        SSQSetState(address, 0x10, ppn);
        sysPutLock(&sqlock);
        break;
    case 0x1c:
        /* set state to 0x14 */
        SSQSetState(address, 0x14, ppn);
        sysPutLock(&sqlock);
        break;
    default:
        printf("ccNackRO: unknown state: 0x%x\n", line_state);
        sysPutLock(&sqlock);
        break;
    }
}

/* ************************************************************
 *
```

```
* NACK(X) - get a back to an exclusive line request
*
*....................................................*/

void ccBackRW(void *address, int header, int opdata,
              int *bufPtr, int node) {
    int line_state;
    int *bufPtr;
    eventBuffer eb;
    int ppn;

    sysGetLock(&sqlock);

    ppn = (header >> 4) & 0xfffff;
    line_state = SSQGetState(address, ppn);

    switch(line_state & 0xff) {
    case 0x8:
        /* resend load request */
        SSQSetState(address, 0x8, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_RESENDLOAD, 4, header, address,
                   opdata, faultCP);
        break;
    case 0x9:
    case 0x10:
        /* resend store request */
        SSQSetState(address, 0x10, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_RESENDSTORE, 4, header, address,
                   opdata, faultCP);
        break;
    case 0xa:
    case 0xe:
        /* signal invalidate */
        SSQSetState(address, 0x0, ppn);
        bufPtr = (int *)SSQDequeue(address, NULL, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_INVALIDATE, 4, node, address, bufPtr, ppn);
        break;
    case 0xc:
        /* signal invalidate and resend load */
        SSQSetState(address, 0x8, ppn);
        bufPtr = (int *)SSQDequeue(address, NULL, ppn);
        sysPutLock(&sqlock);
        eb.header = header;
        eb.opdata = opdata;
        eb.faultCP = faultCP;
        add_eh_job(EVENT_SIGNAL_INV_LOAD, node, address, bufPtr, ppn, &eb);

        /* compiler bug prevents the following code */
        /* add_eh_job(EVENT_SIGNAL_INV_LOAD, 6, node, address, bufPtr,
                   header, opdata, faultCP); */
        break;
    case 0xd:
    case 0x14:
        /* signal invalidate and resend store */
        SSQSetState(address, 0x10, ppn);
        SSQGetFirstW(address, &eb, ppn);
        bufPtr = (int *)SSQDequeue(address, NULL, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_INV_STORE, node, address, bufPtr, ppn, &eb);
        break;
    case 0x18:
        /* set state to 0x9 */
```

```
        SSQSetState(address, 0x9, ppn);
        sysPutLock(&sqlock);
        break;
    case 0x1c:
        /* set state to 0xd */
        SSQSetState(address, 0xd, ppn);
        sysPutLock(&sqlock);
        break;
    default:
        kprintf("ccBackRW: unknown state: 0x%x\n", line_state);
        sysPutLock(&sqlock);
        break;
    }
}

/*****************************************************
*
* Invalidate() - home node is asking us to invalidate a line
*
*****************************************************/

void ccInvalidate(void *address, void *bufPtr, int node) {
    /* an invalidate message has come in. It tells us which
       cache line to invalidate and which yankbuffer ptr to
       use when returning the invalidation. Also, the node which
       asked us to do the invalidation in the first place */

    int line_state;
    int new_state;
    int sq_return;
    int ppn;

    Vprintf("ccInvalidate: got msg for %p\n", address);

    sysGetLock(&sqlock);

    /* we need to find what the physical page backing this line is
       because the invalidation message does not tell us */
    ppn = PPM_lookup(address);

    if (ppn == -1) {
        /* then we don't have this line (or any line in this page, for that
           matter) locally. Therefore, can simply ACK this invalidate
           request */
        sysSendInvalidateAck(bufPtr, node);
        return;
    }

    /* otherwise, now given the ppn, we can look in the sq hash table
       to find out what the state of this line is */
    line_state = SSQGetState(address, ppn);

    if (line_state == -1) {
        /* this is a critical error: should never happen because we performed
           a lookup in the LPT and since we hold the sqlock, nothing should
           change the mapping from under us */
        printf("ccInvalidate: vpn does not match ppn\n");
        return;
    } else if (line_state & 0x100) {
        /* if this page is marked, it's OK because we are doing whoever
           locked it a favor by evicting lines */
        printf("ccInvalidate: page is marked for eviction\n");

        /* so keep on going */
    }
```

```c
    Vprintf("ccInvalidate(%p): line_state is 0x%x\n",
            address, line_state);

    switch(line_state & 0xff) {
    case 0x0:
        new_state = 0;
        break;
    case 0x8:
        new_state = 0xc;
        break;
    case 0x9:
        new_state = 0xd;
        break;
    case 0xa:
        new_state = 0xe;
        break;
    case 0x10:
        new_state = 0x14;
        break;
    case 0x18:
        new_state = 0x1c;
        break;
    default:
        new_state = -1;
        printf("ccInvalidate: new_state set to -1\n");
        break;
    }

    SSQSetState(address, new_state, ppn);

    if (new_state == 0) {
        /* we can invalidate right now */
        int *phys_ptr;

        Vprintf("ccInvalidate(%p): invalidating now...\n", address);

        line_state = sysPUTCSTAT(address, USD_INVALID);
        if (line_state == BSD_DIRTY) {
            /* ship back the dirty line */
            phys_ptr = get_offsetptr_into_ppn(ppn, address);

            Vprintf("ccInvalidate(%p): sending ack + line\n", address);

            /* remember, this automatically unlocks sqlock */
            sysSendLine(phys_ptr, bufPtr, node);
        } else {
            /* a simple ack will do */
            Vprintf("ccInvalidate(%p): sending ack\n", address);

            /* remember, this automatically unlocks sqlock */
            sysSendInvalidateAck(bufPtr, node);
        }
    } else {
        /* store the yankbuffer ptr because it will be needed later since
           we cannot do the invalidation right away */
        SSQEnqueue(-1, address, bufPtr, NULL, ppn);
        sysPutLock(&sqlock);
    }
}

void ccreturnLoad(void *address, int header, int node) {
    /* a readonly line has come back for us to install locally */
    int line_state, sq_return;
```

```c
    eventBuffer eb;
    void *bufPtr;
    int ppn;
    update_ombc();

    sysGetLock(&sqlock);

    ppn = (header >> 4) & 0xfffff;

    line_state = SSQGetState(address, ppn);
    if (line_state == -1) {
        /* this should not occur since we don't send out a
           request unless we have backing! */
        printf("ccreturnLoad: vpn does not match ppn (%d)\n", ppn);
        return;
    } else if (line_state & 0x100) {
        /* a marked page is ok -- this request will be fulfilled just like */
        printf("ccreturnLoad: page %d is marked\n", ppn);
    }

    switch(line_state & 0xff) {
    case 0x8:
    case 0xc:
        /* we were waiting for this read to come back so that we can
           install the line */
        MSG_install_RO(address, header, node);
        break;
    case 0x9:
        /* we had gotten a NACK for an X line and were waiting for this
           read ack to come back */
        SSQSetState(address, 0x10, ppn);

        SSQGetFirstW(address, &eb, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_RESENDSTORE, 4, address,
                   eb.header, eb.opdata, eb.faultCP);
        break;
    case 0xa:
        /* just waiting for this ack to come back because we
           already installed the line */
        SSQSetState(address, 0x0, ppn);
        sysPutLock(&sqlock);
        break;
    case 0xd:
        /* need to signal invalidate and resend store */
        SSQSetState(address, 0x10, ppn);
        SSQGetFirstW(address, &eb, ppn);
        bufPtr = SSQDequeue(address, NULL, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_INV_STORE, node, address, bufPtr, ppn, &eb);
        break;
    case 0xe:
        bufPtr = SSQDequeue(address, NULL, ppn);
        SSQSetState(address, 0x0, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_INVALIDATE, 4, node, address, bufPtr, ppn);
        break;
    case 0x18:
        /* supposed to install all RO up till first write */
        /* for now, we just transition to state 0x10 */
        SSQSetState(address, 0x10, ppn);
        sysPutLock(&sqlock);
        break;
    case 0x1c:
```

```c
        /* got a RO response.  Now waiting for the X response. */
        SSQSetState(address, 0x14, ppn);
        sysPutLock(&sqlock);
        break;
    default:
        printf("**** ccreturnload: unknown line state of %p (%d)\n",
               address, line_state);
        sysPutLock(&sqlock);
        break;
    }
}

void ccinstall_RO_Done(void *address, int node, int ppn) {
    int line_state;
    int *bufPtr;

    /* no more events */
    line_state = SSQGetState(address, ppn);

    switch(line_state & 0xff) {
    case 0xc:
        bufPtr = (int *)SSQDequeue(address, NULL, ppn);
        SSQSetState(address, 0x0, ppn);
        sysPutLock(&sqlock);
        add_eh_job(EVENT_SIGNAL_INVALIDATE, 4, node, address, bufPtr, ppn);
        break;
    case 0x8:
    case 0xc:
    case 0x10:
    case 0x14:
    case 0x9:
    case 0xa:
    case 0xe:
    case 0xd:
        printf("error: install_ro: unexpected state (%d) for line %p\n",
               line_state, address);
        sysPutLock(&sqlock);
    default:
        SSQSetState(address, 0x0, ppn);

        /* now check if page is marked for eviction.  If it is, and no
           more events remain to that page, it may be evicted */
        if (line_state & 0x100) {
            kprintf("ccinstall_RO_Done: page %d is marked for eviction.  Checking if no
more software entries...\n", ppn);
            line_state = SSQGetState(NULL, ppn); /* the NULL tells you to get informatio
n
                                                    on just whether any software events
which                                                 target that line remain */
            if (!line_state) {
                /* the softwareq queue pointer for that page entry is NULL, so we
                   can begin evicting the line.  Since we are the PIMH, we can't
                   do this ourselves - the EH has to be told to do this. */
                kprintf("ccinstall_RO_Done: all clear.  Page may be evicted\n");
                sysPutLock(&sqlock);
                add_eh_job(EVENT_SIGNAL_EVICT, 1, ppn,
                           ((int)address >> 12) & 0x3fffffffff);
                return;
            }
        }
        sysPutLock(&sqlock);
        return;
    }
}

MSG_Install_RO(void *address, int header, int node) {
    /* Install the readonly line identified by VA address.
       sqlock is already locked by us. */

    int ppn;
    int *phys_ptr;
    eventbuffer eb;
    int dequeue_status;
    int line_state;
    int *bufPtr;

    ppn = ((header >> 4) & 0xffffff);

    /* get the cache-line which contains the address so that we can
       store the contents of the message into the physical page */
    phys_ptr = get_offsetptr_into_ppn(ppn, address);

    /* perform sms's to physical memory and flush that line */
    sysSMSMessage(phys_ptr);

    /* set status bits to readonly */
    sysPUTCSTAT(address, BSB_READONLY);

    sysPutLock(&sqlock);

    while (1) {
        /* dequeue entries from the software queue and satisfy them, one
           at a time.  Unlocking between each satisfication allows other
           threads to add new events, to keep things moving smoothly and
           minimizing spurious messages */
        sysGetLock(&sqlock);

        /* dequeue the next event to this address */
        dequeue_status = (int)SSQDequeue(address, &eb, ppn);

        switch(dequeue_status) {
        case 0:
            /* no more events remain to be handled for this line */
            /* the following function will also unlock the SQ */
            ccinstall_RO_Done(address, node, ppn);
            return;
        case 2:
            /* this event dequeued was also the last event in the SQ for
               this line */
            switch ((eb.header >> 56) & 0xff) {
                /* extract the opcode from the header and decide which
                   operation to perform.  For a readonly line, only a LD
                   is a valid operation */
            case OPCODE_LD:
                mbufLoadUpdate(eb.header, eb.faultCP, 0,
                               phys_ptr[((int)eb.address >> 3) & 0x7]);

                /* ro_done will also unlock the sq */
                ccinstall_RO_Done(address, node, ppn);
                return;
            default:
                /* otherwise, we have a problem since the line is given to us
                   readonly.  This is where we skip all stores ? */
                printf("MSG_Install_RO: skipping opcode %dn",
                       (eb.header >> 56) & 0xff);
                ccinstall_RO_Done(address, node, ppn);
                return;
            }
            break;
```

```c
  default:
    switch ((eb.header >> 56) & 0xff) {
      /* extract the opcode from the header and decide which
         operation to perform.  For a readonly line, only a LD
         is a valid operation */
      case OPCODE_LD:
        mbarLoadUpdate(eb.header, eb.faultCP, 0,
                       phys_ptr[((int)eb.address >> 3) & 0x7]);
        sysPutLock(&sqlock);
        break;
      default:
        /* otherwise, we have a problem since the line is given to us
           readonly.  This is where we skip all stores ? */
        /* TODO */
        printf("MSG_Install_RO: skipping opcode %d\n",
               (eb.header >> 56) & 0xff);
        break;
    }
    break;
  }
}

void ccreturnStore(void *address, int header, int node) {
  /* Install the exclusive line identified by VA address.
     sqlock is already locked by whoever called us */

  int ppn;
  int *phys_ptr, *bufptr;
  eventbuffer eb;
  int line_state, new_state, dequeue_status;

  /* usually, the first store event was performed remotely at the
     home node as well, so this means that the line is dirty only if
     one more write event is going to be handled by us within this
     installation.  Dirty_count keeps track and helps us set status
     to DIRTY instead of READWRITE.  We COULD have set status to
     dirty regardless of dirty_count, but then the line may be
     considered dirty event when it isn't always so.  So this is
     a bit of an optimization */
  int dirty_count = -1;

  update_ombc();
  sysGetLock(&sqlock);

  /* when a line comes back, try putcstat invalid and see what the
     status bits already are.  If they are read/write or dirty, that
     means that the line is still here, and we simply do the requests
     without needing to recopy the line */
  ppn = ((header >> 4) & 0xfffff);

  line_state = sysPUTCSTAT(address, BSB_INVALID);

  /* this is the checking that we perform to make sure that we don't
     already hold a more recent copy */
  if ((line_state == BSB_INVALID) ||
      (line_state == BSB_READONLY)) {

    phys_ptr = get_of(setptr_into_ppn(ppn, address));

    /* perform smg's to physical memory and flush that line */
    sysSMSHessage(phys_ptr);
    sysPUTCSTAT(address, BSB_EXCLUSIVE);
  } else {
    printf("ccReturnStore: duplicate %p comestr", address);
```

```c
  sysPUTCSTAT(address, line_state);
  dirty_count = 0; /* just to be safe */
}

while (1) {
  /* dequeue entries from the software queue and satisfy them */
  dequeue_status = (int)SSQDequeue(address, &eb, ppn);

  switch(dequeue_status) {
    case 0:
      /* no more event to dequeue: we are done */
      /* must unlock the sq */

      /* if the dirty count is > 0, putcstat dirty */
      if (dirty_count > 0)
        sysPUTCSTAT(address, BSB_DIRTY);

      line_state = SSQGetState(address, ppn);
      if (line_state & 0x4) {
        bufptr = (int *)SSQDequeue(address, NULL, ppn);
        SSQSetState(address, 0x0, ppn);
        if (bufptr == NULL) {
          kprintf("ccreturnstore: bufptr is NULL!\n");
        } else {
          sysPutLock(&sqlock);
          add_eh_job(EVENT_SIGNAL_INVALIDATE, 4,
                     node, address, bufptr, ppn);
        }
      } else {
        SSQSetState(address, line_state & 0xf, ppn);
        if (line_state & 0x100) {
          kprintf("ccreturnstore: page %d is marked for eviction.  Checking if
                   no more software entries...\n", ppn);
          line_state = SSQGetState(NULL, ppn);
          if (line_state) {
            kprintf("ccreturnstore: all clear.  Page may be evicted\n");
            sysPutLock(&sqlock);

            /* again, as in the readonly installation in the previous
               function, we ask the EH to do the eviction. */
            add_eh_job(EVENT_SIGNAL_EVICT, 2, ppn,
                       ((int)address >> 12) & 0x3fffffffff);

            return;
          }
          sysPutLock(&sqlock);
        }
        return;
      }
    default:
      switch ((eb.header >> 56) & 0xff) {
        /* extract the opcode from the header and decide which
           operation to perform. */
        case OPCODE_LD:
          mbarLoadUpdate(eb.header, eb.faultCP, 0, *(int*)eb.address);
          break;
        case OPCODE_ST:
        case OPCODE_PST:
          /* store into the physical address */
          *(int*)eb.address = eb.opdata;
          dirty_count++;
          mbarStoreUpdate(eb.faultCP);
          break;
        default:
          /* otherwise, we have a problem since the line is given to us
             readonly.  This is where we skip all stores ? */
```

```
        /* ToDO */
        kprintf("MSG_install_RW: skipping opcode %d\n",
                (eb.header >> 56) & 0xff);

        break;

    break;
    }

  }
}

void ccEvictPage(int ppn, int vpn) {
  /* this page-eviction procedure should be called by the event
     handler when it is asked to clean a page */

  /* In order to evict this page, we need to do the following:
     unmap the page using ppm_unmap.  Then putcstat all individual
     lines within the page to BSB_INVALID, and ship back any
     dirty ones */

  int i, line_state;
  int *base_address, *base_physaddr;
  int dest_node;

  /* since this is executing within an eventhandler slot, we don't
     have to worry about new events to this page sneaking in */

  PPM_unmap(vpn);
  base_address = sysSetPtr((P_RW << 60) | (vpn << 12));
  dest_node = sysGFRB(base_address);

  base_physaddr = sysSetPtr((P_PHYSICAL << 60) | (ppn << 12));
  for (i = 0; i < 64; i++) {
    /* for each cache line ... */
    line_state = sysPUTCSTAT(base_address + (i * 8), BSB_INVALID);
    if (line_state == BSB_DIRTY) {
      /* send back a dirty line */
      sysPushDirty(base_physaddr + (i * 8),
                   base_address + (i * 8),
                   dest_node);
    }
  }

  /* afterwards, remember to reclaim the page by returning it to
     the page pool with PPM_reclaim_remote(ppn) */

}
```

```
/* ......................................................
 * .                                                    .
 * . General Event handler (running in v4-h0 on each node .
 * . Yevgeny Gurevich - 1/29/95 v1.0                    .
 * .                  - 3/13/95 v1.5                    .
 * .                  - 7/21/95 v2.0                    .
 * .                                                    .
 * ......................................................
 */

#include <stdio.h>
#include <varargs.h>
#include "sysfunc.h"
#include "cc_funcs.h"
#include "sqlefs.h"
#include "ccdefs.h"
#include "signaldefs.h"
#include "opcodes.h"
#include "eh.h"
#include "tmanager.h"

#define VERBOSE 0

#if VERBOSE
#define Vprintf kprintf
#else
#define Vprintf Vprintf
#endif

int event_lockword = 0;
extern struct GlobalThreadState tProcesses;

/* several job buffers are needed although only one is
   employed in the current implementation */

/* this buffer should be used by user thread slots only */
int event_job_buffer[257];
int *event_buf_cur = event_job_buffer;
int *event_buf_free = event_job_buffer;
int *event_buf_end = &event_job_buffer[256];
int *event_buf_start = event_job_buffer;

/* this buffer is used by the EH to recirculate its own
   events */
int event_job_buffer2[129];
int *event_buf_cur2 = event_job_buffer2;
int *event_buf_free2 = event_job_buffer2;
int *event_buf_end2 = &event_job_buffer2[128];
int *event_buf_start2 = event_job_buffer2;

/* must make another buffer for the MH's to use */

/* add a new job entry word into an event buffer.
   Unfortunately, this is tailored specifically for the
   normal user-mode event buffer. May need to modify this
   to take an extra argument to decide which buffer this
   word is being added to */
int *add_eh_entry(int arg, int *temp_ptr) {
    *temp_ptr = arg;
    temp_ptr++;
    if (temp_ptr == event_buf_end)
        temp_ptr = event_job_buffer;

    if (temp_ptr == event_buf_cur) {
```

```
        printf("fatal error: out of event buffer space\n");
        /* do something intelligent here ? */
        return NULL;
    }
    return temp_ptr;
}

/* eventlockword >> sqlock ??? */

/* ......................................................
 * .                                                    .
 * . adding a request to an event handler               .
 * .                                                    .
 * ......................................................
 */

/* again, as noted above, this is tailored specifically to
   user-threads adding events to the first event buffer.
   In the current implementation, event the message handlers
   add events using this proc, but that should change - either
   have separate procs for the different buffers, or make this
   function take an extra argument which tells it which buffer
   to add the event to */
void add_eh_job(va_alist)
   va_dcl
{
    int type;
    int num_args;
    int next_arg;
    int *buffer;
    int *temp_buf_free;

    va_list ap;

    va_start(ap);

    type = va_arg(ap, int);
    sysGetLock(&event_lockword);

    if (event_buf_free) {
        eventBuffer *ebuffer;

        temp_buf_free = event_buf_free;

        temp_buf_free = add_eh_entry(type, temp_buf_free);

        switch(type) {
        case EVENT_SIGNAL_REQUEST:
            buffer = va_arg(ap, int *);
            temp_buf_free = add_eh_entry(buffer[2], temp_buf_free);
            temp_buf_free = add_eh_entry(buffer[3], temp_buf_free);
            temp_buf_free = add_eh_entry(buffer[4], temp_buf_free);
            temp_buf_free = add_eh_entry(buffer[5], temp_buf_free);
            temp_buf_free = add_eh_entry(buffer[1], temp_buf_free);

            event_buf_free = temp_buf_free;
            break;
        case EVENT_SIGNAL_INV_LOAD:
        case EVENT_SIGNAL_INV_STORE:
            next_arg = va_arg(ap, int); /* node */
            temp_buf_free = add_eh_entry(next_arg, temp_buf_free);

            next_arg = va_arg(ap, int); /* address */
            temp_buf_free = add_eh_entry(next_arg, temp_buf_free);

            next_arg = va_arg(ap, int); /* bufptr */
```

```c
    temp_buf_free = add_eh_entry(next_arg, temp_buf_free);

    next_arg = va_arg(ap, int);    /* ppn */
    temp_buf_free = add_eh_entry(next_arg, temp_buf_free);

    ebuffer = va_arg(ap, eventBuffer *);
    temp_buf_free = add_eh_entry(ebuffer->header, temp_buf_free);
    temp_buf_free = add_eh_entry(ebuffer->opdata, temp_buf_free);
    temp_buf_free = add_eh_entry(ebuffer->faultCP, temp_buf_free);

    event_buf_free = temp_buf_free;
    break;
  default:
    num_args = va_arg(ap, int);
    while (num_args) {
      next_arg = va_arg(ap, int);
      temp_buf_free = add_eh_entry(next_arg, temp_buf_free);
      num_args--;
    }
    event_buf_free = temp_buf_free;
    break;
  } else {
    printf("eventbuffree is NULL!\n");
  }

  sysPutLock(&event_lockword);

  va_end(ap);

  /* finally, do the special signalling of the EH to make sure it
     is awake and can take software events */
  sysHardwareSignalEH();
}

/* this is a helper function for reading out events. Again, may need
   to modify this to allow taking events from the different event
   queues */
void event_buf_advance() {
  event_buf_cur++;
  if (event_buf_cur == event_buf_end)
    event_buf_cur = event_job_buffer;
}

/* In this loop, the EH pops software jobs from its queue and
   processes them. This will need to change so that the EH
   reads not one, but all three of its queues.  (Only code
   for one has been written) */
void EH_SoftwareDequeueLoop() {
  /* run through the software queue and do what you are supposed to */
  int function;
  int node_id, line_state, line_status, ppn;
  void *vaddr;
  void *yank_buffer;
  int header;
  int opdat;
  void *faultCP;
  struct ThreadContext *c;

  while (1) {
    if ((!event_buf_cur) ||
        (event_buf_cur == event_buf_free)) {
      /* you are done processing */

      /* TODO:

         before you return, try examining the watermark of
         available remote backing pages and perhaps perform
         selective evictions if things are getting tight.
         This means examining the event table, looking for
         entries whose software queue pointers are null as
         likely candidates for eviction. */

      return;
    }

    /* pop off the next event and dispatch */
    function = (int)event_buf_cur[0];
    event_buf_advance();

    switch (function) {
    case EVENT_SIGNAL_INVALIDATE:
      /* need to invalidate a local cache line */
      node_id = (int)event_buf_cur[0];
      event_buf_advance();

      vaddr = (void *)event_buf_cur[0];
      event_buf_advance();

      yank_buffer = (void*)event_buf_cur[0];
      if (yank_buffer == NULL) {
        kprintf("EH: yank_buffer is NULL!\n");
      }
      event_buf_advance();

      ppn = event_buf_cur[0];
      event_buf_advance();

      Vprintf("EH: performing invalidate of %p\n", vaddr);

      sysGetLock(&sqlock);
      line_state = SSQGetState(vaddr, ppn);

      /* zero out the need-to-invalidate bit */
      SSQSetState(vaddr, line_state & 0x1b, ppn);
      line_status = sysPUTCSTAT(vaddr, BSB_INVALID);

      /* both sysSend's will unlock the sq */
      if (line_status == BSB_DIRTY) {
        int *phys_ptr;

        /* ppn = PPM_lookup(vaddr); */
        phys_ptr = get_offsetptr_into_ppn(ppn, vaddr);

        sysSendLine(phys_ptr, yank_buffer, node_id);
      } else {
        sysSendInvalidateAck(yank_buffer, node_id);
      }

      break;
    case EVENT_SIGNAL_RESENDSTORE:
    case EVENT_SIGNAL_RESENDLOAD:
      /* resend a load or store request for a remote line */
      header = (int)event_buf_cur[0];
      event_buf_advance();

      vaddr = (void *)event_buf_cur[0];
      event_buf_advance();

      opdat = (int)event_buf_cur[0];
      event_buf_advance();
```

```
        faultCP = (void *)event_buf_cur[0];
        event_buf_advance();

    SendCCMessage(header, vaddr, opdat, faultCP);
        break;
    case EVENT_SIGNAL_INV_STORE:
    case EVENT_SIGNAL_INV_LOAD:
        /* perform a line invalidation, and then send a request
            for that line again */
        node_id = (int)event_buf_cur[0];
        event_buf_advance();

        vaddr = (void *)event_buf_cur[0];
        event_buf_advance();

        yank_buffer = (void *)event_buf_cur[0];
        event_buf_advance();

        ppn = event_buf_cur[0];
        event_buf_advance();

    Vprintf("EH: performing invalidate + resend store req of %p\n",
            vaddr);

    sysGetLock(&sqlock);
    line_state = SSQGetState(vaddr, ppn);

        /* zero out the need-to-invalidate bit */
    SSQSetState(vaddr, line_state & 0x1b, ppn);
    line_status = sysPUTCSTAT(vaddr, BSB_INVALID);

        /* both sysSend's will unlock the sq */
    if (line_status == BSB_DIRTY) {
            int *phys_ptr;

            /* ppn = PPM_lookup(vaddr); */
            phys_ptr = get_offsetptr_into_ppn(ppn, vaddr);

            sysSendLine(phys_ptr, yank_buffer, node_id);
        } else {
            sysSendInvalidateAck(yank_buffer, node_id);
        }

        header = (int)event_buf_cur[0];
        event_buf_advance();

        opdat = (int)event_buf_cur[0];
        event_buf_advance();

        faultCP = (void *)event_buf_cur[0];
        event_buf_advance();

    SendCCMessage(header, vaddr, opdat, faultCP);
        break;
    case EVENT_SIGNAL_REQUEST:
        /* retry a request (this is on the home node) after
            invalidations have completed */
    ccEventRequest();
        break;
    case EVENT_SIGNAL_EVICT:
        /* evict a shared virtual page to free up a backing page */
        header = event_buf_cur[0];
        event_buf_advance();
        opdat = event_buf_cur[0];
        event_buf_advance();
        kprintf("eh.c: signal to evict page (%d -- %d)\n", opdat, header);
        ccEvictPage(header, opdat);
        break;
    case EVENT_SCHEDULE:
        /* perform thread scheduling */
        tSchedule();
        break;
    case EVENT_SLEEP:
        /* put a thread to sleep */
        tc = (struct ThreadContext *)event_buf_cur[0];
        event_buf_advance();
        tPutToSleep(tc);
        break;
    case EVENT_FORK:
        /* tfork has requested that a new process be added to the
            pending list ... */
        tAddPending((struct ThreadContext *)event_buf_cur[0]);
        event_buf_advance();

        tSchedule();
        break;
    case EVENT_WAKE:
        /* wake up a signalled thread */
        tc = (struct ThreadContext *)event_buf_cur[0];
        event_buf_advance();

        tc->need_to_wake = TRUE;
        tc->signalData = event_buf_cur[0];

        if (tc->VSlot == -1) {
            struct ThreadContext *cur;
            int found;

            for (cur = tProcesses.Pending, found = 0;
                    cur != NULL; cur = cur->Next) {
                if (cur == tc) {
                    found = 1;
                    break;
                }
            }

            if (!found)
                tAddPending(tc);
        }

        event_buf_advance();
        break;
    case EVENT_KILL:
        /* remove this thread from the running list and add it to the
            kill list */
        tKill((struct ThreadContext *)event_buf_cur[0]);
        event_buf_advance();

        tSchedule();
        break;
    default:
        kprintf("****** EH_SoftwareDequeueLoop: unknown job type (%d)\n",
                function);

        break;
    }
}
```

```c
int EH_handle_hsm(int header, void *address, int opdata, void *cp_ptr) {
    int sq_return;
    int ppn;

    ppn = (header >> 4) & 0xffffff;

    if (ppn == -1) {
        /* no physical pages were available */
        /* this is not a completed implementation.  Basically, we try to
           to recirculate the event in an internal buffer (no checking of this
           buffer is coded in SoftwareEnqueueLoop yet) and try to evict pages
           in the meantime */
        printf("*EH_handle_hsm out of physical pages.  Must recirculate this event\n");
        *event_buf_[free2] = header;
        event_buf_[free2]++;
        if (event_buf_[free2] == event_buf_end2)
            event_buf_[free2] = event_job_buffer2;
        *event_buf_[free2] = (int)address;
        event_buf_[free2]++;
        if (event_buf_[free2] == event_buf_end2)
            event_buf_[free2] = event_job_buffer2;
        *event_buf_[free2] = opdata;
        event_buf_[free2]++;
        if (event_buf_[free2] == event_buf_end2)
            event_buf_[free2] = event_job_buffer2;
        *event_buf_[free2] = (int)cp_ptr;
        event_buf_[free2]++;
        if (event_buf_[free2] == event_buf_end2)
            event_buf_[free2] = event_job_buffer2;

        /* try to include code to look for pages to evict here...*/

        return -1;
    }

    sysGetLock(&sqlock);
    sq_return = SSGEnqueue(header, address, opdata, cp_ptr, ppn);
    if (sq_return == -1) {
        /* if returned -1, then the vpn did not match in the event
           table.  This means we have to perform a page lookup
           locally */
        printf("*EH_handle_hsm: stale mapping found for %p\n", address);

        /* need to find new mapping */
        ppn = PPM_lookup(address);
        if (ppn != -1) {
            /* To Do: */
            /* if there was a real mapping in the page table,
               this means that the event table entry is just stale.
               So this means that we need to overwrite it somehow
               (SSGSetVpn(address, ppn) perhaps) and keep on going since
               the entry is truly valid */

        } else {
            /* if ppn = -1, that means that either there are no pages left,
               or no pages AT THE TIME THE MISS OCCURED */
            /* so we can try to perform a PPM_map(vpn) and see if we get
               a real backing page back.  If not, that means that we are
               truly out of physical pages, so need to recirculate
               and perform page eviction again */

        }
        } else if (sq_return == -2) {
            /* page is locked - cannot add new software events to it */
            /* this means we have to recirculate again ... */
```

```c
        printf("*ppn %d is locked down, recirculating %p\n", ppn, address);
        *event_buf_[free2] = header;
        event_buf_[free2]++;
        if (event_buf_[free2] == event_buf_end2)
            event_buf_[free2] = event_job_buffer2;
        *event_buf_[free2] = (int)address;
        event_buf_[free2]++;
        if (event_buf_[free2] == event_buf_end2)
            event_buf_[free2] = event_job_buffer2;
        *event_buf_[free2] = opdata;
        event_buf_[free2]++;
        if (event_buf_[free2] == event_buf_end2)
            event_buf_[free2] = event_job_buffer2;
        *event_buf_[free2] = (int)cp_ptr;
        event_buf_[free2]++;
        if (event_buf_[free2] == event_buf_end2)
            event_buf_[free2] = event_job_buffer2;
        sysPutLock(&sqlock);
        return -1;
    } else {
        /* this is an optimization which would be nice to perform : to notice
           that if the request was for a code line, to fetch in other lines
           past it since chances are, they will be needed sooner or later.. */
        /* for now, we don't do this */

        /* finally, return the result of enqueueing the latest event */
        return sq_return;
    }
}
```

```c
#include <stdio.h>
#include <stdlib.h>

#include "lpt.h"
#include "pplist.h"
#define myprintf printf

/* ...................................... */
/* * LPT                                * */
/* ...................................... */

LPTEntry LPTEntry_init()
{
    LPTEntry e;

    e.vpn = LTLBHashNl1;
    e.ppn = -1;
    /* "LL's removed!*/
    e.status1 = 0xaaaaaaaaaaaaaaaa;
    e.status2 = 0xaaaaaaaaaaaaaaaa;
    return e;
}

#if IS_ANSI
LPT_init(pLPT table)
#else
LPT_init(table)
pLPT table;
#endif
{
    int i;
    for (i = 0; i < LPTSIZE; i++) {
        table->htab[i] = LPTEntry_init();
        table->table_hits[i] = 0;
        /* printf("LPT_init %d completed\n", i); */
    }
    table->numhits = LPTSIZE;
    table->numdels = 0;
    table->numcleans = 0;
    table->numlookups = 0;
    table->numiters = 0;
}

#if IS_ANSI
int LPT_calchash(int vpn, long l)
#else
int LPT_calchash(vpn, l)
int vpn;
long l;
#endif
{
    int result = 0;
    if (!l) {
        /* initial position probed */
        result = (
            /* vpn ^ 0xb525L; */

            (((vpn >> 24) & 0xf{L} << 16) |
            (((vpn >> 16) & 0xf{L} << 24) |
            (((vpn >>  8) & 0xf{L} |
            ((vpn & 0xf{L} << 8)

            ) ^ 0x134aL;
        )

    if (!(result % 2)) return result;
    else return result;
```

```c
#else
int LPT_insert(table, Vpn, Ppn, status1, status2)
pLPT table;
int Vpn;
int Ppn;
int status1;
int status2;
#endif
{
    /* creates a new vpn-ppn mapping in thisnode's hash table */
    /* returns 1 on success, or -1 for failure */
    long i;
    int j;
    i = 0;

    table->numlookups++;
    while (i != LPTSIZE) {
        table->numiters++;
        j = LPT_calchash(Vpn, i);
        j %= LPTSIZE;
        if (    (table->htab[j].vpn == LTLBHashNil) ||
                (table->htab[j].vpn == LTLBHashDeleted)) {
            /* insert into this slot */
            if (table->htab[j].vpn == LTLBHashNil)
                table->numnils--;
            if (table->htab[j].vpn == LTLBHashDeleted)
                table->numdels--;
            table->htab[j].vpn = Vpn;
            table->htab[j].ppn = Ppn;
            table->htab[j].status1 = status1;
            table->htab[j].status2 = status2;
            table->numable_hits[j]++;
#if 0
            if ((table->numnils < (LPTSIZE/8)) && globalclean) {
                /* repartition table */
                LPT_clean(table);
            }
#else
            if ((table->numdels > (LPTSIZE/16)) && globalclean) {
                /* repartition table */
                LPT_clean(table);
            }
#endif
            return 1;
        } else {
            i++;
        }
    }

    myprintf("LPT_insert: Hash Table Overflow.\n");
    return -1;
}

#if IS_ANSI
int LPT_remove(pLPT table, int Vpn)
#else
int LPT_remove(table, Vpn)
pLPT table;
int Vpn;
#endif
{
    /* remove the vpn-ppn mapping from hash table */
    /* returns the ppn on success, -1 on failure */
    long i;
    int j;
```

```c
    i = 0;

    table->numlookups++;
    while (i != LPTSIZE) {
        table->numiters++;
        j = LPT_calchash(Vpn, i);
        j %= LPTSIZE;

        if (table->htab[j].vpn == LTLBHashNil) {
            /* could not find it! */
            return (int)-1;
        } else if (table->htab[j].vpn == LTLBHashDeleted) {
            i++;
        } else if (table->htab[j].vpn == Vpn) {
            table->htab[j].vpn = LTLBHashDeleted;
            table->numdels++;
            i = table->htab[j].ppn;
            table->htab[j].ppn = -1;
            return i;
        } else {
            i++;
        }
    }

    myprintf("LPT_remove: Hash Table Overflow.\n");
    return (int)-1;
}

#if IS_ANSI
int LPT_lookup(pLPT table, int Vpn)
#else
int LPT_lookup(table, Vpn)
pLPT table;
int Vpn;
#endif
{
    /* lookup a vpn-ppn mapping in LPT for the map */
    /* VPN may NOT be 0 or 1 */
    /* return -1 on error */
    long i;
    int j;

    /* profiling */
    table->numlookups++;
    i = 0;
    while (i != LPTSIZE) {
        table->numiters++;
        j = LPT_calchash(Vpn, i);
        j %= LPTSIZE;
        if (table->htab[j].vpn == LTLBHashNil) {
            /* could not find it! */
            return (int)-1;
        } else if (table->htab[j].vpn == Vpn) {
            return table->htab[j].ppn;
        } else {
            i++;
        }
    }

    /*was an fprintf to sterr*/
    printf("LPT_lookup for %ld: Hash Table Overflow.\n",
            (unsigned long)Vpn);
    return (int)-1;
}

#if IS_ANSI
pLPTEntry LPT_findEntry(pLPT table, int Vpn)
```

```c
#else
pLPTEntry LPT_findEntry(table, Vpn)
pLPT table;
int Vpn;
#endif
{
    /* lookup a vpn-ppn mapping in LPT for the map */
    /* return the ptr to the hash table entry */
    long i;
    int j;

    /* profiling */
    table->numlookups++;
    i = 0;
    while (i != LPTSIZE) {
        table->numiters++;
        j = LPT_calchash(Vpn, i);
        i %= LPTSIZE;
        if (table->htab[j].vpn == (int)LTLBHashNil) {
            /* could not find it! */
            return NULL;
        } else if (table->htab[j].vpn == Vpn) {
            return &(table->htab[j]);
        } else {
            i++;
        }
    }

    /*was an fprintf to stderr*/
    printf("LPT_findEntry for %d: Hash Table Overflow.\n",
            (unsigned long)Vpn);
    return NULL;
}

#if IS_ANSI
LPT_stats(pLPT table)
#else
LPT_stats(table)
pLPT table;
#endif
{
    /*
    int i;
    for (i = 0; i < LPTSIZE; i++) {
        myprintf(("<%3d> %5d   \n",i,table->table_hits[i]); i++;
        myprintf(("<%3d> %5d   \n",i,table->table_hits[i]); i++;
        myprintf(("<%3d> %5d   \n",i,table->table_hits[i]); i++;
        myprintf(("<%3d> %5d   \n",i,table->table_hits[i]);
    }
    */

    if (table->numlookups)
    /*a change here: removing '...'s*/
    /*were between lookups and left parend*/
    myprintf(("\nStats: %ld iterations of %ld lookups (%2.2f looks / iter)\nNum cleans:
%ld\n\n",
            table->numiters, table->numlookups,
            (float)table->numiters/
            (float)table->numlookups),
            table->numcleans);
}

#if IS_ANSI
LPTEntry_unparse(pLPTEntry *entry)
#else
LPTEntry_unparse(entry)
```
```c
LPTEntry *entry;
#endif
{
    if (entry->vpn == LTLBHashNil) {
        myprintf("NIL");
    } else if (entry->vpn == LTLBHashDeleted) {
        myprintf("Deleted");
    } else {
        if (entry->ppn)
            myprintf("0x%5lx -- 0x%2lx ", (unsigned long)entry->vpn,
                    (unsigned long)entry->ppn);
        else myprintf("0x%5lx -- 0x%2lx ", (unsigned long)entry->vpn,
                    (unsigned long)entry->ppn);
    }
}

#if IS_ANSI
LPT_unparse(pLPT table)
#else
LPT_unparse(table)
pLPT table;
#endif
{
    long i;

    /* profiling */
    for (i = 0; i < LPTSIZE; i++) {
        myprintf("<%4ld> ",i);
        LPTEntry_unparse(&table->htab[i]);
    }
/*
    i++;
    myprintf("<%4ld> ",i);
    LPTEntry_unparse(&table->htab[i]);

    i++;
    myprintf("<%4ld> ",i);
    LPTEntry_unparse(&table->htab[i]);

    i++;
    myprintf("<%4ld> ",i);
    LPTEntry_unparse(&table->htab[i]); */ /*only print one*/
    myprintf("\n");
}
```

```c
/* ........................................
 *
 * Body of the LTLB Miss Handler
 *
 * The heart of the Physical Memory Manager
 *
 * ........................................ */

#include <stdio.h>
#include "pointers.h"
#include "ppm.h"
#include "pplist.h"
#include "lpt.h"

#define VERBOSE 0
#if VERBOSE
#define Vprintf printf
#else
#define Vprintf
#endif

/* magic numbers for ltlb miss-response design */
#define PPM_UNMAP_MAGIC          0x80000
#define PPM_LOOKUP_MAGIC         0x80002
#define PPM_MAP_MAGIC            0x80003
#define PPM_RECLAIML_MAGIC       0x80004
#define PPM_RECLAIMR_MAGIC       0x80005
#define PPM_LOCAL2REMOTE_MAGIC   0x80006

int sysGPHR(void *);
int *sysSetPtr(int);

extern int ltlb_data_for_MH;
extern LPTable theLPT;

/* handle a PUTCSTAT operation which has faulted to the ltlb.
   Instead of adding the mapping into the ltlb like we normally do,
   we simply modify the block-status bits in the page table
   directly, and return the old bits, just as the instruction
   expects */
void PPM_handle_putcstat(int header, void *vaddr, int opdata,
                         int *faultCP, int vpn) {

    int status1, status2;
    int cache_line;
    int old_status;
    int mask;
    pLPTEntry entry;

    Vprintf("handling putcstat of %d to %p (vpn 0x%x)\n", opdata, vaddr, vpn);

    /* since this was not in the ltlb, need to get the block status, set
       the block status, and return the block status */

    entry = LPT_findentry(&theLPT, vpn);
    if (entry) {
        status1 = entry->status;
        status2 = entry->status2;
    } else {
        kprintf("critical error: putcstat finds no mapped page\n");
        return;
    }

    cache_line = ((int)vaddr & 0xfc0) >> 6;
    Vprintf("cache_line is %d, status1 = 0x%lx, status2 = 0x%lx\n",
            cache_line, status1, status2);
```

```c
    if (cache_line > 31) {

        Vprintf("cache_line > 31\n");

        cache_line -= 32;
        old_status = (status2 >> (cache_line * 2)) & 0x3;
        mask = (0x3 << (cache_line * 2));
        status2 &= (~mask);

        Vprintf("status2 after masking: 0x%lx\n", status2);

        status2 |= (opdata << (cache_line * 2));

        Vprintf("status2 after setting %d: 0x%lx\n", opdata, status2);

    } else {
        old_status = (status1 >> (cache_line * 2)) & 0x3;
        mask = (0x3 << (cache_line * 2));
        status1 &= (~mask);

        Vprintf("status1 after masking: 0x%lx\n", status1);

        status1 |= (opdata << (cache_line * 2));

        Vprintf("status1 after setting %d: 0x%lx\n", opdata, status1);

    }

    /* now write back into page table */
    entry->status1 = status1;
    entry->status2 = status2;

    Vprintf("returning old_status of %d\n", old_status);
    /* now return old bits to caller */
    /* old_status -= 1; */
    old_status += 1;
    *faultCP = old_status;
}

/* handle a generic ltlb miss. */
/* returns 0 if the operation which caused the fault is NOT to be retried.
   this is the case when the fault-response mechanism is used.  Returns 1
   if the hardware should now retry the faulting op. */
int PPM_handle_miss(int header, void *vaddr, int opdata, int *faultCP) {
    int vpn, ppn;
    int ltlbCP;
    int *ltlbCPptr;
    int status1, status2;

    pLPTEntry entry, old_entry;

    Vprintf("PPM_handle_miss called with:\n\tHEADER = %x\n\tVADDR = %p\n\tOPDAT = %x\n\tfa
ultCP = %p\n",
            header, vaddr, opdata, faultCP);

    vpn = (int)vaddr << 10;
    vpn >>= 22;

    if (vpn >= 0x80000) {
        switch(vpn) {
        case PPM_UNMAP_MAGIC:
            /* unmapping a page involves not only taking it out of the
               page table but also removing the mapping from the LTLB */
            vpn = (opdata << 10) >> 22;
```

```c
        ltlb_data_for_MH = !PPM_unmap(vpn);

        /* this is more tricky: we must remove the mapping from the LTLB */
        ltlbCP = 0x200000;
        ltlbCP |= 0x75c0000000000000;
        ltlbCPptr = sysSetPtr(ltlbCP);
        {
            int ltlb_set;
            int word1, word2, word3;
            int lru;
            int valid;

            ltlb_set = (vpn & 0x3f) << 3;

            word1 = ltlbCPptr[ltlb_set];

            if (((word1 & 0x3fffffffffffffff) >> 20) == vpn) {
                /* this is the entry to set invalid */
                word1 &= 0x7fffffffffffffff;  /* turn off valid bit */
                ltlbCPptr[ltlb_set] = word1;
            } else {
                ltlbCPptr += 4;
                word1 = ltlbCPptr[ltlb_set];

                if (((word1 & 0x3fffffffffffffff) >> 20) == vpn) {
                    /* this is the entry to set invalid */
                    word1 &= 0x7fffffffffffffff;  /* turn off valid bit */
                    ltlbCPptr[ltlb_set] = word1;
                }
            }
        }

        *faultCP = 1;
        return 0;

    case PPM_LOOKUP_MAGIC:
        ltlb_data_for_MH = !PPM_lookup((opdata << 10) >> 22);
        *faultCP = 1;
        return 0;

    case PPM_MAP_MAGIC:
        ltlb_data_for_MH = !PPM_map((opdata << 10) >> 22);
        *faultCP = 1;
        return 0;

    case PPM_RECLAIML_MAGIC:
        ltlb_data_for_MH = !PPM_reclaim_local((opdata << 10) >> 22);
        *faultCP = 1;
        return 0;

    case PPM_RECLAIMR_MAGIC:
        ltlb_data_for_MH = !PPM_reclaim_remote((opdata << 10) >> 22);
        *faultCP = 1;
        return 0;

    case PPM_LOCAL2REMOTE_MAGIC:
        ltlb_data_for_MH = !PPM_local2remote(opdata);
        *faultCP = 1;
        return 0;

    default:
        break;
    }

    /* fall-through to standard ltlb miss */


    if ((header >> 56) == 69) {
        /* If operation was a PUTCSTAT, handle it differently from all others */
        PPM_handle_putcstat(header, vaddr, opdata, faultCP, vpn);
        return 0;
    }

    /* find the current virtual-physical mapping in the page table */
    entry = LPT_findEntry(&theLPT, vpn);
    if (!entry) {
        /* if no entry found, try to create a new mapping
           this is the on-demand allocation */
        ppn = !PPM_map(vpn);
        entry = LPT_findEntry(&theLPT, vpn);
        if (!entry) {
            /* if still couldn't make a mapping, this is a critical error! */
            /* ususally, even if out of pages, the ppn will be -1 but a
               mapping WILL exist */
            printf("critical error: could not make new mapping!\n");
            return 0;
        }
    } else {
        ppn = entry->ppn;
    }

    /* status bits from existing entry or an entry you just added */
    status1 = entry->status1;
    status2 = entry->status2;

    /* for now, all cases drop through so that PPN is set and
       we can look up a mapping in the LPT to place into the LTLB */
    /* here's where the fun really begins.  First, we need a
       configuration-space pointer to the ltlb */
    ltlbCP = 0x200000;
    ltlbCP |= 0x75c0000000000000;
    ltlbCPptr = sysSetPtr(ltlbCP);
    {
        int ltlb_set;
        int word1, word2, word3;
        int lru;
        int valid;

        ltlb_set = (vpn & 0x3f) << 3;

        word1 = ltlbCPptr[ltlb_set];
        lru = (word1 >> 62) & 0x1;
        valid = (word1 >> 63) & 0x1;

        Vprintf("ltlbCPptr = %p, ltlb_set = %x, word1 = %x, lru = %d, valid = %d\n",
                ltlbCPptr, ltlb_set, word1, lru, valid);

        /* if this ltlb entry is not least-recently-used, go to the second
           entry in the set */
        if (!lru && valid) {
            ltlbCPptr += 4;
            word1 = ltlbCPptr[ltlb_set];
            lru = (word1 >> 62) & 0x1;
            valid = (word1 >> 63) & 0x1;
        }

        /* read out current entry if it is valid */
        /* valid entries need to have their block-status bits written
           back into the page table, before the entried are evicted from
           the LTLB */
        if (valid) {
            int old_ppn, old_vpn;
```

```c
    word2 = ltlbCPptr[ltlb_set + 1];
    word3 = ltlbCPptr[ltlb_set + 2];

    Vprintf("read out: %x %x %x\n", word1, word2, word3);

    old_ppn = word1 & 0xfff;
    old_vpn = (word1 & 0x3fffffffffffffff) >> 20;

    Vprintf("old ppn = %x, old vpn = %x, oldstatus1 = %x, oldstatus2 = %x\n",
        old_ppn, old_vpn, word2, word3);

    old_entry = LPT_findEntry(&theLPT, old_vpn);
    if (old_entry) {
        old_entry->status1 = word2;
        old_entry->status2 = word3;
    } else {
        printf("error: evicted LTLB entry not in page table!\n");
    }

    /* now we can set the new mapping, overwriting the old one because
       that has been safely copied into the page table */
    word1 = (1 << 61);
    word1 |= (vpn << 20);
    word1 |= ppn;

    Vprintf("setting %p to 0x%x\n", ltlbCPptr + ltlb_set, word1);
    Vprintf("setting %p to 0x%x\n", ltlbCPptr + ltlb_set + 1, status1);
    Vprintf("setting %p to 0x%x\n", ltlbCPptr + ltlb_set + 2, status2);

    ltlbCPptr[ltlb_set] = word1;
    ltlbCPptr[ltlb_set + 1] = status1;
    ltlbCPptr[ltlb_set + 2] = status2;
}

/* return 1 to restart the faulting op since the new data is now
   in the LTLB */
return 1;
}
```

```c
/*****************************************
 ...
 A new version of the physical page manager for the page
 table stuff in lpt.c.

 Andrew Shultz 6/22/95
 ...
 *****************************************/

#include "/home/mm/apps/runtime/v2.0/yev/pointers.h"
#include "pplist.h"

#define PPLISTTEST 0

#if PPLISTTEST
void * fake;
#endif

#if IS_ANSI
PPList_init(pplist thelist, int start, end)
#else
PPList_init(thelist, start, end)
pplist thelist;
int start, end;
#endif
{
    int *workingPtr;
    int workingNum = start;

#if PPLISTTEST
if(start != -1) /*for initing the second, null list*/
{
    fake = alloc(PAGE_SIZE*30);
    printf("fake: %p\n", fake);
    initFake(fake);
}
#endif

thelist->first = start; /*set the list to point to the start*/
thelist->last = end; /*set up the last*/

/*for empty starts*/
if(start == PPNULL)
{
    return;
}

/*now connect the pages in between in a chain*/
while(workingNum < end)
{
    /*write to the physical page*/
    workingPtr = (int *)createPointer(workingNum*PAGE_SIZE,
                                      17,
                                      P_PHYSICAL);

#if VERBOSE
    printf("workingPtr: %p\n", workingPtr);
#endif
    workingNum++;
    *workingPtr = workingNum;
}

/*last one, workingNum = end*/

    workingPtr = (int *)createPointer(workingNum*PAGE_SIZE,
                                      17,
                                      P_PHYSICAL);

    *workingPtr = PPNULL;
}

int PPList_getpage(thelist)
pplist thelist;
{
    int ret = thelist->first;
    int *workingPtr;

    if(thelist->first != PPNULL)
    {
        workingPtr = (int *)createPointer(thelist->first*PAGE_SIZE,
                                          17,
                                          P_PHYSICAL);

        thelist->first= *workingPtr;
        *workingPtr = PPNULL;
    }

    return(ret);
}

PPList_putpage(thelist, thePage)
pplist thelist;
int thePage;
{
    int *workingPtr;

    if(thelist->first == PPNULL)
    {
        thelist->last = thelist->first = thePage;
        workingPtr = (int *)createPointer(thePage*PAGE_SIZE,
                                          17,
                                          P_PHYSICAL);

        *workingPtr = PPNULL;
    }
    else
    {
        workingPtr = (int *)createPointer(thelist->last*PAGE_SIZE,
                                          17,
                                          P_PHYSICAL);

        *workingPtr = thePage;

        workingPtr = (int *)createPointer(thePage*PAGE_SIZE,
                                          17,
                                          P_PHYSICAL);

        *workingPtr = PPNULL;

        thelist->last = thePage;
    }
}

PPList_unparse(thelist)
pplist thelist;
{
    int *workingPtr;
    int temp = thelist->first;

    printf("Printing PP List:\n");

    while(temp != -1)
```

```
    {
        print(("%d, ", temp);
        workingPtr = (int *)createPointer(temp*PAGE_SIZE,
                17,
                P_PHYSICAL);

        temp = *workingPtr;
    }

    print(("End.\n");
}
```

```
/* ..................................................
 * ..................................................

 A physical page manager.
 Uses LPT code in lpt.c and pplist code in pplist.c

 Andrew Shultz 7/29/95
 ................................................. */

#include <stdlib.h>
#include <stdio.h>

#include "lpt.h"
#include "pplist.h"
#include "ppm.h"
#include "pointers.h"

/*global variables*/
PPList remote;
PPList local;
int remotesize;
int localsize;
LPTable theLPT;

#if IS_ANSI
int PPM_Init(int start)
#else
int PPM_Init(start)
int start;
#endif
{
    PPList_Init(&local, start, NUMPHYSPAGES - 1);
    localsize = NUMPHYSPAGES - start;
    PPList_Init(&remote, PPNULL, PPNULL);
    remotesize = 0;
    LPT_Init(&theLPT);

    SHashCreate(&(local.first), &theLPT);
    printf(("PPM Initialized\n");
}

#if IS_ANSI
int PPM_local2remote(int num)
#else
int PPM_local2remote(num)
int num;
#endif
{
    int i, hold;

    for(i=0; i<num; i++)
    {
        hold = PPList_getpage(&local);
        if(hold == PPNULL)
            break;
        else
            PPList_putpage(&remote, hold);
    }

    localsize -= i;
    remotesize += i;

    return i;
}
```

```
#if IS_ANSI
int PPM_remote2local(int num)
#else
int PPM_remote2local(num)
int num;
#endif
{
    int i, hold;

    for(i=0; i<num; i++)
    {
        hold = PPList_getpage(&remote);
        if(hold == PPNULL)
            break;
        else
            PPList_putpage(&local, hold);
    }

    remotesize -= i;
    localsize += i;

    return i;
}

#if IS_ANSI
int PPM_lookup(int vpn)
#else
int PPM_lookup(vpn)
int vpn;
#endif
{
    return LPT_lookup(&theLPT, vpn);
}

#if IS_ANSI
int PPM_map(int vpn)
#else
int PPM_map(vpn)
int vpn;
#endif
{
    int page;

    /*I think this vpn may need to become a virtual address*/

    if(sysGPRB(sysSetPtr((vpn << 12) | (P_RW << 60)]) == sysGetNodeId()) {
        page = PPList_getpage(&local);
        if (page != PPNULL) {
            LPT_insert(&theLPT, vpn, page,
                       0xaaaaaaaaaaaaaaaa,
                       0xaaaaaaaaaaaaaaaa);
            localsize --;
        } else {
            printf("PPM: No more local pages to give out.\n");
            return -1;
        }
    } else {
        page = PPList_getpage(/*&remote*/ &local);
        if(page != PPNULL) {
            LPT_insert(&theLPT, vpn, page, 0, 0);
            remotesize --;
        } else {
            printf("PPM: No more remote pages to give out.\n");
            LPT_insert(&theLPT, vpn, -1, 0, 0);
        }
    }
}
```

```c
    return page;
}

#if IS_ANSI
int IPPM_unmap(int vpn)
#else
int IPPM_unmap(vpn)
#endif
{
    return LPT_remove(&thePT, vpn);
}

#if IS_ANSI
int IPPM_reclaim_local(int ppn)
#else
int IPPM_reclaim_local(ppn)
#endif
{
    if(ppn != PPNULL)   /*allows careless use*/
    {
        PPlist_putpage(&local, ppn);
        localsize ++;
    }
    return ppn;   /*what should this return*/
}

#if IS_ANSI
int IPPM_reclaim_remote(int ppn)
#else
int IPPM_reclaim_remote(ppn)
#endif
{
    if(ppn != PPNULL)
    {
        PPlist_putpage(/* &remote */ &local, ppn);
        remotesize ++;
    }
    return ppn;   /*what should this return*/
}

#if IS_ANSI
int IPPM_remote_left()
#else
int IPPM_remote_left()
#endif
{
    return remotesize;
}

#if IS_ANSI
int IPPM_local_left()
#else
int IPPM_local_left()
#endif
{
    return localsize;
}
```

```
/*
 * ...............................................................
 *
 * Code for managing the memory-coherence software-event queue
 * This is the event table code, the queue node code, and the
 * event node code
 *
 * Please note that this is meant to run within the M-Machine simulator and
 * will need modifications once it's moved into plain C to be run within the
 * runtime system itself
 *
 * ...............................................................
 */

#include <stdio.h>
#include "stddefs.h"
#include "classdefs.h"
#include "msim_lib.h"
#include "ccmagic.h"
#include "sq.h"

int fprintf(FILE *fp, const char *format, ...);

/* the enodelist contains spare, statically-allocated event nodes */
ENode ENodeList[200];
ENode *freeENode = ENodeList;

void ENodeList_init() {
    int i;
    for (i = 0; i < 199; i++)
        ENodeList[i].next = &(ENodeList[i+1]);
    ENodeList[i].next = NULL;
}

ENode* ENodeList_pop() {
    ENode *retval;
    if (freeENode) {
        retval = freeENode;
        freeENode = freeENode->next;

        retval->address.wval.ival = -1LL;
        retval->data = 0LL;
        retval->CP.wval.ival = -1LL;
        retval->next = NULL;

        return retval;
    } else
        return NULL;
}

void ENodeList_push(ENode *node) {
    node->next = freeENode;
    freeENode = node;
}

void ENode_unparse(ENode *node) {
    fprintf(stdout, "0x%slt", CvtXL64(node->header));
    Pointer_print(stdout, node->address.wval.pval);
    fprintf(stdout, "\t0x%slt", CvtXL64(node->data));
    Pointer_print(stdout, node->CP.wval.pval);
    fprintf(stdout, "\n");
}

/* the SQNodeList is a statically-allocated list of available SQNode's
   or the software queue nodes which contain collections of event nodes
   for a particular cache line */
```

```
SQNode SQNodeList[200];
SQNode *freeSQNode = SQNodeList;

void SQNodeList_init() {
    int i;
    for (i = 0; i < 199; i++)
        SQNodeList[i].next = &(SQNodeList[i+1]);
    SQNodeList[i].next = NULL;
}

SQNode* SQNodeList_pop() {
    SQNode *retval;
    if (freeSQNode) {
        retval = freeSQNode;
        freeSQNode = freeSQNode->next;

        retval->next = NULL;
        retval->address = -1LL;
        retval->events = NULL;
        retval->tail = NULL;
        retval->state.px = 0;
        retval->state.pr = 0;
        retval->state.inv = 0;
        retval->state.ax = 0;
        retval->state.nx = 0;
        retval->invalidate_ptr = -1LL;
        retval->emptied = TRUE;

        return retval;
    } else
        return NULL;
}

void SQNodeList_push(SQNode *node) {
    node->next = freeSQNode;
    freeSQNode = node;
}

/* return the cache-line state for a particular line */
SQ_STATE SQNode_getState(SQNode *node, ULong64 address) {
    SQ_STATE retstate;

    retstate.px = 0;
    retstate.pr = 0;
    retstate.inv= 0;
    retstate.ax = 0;
    retstate.nx = 0;

    while(node != NULL) {
        if (node->address == (address & 0x3fffffffffffc0LL))
            return node->state;
        if (node->address > (address & 0x3fffffffffffc0LL))
            return retstate;
        node = node->next;
    }

    return retstate;
}

/* set the state of a particular line */
int BackLog_setState(HEntry *entry, ULong64 address, int newState) {
    /* returns 0 on failure, 1 on success, and 2 if sq can now be
       deleted */
    SQNode *node, *prev = NULL;
```

```c
  node = entry->head;

  for (prev = node; node != NULL; prev = node, node = node->next) {
    if (node->address == (address & 0x3ffffffff{{{c0LL)) {
      node->state.px = (newState >> 4) & 0x1;
      node->state.pr = (newState >> 3) & 0x1;
      node->state.inv= (newState >> 2) & 0x1;
      node->state.ax = (newState >> 1) & 0x1;
      node->state.nx = (newState >> 0) & 0x1;
      if ((newState == 0) && (node->invalidate_ptr == -1LL)) {
        if (node->events) {
          printf("problem: set state to 0 but events remain.\n");
        } else {
          if (prev == node)
            entry->head = node->next;
          else {
            prev->next = node->next;
          }
          SQNodeList_push(node);
        }
        return 2;
      }
      return 1;
    }
  }
  return 0;
}

/* statically-allocated hash table */
REntry BTable[NUM_NODES][BACKING_SIZE];

void REntry_unparse(REntry *entry) {
  SQNode *curnode;
  ENode *enode;

  for (curnode = entry->head; curnode != NULL; curnode = curnode->next) {
    printf(stdout, "      0x%s [%1d%1d%1d%1d%1d%1d] [%s] ===>\n",
           CvtXL64(curnode->address),
           curnode->state.px,
           curnode->state.pr,
           curnode->state.inv,
           curnode->state.ax,
           curnode->state.nx,
           CvtXL64(curnode->invalidate_ptr));
    for (enode = curnode->events; enode != NULL; enode=enode->next) {
      printf(stdout, "        <>     ");
      ENode_unparse(enode);
    }
  }
}

void Backing_unparse(REntry *entry) {
  int i;
  for (i = 0; i < BACKING_SIZE; i++) {
    if (entry[i].vpn != -1LL) {
      printf(stdout, "PN <%3d> -- VPN <0x%s> -- [status %d]\n",
             i,
             CvtXL64(entry[i].vpn),
             entry[i].status);
      REntry_unparse(entry + i);
    }
  }
}

void Backing_init(REntry *BackingHashTable) {
```

```c
  int i;
  for (i = 0; i < BACKING_SIZE; i++) {
    BackingHashTable[i].head = NULL;
    BackingHashTable[i].vpn = -1LL;
    BackingHashTable[i].status = 0;
  }
}

/* whenever we operate on an event table entry, we start at its
   head, and look for the correct cache line. The queue nodes are
   sorted by address, so improve lookup speed */
int Backing_addInvalidate(REntry *entry, ULong64 address, ULong64 ptr) {
  SQNode *curQnode = NULL, *prevQnode = NULL;

  curQnode = entry->head;

  for (prevQnode = curQnode; curQnode != NULL;
       prevQnode = curQnode, curQnode = curQnode->next) {
    /* try to find matching cache line address */
    if (curQnode->address == (address & 0x3ffffffff{{{c0LL)) {
      /* cache-line match */
      if (curQnode->invalidate_ptr != -1LL) {
        fprintf(stderr,
                "addInvalidate: tried to overwrite existing ptr\n");
        return 0;
      } else {
        curQnode->invalidate_ptr = ptr;
        return 1;
      }
    } else if (curQnode->address > (address & 0x3ffffffff{{{c0LL)) {
      /* the list is sorted, so if we didn't find the address,
         too bad. */
      return 0;
    }
  }
  return 0;
}

/* return the invalidate pointer stored for the cache line
   identified by address */
ULong64 Backing_getInvalidate(REntry *entry, ULong64 address) {
  SQNode *curQnode = NULL, *prevQnode = NULL;
  ULong64 retval;

  curQnode = entry->head;

  for (prevQnode = curQnode; curQnode != NULL;
       prevQnode = curQnode, curQnode = curQnode->next) {
    /* try to find matching cache line address */
    if (curQnode->address == (address & 0x3ffffffff{{{c0LL)) {
      /* cache-lines match */
      if (curQnode->invalidate_ptr == -1LL) {
        fprintf(stderr, "getInvalidate: no such pointer set\n");
        return 0LL;
      } else {
        retval = curQnode->invalidate_ptr;
        curQnode->invalidate_ptr = -1LL;
        return retval;
      }
    } else if (curQnode->address > (address & 0x3ffffffff{{{c0LL)) {
      return 0LL;
    }
  }
  return 0LL;
}
```

```c
/* add a new event to the entry.  The enode contains all of the
   necessary information, including virtual address.
   This automatically modifies line-state bits. */
int Backing_addEvent(REntry *entry, ENode *node) {
  /* return codes:
     0 - error
     1 - ok.  No need to send message
     2 - ok.  send a message as well.
     4 - ok.  Stop thread.  no need to send message.
     6 - ok.  Stop thread.  Send a message as well.
     8 - need to recirculate because page is 'locked'
  */

  SQNode *curQnode = NULL, *prevQnode = NULL, *newQnode = NULL;
  int retval = 0;

  if ((node->header >> 52) & 0xfLL == 0xfLL) {
    /* this is an lcache request */
    retval |= 0x4;
  }

  curQnode = entry->head;

  for (prevQnode = curQnode; curQnode; curQnode != NULL;
    prevQnode = curQnode, curQnode = curQnode->next) {

    /* try to find matching cache line address */
    if ((curQnode->address ==
      (node->address.wval.pval.address & 0x3fffffffc0LL))) {
      /* cache-lines match */
      if (curQnode->tail) {
        curQnode->tail->next = node;
        curQnode->tail = node;

        /* if adding a new event - if state was a non-X state and
           got a WR or WI event, set px bit and message must be sent */
        if ((node->header & 0xf) != FAULT_BLOCK_RI) {
          if ((curQnode->state.px == 0) &&
            (curQnode->state.ax == 0) &&
            (curQnode->state.nx == 0)) {
            curQnode->state.px = 1;
            if (curQnode->emptied)
              retval |= 0x8;
            curQnode->emptied = FALSE;
            return (retval | 0x2);
          } else {
            if (curQnode->emptied)
              retval |= 0x8;
            curQnode->emptied = FALSE;
            return (retval | 0x1);
          }
        } else {
          if (curQnode->emptied)
            retval |= 0x8;
          curQnode->emptied = FALSE;
          return (retval | 0x1);
        }
      } else {
        /* first event of its kind */
        curQnode->tail = node;
        curQnode->events = node;

        /* set the state to either px or pr depending on
           type of miss event */
        if ((node->header & 0xf) == FAULT_BLOCK_RI)
```

```c
          curQnode->state.pr = 1;
        else if ((node->header & 0xf) == FAULT_LTLB_MISS) {
          /* case out on operation */
          int opcode = ((int)(node->header >> 56) & 0x3fLL);
          if (opcode == MEMU_LD_ACTION)
            curQnode->state.pr = 1;
          else curQnode->state.px = 1;
        } else
          curQnode->state.px = 1;

        if (curQnode->emptied)
          retval |= 0x8;
        curQnode->emptied = FALSE;
        return (retval | 0x2);
      }
    } else if (curQnode->address >
      (node->address.wval.pval.address & 0x3fffffffc0LL)) {
      break;
    }
  }

  /* no lines there, yet, so we need to add a new software queue
     node */
  newQnode = SQNodeList_pop();
  if (newQnode) {
    newQnode->address =
      (node->address.wval.pval.address & 0x3fffffffc0LL);
    newQnode->next = curQnode;
    newQnode->events = node;
    newQnode->tail = node;

    if ((node->header & 0xf) == FAULT_BLOCK_RI)
      newQnode->state.pr = 1;
    else if ((node->header & 0xf) == FAULT_LTLB_MISS) {
      /* case out on operation */
      int opcode = ((int)(node->header >> 56) & 0x3fLL);
      if (opcode == MEMU_LD_ACTION)
        newQnode->state.pr = 1;
      else newQnode->state.px = 1;
    } else
      newQnode->state.px = 1;
    if (newQnode->emptied)
      retval |= 0x8;
    newQnode->emptied = FALSE;

    if ((prevQnode) && (prevQnode != curQnode)) {
      prevQnode->next = newQnode;
      return (retval | 0x2);
    } else if (prevQnode == curQnode) {
      entry->head = newQnode;
      return (retval | 0x2);
    } else if (curQnode) {
      /* no mapping to a physical page found */
      fprintf(stderr, "error: no backing for VPN %s exists\n",
        CvtXL64((node->address.wval.pval.address >> 12) &
          0x3fffffffLL));
      return 0;
    } else {
      fprintf(stderr, "sg.c: unknown condition for backing_addEvent\n");
      return 0;
    }
  } else {
    fprintf(stderr, "sg.c: can't get new sqnode!\n");
    return 0;
  }
}
```

```c
/* pop off the next event in the entry targeting the cache-line
   identified by address */
ENode *Backing_popEvent(REntry *entry, ULong64 address) {
    SQNode *curnode = NULL, *prevnode = NULL;
    ENode *retval;

    curnode = entry->head;

    for (prevnode = curnode; curnode != NULL;
         [prevnode = curnode, curnode = curnode->next) {
        /* try to find matching cache line address */
        if (curnode->address == (address & 0x3fffffffffffc0LL)) {
            /* cache-lines match */
            retval = curnode->events;

            if (!curnode->events) {
                curnode->emptied = TRUE;
            }

            if (curnode->events == curnode->tail) {
                curnode->tail = NULL;
                curnode->events = NULL;
                /* SQNodeList_push(curnode); */
                entry->head = NULL; */
            } else {
                curnode->events = curnode->events->next;
            }

            return retval;
        }
    }

    /* line not found */
    return NULL;
}

int header_isWriteEvent(ULong64 header) {
    int opcode = ((int)((ULong64)header >> 56));
    return (opcode != 47);
}

ENode *Backing_firstWriteEvent(REntry *entry, ULong64 address) {
    SQNode *curnode = NULL, *prevnode = NULL;
    ENode *retval;

    curnode = entry->head;

    for (prevnode = curnode; curnode != NULL;
         prevnode = curnode, curnode = curnode->next) {
        /* try to find matching cache line address */
        if (curnode->address == (address & 0x3fffffffffffc0LL)) {
            /* cache-lines match */
            /* run through nodes until reach event whose header is
               the first 'write' event */
            for (retval = curnode->events; retval != NULL;
                 retval = retval->next) {
                if (header_isWriteEvent(retval->header))
                    return retval;
            }
        }
    }

    /* line not found */
    return NULL;
}
```

```c
void BackingTable_unparse(int n) {
    Backing_unparse(BTable[n]);
}
```

```c
/*******************************************
 * tmanager.c: Core M-Machine runtime system thread management code  *
 * Written by:                  Yevgeny Gurevich                     *
 *                              June 1995 - August 1995              *
 *******************************************/

#include <stdio.h>
#include <varargs.h>
#include "sysfunc.h"
#include "tmanager.h"
#include "signaldefs.h"
#include "signal.h"
#include "pointers.h"

#define VERBOSE_TMANAGER  0
#define VERBOSE_TSIGNALLER 0

#if VERBOSE_TMANAGER
#define Vprintf kprintf
#else
#define Vprintf Vprintf
#endif

void *malloc(int);
int buddyFree(void *);
void tAddChild(struct ThreadContext *, void *);

struct GlobalThreadState tProcesses;

/* Initialize the HContext data structure */
void HContext_Init(struct HContext *context) {
    int i;

    for (i = 0; i < 16; i++) {
        context->int_reg_file[i] = 0;
        context->fp_reg_file[i] = (float)0;
    }

    context->local_cc[0] = 0;
    context->local_cc[1] = 0;
    context->local_cc[2] = 0;
    context->local_cc[3] = 0;

    context->empty_scoreboard = 0;
    context->hardware_mbar_counter = 0;
    context->software_mbar_counter = 0;

    context->restartIPvector[0] = NULL;
    context->restartIPvector[1] = NULL;
    context->restartIPvector[2] = NULL;
    context->restartIPvector[3] = NULL;
}

/* Initialize the threadcontext data structure */
void ThreadContext_Init(struct ThreadContext *context) {
    context->Next = NULL;
    context->VSlot = -1;
    context->Parent = NULL;
    context->Sibling = NULL;
    context->Children = NULL;
    HContext_Init(&(context->hthreads[0]));
    HContext_Init(&(context->hthreads[1]));
    HContext_Init(&(context->hthreads[2]));
    HContext_Init(&(context->hthreads[3]));
    context->flags = 0;
```

```c
    context->SCC = 1024;
    context->SCL = 1024;

    context->signalData = 0xdeadbeef;
    context->need_to_wake = FALSE;
    context->need_to_sleep = FALSE;
    context->need_to_block = FALSE;
}

/* print out global thread information for debugging purposes */
void ThreadContext_unparse(struct ThreadContext *context) {
    struct ThreadContext *cur = NULL;

    for (cur = context; cur != NULL; cur = cur->Next) {
        printf("TC %p:\n", cur);
        printf("\tNext        %p | signalData    %lx (swb %d%d%d)\n",
               cur->Next, cur->signalData,
               cur->need_to_sleep,
               cur->need_to_wake,
               cur->need_to_block);
        printf("\tVSlot       %lx         | flags        %lx\n",
               cur->VSlot, cur->flags);
        printf("\t-------------------------------------------------\n");
        printf("\tH0 IP0      %p\n", cur->hthreads[0].restartIPvector[0]);
        printf("\tH0 IP1      %p\n", cur->hthreads[0].restartIPvector[1]);
        printf("\tH0 IP2      %p\n", cur->hthreads[0].restartIPvector[2]);
        printf("\tH0 IP3      %p\n", cur->hthreads[0].restartIPvector[3]);
        printf("\tH0 SP       %p\n", cur->hthreads[0].int_reg_file[2]);
        printf("\tH0 DP       %p\n", cur->hthreads[0].int_reg_file[5]);
        printf("\tH0 AP       %p\n", cur->hthreads[0].int_reg_file[12]);
        printf("\t-------------------------------------------------\n");
        printf("\tMbar Counters  %d %d %d %d\n",
               cur->hthreads[0].hardware_mbar_counter,
               cur->hthreads[1].hardware_mbar_counter,
               cur->hthreads[2].hardware_mbar_counter,
               cur->hthreads[3].hardware_mbar_counter);
        printf("\t-------------------------------------------------\n");
        printf("\tParent      %p\n", cur->Parent);
        if (cur->Children) {
            struct ThreadContext *cur2;
            printf("\tChildren:\n\t\t");
            for (cur2 = cur->Children; cur2 != NULL; cur2 = cur2->Sibling) {
                printf("%p ", cur2);
            }
            printf("\n");
        }
        printf("\t-------------------------------------------------\n");
    }
}

/* deallocate a thread context - careful in terms of which locks are held */
int ThreadContext_free(struct ThreadContext *tc) {
    /* for now we will do nothing.  In the future, a REQUEST for a
       free operation to occur needs to be given to the buddylist
       manager */
    /* return buddyFree((void*)tc); */
}

/* allocate a new thread context */
struct ThreadContext *tAlloc() {
    /* first allocate new memory for a threadcontext, and then
       initialize the data strucure */
    struct ThreadContext *tc;
    tc = (struct ThreadContext*)malloc(sizeof(struct ThreadContext));
    ThreadContext_Init(tc);
```

```
    return tc;
}

/* Initialize all on-node data structures dealing with
   thread management.  This code is to be executed as
   part of the boot sequence. */
void tInit(void *DataPtr) {
    /* Initialize processes */
    struct ThreadContext *self;
    int **CP;
    int *myCP;

    tProcesses.occupied = 0x1;
    tProcesses.Running = NULL;
    tProcesses.RunningEnd = NULL;
    tProcesses.Pending = NULL;
    tProcesses.PendingEnd = NULL;
    tProcesses.Kill = NULL;
    tProcesses.KillEnd = NULL;

    /* set up own processes properly */
    self = tAlloc();
    self->hthreads[0].int_reg_file[5] = (int)DataPtr;
    self->flags = 0x11;  /* hthread 0 is both full and may issue */
    self->vSlot = 0;       /* I am running in slot 0 */

    /* now write this CP into own context ptr */
    CP = (int**)sysSetPtr(0x77c000000510000);
    myCP = *CP;

    myCP[0x2000] = (int)self;
    tAddRunning(self);

    /* Initialize the signal table */
    sigTab_init();
}

void sysUnparse() {
    print("Pending:\n");
    ThreadContext_unparse(tProcesses.Pending);
    print("Kill:\n");
    ThreadContext_unparse(tProcesses.Kill);
    print("Running:\n");
    ThreadContext_unparse(tProcesses.Running);
}

/* add a thread context to the pending threads list */
void tAddPending(struct ThreadContext *tc) {
    if (tProcesses.PendingEnd) {
        tProcesses.PendingEnd->Next = tc;
        tProcesses.PendingEnd = tc;
        tProcesses.PendingEnd->Next = NULL;
    } else {
        tProcesses.Pending = tc;
        tProcesses.PendingEnd = tc;
    }
    tc->Next = NULL;
}

/* add a thread context to the to-be-killed threads list */
void tAddKill(struct ThreadContext *tc) {
    if (tc) return;
    if (tProcesses.KillEnd) {
        tProcesses.KillEnd->Next = tc;
```

```
        tProcesses.KillEnd = tc;
        tProcesses.KillEnd->Next = NULL;
    } else {
        tProcesses.Kill = tc;
        tProcesses.KillEnd = tc;
    }
    tc->Next = NULL;
}

/* pop the first thread context from the pending threads list */
struct ThreadContext *tPopPending() {
    struct ThreadContext *retval;

    if (tProcesses.Pending) {
        retval = tProcesses.Pending;
        if (tProcesses.PendingEnd == tProcesses.Pending)
            tProcesses.PendingEnd = NULL;
        tProcesses.Pending = tProcesses.Pending->Next;
        return retval;
    } else {
        return NULL;
    }
}

/* add a thread context to the running threads list */
void tAddRunning(struct ThreadContext *tc) {
    if (!tc) return;
    if (tProcesses.RunningEnd) {
        tProcesses.RunningEnd->Next = tc;
        tProcesses.RunningEnd = tc;
    } else {
        tProcesses.Running = tc;
        tProcesses.RunningEnd = tc;
    }

    tc->Next = NULL;
    tProcesses.occupied |= (1 << tc->vSlot);
}

/* pop the first thread context from the running threads list */
struct ThreadContext *tPopRunning() {
    struct ThreadContext *retval;

    if (tProcesses.Running) {
        retval = tProcesses.Running;
        if (tProcesses.RunningEnd == tProcesses.Running)
            tProcesses.RunningEnd = NULL;
        tProcesses.Running = tProcesses.Running->Next;
        return retval;
    } else {
        return NULL;
    }
}

/* pop the first thread context from the to-be-killed threads list */
struct ThreadContext *tPopKill() {
    struct ThreadContext *retval;

    if (tProcesses.Kill) {
        retval = tProcesses.Kill;
        if (tProcesses.KillEnd == tProcesses.Kill)
            tProcesses.KillEnd = NULL;
        tProcesses.Kill = tProcesses.Kill->Next;
        return retval;
    } else {
        return NULL;
    }
}
```

```
}

/* Fork off a new thread
Interface is:
(void *threadIP, void *DataPtr, void *returnPtr, int numargs,
void *parentTC, ... ) */
void* SYSFork(va_alist)
va_dcl
{
struct ThreadContext *newContext;
int *temp;
char *threadIP;
void *DataPtr, *returnPtr, *returnval, *passed_parent;
int numargs, i;

va_list ap;
va_start(ap);
threadIP = va_arg(ap, char *);
DataPtr = va_arg(ap, void *);
returnPtr = va_arg(ap, void *);
numargs = va_arg(ap, int);
passed_parent = va_arg(ap, void *);

/* allocate a new thread context */
newContext = tAlloc();
newContext->hthreads[0].int_reg_file[5] = (int)DataPtr;
newContext->hthreads[0].restartIPvector[0] = threadIP;
newContext->hthreads[0].restartIPvector[1] = threadIP + 4;
newContext->hthreads[0].restartIPvector[2] = threadIP + 8;
newContext->hthreads[0].restartIPvector[3] = threadIP + 12;

newContext->flags = 0x11; /* hthread 0 is both full and may issue */

/* allocate stack/AP for the new thread */
temp = (int *)malloc(0x4000);  /* user stack */
if (temp) {
    temp += 2047;  /* end of stack */
    newContext->hthreads[0].int_reg_file[1] = (int)threadIP;
    newContext->hthreads[0].int_reg_file[2] = (int)temp;
    newContext->hthreads[0].int_reg_file[4] = (int)returnPtr;
    newContext->hthreads[0].int_reg_file[12] = (int)temp;

    newContext->hthreads[0].empty_scoreboard = 0x1034;

    /* put arguments in new thread's stack. The first six also
       go into the thread's register file */
    /* problem is how to deal with differences between FP and INT
       arguments. For now, all are assumed to be int arguments */
    for (i = 0; i < numargs; i++) {
        int arg;

        arg = va_arg(ap,int);
        *temp = arg;
        if (i < 6) {
            newContext->hthreads[0].int_reg_file[6 + i] = arg;
            newContext->hthreads[0].empty_scoreboard |=
                1 << (6 + i);
        }
        temp -= 1;
    }

    va_end(ap);

    /* add the newly-created context as a child of its parent. */
```

```
    tAddChild(newContext, passed_parent);

    /* add the newly-created context to the pending list of processes
       signifying that it is ready to issue and may be scheduled. */
    /* tSafeAddPending(newContext); */
    add_eh_job(EVENT_FORK, 1, newContext);

    /* signal the event handler to perform scheduling to hopefully
       let this context run */
    add_eh_job(EVENT_SCHEDULE, 0);

    /* return the child context pointer to the parent */
    returnval = (void*)sysSetPtr(((int)newContext & 0x0fffffffffffffff) |
                                 (P_KEY << 60));

    return returnval;
} else {
    printf("tfork: malloc failed\n");
    va_end(ap);
    return NULL;
}
}

int ConvertHCflags(int flags, int slot) {
/* convert from relative-hthread to absolute cluster */
int new_flags;

switch (slot) {
case 0:
    return flags;
case 1:
    new_flags = ((((flags >> 3) & 0x1) |
                 (((flags & 0x7) << 1) |
                 ((((flags >> 7) & 0x1) << 4) |
                 ((((flags >> 4) & 0x7) << 5));

    return new_flags;
case 2:
    new_flags = ((((flags >> 2) & 0x3) |
                 (((flags & 0x3) << 2) |
                 ((((flags >> 6) & 0x3) << 4) |
                 ((((flags >> 4) & 0x3) << 6));

    return new_flags;
case 3:
    new_flags = ((((flags >> 1) & 0x7) |
                 (((flags & 0x1) << 3) |
                 ((((flags >> 5) & 0x7) << 4) |
                 ((((flags >> 4) & 0x1) << 7));

    return new_flags;
default:
    return 0;
}
}

int ConvertCHflags(int flags, int slot) {
/* convert from absolute cluster to HO relative notation */
int new_flags;

switch (slot) {
case 0:
    return flags;
case 1:
    new_flags = ((((flags & 0xe) >> 1) |
                 (((flags & 0x1) << 3) |
                 ((((flags & 0xe0) >> 1) |
                 ((((flags & 0x10) << 3));

    return new_flags;
```

```
    case 2:
        new_flags = (((flags & 0xc) >> 2) |
                     ((flags & 0x1) << 2) |
                     ((flags & 0xc0) >> 2) |
                     ((flags & 0x10) << 2));

        return new_flags;
    case 3:
        new_flags = (((flags & 0x8) >> 3) |
                     ((flags & 0x7) << 1) |
                     ((flags & 0x80) >> 3) |
                     ((flags & 0x70) << 1));

        return new_flags;
    default:
        return 0;
    }
}

/* read hardware hthread state from vslot, slot, and store it into the
   hthread, hc. Returns the state of the membar counter of the hthread */
int tEvictHThread(int *CP, struct HContext *hc, int slot) {
    int i;

    for (i = 0; i < 16; i++) {
        hc->int_reg_file[i] = CP[i];
    }
    CP += 32;
    for (i = 0; i < 16; i++) {
        hc->fp_reg_file[i] = CP[i];
    }
    CP += 32;
    for (i = 0; i < 4; i++) {
        hc->local_cc[i] = CP[i];
    }
    CP += 160;
    for (i = 0; i < 4; i++) {
        hc->restart_ifvector[i] = (char *)CP[i];
    }
    CP += 4;
    hc->hardware_mbar_counter = CP[0];
    hc->empty_scoreboard = CP[1];
    return hc->hardware_mbar_counter;
}

/* attempt to evict the threadcontext from its slot
   returns the slot which was freed up due to the eviction on success,
   -1 on failure */
int tEvict(struct ThreadContext *tc,
           int slot, evict_flags, cluster;
           int *CP;
           int mbar_stat;

    slot = tc->vslot;

/* In order to be able to evict a thread, its hardware mbar counter
   should be zero, or its software mbar counter should equal its
   hardware mbar counter (which means that all events now use the
   threadcontext pointer and not the absolute configspace pointer for the
   thread */
```

```
    /* first thing is to cons up a configspace ptr to the thread,
       read out the current hissue status and set it to 0, waiting
       for pipe to hopefully drain */

    CP = sysMakeCP(slot);
    evict_flags = CP[0x20001];

    /* kill hissue bits to stop the thread from sending ops into the pipe */
    CP[0x20001] = evict_flags & 0x0f;

    /* now that we have the context, start storing stuff away into it.
       Remember that the CP points to physical clusters and we need
       logical units here. */
    Vprintf("Evicting tc %p from slot %d [flags %x]\n", tc, slot, evict_flags);
    tc->flags = ConvertCHflags(evict_flags, slot);

    for (cluster = 0; cluster < 4; cluster++) {
        if (evict_flags & (1 << cluster)) {
            mbar_stat = tEvictHThread(CP + (0x200 * cluster),
                                      &(tc->hthread[(4 * cluster) - slot)),
                                      & 4], slot);

            if (mbar_stat) {
                /* can't evict after all */
                /* so reset to issuing */
                CP[0x20001] = evict_flags;
                printf("mbar for %d of %p was not zero. Can't evict\n",
                       tc, cluster);
                return -1;
            }
        }
    }

    tc->vslot = -1;
    tc->Next = NULL;

#if VERBOSE_TMANAGER
    kprintf("evicted context %p:\n", tc);
    ThreadContext_unparse(tc);
#endif

    /* now set the occupied state of the slot to zero and also
       set the occupied hthread bits of the slot to unoccupied */
    CP[0x20001] = 0x0;
    tProcesses.occupied &= ~(1 << slot);
    return slot;
}

int tGetSlot() {
    /* returns a thread slot available for a thread to be set
       running in. If no slots are free, evicts an existing
       thread if possible */

    if ((tProcesses.occupied & 0x1) == 0) {
        /* vslot 0 available */
        return 0;
    } else if ((tProcesses.occupied & 0x2) == 0) {
        /* vslot 1 available */
        return 1;
    } else if ((tProcesses.occupied & 0x4) == 0) {
        /* vslot 2 available */
        return 2;
    } else if ((tProcesses.occupied & 0x8) == 0) {
        /* vslot 3 available */
        return 3;
```

```
} else {
    int evict_attempts, evict_slot;
    struct ThreadContext *evict;
    Vprintf("No slots available.  Trying to evict a thread.\n");
    for (evict_attempts = 0; evict_attempts < 4; evict_attempts++) {
        if ((evict = tPopRunning()) != -1) {
            if ((evict_slot = tEvict(evict)) != -1) {
                /* evicted the thread, which means we attach it to the
                   pending list */

                /* If the thread needs to block, it should not be
                   added to the pending list, since it is already
                   somewhere in the signal list. */
                if (!evict->need_to_block)
                    tAddPending(evict);
                return evict_slot;
            } else {
                /* If we couldn't evict the thread, we return it to
                   the back of the running list */
                tAddRunning(evict);
            }
        } else {
            return -1;
        }
    }

    /* could not evict any running thread, and all thread slots
       are taken */
    return -1;
}
```

```
void tInstallHThread(int *CP, struct HContext *hc, int slot) {
    /* set up hthread state by copying it into hardware through the
       configspace pointer (CP) */
    int i;

    CP[0] = hc->int_reg_file[0];

    for (i = 2; i < 16; i++) {
        CP[i] = hc->int_reg_file[i];
    }

    CP += 32;
    for (i = 0; i < 16; i++) {
        CP[i] = hc->fp_reg_file[i];
    }

    CP += 32;

    for (i = 0; i < 4; i++) {
        CP[i] = hc->local_cc[i];
    }

    CP += 164;

    i = hc->hardware_mbar_counter;
    CP[0] = i;
    CP[1] = hc->empty_scratchpad;

    CP += 4;

    /* now let things fly for this hthread */
    for (i = 0; i < 4; i++) {
        CP[i] = (int)hc->instantiflyvector[i];
```

```
int tInstall(struct ThreadContext *tc, int slot) {
    /* install the thread context into a free slot */
    int flags;
    int *CP;
    int cluster;
    CP = sysMakeCP(slot);

    /* the first thing is to set all of the proper hFull bits */
    flags = ConvertHCflags(tc->flags & 0x0f, slot);

    CP[0x200] = flags;
    CP[0x200] = (int)tc;
    CP[0x200] = tc->SCC;
    CP[0x204] = tc->SCL;

    for (cluster = 0; cluster < 4; cluster++) {
        if (flags & (1 << cluster)) {
            tInstallHThread(CP + (0x200 * cluster),
                            &(tc->hthreads[((4 * cluster) -slot) & 4]), slot);
        }
    }

    tc->VSlot = slot;

    /* add this thread to the list of running threads */
    tAddRunning(tc);
    tProcesses.occupied |= (1 << slot);
    return 1;
}
```

```
void tCleanThreads() {
    /* kill dead threads and take them out of commission */
    struct ThreadContext *candidate = NULL, *recirculated = NULL;
    int hardware_flags, i;
    int *CP;
    int slot;

    /* no need to do lots of locking here since a thread on
       the kill list is not issuing, and no one else will change
       its thread slot.  So no need to lock things as we
       create the CP to that slot and read out the mbar counters */
    while ((candidate = tPopKill()) != NULL) {
        if (candidate == recirculated)
            /* break out of loop */
            break;

        slot = candidate->VSlot;
        Vprintf("tCleanThreads: taking %p (slot %d) out.\n", candidate, slot);
        if (slot != -1) {
            CP = sysMakeCP(slot);
            /* read out the mbar counters for the thread to be killed.
               If the thread still has positive mbar values, it cannot
               be killed until everything comes back. */
            if (CP[0xe4] || CP[0x2e4] || CP[0x4e4] || CP[0x6e4]) {
                Vprintf("tCleanThreads: can't take thread out yet.\n");
                tAddKill(candidate);
                if (recirculated == NULL)
                    recirculated = candidate;
            } else {
                CP[0x200] &= 0xf;     /* kill issue bits */
                CP[0x200] = 0;        /* kill occupied bits */
```

```
          /* this is a little questionable -
             interface to buddylist code */
          ThreadContext_free(candidate);
          /* finally, modify the occupied information, demonstrating that
             the thread slot is no longer occupied by a running thread */
          tProcesses.occupied &= ~(1 << slot);
      }
  } else {
      if (candidate->hthreads[0].hardware_mbar_counter ||
          candidate->hthreads[1].hardware_mbar_counter ||
          candidate->hthreads[2].hardware_mbar_counter ||
          candidate->hthreads[3].hardware_mbar_counter) {
          vprintf("tCleanThreads: can't take thread out yet\n");
          tAddKill(candidate);
          if (recirculated == NULL)
              recirculated = candidate;
      } else {
          ThreadContext_free(candidate);
      }
  }
}

void tHandleSignals() {
  /* for any threads which have their need_to_wake flag set,
     fill in their blocked register with the signal data         */
  /* must do this for all pending and running threads            */
  struct ThreadContext *cur;
  int *CP;

  for (cur = tProcesses.Pending; cur != NULL; cur = cur->Next)
      if (cur->need_to_wake && !cur->need_to_sleep) {
          cur->hthreads[0].int_reg_file[10] = cur->signalData;
          cur->hthreads[0].empty_scoreboard |= 0x40;
          cur->need_to_wake = FALSE;
          cur->need_to_block = FALSE;
          cur->signalData = 0xdeadbeef;
          vprintf("tHandleSignals: woke pending %p\n", cur);
      }

  for (cur = tProcesses.Running; cur != NULL; cur = cur->Next)
      if (cur->need_to_wake && !cur->need_to_sleep) {
          /* need to use actual context pointer to the slot */
          CP = sysMakeCP(cur->vSlot);

          /* assume to wake up the thread at cluster 0            */
          /* cluster depends on the vSlot we are in!             */
          /* the offset 10 is used because 110 is the return
             register to be filled with the signal data          */
          CP[(0x200 * cur->vSlot) + 10] = cur->signalData;   /* write data */
          cur->need_to_wake = FALSE;
          cur->need_to_block = FALSE;
          cur->signalData = 0xdeadbeef;
          vprintf("tHandleSignals: woke running %p\n", cur);
      }
}

int tSchedule() {
  /* performs thread management                                  */

  /* jobs include taking threads off pending lists and setting them
     to run, killing threads which are on the kill list, and
     examining running threads to determine whether they have asked
     to sleep on signals.                                        */
```

```
  struct ThreadContext *candidate = NULL;
  int hardware_flags, i;
  int install_slot;

  vprintf("tSchedule ............ \n");

  /* first, kill all threads waiting to be killed */
  tCleanThreads();

  /* examine all running threads to determine whether they need to be
     put to sleep, or wakened */
  tHandleSignals();

  /* takes a thread off the pending list, and installs it in
     an available user-threadslot.  If no threadslots are available,
     attempts to evict a thread from a user-level threadslot
     and add it to the pending list. */
  candidate = tPopPending();

  if (!candidate) {
      vprintf("Nothing to do: pending list is empty\n");
      return 1;
  }

  /* we got a candidate to schedule.  First, get a free thread slot */
  install_slot = tGetSlot();
  if (install_slot == -1) {
      /* if no slots are available, we must place the candidate
         back on the pending list */
      printf("Installation not possible at this time.\n");
      printf("Waiting for a later chance.\n");
      tAddPending(candidate);
      return -1;
  } else {
      /* we should now install the candidate into the thread slot */
      vprintf("Installing context %p into vslot %d\n",
              candidate, install_slot);
      tInstall(candidate, install_slot);
      return 1;
  }
}

/* put the thread to sleep, letting it wait for a signal
   to arrive.  If need to wake is already set for the thread, it
   will be awakened immediately */
void tPutToSleep(struct ThreadContext *tc) {
  struct ThreadContext *cur, *prev;
  int status = 0;
  vprintf("tPutToSleep: putting %p to sleep\n", tc);

  if (tc->need_to_wake) {
      vprintf("tPutToSleep: %p needs to be waked already\n", tc);

      tc->need_to_sleep = FALSE;

      if (tc->vSlot == -1) {
          status = 0;
          /* if thread is NOT running, search through the pending list.
             If this thread is not there, then add it */
          for (cur = tProcesses.Pending, prev = tProcesses.Pending;
               cur != NULL; prev = cur, cur = cur->Next) {
              if (cur == tc) {
                  status = 1;
                  break;
              }
          }
```

```c
    }

    if (!status)
        tAddPending(tc);

    /* !!! In the sleep register (!!0 is the return register) */
    tc->hthreads[0].int_reg_file[10] = tc->signalData;
    tc->hthreads[0].empty_scoreboard |= 0x40; /* set !!0 full */
    tc->need_to_wake = FALSE;
    tc->need_to_block = FALSE;
    tc->signalData = 0xdeadbeef;
    Vprintf("tPutToSleep: woke non-running %p while puttosleeping\n",
            tc);

} else {
    int *cr;

    /* need to use actual context pointer to the slot */
    cr = sysMakeCP(tc->VSlot);

    /* assume to wake up the thread at cluster 0       */
    /* cluster depends on the VSlot we are in!         */
    CP((0x200 * tc->VSlot) + 10) = tc->signalData; /* write data */
    tc->need_to_wake = FALSE;
    tc->need_to_block = FALSE;
    tc->signalData = 0xdeadbeef;
    Vprintf("tPutToSleep: woke running %p while putsleeping\n", tc);

    return;
} else {
    tc->need_to_sleep = FALSE;
    tc->need_to_block = TRUE;
}

/* find and remove from the running list or pending list, since this
   thread is now considered "blocked".  The only exception is when it
   ALSO has the need_to_wake set.  In which case it is woken
   automatically (above) */

if (tc->VSlot != -1) {
    for (cur = tProcesses.Running, prev = tProcesses.Running;
         cur != NULL; prev = cur, cur = cur->Next) {
        if (cur == tc) {
            /* found the process we want to remove */
            if (tProcesses.Running == cur) {
                /* if we are already the first to be removed, stay there */
                status = 1;
            } else {
                /* if we are somewhere in the middle, remove ourselves */
                prev->Next = cur->Next;
                status = 2;
            }

            if ((tProcesses.RunningEnd == cur) && (cur != prev))
                tProcesses.RunningEnd = prev;

            /* now that we found it, instead of evicting this thread,
               which is quite expensive, we simply add it to the front of
               the running queue so it is the first popped off if an
               eviction is necessary */
            if (status == 2) {
                tc->Next = tProcesses.Running;
                tProcesses.Running = tc;
            }
            break;
        }
    }
} else {
    for (cur = tProcesses.Pending,
         prev = tProcesses.Pending;
         cur != NULL; prev = cur, cur = cur->Next) {
        if (cur == tc) {
            /* found the process we want to remove */
            if (tProcesses.Pending == cur) {
                tProcesses.Pending = cur->Next;
            } else {
                prev->Next = cur->Next;
            }

            if (tProcesses.PendingEnd == cur) {
                if (cur == prev)
                    tProcesses.PendingEnd = NULL;
                else
                    tProcesses.PendingEnd = prev;
            }
            break;
        }
    }
}

if (!status) {
    /* then the thread which is to be put to sleep is not running */
    printf("thread to be put to sleep is not running\n");
}
}
```

```c
#include <stdio.h>
#include <varargs.h>
#include "sysfunc.h"
#include "tmanager.h"
#include "tsignal.h"
#include "signaldefs.h"
#include "pointers.h"
#include "eh.h"

void kprintf(char *, ...);
void *malloc(int);
int buddyFree(void *);

int userthreadLock;
extern struct GlobalThreadState tProcesses;

#define VERBOSE_TMANAGER 0

#if VERBOSE_TMANAGER
#define Vprintf kprintf
#else
#define Vprintf
#endif

int *SYSgetSelfTC() {
    int **CP;
    int *myCP;
    struct ThreadContext *myTC;

    CP = (int **)sysSetPtr(0x77c0000005l0000);

    /* get 'own' configspace pointer */
    myCP = *CP;

    /* given my configspace ptr, I can get my context ptr */
    myTC = (struct ThreadContext *)(myCP[0x20001]);
    return (int *)sysSetPtr(((int)myTC & 0x0fffffffffffffff) | (P_KEY << 60));
}

struct ThreadContext *tSelfTC() {
    int **CP;
    int *myCP;
    struct ThreadContext *myTC;

    CP = (int **)sysSetPtr(0x77c0000005l0000);

    /* get 'own' configspace pointer */
    myCP = *CP;

    /* given my configspace ptr, I can get my context ptr */
    return (struct ThreadContext *)(myCP[0x20001]);
}

int *SYSgetParent(int *tc) {
    int *ptc;
    struct ThreadContext *retval;

    sysGetLock(&userthreadLock);
    ptc = ((int *)tc & 0x0fffffffffffffff) | (P_RW << 60);
    retval = (struct ThreadContext *)sysSetPtr(ptc);
    retval = retval->Parent;
    sysPutLock(&userthreadLock);
    return (int *)sysSetPtr(((int)retval & 0x0fffffffffffffff) | (P_KEY << 60));
}
```

```c
/* calling thread wishes to exit with returnvalue retval */
/* must
   - signal the parent thread that the child has died
   - place thread on kill list
   - signal the event handler
 */

int SYSExit(int retval) {
    struct ThreadContext *myTC;
    struct ThreadContext *prev, *cur;

    myTC = tSelfTC();
    Vprintf("tExit called by %p with returnval 0x%x\n", myTC, retval);

    /* we wish to evict ourselves from our own threadslot    */
    /* and signal a parent that we have exited               */
    /* signalling our own exit can be done first:            */

    SYStSignal(sysSetPtr(((int)myTC & 0x0fffffffffffffff) | (P_KEY << 60)),
               T_CHILD_EXIT);

    /* after we have signalled an exit to our parent, we need to inform
       the event handler that we would like to be removed.  This is done
       by removing ourselves from the running list, adding ourselves to
       the KILL list, and then signalling for a schedule */
    sysGetLock(&userthreadLock);
    if (myTC->Parent != NULL) {
        if (sysGPRB((void *)myTC->Parent) != sysGetNodeId()) {
            /* all of this removal may have to occur remotely */
            Vprintf("tExit: myTC->Parent (%p) is remote.\n", myTC->Parent);
        } else {
            /* remove self from list of siblings */
            if (myTC->Parent->Children) {
                for (cur = myTC->Parent->Children,
                     prev = myTC->Parent->Children;
                     cur != NULL; prev = cur, cur = cur->Sibling) {
                    if (cur == myTC) {
                        /* found the process we want to remove */
                        if (cur == prev) {
                            myTC->Parent->Children = cur->Sibling;
                        } else {
                            prev->Sibling = cur->Sibling;
                        }
                    }
                }
            }
        }
    }

    sysPutLock(&userthreadLock);
    add_eh_job(EVENT_KILL, 1, myTC);

    /* at this point we are all done and will block until killed */
}

void tAddChild(struct ThreadContext *tc, void *passed_parent) {
    /* first, acquire a lock and get own thread slot */
    struct ThreadContext *parentTC;
    struct ThreadContext *temp, *cur;

    if (!passed_parent) {
        /* if passed_parent is null, this means that the parent of
           the child is the 'current' tc, whatever that happens to be */
        parentTC = tSelfTC();

        sysGetLock(&userthreadLock);
        if (parentTC->Children == NULL) {
            parentTC->Children = tc;
```

```c
    } else {
        for (cur = parentTC->Children, temp = parentTC->Children;
             cur != NULL; cur = cur->Sibling)
            temp = cur;
        temp->Sibling = tc;
    }

    tc->Parent = parentTC;
} else {
    sysGetLock(&userthreadLock);
    tc->Parent = (struct ThreadContext *)sysSetPtr(((int)passed_parent & 0x0fffffff
[[[[[[[ | (P_RW << 60]);
    }
    Vprintf("tAddChild: set tc->Parent to %p\n", tc->Parent);
    sysPutLock(&userthreadLock);
}

void tKill(struct ThreadContext *target) {
    struct ThreadContext *cur, *prev;
    int found;

    for (cur = tProcesses.Running,
         prev = tProcesses.Running,
         found = 0;
         cur != NULL; prev = cur, cur = cur->Next) {
        if (cur == target) {
            found = 1;
            /* found the process we want to remove */
            if (tProcesses.Running == cur) {
                tProcesses.Running = cur->Next;
            } else {
                prev->Next = cur->Next;
            }

            if (tProcesses.RunningEnd == cur) {
                if (cur == prev)
                    tProcesses.RunningEnd = NULL;
                else
                    tProcesses.RunningEnd = prev;
            }
        }
    }

    if (!found) {
        for (cur = tProcesses.Pending,
             prev = tProcesses.Pending,
             found = 0;
             cur != NULL; prev = cur, cur = cur->Next) {
            if (cur == target) {
                found = 1;
                /* found the process we want to remove */
                if (tProcesses.Pending == cur) {
                    tProcesses.Pending = cur->Next;
                } else {
                    prev->Next = cur->Next;
                }

                if (tProcesses.PendingEnd == cur) {
                    if (cur == prev)
                        tProcesses.PendingEnd = NULL;
                    else
                        tProcesses.PendingEnd = prev;
                }
            }
        }
    }

    if (!found) {
        kprintf("ERR: critical error (KILL %p not found running or pending\n", target);
    } else {
        tAddKill(target);
    }
}
```

```c
/* ...................................................
 * tsignal.c: Core M-Machine runtime system thread signalling code  *
 * Written by:                                     Yevgeny Gurevich  *
 *                                                      7/1995       *
 * ...................................................*/

#include <stdio.h>
#include <varargs.h>
#include "manager.h"
#include "tsignal.h"
#include "signaldefs.h"
#include "pointers.h"
#include "sysfunc.h"
#include "sh.h"

void kprintf(char *, ...);
int sysSendSignal(int, int *, int);
int sysSendWakeup(int, struct ThreadContext *, int);
int sysSendSleep(void *, struct ThreadContext *, int);
int sysCPHR(void *);
int sysGetModeId();

#define SIGTAB_SIZE 8
#define VERBOSE_SIGNALLER 0

#if VERBOSE_SIGNALLER
#define Vprintf kprintf
#else
#define Vprintf
#endif

/* lock for concurrent access to signal table */
int SigTab_lock = 0;

typedef struct se {
    struct se *next;
    int signal_word;
    int signal_data;
    struct ThreadContext *sleeper;
} signal_entry;

typedef struct signal_top_entry {
    signal_entry *entry;
    signal_entry *last;
} STE;

STE signal_hash_table[SIGTAB_SIZE];

/* Initialize the on-node signal table */
void SigTab_init() {
    int i;

    for (i = 0, i < SIGTAB_SIZE; i++) {
        signal_hash_table[i].entry = NULL;
        signal_hash_table[i].last = NULL;
    }
}

void signal_entry_unparse(entry)
signal_entry *entry;
{
    printf("\t[<?p> %lx %lx]",
        entry->sleeper,
        entry->signal_word,
        entry->signal_data);

void SigTab_unparse() {
    int i;
    signal_entry *entry;

    for (i = 0; i < SIGTAB_SIZE; i++) {
        if(signal_hash_table[i].entry)
            print("<?d>: ", i);
        for (entry = signal_hash_table[i].entry; entry != NULL;
            entry = entry->next) {
            signal_entry_unparse(entry);
            print("\n");
        }
    }
}

/* calculate a hash function */
int SigTab_calchash(int value) {
    int result = 0;
    value >>= 10;
    result = (value ^ 52);
    return result;
}

/* add a new entry to the signal hash table */
/* if sleeper is NULL, this is a dormant signal, otherwise it's a sleeper */
/* returns 1 on success, 0 on failure */
int SigTab_add_entry(int signal_word, int signal_data,
                     struct ThreadContext *sleeper) {
    int j;
    signal_entry *new;

    j = SigTab_calchash(signal_word);
    j &= SIGTAB_SIZE;

    new = (signal_entry *)malloc(sizeof(signal_entry));

    if (new) {
        new->next = NULL;
        new->signal_word = signal_word;
        new->signal_data = signal_data;
        new->sleeper = sleeper;
    } else {
        kprintf("signal_table_insert: could not allocate new entry\n");
        return 0;
    }

    if (signal_hash_table[j].entry) {
        /* add to existing chain */
        signal_hash_table[j].last->next = new;
        signal_hash_table[j].last = new;
    } else {
        signal_hash_table[j].entry = new;
        signal_hash_table[j].last = new;
    }

    return 1;
}

signal_entry* SigTab_pop_entry(int signal_word, int signal_data,
                     int flag) {
    /* if flag is 0, means looking for a real sleeper,
       if flag is 1, means looking for a dormant signal */
    int j;
```

```c
    signal_entry *cur, *prev;

    j = SigTab_calchash(signal_word);
    i = SIGTAB_SIZE;

    for (cur = signal_hash_table[j].entry, prev = cur;
         cur != NULL; prev = cur, cur = cur->next) {
        if ((cur->signal_word == signal_word) &&
            ((flag && (cur->sleeper != NULL) &&      /* looking for sleeper */
              ((cur->signal_data & signal_data) ||   /* mask match */
               !cur->signal_data)) ||                /* 0 - all match */
             (flag && (cur->sleeper == NULL) &&      /* looking for dormant sig */
              ((cur->signal_data & signal_data) ||   /* mask match */
               !signal_data)))) {                    /* 0 - all match */
            if (cur == prev) {
                signal_hash_table[j].entry = cur->next;
                cur->next = NULL;
            } else {
                prev->next = cur->next;
                cur->next = NULL;
            }

            if (signal_hash_table[j].last == cur) {
                if (cur == prev) {
                    signal_hash_table[j].last = signal_hash_table[j].entry;
                } else {
                    signal_hash_table[j].last = prev;
                }
            }

            return cur;
        }
    }
    return NULL;
}

int SYSSignal(void *signal_word, int signal_data) {
    /* perform a signal on signal_word with the data signal_data */
    /* If no sleepers on signal_data are found, make a dormant entry.
       Otherwise, wake all sleepers */
    signal_entry *entry;
    int signal_int, found_one = FALSE;
    int signal_word_home;

    /* need to check here for the protections of signal_word */
    if (((((int)signal_word >> 60) & 0xf) != P_KEY) {
        kprintf("**** tSignal: invalid (non-key) signal word passed (%p)\n",
                signal_word);
        return -1;
    }

    /* If the signal_word does not map to our own node, then we
       must send a signal message to the home node of this signal. */
    if ((signal_word_home = sysGPRH((void*)signal_word)) != sysGetNodeId()) {
        /* not our node, so send a signal message to it */
        sysSendSignal(signal_word_home, signal_word, signal_data);
        return 1;
    }

    sysGetLock(&SigTab_lock);

    signal_int = (int)signal_word & 0x003fffffffffffff;

    /* first, find if any targets are already sleeping and wake each
       one with the signal - that is LOOK FOR SLEEPERS */
    while ((entry = SigTab_pop_entry(signal_int, signal_data, 0)) != NULL) {
        /* found the context which was sleeping.
           This means that we need to wake it
           - copy appropriate data into its context
           - set its status to 'woken'
           - move it to pending list. */
        Vprintf("tSignal(%p): found sleeper %p\n",
                signal_word, entry->sleeper);

        /* now, if the sleeper is a remote thread, we need to wake it
           by sending a message */
        if ((signal_word_home = sysGPRH((void*)entry->sleeper)) !=
            sysGetNodeId()) {
            /* not our node, so send a signal message to it */
            /* the danger here is that we get stuck sending messages
               with the threadlock locked.  It is possible to lock up
               the event system this way... */
            sysPutLock(&SigTab_lock);
            sysSendWake(signal_word_home, entry->sleeper, signal_data);
            sysGetLock(&SigTab_lock);
            buddyFree(entry);
            found_one = TRUE;
        } else {
            add_eh_job(EVENT_WAKE, 2, entry->sleeper, signal_data);
            buddyFree(entry);
            found_one = TRUE;
        }
    }

    if (!found_one) {
        /* insert a new 'dormant' entry */
        Vprintf("no sleepers yet.  Adding dormant signal\n");
        SigTab_add_entry(signal_int, signal_data, NULL);
        Vprintf("tSignal: signal placed for %p\n", signal_word);
    }

#if VERBOSE_SIGNALLER
    SigTab_unparse();
#endif
    sysPutLock(&SigTab_lock);

    /* add a schedule event ? */
    add_eh_job(EVENT_SCHEDULE, 0);
}

int SYSWake(struct ThreadContext *tc, int signal_data) {
    Vprintf("called SYSWake for %p with 0x%lx\n", tc, signal_data);

    /* this goes in the MH->EH job queue */
    add_eh_job(EVENT_WAKE, 2, tc, signal_data);
}

int SYSSleepRemote(void *signal_word,
                   struct ThreadContext *tc,
                   int data_mask,
                   int *return_val_location) {
    /* a sleep request from a remote node has come in */
    /* return 1 if found a dormant signal, return 0 if no signal found yet */
    int signal_int;
    signal_entry *entry;

    Vprintf("tSleepRemote(%p, %p, %lx)\n", signal_word, tc, data_mask);

    sysGetLock(&SigTab_lock);

    signal_int = (int)signal_word & 0x003fffffffffffff;
```

```
    /* first, find if a dormant signal exists */
    entry = SigTab_pop_entry(signal_int, data_mask, 1);
    if (entry) {
        /* found the dormant signal */

        *return_val_location = entry->signal_data;
        buddyFree(entry);

        Vprintf("tSleepRemote: found dormant signal (word 0x%lx)\n",
                signal_int);
#if VERBOSE_SIGNALLER
        SigTab_unparse();
#endif

        sysPutLock(&SigTab_lock);
        add_eh_job(EVENT_SCHEDULE, 0);
        return 1;
    } else {
        /* we decide to sleep on the signal_word until we are wakened */
        Vprintf("tSleepRemote: adding new signal entry for self (%p)\n", tc);
        SigTab_add_entry(signal_int, data_mask, tc);

        /* now we must do something difficult - ask to be put to sleep */
        /* since we are remote, we have already asked to be put to sleep,
           so we don't do this */
        sysPutLock(&SigTab_lock);
        return 0;
    }
}

/* put self to sleep, sleeping on signal_word, accepting signals which
   intersect with our mask, and returning the signal data, if any */
int SysSleep(void *signal_word, int data_mask) {
    signal_entry *entry;
    int signal_int;
    int return_val;
    int *CP;
    int *myCP;
    struct ThreadContext *tc;

    /* need to check here for the protections of signal_word */
    if ((((int)signal_word >> 60) & 0xf) != P_KEY) {
        Vprintf("*** tSleep: invalid (non-key) signal word passed (%p)\n",
                signal_word);
        return -1;
    }

    tc = tSelfTC();

    Vprintf("tSleep(%p, %lx) called by %p\n",
            signal_word, data_mask, tc);

    /* need to check if the signal_word is remote.  If it is, need to
       send a sleep signal to the home node of the signal word */
    if (sysGPRNB(signal_word) != sysGetNodeId()) {
        tc->need_to_sleep = TRUE;
        sysSendSleep(signal_word, tc, data_mask);

        /* block ?? */
        return sysSignalSleep(EVENT_SLEEP, tc);
    }

    sysGetLock(&SigTab_lock);
```

```
    signal_int = (int)signal_word & 0x003fffffffffffff;

    /* first, find if a dormant signal exists */
    entry = SigTab_pop_entry(signal_int, data_mask, 1);
    if (entry) {
        /* found the dormant signal */

        return_val = entry->signal_data;
        buddyFree(entry);

        Vprintf("tSleep: found dormant signal (word 0x%lx)\n", signal_int);
#if VERBOSE_SIGNALLER
        SigTab_unparse();
#endif

        sysPutLock(&SigTab_lock);
        add_eh_job(EVENT_SCHEDULE, 0);
        return return_val; /* no blocking is necessary since the signal
                              data is already available to us.  We simply
                              return with it */
    } else {
        /* we decide to sleep on the signal_word until we are wakened */
        Vprintf("tSleep: adding new signal entry for self (%p)\n", tc);
        SigTab_add_entry(signal_int, data_mask, tc);

        /* now we must do something difficult - ask to be put to sleep */
        tc->need_to_sleep = TRUE;
        sysPutLock(&SigTab_lock);

        /* when sysSignalSleep returns, we WILL have been wakened */
        /* also have to worry about locking somewhat */
        /* since we have unlocked the sigtable and threadlock, it is
           possible for someone to already have set tc->need_to_wake
           at the time we go to sleep.  In fact, since at this point
           someone may have already set need_to_wake AND the scheduler
           have come in and automatically filled our 'sleep' register,
           we cannot simply signal right now.  What we have to prepare
           is a flag which tells the handler that we WILL signal them
           in the future so don't touch us just yet.          */
        return sysSignalSleep(EVENT_SLEEP, tc);
    }
}
```

# Appendix E

# Sample User Programs

This chapter contains the source code for two user-level programs which make calls on MARS primitives. Matmul1.c is a parallel-matrix-multiply program. Jacoby6.c is a version of an iterative jacobian-matrix relaxation.

```c
#include <stdio.h>

#define MATSIZE 4
#define HSPAWN  0
#define TSPAWN  1
#define GCC_TEST 0

static int matrix1[4][4] = {
  { 0, 3, 4, 5 },
  { 5, 6, 2, 4 },
  { 6, 3, 87, 46 },
  { 5, 8, 33, 64)
};

static int matrix2[4][4] = {
  { 12, 45, 92, 4 },
  { 6, 82, 36, 75 },
  { 9, 61, 11, 6 },
  { 5, 2, 4, 3 }
};
static int final[4][4];

/* vectors for synchronization */
static int donevec[4] = { 0, 0, 0, 0};
static int donevec1[4] = { 0, 0, 0, 0};

int tspawn(int numargs, void *threadip, int dest_node, ...);
int hspawn(int numargs, void *threadip, int dest_node, ...);

void runthread(int row[4], int column) {
  /* performs an individual rox x column calculation */
  int i,j, total;
  for (j = 0; j < MATSIZE; j++) {
    total = 0;
    for (i = 0; i < MATSIZE; i++)
      total += row[i]*matrix2[i][j];
    printf("total = %d\n", total);
    final[column][j] = total;
  }
  donevec[column] = 1;
}

void print_matrix(int m[4][4]) {
  int i,j;
  for (i = 0; i < MATSIZE; i++) {
    for (j = 0; j < MATSIZE; j++)
      printf("%5d", m[i][j]);
    printf("\n");
  }
  printf("\n\n");
}

int main() {
  int i, j;

  print_matrix(matrix1);
  print_matrix(matrix2);

#if GCC_TEST
  /* testing linear code without threads */
  for (i = 0; i < MATSIZE; i++) {
    runthread(matrix1[i], i);
  }

  for (i = 0; i < MATSIZE; i++) {
    runthread(matrix1[i], i);
  }

  for (i = 0; i < MATSIZE; i++) {
    runthread(matrix1[i], i);
  }

  for (i = 0; i < MATSIZE; i++) {
    runthread(matrix1[i], i);
  }

#else
#if (HSPAWN)
  /* spawn H-Threads on local node */
  i = hspawn(2, runthread, 1, matrix1[0], 0);
  if (!i) printf("hspawn 1 failed\n");
  i = hspawn(2, runthread, 2, matrix1[1], 1);
  if (!i) printf("hspawn 1 failed\n");
  i = hspawn(2, runthread, 3, matrix1[2], 2);
  if (!i) printf("hspawn 1 failed\n");
  runthread(matrix1[3], 3);

#else
  /* spawn V-Threads over three nodes */
  tspawn(2, runthread, 1, matrix1[0], 0);
  tspawn(2, runthread, 1, matrix1[1], 1);
  tspawn(2, runthread, 2, matrix1[2], 2);
  tspawn(2, runthread, 3, matrix1[3], 3);
#endif
#endif

  /* perform in-memory barrier synchronization */
  while(!donevec[0] ||
        !donevec[1] ||
        !donevec[2] ||
        !donevec[3]) ;
  print_matrix(final);
  return 1;
}
```

```c
#include <stdio.h>
#include "syscalls.h"
#include "tsignal.h"

#ifdef SMALL_TEST
#define MAXX 4
#define MAXY 4
#else
#define MAXX 10
#define MAXY 10
#endif

#define TSPAWN 1

#ifdef SMALL_TEST
int matrix[MAXX][MAXY] = {
  { 91, 81, 81, 54 },
  { 57, 51, 98, 95 },
  { 54, 96, 55, 98 },
  { 77, 95, 56, 87 }};
#else
int matrix[MAXX][MAXY] = {
  { 91, 81, 54, 98, 97, 55, 86, 83, 69 },
  { 57, 51, 98, 95, 69, 57, 84, 85, 69, 50 },
  { 54, 96, 55, 98, 87, 77, 86, 69, 47, 25 },
  { 77, 95, 56, 87, 63, 87, 69, 44, 27, 39 },
  { 85, 98, 98, 57, 88, 69, 40, 31, 46, 57 },
  { 81, 54, 97, 86, 69, 34, 37, 54, 58, 39 },
  { 96, 98, 77, 69, 25, 46, 59, 39, 22, 53 },
  { 98, 57, 69, 31, 57, 39, 37, 36, 47, 7 },
  { 98, 69, 46, 39, 53, 47, 49, 9, 50 },
  { 69, 39, 47, 49, 50, 7, 15, 20, 39, 43 }};
#endif

int matrix2[MAXX][MAXY];

void print_matrix(int mat[MAXX][MAXY]) {
  int i, j;

  for (i = 0; i < MAXX; i++) {
    for (j = 0; j < MAXY; j++) {
      printf("%3d", mat[i][j]);
    }
    printf("\n");
  }
}

void calc_row(int row, int mat[MAXX][MAXY], int dest[MAXX][MAXY]) {
  int i;

  for (i = 0; i < MAXY; i++) {
    if (!i) {
      if (!row) {
        dest[row][i] =   ((0 +
                           100 +
                           mat[row][i + 1] +
                           mat[row + 1][i]) / 4.0);
      } else if (row == (MAXX - 1)) {
        dest[row][i] =   ((0 +
                           mat[row - 1][i] +
                           mat[row][i + 1] +
                           100) / 4.0);
      } else {
        dest[row][i] =   ((0 +
                           mat[row - 1][i] +
                           mat[row][i + 1] +
                           mat[row + 1][i]) / 4.0);
      }
    } else if (i == (MAXY - 1)) {
      if (!row) {
        dest[row][i] =   ((mat[row][i - 1] +
                           100 +
                           100 +
                           mat[row + 1][i]) / 4.0);
      } else if (row == (MAXX - 1)) {
        dest[row][i] =   ((mat[row][i - 1] +
                           mat[row - 1][i] +
                           100 +
                           100) / 4.0);
      } else {
        dest[row][i] =   ((mat[row][i - 1] +
                           mat[row - 1][i] +
                           100 +
                           mat[row + 1][i]) / 4.0);
      }
    } else {
      if (!row) {
        dest[row][i] =   ((mat[row][i - 1] +
                           100 +
                           mat[row][i + 1] +
                           mat[row + 1][i]) / 4.0);
      } else if (row == (MAXX - 1)) {
        dest[row][i] =   ((mat[row][i - 1] +
                           mat[row - 1][i] +
                           mat[row][i + 1] +
                           100) / 4.0);
      } else {
        dest[row][i] =   ((mat[row][i - 1] +
                           mat[row - 1][i] +
                           mat[row][i + 1] +
                           mat[row + 1][i]) / 4.0);
      }
    }
  }
}

/* child thread code which calculates values for an individual
   matrix row */
void run_thread(int row) {
  int i;
  int avg;
  int stage;

  for (stage = 0; stage < 15; stage++) {
    /* depending on which iteration we are in, we use either
       matrix or matrix2 as the source, and the other matrix as
       the destination */
    if (stage % 2)
      calc_row(row, matrix2, matrix);
    else
      calc_row(row, matrix, matrix2);

    /* perform barrier, and then wait for flying off again */
    tSignal(_getParent(_getSelfTC()), T_ALL_SIGNALS);

    /* wait for parent to tell us to go off again... */
    tSleep(_getParent(_getSelfTC()), T_ALL_SIGNALS);
    printf("child %d waking...\n", row);
  }
```

```
}

/* vector of children which are spawned */
void *childvector[] = ( NULL, NULL, NULL, NULL, NULL,
                        NULL, NULL, NULL, NULL, NULL);

int main() {
    int iteration;
    printf("Jacoby Starting\n");
    printf("Matrix = %p, matrix2 = %p\n", matrix, matrix2);

/* in the small test, we spawn off 4 threads.  In the large, 10 */
#ifdef SMALL_TEST
    childvector[0]  =  tspawn(1, run_thread, 1, 0);
    childvector[1]  =  tspawn(1, run_thread, 1, 1);
    childvector[2]  =  tspawn(1, run_thread, 1, 2);
    childvector[3]  =  tspawn(1, run_thread, 1, 3);
#else
    childvector[0]  =  tspawn(1, run_thread, 1, 0);
    childvector[1]  =  tspawn(1, run_thread, 1, 1);
    childvector[2]  =  tspawn(1, run_thread, 1, 2);
    childvector[3]  =  tspawn(1, run_thread, 2, 3);
    childvector[4]  =  tspawn(1, run_thread, 2, 4);
    childvector[5]  =  tspawn(1, run_thread, 2, 5);
    childvector[6]  =  tspawn(1, run_thread, 3, 6);
    childvector[7]  =  tspawn(1, run_thread, 3, 7);
    childvector[8]  =  tspawn(1, run_thread, 3, 8);
    childvector[9]  =  tspawn(1, run_thread, 1, 9);
#endif
    for (iteration = 1; iteration < 15; iteration++) {

        /* for each iteration, the parent sleeps until all
           of the children have completed their calculation.
           Then it prints out the result, and tells the
           children to go on, by signalling them */

        printf("Iteration ----- %d -----\n", iteration);
#ifdef SMALL_TEST
        tSleep(childvector[0],  T_ALL_SIGNALS);
        tSleep(childvector[1],  T_ALL_SIGNALS);
        tSleep(childvector[2],  T_ALL_SIGNALS);
        tSleep(childvector[3],  T_ALL_SIGNALS);
#else
        tSleep(childvector[0],  T_ALL_SIGNALS);
        tSleep(childvector[1],  T_ALL_SIGNALS);
        tSleep(childvector[2],  T_ALL_SIGNALS);
        tSleep(childvector[3],  T_ALL_SIGNALS);
        tSleep(childvector[4],  T_ALL_SIGNALS);
        tSleep(childvector[5],  T_ALL_SIGNALS);
        tSleep(childvector[6],  T_ALL_SIGNALS);
        tSleep(childvector[7],  T_ALL_SIGNALS);
        tSleep(childvector[8],  T_ALL_SIGNALS);
        tSleep(childvector[9],  T_ALL_SIGNALS);
#endif
        printf("barrier reached\n");
        if (iteration % 2)
            print_matrix(matrix2);
        else print_matrix(matrix);

        tSignal(_getSelfTC(), 1);

    }
}
```

# Bibliography

[1] Nicholas P. Carter, Stephen W. Keckler, and William J. Dally. Hardware support for fast capability-based addressing. In *Proceedings of the Sixth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS VI)*, pages 319–327. Association for Computing Machinery Press, October 1994.

[2] Jeffrey S. Chase, Henry M. Levy, Miche Baker-Harvey, and Edward D. Lazowska. How to use a 64-bit virtual address space. Technical Report 92-03-12, University of Washington, 1992.

[3] Thomas H. Cormen, Charles E. Leiserson, and Ronald L. Rivest. *Introduction to Algorithms*. MIT Press, Cambridge, Massachusetts, 1993.

[4] William J. Dally, Stephen W. Keckler, Nick Carter, Andrew Chang, Marco Fillo, and Whay S. Lee. M-Machine architecture v1.0. Concurrent VLSI Architecture Memo 58, Massachusetts Institute of Technology, Artificial Intelligence Laboratory, January 1994.

[5] William J. Dally, Stephen W. Keckler, Nick Carter, Andrew Chang, Marco Fillo, and Whay S. Lee. The MAP instruction set reference manual v1.3. Concurrent VLSI Architecture Memo 59, Massachusetts Institute of Technology, Artificial Intelligence Laboratory, February 1995.

[6] Abraham Silberschatz, James L. Peterson, , and Peter B. Galvin. *Operating System Concepts*. Addison-Wesley, Reading, Massachusetts, third edition, 1992.

[7] Andrew S. Tannenbaum. *Modern Operating Systems.* Prentice Hall, Englewood Cliffs, NJ, 1992.