

# The Yalnix Kernel

## CS230/330 Operating Systems, Winter 2002

Handed out : January 7, 2002  
Due : March 8, 2002

## 1 Preface

In this project, you will be designing and implementing a simple operating system kernel. This document provides an overview of the entire kernel project. At first glance, this is going to appear very daunting if not impossible. However, keep in mind that the project is divided into several well-defined sub-projects that follow the material presented in the lectures. Furthermore, your instructor and TAs all have experience working on prior versions of the kernel project. Therefore, do not panic if you do not understand all of this document after a first reading—it will all become clear in due course. In the meantime, you might just want to sleep with this document under your pillow.

## 2 Introduction

Having lost 99% of its stock price, the Engone Gigacorp of Houston, Texas has decided that it will be able to regain its former glory by introducing a new handheld game console and a whole lineup of games centered around the lovable characters from the hit game “Extreme Texas-42.” However, having blown 147 million dollars on track lighting, Aeron chairs, latte machines, and CASE tools, only \$1635.23 is actually available to build a prototype. Therefore, the contract to build the prototype, was awarded to iBubba.com consulting (ISO-9000 certified) from Amarillo. Not surprisingly, the name of this prototype has become known as the “GameBubba.”

Your task is to write the operating system for the GameBubba prototype, which in honor of the President you have decided to call “Yalnix” (a horrible Texas-based pun on “Unix”). Since the real hardware has yet to be fabricated, your operating system will run a machine simulator that has been developed for Sun SPARC Solaris systems. Yalnix includes the main process and memory management features of Unix, plus support for a variety of I/O devices including a serial port, a graphical display device, a pen-based input device, and a disk drive. Although your kernel needs to talk to a variety of low-level devices, it will not be necessary to program any portions of your project in assembly language nor will you need to worry about any low-level CPU-specific details. However, you will need to implement process management, memory management, process scheduling, I/O, demand paging, and interprocess communication.

## 3 Project Overview

You will be implementing a fully functional operating system kernel from scratch. Although your kernel will run on simulated hardware, you will find that environment provides the same challenges as you would find on a real machine. Furthermore, you will be responsible for the design and implementation of the *entire* kernel (i.e., you will not be given any skeleton code, nor will you be writing small components for an existing kernel). Because of this, you have the freedom to implement your kernel however you see fit as long as it runs on the emulated hardware and provides the functionality described in this document.

Your kernel will have to address three general problems. First, it needs to control the low-level hardware. To do this, you will have to write a collection of low-level device drivers. Second, your kernel needs to provide the abstraction of a “process.” A process is simply a running program. Therefore, your kernel will need to figure out how to allocate memory between processes and to schedule their execution on the CPU. Finally, your kernel must provide a system call interface to user programs so that they can request the services of your operating system.

The rest of this document describes various aspects of the system in detail. As the project is completed, you will receive more handouts containing additional information and implementation tips.

## 4 The GameBubba Hardware

The GameBubba system currently consists of 4 Mbytes of DRAM, an 16 Mbyte disk drive, a serial port, an LCD display, and a pen input device. All of these devices are controlled by special memory-mapped control ports. Furthermore, some of the devices use DMA (Direct Memory Access) to transfer data. In addition, the hardware generates a number of traps and interrupts that your kernel must handle in response to various events.

The file `/yalnix/include/hardware.h` contains all of the constants and data structures needed to control the hardware in your kernel implementation. You may want to look at this file as you read through this section.

### 4.1 Access to the Raw Hardware

The hardware is controlled by reading and writing values to a collection of memory-mapped control ports. All this really means is that each hardware device (e.g., serial port, disk, etc.) monitors the memory bus and responds to load/store operations that match a certain range of memory addresses. The `hardware.h` file defines symbolic names for these special memory addresses as shown in the following table:

Control Port	Address	Description	Readable	Writable
MEM_VECTOR_BASE	0x3000000	Interrupt vector Address	Yes	Yes
MEM_PHYSM	0x3000004	Size of physical memory	Yes	No
MEM_CONSOLE	0x3000008	Write character to console	Yes	Yes
MEM_VM_ENABLE	0x3002000	Enable virtual memory	Yes	Yes
MEM_TLB_FLUSH	0x3002004	Flush the TLB	No	Yes
MEM_PT0_BASE	0x3002008	Region 0 page table base	Yes	Yes
MEM_PT0_SIZE	0x300200c	Region 0 page table size	Yes	Yes
MEM_PT1_BASE	0x3002010	Region 1 page table base	Yes	Yes
MEM_PT1_SIZE	0x3002014	Region 1 page table length	Yes	Yes
MEM_PEN_EVENT	0x3004000	Pen event port	Yes	No
MEM_SERIAL_CS	0x3006000	Serial status and control port	Yes	Yes
MEM_SERIAL_ADDR	0x3006004	Serial address	No	Yes
MEM_DISK_CS	0x3008000	Disk control port	No	Yes
MEM_DISK_ADDR	0x3008004	Disk address port	No	Yes
MEM_DISK_NBLK	0x3008008	Number of blocks on disk device	Yes	No
MEM_DISPLAY_CTRL	0x300a000	Display control port	No	Yes

In `hardware.h`, constants such as `MEM_VECTOR_BASE` are defined as follows:

```
#define MEM_VECTOR_BASE 0x3000000
```

To read or write from a port, you need to use the C assignment operator. For example, this code reads a value from a port:

```
unsigned int value = *((unsigned int *)MEM_PEN_EVENT);
```

To write to a port, you use assignment in the other direction like this:

```
*((unsigned int *) MEM_VECTOR_BASE) = (unsigned int) value;
```

All of the memory-mapped ports are 32 bits wide. In C, the value is represented as a type of `unsigned int`. You must use a C cast to convert other data types (such as pointers) to type `unsigned int` when writing to a port and you must also use a cast to convert the value read from a port to the desired type if you need to interpret the value as anything other than an `unsigned int`.

To reduce the amount of typing and make your code easier to read, `hardware.h` defines a few macros that you can use to access the ports. For example:

```
unsigned int value = ReadPort(MEM_PEN_EVENT);
...
WritePort(MEM_VECTOR_BASE,value);
```

Note: `ReadPort()` and `WritePort()` are macros that expand to the code above (with all of the casts). They are not library functions.

Access to the HW ports is only allowed when the CPU is running in supervisor (kernel) mode. Your kernel will be able to access the ports, but user applications will not.

## 4.2 Traps

When a device needs attention or a program generates an error, a trap (or interrupt) is generated. The following list briefly describes all of the traps that can be generated by the hardware:

- **TRAP\_KERNEL**

This trap is generated when a user program executes a system call. See section 7.3.

- **TRAP\_CLOCK**

This is a periodic timer interrupt generated by the machine's hardware clock. Your kernel will use this to implement time-sharing and to context-switch between different programs. See Section 7.2.

- **TRAP\_ILLEGAL**

This trap occurs when a user program tries to execute an illegal CPU instruction. When it occurs, the offending process should be terminated.

- **TRAP\_MEMORY**

This trap results when a user process tries to access a memory address that is not mapped in any of the hardware page tables. It can also occur if a user process tries to access a portion of memory for which it has insufficient access permissions (if a user program tries to access a privileged I/O port for instance). See Section 5.1.

- **TRAP\_MATH**

This trap occurs when a math operation executed by a user program generates a fault. For example, division by zero.

- **TRAP\_SERIAL**

This interrupt is generated by the serial device when it receives data or when it has finished sending data.

- **TRAP\_PEN**

This interrupt occurs whenever there is activity on the pen input device. After this interrupt has occurred, the status of the pen device should be read by reading from `MEM_PEN_EVENT`.

- **TRAP\_DISK**

This interrupt occurs whenever the last I/O operation performed on the disk storage device has been completed.

When a trap occurs, the CPU looks in an interrupt vector table to find out the location of an interrupt handler function. The above constants (e.g., `TRAP_CLOCK`) represent integer offsets into the vector table. The interrupt vector table is nothing more than an array of

pointers to functions that your kernel must initialize. To do this, the privileged machine control port `MEM_VECTOR_BASE` must be set to point to the vector table within your kernel. The vector table must be statically dimensioned by your kernel with a size of `TRAP_VECTOR_SIZE`. Only the entries defined above are used by the hardware. All other entries in the vector table must be initialized by your kernel to `NULL`.

The interrupt handler functions (the entries in the vector table) must have the following prototype:

```
void Trap_Handler(cpu_context_frame *)
```

(See section 5.3 for a description of the type `cpu_context_frame`).

Finally, it should be noted that the hardware blocks the delivery of all other interrupts until your kernel returns from the current trap (e.g., your trap handling functions will not be interrupted by other traps). This means that your kernel is generally not interruptible and that you will not need to implement kernel locking and synchronization primitives (this is good, be thankful).

## 4.3 Memory Management Hardware

The hardware features a virtual memory system that is managed by a special MMU controller. The MMU manages physical memory and is used to provide programs with the illusion of a large address space.

### 4.3.1 Physical Memory (DRAM)

Physical memory begins at the address defined by `PMEM_BASE` which is found in `hardware.h`. Although the current version of the hardware provides only 4 Mbytes of DRAM, the actual amount of memory is supplied to the kernel during the boot process and may vary in future revisions of the hardware. The size of physical memory can be read from `MEM_PHYSM`.

Logically, the physical memory is divided into a collection of pages of size `PAGESIZE`. These pages are referred to by numbers known as *page frame numbers*. To better illustrate the idea, the physical layout of the machine looks something like that shown in Figure 1.

In the figure, the 4 Mbytes of physical memory begins at physical memory address `0x1000000` (page frame 2048) and extends to address `0x13fffff` (page frame 2559). Physical memory addresses lower than `PMEM_BASE` are used by the boot ROM other hardware devices. These addresses should not be used by your kernel. Similarly, attempts to access physical memory above the top of available memory will result in undefined hardware behavior (and a possible hardware fault).

In order to run concurrent user programs, your kernel needs to keep track of the physical memory pages that are in use and those that are available. To do this, you might just create an array of 512 entries and record a 0 or 1 to indicate usage. Of course, there are more elegant (and more preferred) ways to do this. Also, the amount of physical memory may change as the hardware is improved so it would be a bad idea to hard-wire a fixed sized array into your kernel except for initial testing.

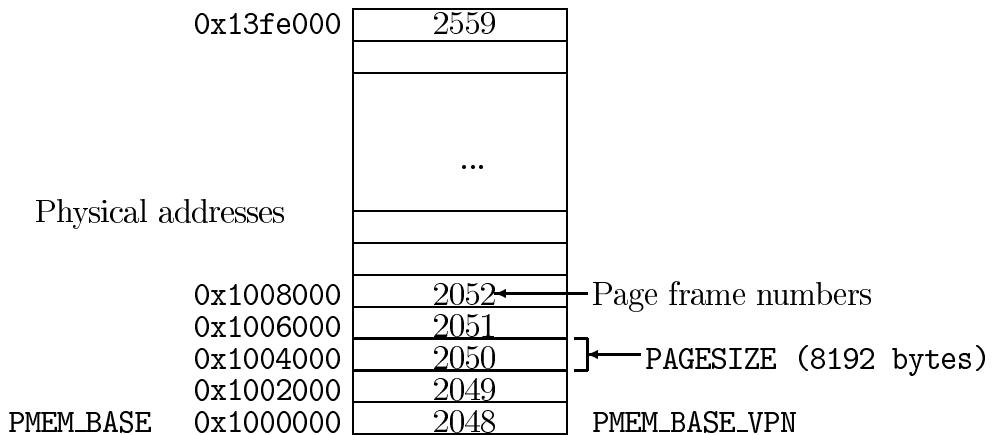


Figure 1: Physical memory layout

#### 4.3.2 Virtual Memory

Although the kernel, user programs, and data are stored in physical memory, every program runs within a virtual memory space that is managed by a special MMU controller. This controller provides a simple paged memory scheme that relies upon the use of page tables.

The virtual address space of the machine is divided into two *regions* called region 0 and region 1. By convention, region 0 is used by the operating system kernel and region 1 is used by user processes executing on the Yalnix system. Region 0 will be managed by your kernel so that it is shared by all processes. Region 1 will be used to map the currently running user process and will be changed by the OS whenever it performs a context switch. The hardware provides these two regions so that it can protect kernel data structures from illegal tampering by user applications. When the CPU is executing in privileged mode, references to both regions are allowed. However, when the CPU is not executing in privileged mode, references to region 0 are forbidden. The hardware does all of this automatically without any intervention from your kernel.

The virtual address space and its regions are defined by a number of constants located in the `hardware.h` file.

- **VMEM\_0\_BASE:** The base address of virtual memory region 0.
- **VMEM\_0\_SIZE:** The size of virtual memory region 0.
- **VMEM\_0\_LIMIT:** The lowest address not part of virtual memory region 0.
- **VMEM\_1\_BASE:** The base address of virtual memory region 1.
- **VMEM\_1\_SIZE:** The size of virtual memory region 1.
- **VMEM\_1\_LIMIT:** The lowest address not part of virtual memory region 1.

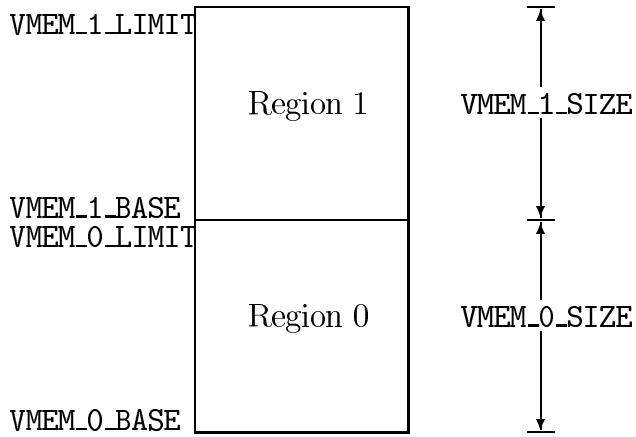


Figure 2: Virtual Memory Layout

Virtual addresses greater than or equal to `VMEM_0_BASE` and less than `VMEM_0_LIMIT` are in region 0, and virtual addresses greater than or equal to `VMEM_1_BASE` and less than `VMEM_1_LIMIT` are in region 1. All other virtual addresses are illegal. Graphically, the structure of the virtual address space is shown in Figure 2.

Like physical memory, the virtual memory regions are divided into pages, each of size `PAGESIZE`. These pages are accessed by numbers known as *virtual page numbers*. The size of the virtual address space is not necessarily the same as physical memory. In fact, in the current implementation, each virtual memory region provides an address space of 16 Mbytes (2048 pages) which is more than enough to accomodate the 4 Mbytes of physical memory.

It is important to note that the division of virtual memory into regions has nothing to do with the organization of the physical memory nor are the regions considered to be “segments.” Thus, you are not going to be dividing the physical memory into separate contiguous regions like this. In fact, each page of physical memory can be assigned arbitrarily to either virtual memory region as needed (and a physical memory page might even be mapped to multiple locations within the virtual address space in special circumstances).

#### 4.3.3 Address Manipulation

Throughout the implementation of your kernel, you will need to manipulate memory addresses, virtual page numbers, and physical page frames. Furthermore, you may need to convert between the various representations. To do this, a number of constants and macros are defined in `hardware.h` for your use.

- `PAGESIZE`

This is the size (in bytes) of a memory page in both virtual and physical memory.

- `PAGEOFFSET`

This is a bit mask that can be used to extract the offset of a memory address within a page. To use it, simply compute `addr & PAGEOFFSET`.

- **PAGESHIFT**

This is a constant that contains the number of bits that an address must be shifted to the right in order to obtain its page number or the number of bits that a page number must be shifted to the left to obtain its base address. For example, if  $vp1$  is a virtual address located in Region 1,  $(vp1 - \text{VMEM\_1\_BASE}) \gg \text{PAGESHIFT}$  would compute the page number within that region.

- **PAGEMASK**

This is a bit mask that can be used to extract the base-address of a page given a memory address. Use `addr & PAGEMASK` to get the base address.

- **UP\_TO\_PAGE( $x$ )**

Rounds the address  $x$  up the next highest page boundary. If  $x$  is on a page boundary, it returns  $x$ .

- **DOWN\_TO\_PAGE( $x$ )**

Rounds the address  $x$  down to the next lowest page boundary. If  $x$  is on a page boundary, it returns  $x$ .

All of the above constants and macros work perfectly well with both virtual and physical addresses. However, the interpretation of the results may vary depending on the region. For example, address computations often need to include the constants `VMEM_0_BASE` or `VMEM_1_BASE` depending on the memory region of interest.

#### 4.3.4 Page Tables

Whenever a reference to a virtual address is made, the MMU hardware translates it into a physical memory address. This translation occurs automatically without intervention from the kernel. However, the mapping that specifies the translation is controlled by the Yalnix kernel which is responsible for filling in the values of a single-level direct mapped page table for **each** of the virtual memory regions.

A page table is created by allocating a contiguous array of page table entries that define the virtual memory mapping for every single page in a virtual memory region. Internally, the kernel is allowed to put the page table wherever it wants as long as it tells the MMU where to find it. This is done by writing to the following special hardware ports:

- **MEM\_PT0\_BASE**: Contains the *physical* memory address of the page table for region 0 of virtual memory.
- **MEM\_PT0\_SIZE**: Contains the size of the page table for virtual memory region 0, i.e, the number of virtual pages in region 0.

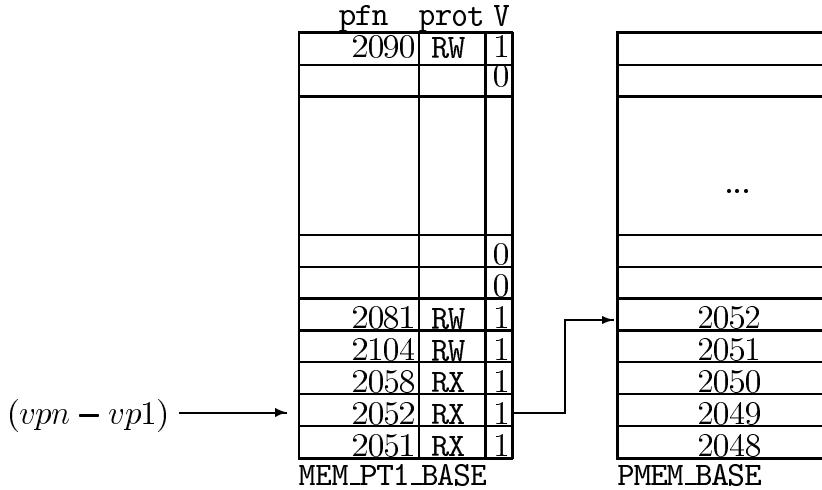


Figure 3: Page Table Lookup for Region 1

- **MEM\_PT1\_BASE:** Contains the *physical* memory base address of the page table for region 1 of virtual memory.
- **MEM\_PT1\_SIZE:** Contains the length of the page table for virtual memory region 1, i.e., the number of virtual pages in region 1.

The MMU hardware never updates these settings: only the kernel can change the virtual-to-physical address mapping.

When a program (or the kernel itself) makes a memory reference to a virtual page  $vpn$ , the MMU finds the corresponding page frame number  $pfn$  by looking in the page table for the appropriate region of virtual memory. The page tables are indexed by virtual page number, and contain entries that specify the corresponding physical page frame numbers.

For example, if the reference is to region 0 and the first virtual page number of region 0 is  $vp0$ , the MMU looks at the page table entry that is  $vpn - vp0$  page table entries above the address in **MEM\_PT0\_BASE**. Likewise, if the reference is to region 1 and the first virtual page number of region 1 is  $vp1$ , the MMU looks at the page table entry that is  $vpn - vp1$  page table entries above the address in **MEM\_PT1\_BASE**. This process is illustrated in Figure 3.

The lookup to translate  $vpn$  to its currently mapped  $pfn$  is carried out wholly in hardware. As a result, the hardware imposes a strict format on the structure of the page table entries. The file **hardware.h** contains the definition of a C data structure (**struct pte**) that has the memory layout required by the MMU. You will need to use this data structure to manipulate page table entries from your kernel.

Each page table entry is 32-bits wide, and contains the following fields:

- **valid:** (1 bit) If this bit is set, the page table entry is valid; otherwise, a memory exception is generated when/if this virtual memory page is accessed.

- **prot:** (3 bits) This field defines the memory protection applied by the hardware to this virtual memory page. The three protection bits are interpreted independently as follows:
  - PROT\_READ:Memory within the page may be read.
  - PROT\_WRITE:Memory within the page may be written.
  - PROT\_EXEC:Memory within the page may be executed as machine instructions.

Execution of instructions requires that their pages be mapped with both PROT\_READ and PROT\_EXEC.

- **access:** (1 bit) This bit is set when the hardware accesses a memory page (reading or writing). It is never cleared by the hardware—your kernel will need to do that.
- **dirty:** (1 bit) This bit is set whenever the MMU detects a write to a page. Your kernel can examine this bit to determine which pages have been modified. This bit is never cleared by the hardware.
- **user1:** (1 bit) User definable. Not used by the hardware. Your kernel can use this for its own purposes.
- **user2:** (1 bit) User definable. Not used by the hardware.
- **pfn:** (24 bits) This field contains the page frame number of the page of physical memory to which this virtual memory page is mapped by this page table entry. This field is ignored if the valid bit is off. You might also want to use this field when implementing demand paging from the disk drive.

`MEM_PT0_BASE` and `MEM_PT1_BASE` contain *physical* memory addresses. As it turns out, the page table for each virtual memory region exactly fits into a single page of physical memory. Therefore, page tables can easily be allocated by your kernel by simply grabbing a free page of physical memory and placing its address in the appropriate page table register. However, in order to modify the contents of the page tables, you will also need to map these pages someplace with in the kernel memory map.

Finally, it is important to realize that your kernel will have to manipulate a variety of different page tables as it runs. In particular, **each** user process will have its own page table that is mapped to region 1 of the virtual memory space. Whenever the kernel performs context switch to another program, it will need to change the region 1 page table by writing a new value to `MEM_PT1_BASE` register. In addition, your kernel will need to manage a region 0 page table for its own memory needs.

#### 4.3.5 The Translation Lookaside Buffer (TLB)

Most architectures contain a TLB to speed up address translation. The TLB caches address translations so that subsequent references to the same virtual page do not have to retrieve

the corresponding page table entry from memory. This is important because it can take several memory accesses to retrieve a page table entry<sup>1</sup>.

The yalnix hardware also contains a TLB that is managed by the hardware MMU. Whenever a program accesses a virtual page, the mapping of that page to a physical frame number is automatically stored in the TLB. However, in certain cases, your kernel will need to flush all or part of the TLB. In particular, after changing a page table entry in memory, the TLB might contain a stale mapping for the virtual addresses corresponding to the modified page table entry. Also, when carrying out a context switch from one process to another, the MMU will continue to map the cached entries of the previous process. Unless the TLB is flushed, this will cause virtual addresses in the new address space to be mapped to the physical pages of the old address space.

The hardware provides the `MEM_TLB_FLUSH` port to control the flushing of all or part of the TLB. When writing a value into `MEM_TLB_FLUSH`, the MMU interprets the value as follows:

- `TLB_FLUSH_ALL`: Flush the entire TLB.
- `TLB_FLUSH_0`: Flush all mappings for virtual addresses in region 0 from the TLB.
- `TLB_FLUSH_1`: Flush all mappings for virtual addresses in region 1 from the TLB.
- `addr`: Flush only the mapping for virtual address `addr` from the TLB, if present.

The symbolic constants above are all defined in `hardware.h`.

#### 4.3.6 Initializing Virtual Memory

When the machine is booted, a program in the boot PROM loads the kernel into physical memory and begins executing it. The boot PROM loads the kernel into *physical* addresses starting at `PMEM_BASE`.

At this point, virtual memory cannot yet be active, since it is the kernel's responsibility to initialize the location and values of the page tables. In order to initialize virtual memory, the hardware initially begins execution of the kernel with virtual memory disabled. The boot ROM loads the kernel into the beginning of *physical* memory, and the machine uses only *physical* addresses when accessing memory until virtual memory is explicitly enabled by the kernel. To enable virtual memory, the kernel must first write the locations of the region 0 and region 1 page tables into `MEM_PT0_BASE` and `MEM_PT1_BASE`. Then, it should store a 1 in `MEM_VM_ENABLE`. After doing this, *all* values used to address memory are interpreted *only* as *virtual* memory addresses by the MMU. Virtual memory cannot be disabled after it has been enabled.

Further details about kernel memory initialization can be found in the Section 5.2.

---

<sup>1</sup>In the case of multi-level page tables and inverse page tables, it might take many memory accesses; in the case of a demand-paged memory system, it might even take a disk access to retrieve a page table entry.

## 4.4 Disk Access

The hardware features a small disk drive that is used to store the kernel itself, user programs, and data. Internally, the disk is organized as a collection of blocks, each of size DISK\_BLOCKSIZE. Currently the size of each block is 1024 bytes (note that this is *different* than the pagesize of DRAM). The total number of disk blocks can be obtained by reading the value of `MEM_DISK_NBLK`. This value may change depending on the size of the disk a user has installed.

Access to the disk is managed through a pair of hardware ports and DMA memory transfer. Specifically, to initiate a memory transfer, it is first necessary to write a *physical memory address* into `MEM_DISK_ADDR`. Then, a disk operation is written into `MEM_DISK_CS`. The value written to `MEM_DISK_CS` has the following format:

Bits 28-31	Bits 0-27
opcode	block number

The opcode field specifies the type of transfer and is one of the following values:

- `DISK_READ`. Transfer a data block from disk to physical memory.
- `DISK_WRITE`. Transfer a data block from physical memory to disk.

For example, the following code shows how one would initiate a disk block transfer.

```
/* Transfer data from the disk drive */
char    *paddr;          /* Raw physical memory address */
int     blocknum;        /* Starting block number on the disk */

*((unsigned int *) MEM_DISK_ADDR) = (unsigned ) paddr;
*((unsigned int *) MEM_DISK_CS)   = (unsigned) (DISK_READ<<28) | blocknum;
...
```

Memory transfers involving the disk storage device must observe proper memory alignment. Specifically, the physical memory address written to `MEM_DISK_ADDR` must be a multiple of `DISK_BLOCKSIZE`. Upon completion of the memory transfer, a `TRAP_DISK` trap is generated.

## 4.5 The Serial Port

The hardware features a serial port that is ordinarily used to communicate with a personal computer and other add-on devices. However, for the purposes of debugging and development of your kernel, you are going to use the serial port as a debugging console where you will be able to launch programs using a simple text-based shell.

The serial port is a line-oriented device that uses DMA to transfer data to and from physical memory. Two special memory ports are used to control the serial device:

- **MEM\_SERIAL\_ADDR.** Writing a value to this memory location sets the *physical* memory address that will be used in the next write to **MEM\_SERIAL\_CS**. Typically, this register is used to set the location of a buffer from which data will be sent or into which data will be received.
- **MEM\_SERIAL\_CS.** This port is used to issue commands to the serial device and to query the device about the last serial interrupt. The format of the value written to this port is as follows:

Bits 24-31	Bits 16-23	Bits 0-15
opcode	unused	Length (in bytes)

where the **opcode** field is one of **SERIAL\_RECEIVE** or **SERIAL\_TRANSMIT**.

To write data to the serial device, the *physical* memory address of the data is first placed in **MEM\_SERIAL\_ADDR**. Afterwards, a **SERIAL\_TRANSMIT** operation is written to **MEM\_SERIAL\_CS**. The length field of this operation should contain the number of bytes of data to send. Upon completion of the write, a **TRAP\_SERIAL** interrupt will be generated. A maximum of 65536 bytes can be sent with a single write operation. However, your kernel will probably want to send data in smaller segments than that.

To receive data from the serial device, your kernel must first create an incoming data buffer and write its *physical* memory address to **MEM\_SERIAL\_ADDR**. Then, a **SERIAL\_RECEIVE** operation must be written to **MEM\_SERIAL\_CS**. In this case, the length field should contain the maximum number of bytes that can be stored in your kernel buffer. When a new line of input is received, a **TRAP\_SERIAL** interrupt will be generated. The constant **SERIAL\_MAX\_LINE** specifies the maximum number of bytes that the hardware can receive on a single line. To simplify programming, your incoming data buffer should always be at least that large.

It is important to note that the serial device can support simultaneous read and transmit operations. That is, a transmit command can be issued while the device is waiting for input and input can be received while the device is transmitting. Whenever either command is issued to the device, the completion of that command is signalled by a **TRAP\_SERIAL** interrupt. However, when this interrupt occurs, your kernel will have to read the value of **MEM\_SERIAL\_CS** to find out more information. When this register is read, the **opcode** field contains either **SERIAL\_RECEIVE** or **SERIAL\_TRANSMIT** to indicate the interrupt type. In addition, for receives, the length field will contain the number of bytes that were actually read by the device.

Although the device can support simultaneous read/transmit operations, *it is illegal to issue a new operation to the device until the previous operation of the same type has completed*. In other words, a new transmit command can not be issued until the previous transmit has completed. In addition, please note that the serial device operates on *physical memory addresses*. Writing a virtual memory address into **MEM\_SERIAL\_ADDR** is likely to lead to extremely erratic and peculiar behavior.

## 4.6 The Display Device

The hardware features a 256x256 pixel LCD display. The display has its own microcontroller that understands simple commands written to a special `MEM_DISPLAY_CTRL` port. Each display command is a 32 bit integer encoded as follows:

Bits 24-31	Bits 16-23	Bits 8-15	Bits 0-7
opcode	argument	x coordinate	y coordinate

The different display opcodes are as follows:

- `DISPLAY_POINT` Set a point. `arg(1) = black, arg(0) = white`
- `DISPLAY_LINE` Draw line from last drawn point.
- `DISPLAY_CLEAR` Clear the display.
- `DISPLAY_CHAR` Draw a character. `arg = ASCII code.`
- `DISPLAY_BOX` Draw solid box from last drawn point.

Your kernel does not have to perform any special decoding of display operations. However, your kernel must allow user programs to write display codes to the display device using a special system call (see Section 7.3.4).

## 4.7 The Pen Device

The hardware features a touch screen that detects the presence of a pen. Whenever the pen is touching the display, the hardware generates a `TRAP_PEN` interrupt. When generated, your kernel should read the status of the pen device by reading from `MEM_PEN_EVENT`. The value read from this location is encoded as follows:

Bits 24-31	Bits 16-23	Bits 8-15	Bits 0-7
event	unused	x coordinate	y coordinate

As with the display device, the event field specifies the type of event that has occurred. It is one of the following values

- `PEN_PRESS` The pen has touched the display
- `PEN_RELEASE` The pen has been lifted from the display
- `PEN_MOTION` The pen has moved.

Your kernel will generally not have to interpret the meaning of `MEM_PEN_EVENT` itself. Rather, it will simply have to queue up events and pass them on to user processes that request to read from the Pen device (see section 7.3.4).

## 4.8 Bootstrapping the kernel

On boot, the hardware starts executing code contained in the *Boot ROM*. This firmware loads the master boot record (block 0) from the disk device into the first page of physical memory and executes the first instruction found starting at `PMEM_BASE`.

The code contained in the master boot record is responsible for loading the rest of the kernel into memory. Specifically, it interacts with the disk controller to load the rest of the kernel and transfers control to the kernel's `main()` procedure. You will have to write the boot loader in order to get the rest of your kernel to run. Further details can be found in section 8.2.

## 4.9 Miscellaneous Hardware Operations

Normally, the CPU continues to execute instructions even when there is no useful work available for it to do. As a result, operating systems usually provide an *idle process* that is executed in this situation. Typically, this is just an empty infinite loop. However, in order to conserve power, the Yalnix hardware provides the following special instruction:

- `void Pause(void)`

Puts the CPU in standby mode until the next trap or interrupt occurs.

The hardware also provides another instruction to stop the CPU.

- `void Halt(void)`

Completely halts the CPU and does not begin execution until rebooted (i.e., until the Yalnix process is started again from the Unix command line). You should use this instruction to exit from your kernel.

Finally, to assist in the debugging of your kernel, the hardware provides a special character output port (the console). To write a character to this device, the following code can be used:

```
/* c is a single ASCII character */
WritePort(MEM_CONSOLE, c);
```

Unlike other devices, the console port does not generate interrupts. Once a character has been written, the value of `MEM_CONSOLE` remains non-zero until the port is ready to write another character. Therefore, to use the console, your kernel will have to use polling to determine when it is safe to write a character. For example:

```
while (ReadPort(MEM_CONSOLE)); /* Ready? */
/* Write next character */
WritePort(MEM_CONSOLE,c);
```

If characters are written to the console while the device is busy, they are silently discarded.

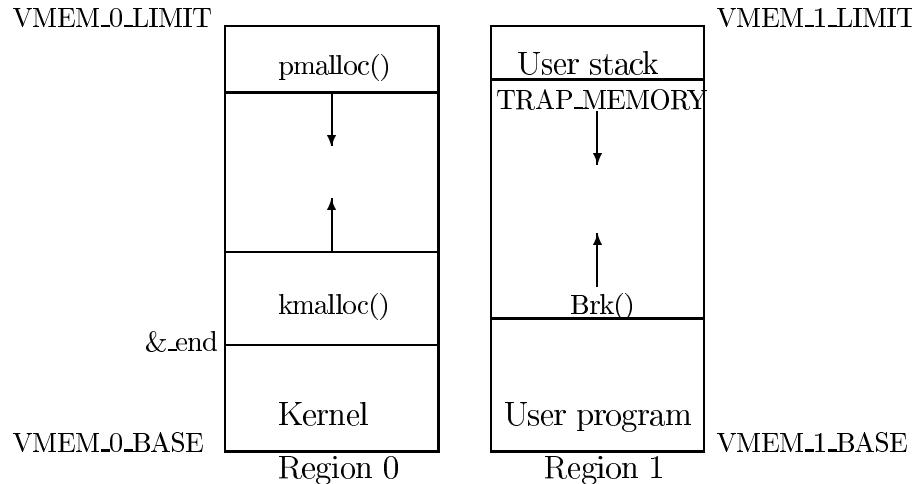


Figure 4: Process and Kernel Memory Management

## 5 Memory Management and CPU Context

When a user program is running on the hardware, its state of execution is encapsulated within a *process*. Associated with each process is a page table defining its memory use and a CPU context frame that contains the state of various CPU registers. In addition to managing the resources assigned to each processes, the kernel must manage memory and resources for its own use. This section provides an overview of how the kernel is supposed to manage memory and the CPU context for user processes and itself.

### 5.1 User Processes

As previously mentioned, region 1 of the virtual address space is reserved for user processes. As these processes run, they consume memory in two ways. First, a program may dynamically allocate memory on the heap using `malloc()` or `new`. Second, a program may require memory to be allocated on the program stack. The stack is where local variables to functions are stored as your program executes. Thus, if your program was highly recursive or involved a deeply nested collection of function calls, you would expect to see the stack grow significantly in size.

The conventional way to handle both types of memory allocations is to allocate memory from opposite ends of the virtual address space and for the heap and the stack to grow towards each other. In this case, the program and data area of a process would start at `VMEM_1.BASE` and grow upwards. The stack would be placed at the top of the virtual address space `VMEM_1.LIMIT` and grow downwards. Figure 4 illustrates this arrangement.

In general, the area of memory containing the program, data, and heap of a program is called the *program area* and the area containing the stack is called the *stack area*. The

limit of the program area (the lowest address not part of the program area) is called the process' *break* and is set from Yalnix user processes with `Brk()` system call. You will need to implement this system call as described in section 7.3.

The library with which your user programs are linked contain versions of `malloc()` and similar procedures that will call your kernel's `Brk()` syscall. Thus, while your kernel must manage user processes' via the `Brk()` syscall, your user programs can call `malloc()`, `realloc()`, etc, as if you were writing regular UNIX programs.

The stack area of a process is automatically grown by your Yalnix kernel in response to a memory exception resulting from an access to a page between the current break and the current stack pointer which has not yet been allocated. When the user process grows the stack beyond its currently allocated stack, the kernel will receive a `TRAP_MEMORY` exception. If the memory address that caused the exception is between the current break for the executing process and below the currently allocated memory for the stack, your kernel should attempt to grow the stack downwards to "cover" this address, if possible. In general, this will involve grabbing as many free frames of physical memory as are necessary. If there is insufficient memory, your kernel will have to kill the user process that caused the exception, but should continue to run all other user processes.

## 5.2 Kernel Memory

In the previous section we explained that at boot time, the hardware loads the kernel into a range of *physical* addresses starting at `PMEM_BASE`. To make it easier for the kernel to run after virtual memory is enabled, `PMEM_BASE` is exactly the same as `VMEM_0_BASE`. Therefore, when initializing the virtual memory system, you should map all of the physical pages used by the kernel to the exact same set of pages in region 0 by making sure that `vpn == pfn` in that range of memory addresses at the bottom of region 0.

When making this initial memory mapping, you will need to know where the kernel ends in memory. To find out, you can use the following values:

- `&_end` is the last address in the kernel's data segment.
- `&_etext` is the lowest address that does not contain executable instructions.

Thus, when your kernel boots, one of the first things it will do is examine the value of `&_end` and `&_etext` to find out how much memory it is using. All of this memory should be marked as "in use" and initialized in the kernel page table. The remaining physical memory pages should be placed on a free memory list before virtual memory is enabled.

To implement dynamic memory management, your kernel must implement it's own memory allocation functions. In fact, the standard C `malloc()`, `realloc()`, and `free` functions are **not** available to your kernel! Therefore, it is up to you to devise a suitable replacement.

There are two types of kernel memory management that your kernel is likely to require. The first is a page allocator that allocates memory in page-sized chunks. To do this, you should implement two functions:

- `void *pmalloc()`. Allocate a single page of memory and return a pointer to its virtual memory address. Returns 0 if no more physical memory is available.
- `void pfree(void *ptr)`. Release a page of memory starting at address `start`.

The page allocator is particularly useful for creating large data structures such as page tables and PCBs. Since the virtual address space for the kernel is relatively large (16 Mbtyes), an easy way to implement the page allocator is to simply start at the top of the region 0 page table and to allocate the first free page found in a descending search of the kernel page table.

In the unlikely event that your kernel needs to allocate a data structure larger than a page, you could extend the `pmalloc()` and `pfree()` functions to accept a page count (e.g., allocate 5 pages of memory). However, you probably won't need to make any allocations larger than a page in your implementation.

The second type of kernel memory allocator you will need is one optimized for small memory allocations (less than a page). To do this, implement two functions similar to `malloc()` and `free()` as follows:

- `void *kmalloc(int nbytes)`. Allocate an object of `nbytes` and return a pointer to its virtual memory location. Returns 0 if no more memory is available.
- `void kfree(void *ptr)`. Free an object previously created with `kmalloc()`.

The implementation of these functions is entirely up to you. One approach, is to implement a simple memory management scheme that utilizes a “kernel break” like the normal C `malloc()` function. This approach is described in the K&R C Programming Language book on page 185. Alternatively, you could try to divide virtual memory up into a collection of memory pools for different sizes of objects. Or, as a worst case scenario, you could just use the page allocator for everything and waste a lot of memory. In any case, a certain amount of creativity on your part will make the difference between an efficient kernel and a bloated kernel.

Finally, if running on real hardware, your kernel would also have to manage its own stack, but we have eliminated that concern for the sake of simplicity. You may assume that the kernel has no stack in the time-space continuum (a feat of the Yalnix design committee).

### 5.3 CPU Context

In order to switch execution from one user process to another, you need to be able to save the state of the running process, and restore the state in which some other process found itself at a previous time. Most of this state is active in physical memory, and can be made unavailable or available by changing the page tables. The rest of the state of the running process (program counter, stack pointer, etc) is passed to your kernel when a trap is executed through the trap vector, in the form of an argument to the function in your kernel that receives the trap.

This argument is of type `cpu_context_frame`, which is defined in `hardware.h`. The following fields are defined within a cpu context frame:

- `int frame_len`: The length (in bytes) of the entire exception frame. The fields defined here form a common header on the exception frame, and are followed by a variable number of bytes of additional hardware information defining the state at the time of the exception. The `frame_len` field defines the total length (including this header) of the cpu context frame. The maximum length of a cpu context frame is defined in `hardware.h` as `MAX_CPU_CONTEXT_FRAME`.
- `int vector`: The vector number of the particular trap, e.g, `TRAP_ILLEGAL`.
- `int code`: A code value giving more information on the particular trap. Its meaning varies depending on the type of trap. In particular, for `TRAP_KERNEL`, it specifies the type of syscall that produced the current trap (see section 7.3).
- `void * addr`: This field is only meaningful for the `TRAP_MEMORY` exception. It contains the memory address being referenced that caused the exception.
- `void * pc`: The *program counter* value at the time of the trap.
- `void * sp`: The *stack pointer* value at the time of the exception.
- `unsigned long regs[8]`: The contents of eight general purpose CPU registers at the time of the exception. In particular, for a `TRAP_KERNEL` syscall, these values give the arguments passed by the user process to the syscall and are used to return the result value from the syscall to the user process. This usage is defined further in section 7.3.

In order to switch contexts from one user process to another, you need to save the cpu context of the running process into its process control block, select another process to run, and restore that process' previously stored cpu context frame into the cpu context argument passed by reference to your trap handler. The hardware takes care of extracting and restoring the actual hardware state into and from the trap handler's cpu context argument. When copying a cpu context back and forth, you absolutely *must* copy the number of bytes indicated by its `frame_len` field, and **not** the number of bytes returned by `sizeof(cpu_context_frame)`.

The values of any privileged memory locations (i.e., `MEM_*` ) are **not** part of the CPU context. These values are associated with the current process by your kernel, not by the hardware, and must be changed by your kernel on a context switch when/if needed.

## 6 The Yalnix Filesystem

In order to run programs, your kernel needs to be able to load programs from disk drive. Not only that, the kernel itself needs to be loaded by the master boot record loader. To do

this, your kernel will need to interact with the disk controller to bring selected disk blocks into memory. However, it also needs to know how data is actually stored on the disk.

Currently, the disk is organized in an extremely lame, but very simple format. In a real operating system, the disk is *much* more complicated than this. However, this should be enough to give you a taste of the trouble associated with disks—besides, we only have 7 weeks to finish the project.

## 6.1 Disk Organization

The disk is organized as follows:

- Block 0. The Yalnix Master Boot Record (MBR). This block is loaded and executed automatically by the hardware during power-up. The code contained in the MBR (which you must write) loads the rest of the yalnix kernel.
- Blocks 1 to n. The Yalnix kernel image. This is the Yalnix kernel itself. The kernel is stored in a special format known as YOFF (Yalnix Object File Format). This is described shortly. The size of the kernel image will depend on how big your kernel is.
- The last 8 blocks on the disk (starting at \*MEM\_DISK\_NBLK-8) contain the disk directory. These blocks contain information about where other files are located on the disk.
- All other blocks are available for use by user programs.

## 6.2 YOFF Files

All user programs and the Yalnix kernel itself are stored in a format known as YOFF (Yalnix Object File Format). YOFF is a special file format that contains information about executable programs. The first *block* in a YOFF file contains a special header described by the following C structure, which is defined in the header file `yoff.h`:

```
typedef struct {
    char        magic[4];      /* Magic bytes "YOFF"          */
    unsigned int length;       /* Length of file in bytes    */
    unsigned int text_base;    /* Text base address          */
    unsigned int text_size;    /* Text segment size          */
    unsigned int text_entry;   /* Text segment entry point   */
    unsigned int data_base;    /* Base address of data segment */
    unsigned int data_size;    /* Size of data segment       */
    unsigned int bss_base;     /* Base address of bss segment */
    unsigned int bss_size;     /* Size of bss section        */
} yoff_header;
```

Immediately following the header block are the raw data blocks that make up the file. The `length` field of the header contains the length of the data to follow.

The `text_base`, `text_size`, and `text_entry` fields describe the program's text-segment. `text_entry` is the memory address of the first instruction to execute in the program (usually corresponding to `_start()` or `main()`). The `data_base` and `data_size` fields describe the location and size of the program's data segment. The `bss_base` and `bss_size` fields describe the bss segment.

To actually load a program, your kernel will first load the program's header block. Then, using the information in the header, the kernel will arrange to load the rest of the program, it will initialize memory, and it will pass control to the program's entry point stored in `text_entry`.

## 6.3 Disk directory

The yalnix kernel is always stored in a sequence of consecutive blocks starting at block 1. However, for Yalnix user programs, your kernel will need to consult the disk directory to find out where the programs are located. The directory is located in the last 8 blocks of the disk—which conveniently fits into a single page of physical memory. The directory is an array of directory entries defined by the following structure (found in `yalnixdir.h`):

```
typedef struct {
    unsigned int      first_block;      /* Starting block */
    unsigned int      length;          /* Length in bytes */
    char              filename[24];     /* Filename */
} yalnix_dirent;
```

Files are always stored in a collection of consecutive disk blocks starting at `first_block`. `filename` is a NULL-terminated string. A zero-length filename indicates an empty directory entry. Since the disk directory is limited to 8 blocks, a maximum of 256 files can be created.

## 6.4 How do files get created anyways?

Your kernel will not need to worry about the creation of files. Special Unix-based tools are available to format and place files on the Yalnix disk device prior to the execution of your kernel. Your kernel will only need to worry about loading this information from the disk as it runs.

Note: since your kernel won't create files, you might find it easier to cache the disk directory in memory so you don't have to keep reloading it.

## 7 The Yalnix Kernel

So far, the discussion has focused exclusively on the underlying hardware of the system and a few data formats. In this section, we describe the major components that will form the implementation of your Yalnix kernel. These are as follows:

- User processes.
- Process scheduling.
- System calls.
- Device drivers.

### 7.1 User Processes

In your kernel, you will need to implement a process control block (PCB) structure that encapsulates the notion of a “user process.” This data structure minimally needs to contain the following:

- An integer process ID that uniquely identifies the process.
- A CPU context (`cpu_context_frame`).
- A memory map specifying the virtual to physical memory mapping. For example, you may need to store the region 1 page table for the process.
- A process state (e.g., whether or not it is running or not).

In addition to this, you will probably need to create a few general purpose data structures such as lists and priority queues that can be used to contain PCBs. This is because processes spend much of their time waiting for I/O devices, waiting to be scheduled on the CPU, and so forth.

### 7.2 Process Scheduling

Your kernel should use the `TRAP_CLOCK` interrupt to implement a round-robin process scheduling scheme. Whenever the `TRAP_CLOCK` interrupt is received, your kernel should save the context for the currently running process and switch to the next process on a ready process queue. The only exception to this rule is the case in which only a single process (other than the kernel idle process) is in a runnable state. In this case, the currently running process should be allowed to continue (i.e., your kernel should never switch to the idle process if there is useful work to do).

It is important to note that the `TRAP_CLOCK` interrupt handler also has to check for processes that are sleeping due to a `Delay()` system call. For example, your kernel will have to check the sleep queue for processes that need to be awakened and put them back on to the ready queue if their time has expired.

## 7.3 Yalnix System Calls

### 7.3.1 How to handle TRAP\_KERNEL

When a user process wants to call the kernel, it executes a trap instruction that results in a `TRAP_KERNEL` interrupt in the kernel. Upon entry to your kernel, the `code` field of the `cpu` context frame indicates *which* kernel call is being invoked (as defined with symbolic constants in `yalnixcall.h`.)

To simplify the interface, a library of assembly routines is provided that performs the trap from a user process. This library provides a standard C procedure call interface for the syscalls described shortly. The arguments supplied to each of the library calls are available in the `regs` fields of the `cpu` context frame received by your `TRAP_KERNEL` handler. Each argument passed to the library procedure is available in a separate `regs` register, in the order passed to the procedure, beginning with `regs[0]`.

Each syscall returns a single integer value, which becomes the return value from the C library procedure call interface for the syscall. When returning from a `TRAP_KERNEL`, the value to be returned should be placed in the `regs[0]` register by your kernel. If any error is encountered for a syscall, one of the error codes of the form `ERROR_type` defined in `yalnix.h` should be returned in `regs[0]`.

The system calls that you will be implementing fall into these general categories:

- Process management.
- Memory management.
- I/O.
- Synchronization.
- Interprocess communication.

Each of these are now described.

**Note:** The system call names listed in this section are the names used by Yalnix *user* processes. They are not necessarily the same function names or prototypes that you will use in the implementation of your kernel.

### 7.3.2 Process management

The following syscalls pertain to the creation, destruction, and management of user processes.

- `int Fork()`

Create a new child process as a copy of the parent (calling) process. On success, in the parent, the new process ID of the child is returned. In the child, the return value should be 0. If there is not enough memory to perform the requested fork operation, the value `ERROR_NOMEM` should be returned.

- `int Exec(char * filename, char ** argvec)`

Replace the currently running program in the calling process with the program stored in the file named by *filename*. The argument *argvec* points to a vector of arguments to pass to the new program as its argument list. The new program receives these arguments as arguments to its `main` procedure. The *argvec* is formatted the same as any C program's *argv* vector. The last entry in the *argvec* vector must be a NULL pointer to indicate the end of the list. By convention, *argvec[0]* is the name of the program to be run, but this is controlled entirely by the calling program. The first argument to the new program's `main` procedure is the number of arguments passed in this vector. On success, there is no return from this call in the calling program, but, rather, the new program begins executing at its entry point, and its conventional `main(argc, argv)` routine is called. If the call fails, one of the following error codes should be returned (depending on the nature of the failure):

- `ERROR_NOMEM`. The new process requires more memory than is available.
- `ERROR_NOEXEC`. The filename specified is not an executable file in the right format.
- `ERROR_NOENT`. The filename specified does not exist.
- `ERROR_FAULT`. One of the arguments points to an illegal address.

To load an executable program, your kernel will have to search for the filename in the disk directory (described in the previous section). You will then need to load the program by looking at its YOFF header and initializing all of the proper data structures.

- `int Exit(int status)`

Terminate execution of the current process and save the integer *status* for possible later collection by the parent process on a call to `Wait`. All resources used by the calling process are freed, except for the saved *status*. This system call should always succeed so the return value is meaningless. When a process exits or is aborted, its children should continue to run normally even though they will no longer have a parent. When the orphans exit at a later time, you need not save or report their exit status since there is no longer anybody around to care.

- `int Wait(int pid, int * status_ptr)`

Collect the process ID and exit status returned by a child process of the calling program. *pid* indicates a specific child PID or can be set to 0 to indicate any child. When a child process exits, it should be added to a queue of child processes not yet collected by its specific parent. After the `Wait` call, the child process is removed from this queue. If none of the child processes have called `Exit()`, the calling process is blocked until its next child calls `Exit`. The process ID of the child process is returned on success, and its exit status is copied to the integer referenced by the *status\_ptr* argument. If the calling process has no children, this function should return the value `ERROR_CHILD`. If *status\_ptr* points to an invalid memory address, `ERROR_FAULT` should be returned. If

`pid` is not a valid process ID or does not correspond to a child of the calling process, `ERROR_CHILD` should be returned.

- `int GetPid(void)`

Returns the process ID of the calling process.

- `int Delay(int clock_ticks)`

The calling process is blocked until *at least* `clock_ticks` clock interrupts have occurred after the call. Upon completion of the delay, the value 0 is returned. If `clock_ticks` is 0, return is immediate. If `clock_ticks` is less than 0, time travel is not carried out and `ERROR_INVAL` is returned.

- `int Yield()`

The calling process is descheduled and returned to the ready state (i.e., placed back on the ready queue). The scheduler then switches to the next runnable process. If no other runnable processes are available, this system call does nothing and control is returned. The function always succeeds and returns 0.

### 7.3.3 Memory management

- `int Brk(char * addr)`

Allocate memory to the calling process' program area, enlarging or shrinking the program's heap storage so that the value `addr` is the new break value of the calling process. The break value of a process is the address immediately above the last address used for its program instructions and data. This call has the effect of allocating or deallocating enough memory to cover only up to the specified address, rounded up to an integer multiple of `PAGESIZE`. The value 0 is returned on success. On failure, the function should return `ERROR_NOMEM` to indicate insufficient memory.

### 7.3.4 Device I/O

Each I/O device is controlled by a dedicated set of system calls.

- `int DisplayWrite(int *buf, int nitems)`

Writes the contents of the buffer referenced by `buf` to the LCD display device. `buf` is a buffer containing integer display opcodes as described in section 4.6. `nitems` is the number of items in the buffer. The return value should be `nitems` if the Write is successful. Otherwise, one of the following error codes should be returned:

- `ERROR_FAULT`. `buf` refers to an illegal address.
- `ERROR_INVAL`. `len` is invalid (a negative number perhaps).

- **int PenRead(int \*buf, int maxevents, int noblock)**

Reads up to *maxevents* pen event codes into the buffer specified by *buf*. Each event is an integer encoded in the format described in section 4.7. Returns the number of events written into the buffer. If *noblock* is set to 0, this call blocks until a pen event is received. If set to 1, the call returns immediately even if no events are pending (in which case 0 is returned). The following error codes should be returned:

- **ERROR\_FAULT.** *buf* refers to an illegal address.
- **ERROR\_INVAL.** *maxevents* or *block* is invalid (a negative number perhaps).

- **int SerialWrite(char \*buf, int len)**

Writes *len* bytes of data referenced by *buf* through the serial port. The call blocks until all of the data has been successfully sent. On successful completion, the return value should be *len*. Otherwise, one of the following error codes should be returned:

- **ERROR\_FAULT.** *buf* refers to an illegal address.
- **ERROR\_INVAL.** *len* is invalid (a negative number perhaps).

- **int SerialRead(char \*buf, int maxlen, int noblock)**

Receive a maximum of *maxlen* bytes from the serial port and store them in *buf*. If *noblock* is set to zero, the call blocks until data is actually received. Otherwise, the call returns immediately even if no data is available. On success, the number of bytes received are returned. Otherwise, one of the following error codes are returned:

- **ERROR\_FAULT.** *buf* refers to an illegal address.
- **ERROR\_INVAL.** *len* is invalid (a negative number perhaps).

- **int ConsoleWrite(char \*buf, int len)**

Writes *len* bytes of data referenced by *buf* on the Yalnix debugging console. This call never blocks. On success, *len* is returned. Otherwise, one of the following error codes should be returned:

- **ERROR\_FAULT.** *buf* refers to an illegal address.
- **ERROR\_INVAL.** *len* is invalid (a negative number perhaps).

### 7.3.5 Interprocess Communication

Your kernel needs to support a simple form of interprocess communication in the form of pipes. The following system calls are used:

- **int PipeCreate(int *pipeid*)**

Creates a named pipe with integer identifier *pipeid*. Data written to *pipeid* using `PipeWrite()` is read using `PipeRead()`. Returns `ERROR_NOMEM` if no memory is available to create a pipe. Returns `ERROR_PIPE` if *pipeid* is a bad pipe descriptor or if a pipe with the same descriptor is already in use.

- **int PipeWrite(int *pipeid*, void \**buffer*, int *len*)**

Writes *len* bytes of data to the pipe with descriptor *pipeid*. The calling process is blocked until the data is successfully read by another process using `PipeRead()`. Returns the number of bytes sent on success. On error, one of the following values should be returned:

- `ERROR_FAULT`. *buffer* refers to an illegal address.
- `ERROR_INVAL`. *len* is invalid (a negative number perhaps).
- `ERROR_PIPE`. *pipeid* is an invalid destination (a non-existent pipe descriptor).

- **int PipeRead(int *pipeid*, void \**buffer*, int *len*, int *noblock*)**

Reads up to *len* bytes from a pipe with descriptor *pipeid*. The *noblock* flag indicates whether or not the read operation should block. On success, it returns the number of bytes that were received. On error, the following codes should be returned:

- `ERROR_FAULT`. *buffer* refers to an illegal address.
- `ERROR_INVAL`. *len* is invalid (a negative number perhaps).
- `ERROR AGAIN`. No data is available. Only returned in non-blocking mode.
- `ERROR_PIPE`. *pipeid* is an invalid pipe descriptor.

```
int PipeClose(int pipeid)
```

Closes the pipe in the calling process. This prevents data from being written or read in that process, but would not prevent the use of the pipe in other processes. A pipe is not actually destroyed until all other processes have close the pipe. `PipeClose()` returns `ERROR_PIPE` if *pipeid* is an invalid pipe descriptor.

Note: when a pipe is created, it is shared by all processes created using `Fork()`. Thus, a parent can communicate with a child process by first creating a pipe and then calling `Fork()`. The `Exec()` system call preserves all open pipes in the calling process—allowing interprocess communication between two separate programs. A pipe is not destroyed until all processes close the associated pipe descriptor. More than one process can read and write from a pipe at once, but the results aren't defined.

### 7.3.6 Synchronization

In order to synchronize processes, your kernel needs to provide semaphores. The following system calls should be implemented:

- `int SemaCreate(int semaid, unsigned int value)` Create a new semaphore identified by the integer *semaid*. The initial value of the semaphore is *value* which must be a non-negative integer. Returns 0 on success or `ERROR_SEMA` if the semaphore identifier *semaid* is already in use.

- `int SemaDestroy(int semaid)`

Destroys the semaphore identified by *semaid*. Returns 0 on success or `ERROR_SEMA` if the semaphore identifier doesn't exist. When the semaphore is destroyed, all processes waiting on the semaphore should be unblocked and the `SemaWait()` syscall should return `ERROR_SEMA`.

- `int SemaSignal(int semaid)`

Signals the semaphore identified by *semaid*. This increments the value of the semaphore and awakens at most one process that might be waiting on the semaphore (if applicable). Returns the value of the semaphore (before the increment) on success or `ERROR_SEMA` if the semaphore identifier is invalid.

- `int SemaWait(int semaid, int noblock)`

Waits on the semaphore identified by *semaid*. If the value of the semaphore is non zero, the value is decremented and the system call returns. If the value is zero, the calling process is blocked until it is awakened by a process calling `SemaSignal()`. If *noblock* is set to 1, the function returns an error instead of blocking when the semaphore is zero. Returns the value of the semaphore (before the decrement) on success or `ERROR_SEMA` if the semaphore identifier is invalid. When *noblock* is set, `ERROR AGAIN` should be returned if waiting on the semaphore would cause the process to block.

### 7.3.7 Argument checking

Your kernel should **verify all arguments passed to a syscall**, and should return an appropriate error if any arguments are invalid. In particular, you need to verify that a pointer is valid before you use it. This means you need to look in the page table in your kernel to make sure that the entire area (such as a pointer and a specified length) are readable and/or writable (as appropriate) before your kernel actually tries to read or write there. For C-style character strings (null-terminated) you will need to check the pointer to each byte as you go (C strings like this are passed to `Exec()`) You should write a common routine to check a buffer with a specified pointer and length for read, write and/or read/write access; and a separate routine to verify a string pointer for read access. The string verify routine would check access to each byte, checking each until it found the '\0' at the end. Insert calls

to these two routines as needed at the top of each syscall to verify the pointer arguments before you use them.

Such checking of arguments is important for two reasons: security and reliability. An unchecked `Read`, for instance, might well overwrite crucial parts of the operating system which might, in some clever way, gain an intruder access as a privileged user. Also, a pointer to memory that is not correctly mapped or not mapped at all would generate a `TRAP_MEMORY` which would never be serviced because the kernel is not interruptible.

## 7.4 Device Driver Implementation

A critical part of your kernel is going to involve the use of the serial device and the disk drive. In the past, students have experienced a considerable amount of unwarranted trouble supporting hardware devices. This section offers a few suggestions on how to structure your kernel to make it easier.

First of all, realize that your kernel is really nothing more than a glorified queuing system. For example, consider the handling of the `TRAP_CLOCK` interrupt. Inside your kernel, you are going to keep a queue of processes that are ready to run (the ready queue). Then, on each clock interrupt, your kernel is going to remove the currently running process and place it at the end of the ready queue. Then, the first process on the ready queue will be popped off the queue and placed on the CPU. This process will run until the next clock interrupt occurs at which point it will be removed from the CPU and replaced with another process. Conceptually, this procedure should be pretty obvious.

The handling of other hardware devices is very similar to the handling of the clock interrupt. The only difference is that unlike the clock interrupt (which occurs automatically), the I/O devices operate using a command-interrupt cycle. That is, you first give the device a command. Then, some time later, an interrupt is generated to indicate completion of the command. Furthermore, once a command has been issued, no other operations can be scheduled on the device until an interrupt is received to signal its completion.

The implementation of your I/O device drivers is going to be very similar to that of the clock interrupt handler. First of all, you will need a queue that holds all of the processes that are waiting to use a particular device. Then, you will need to implement an interrupt handler that responds to the interrupt that is generated when a device operation has been completed (e.g., `TRAP_SERIAL` or `TRAP_DISK`). In this interrupt handler, your kernel will take the process for which the operation was performed and move it to another queue (perhaps the ready queue). Then, the handler will pop off the next process from its device queue and issue a command to the hardware device on its behalf. The final piece of the puzzle is how processes actually move from the ready-queue to the device queue. This usually occurs when a process executes a system call such as `SerialWrite()`. Figure 5 shows the sequence of events.

Another way to view the overall control flow is as a large state machine. Specifically, at any given time, each user process is only allowed to be in one particular state. System calls and interrupts move the process between different states of execution as shown in Figure 6.

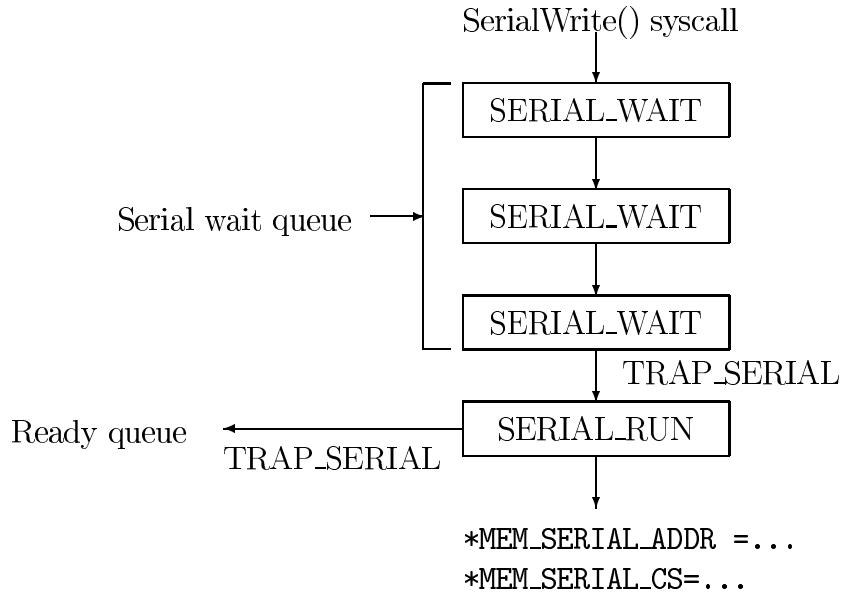


Figure 5: Control flow of serial device

## 8 Implementation hints

### 8.1 Process Initialization

#### Process IDs

Each process in the Yalnix system must have a unique integer process ID. Since an integer allows for over 2 billion processes to be created before overflowing its range, you should simply assign sequential numbers to each process created and not worry about the possibility of wrap-around (though real operating systems do worry about this.) The process ID must be a small, unique integer, not a pointer. The data structure used to map a process ID to its process control block should not be grossly inefficient in space nor time. This probably means that you'll end up using a hash table indexed by process ID to store the process control blocks.

#### Fork

On a `Fork`, a new process ID is assigned, the cpu context from the running parent is copied to the child, and new physical memory is allocated into which to copy the parent's memory space contents. Since the CPU can only access memory by virtual addresses (using the page tables) you must map both the source and the destination of this copy into the virtual address space at the same time. You need not map all of both address spaces at the

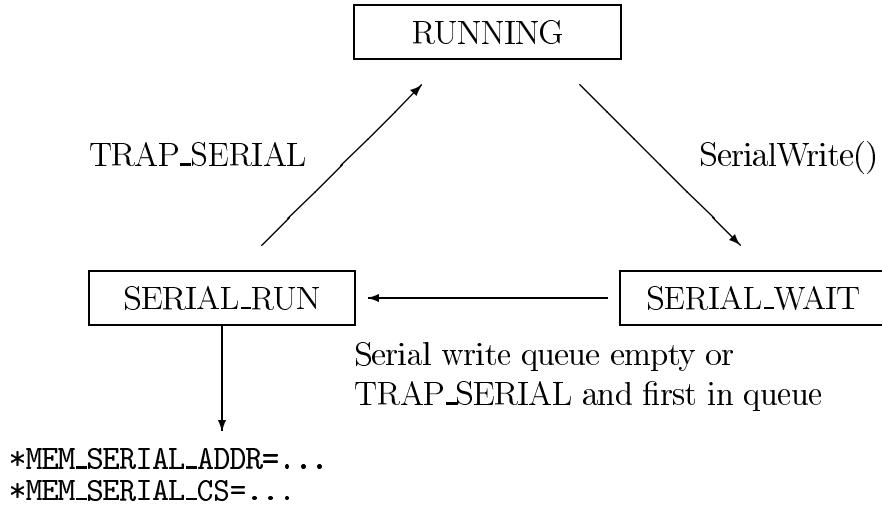


Figure 6: State transition diagram for a process calling `SerialWrite()`

same time, however: the copy may be done piece-meal, since the address spaces are already naturally divided into pages.

### Exec

On an `Exec`, you must load the new program from the specified file into the program region of the calling process, which will in general also require changing the process' page tables. The program counter in the `pc` field of the `cpu` context frame must also be initialized to the new program's entry point.

On an `Exec`, you must also initialize the stack area for the new program. The stack for a new program starts with only one page of physical memory allocated, which should be mapped to virtual address `VMEM_1_LIMIT - PAGESIZE`. As the process executes, it may require more stack space, which can then be allocated as usual for the normal case of a running program.

To complete initialization of the stack for a new program, the argument list from the `Exec` must first be copied onto the stack. The SPARC architecture also requires that an additional number of bytes (defined in `hardware.h` as `INITIAL_STACK_FRAME_SIZE`) be reserved immediately below the argument list. The stack pointer in the `sp` field of the `cpu` context frame should then be initialized to the lowest address of this space reserved for the initial stack frame. Like all addresses that you use, this must be a virtual address.

### Pipes

Pipes can be shared by different processes (pipe descriptors are duplicated by `Fork()`). Your kernel will need to keep a reference count of all pipe descriptors and only destroy a pipe when all references go away. The `PipeClose()` system decreases the reference count on a pipe descriptor and makes it unavailable in the calling process. Pipes are also closed when a process calls `Exit()`.

The pipe descriptor given to `PipeCreate()` is local to the calling process. Thus, if two entirely different processes call `PipeCreate()` with the same descriptor value, that creates two different pipes.

## 8.2 Kernel Initialization

Your kernel runs entirely off the emulated hardware. When the hardware starts, a boot ROM loads the master boot record (block 0) from the disk into the first page of physical memory. Execution then begins by simply setting the PC to first memory location in physical memory (`PMEM_BASE`).

The master boot record contains code that you have to write. It's also quite tiny (less than 1024 bytes). However, this is enough space to write a small kernel loader. The role of this loader is to bring the rest of the kernel into physical memory and to invoke its `main()` function.

The `main()` procedure invoked by your boot loader should conform to the following prototype:

```
int  
main(cpu_context_frame *context)
```

The `context` argument is a pointer to the initial cpu context frame (see section 5.3). Your kernel should use this cpu context frame as the basis for other cpu context frames, notably the cpu context frame that starts the initial process at boot time. Processes created from a `Fork`, instead, copy the cpu context frame from the parent process. (Note that in Yalnix, all processes except for the first process started at boot time are created by `Fork`.)

Before allowing the executing of user processes, the `main` routine should perform any initialization necessary for your kernel or required by the hardware. In addition to any initialization you may need to do for your own data structures, your kernel initialization should probably include the following steps:

- Initialize the interrupt vector entries for each trap, by making them point to the correct subroutines in your kernel.
- Initialize `MEM_VECTOR_BASE` to point to your interrupt vector array.
- Build a structure to keep track of what page frames are free in physical memory. For this purpose, you might be able to use a linked list of physical frames, implemented in the frames themselves. Or you can have a separate structure, which is probably easier, though slightly less efficient. This list of free pages should be based on the value of `MEM_PHYSM` which should be read in your `main()` procedure. The list should obviously not include any memory that is already in use by your kernel.
- Enable virtual memory

- Create the idle process based on the cpu context passed to your `main` routine. The idle process is the process that gets scheduled onto the CPU when there is no other process ready to run. See `Pause` in section 4.9.
- Create the first process and load the initial program into it. The process will serve the role of the `init` process in UNIX<sup>2</sup>.
- Return from your `main` routine. The machine will begin running the program defined by the current page tables and by the values returned in the cpu context frame (values which you have presumably modified to point to the initial context of the initial process).

### 8.3 Plan of Attack

Although the kernel project is being divided into distinct sub-projects, here is a feel for what is involved in the entire implementation effort.

1. “Wrap your brain” around the assignment. Read it carefully and understand first the hardware, then the operations you will need to implement. Understand all the points in the code at which your kernel can be executed, i.e, make a comprehensive list of kernel entry points (hint: the kernel runs in privileged CPU mode, which can only be set by the hardware.)
2. Think about all of the queues and process states that your kernel will need. When in doubt, remember that your kernel is really nothing more than a huge queueing system.
3. Write high-level pseudo-code for each syscall, interrupt and exception. Then decide on the things you need to put in what data structures (notably in the Process Control Block) to make it all work. Iterate until the pseudo-code and the main prototype data structures mesh well together.
4. Implement the Yalnix kernel loader. This code serves as the master boot record and is responsible for bringing the rest of the kernel into memory. Details instructions on this stage of the project will be given in another project handout.
5. Implement a skeleton of your kernel in which there are simple functions for every trap and system call.
6. Write some debugging support and utility functions to help you with the development of the kernel. For example, you might implement some function to write debugging information to the hardware console.
7. Write a simple in-kernel idle “process” that prints a message whenever it runs. Modify your kernel so that it boots up and starts running the idle process. Then you should experiment with generating different types of interrupts as the idle process runs (`TRAP_CLOCK`, `TRAP_PEN`, `TRAP_SERIAL`, etc.).

---

<sup>2</sup>The init process in UNIX forks all other processes; in particular, it forks all the `getty` programs that display the login prompt on all terminals, and it forks off all the network daemons.

8. Implement kernel memory management and the `TRAP_MEMORY` handler.
9. Create some functions for creating processes and define your PCB structure.
10. Implement the `ConsoleWrite()` system call so that programs can generate some output.
11. Write code that loads a user program from the disk drive into memory and starts it as a process.
12. Write an idle program (a single file that simply loops forever and prints some messages using `ConsoleWrite()`). Modify your kernel so that it loads the idle program into memory and starts its execution. If you can get the idle process to run, you have successfully bootstrapped virtual memory.
13. Write an “init” program. The simplest init program would just loop forever and print messages. Modify your kernel `main()` function to load “init” in addition to the idle program. You will now have two separate programs loaded at the same time.
14. Modify the `TRAP_CLOCK` handler to context switch between the idle and init programs.
15. Implement the `Brk` syscall and call it from the init program. At this point you have a substantial part of the memory management code working.
16. Implement the `Delay` syscall and call it from the init program. Make sure your idle program then runs for several clock ticks uninterrupted. At this point you have your process queues working smoothly and can context switch in at least two places.
17. Implement a few additional I/O syscalls so that user programs can start to generate some output. For example, `DisplayWrite()`, `SerialWrite()`.
18. Implement the `PenRead()` and `SerialRead()` syscalls. Now your kernel should have fairly good control over of the I/O devices.
19. Implement the `Fork` syscall. If you get this to work you are almost done with the memory system.
20. Implement the `Exec` syscall. You have already done something similar by writing `LoadProgram`.
21. Write another small program that does not do much (print messages and loop for instance). Call `Fork` and `Exec` from your init program to get this third program running. Watch for context switches.
22. Implement and test the `Exit` and `Wait` syscalls.
23. By now, you should be able to run a simple shell process on your kernel.
24. Implement semaphores.

25. Implement the `PipeCreate()`, `PipeWrite()`, `PipeRead()`, and `PipeClose()` syscalls.
26. Implement and test anything else that is remaining. This should be a piece of cake unless you are still working on process management the night before the deadline.
27. Look at your work and wonder in amazement.

One final note about the schedule is that time management is absolutely critical to success. It **is** possible to complete the kernel by making slow and steady progress. On the other hand, if you wait until the last minute to start, you may find yourself spending the night in the lab!

## 9 Your Assignment

To complete this assignment, you need to implement a Yalnix kernel that runs on the hardware defined in this document and handles all of the traps and system calls defined in this document. Specifically, you need to provide:

- A kernel loader (the Master Boot Record).
- A `main` routine to perform kernel initialization,
- A procedure to handle each defined trap.
- A procedure (called by your `TRAP_KERNEL` handler) to implement each defined Yalnix syscall.
- Implementations of the required I/O functions for each of the I/O devices.

Minimally, to receive a passing grade for the project, your kernel must be able to boot up to a simple shell and execute simple commands. In addition, your kernel must support the basic I/O devices (display, pen, serial).

## 10 Logistics

### 10.1 Compiling your kernel

The file `/yalnix/sample/Makefile.template` contains a basis for the `Makefile` we recommend you use for this project. The comments within it should be self-explanatory. There is some magic involved in the make macros with names following the pattern `*_LINK_FLAGS` which you may safely ignore but which may not be safely taken out.

The GNU C compiler (`gcc`) can be used to compile your kernel and user programs. Although `g++` has been installed on the machine, your kernel must be implemented in ANSI C. In part this is due to the fact that your kernel does not have access to the standard C or C++ libraries—making it extremely difficult to utilize basic languages features like the “new” and “delete” operators.

## 10.2 Running your kernel

Your kernel will be compiled and linked into a few executable programs. The program 'ylnixmbr' is the ylnix master boot record. The program 'ylnix' is the kernel image. To run your kernel, these programs must first be copied to the ylnix disk drive. A special tool, ydisk is used to do this. For example:

```
% ydisk mbr ylnixmbr
*** ydisk: installing master boot record.
    Boot record size : 940 bytes
    Text segment      : 820 bytes
    Master boot record successfully installed!
% ydisk kernel ylnix
*** ydisk: installing kernel image.
    Kernel image size : 24612 bytes
    Text segment      : 5040 bytes
    Kernel image successfully installed!
```

The supplied Makefile automatically tries to install these programs on your ylnix disk device—so you will not have to type these commands yourself. However, you will see the above output when compiling your kernel.

To boot your kernel, the yboot program is used as follows:

```
yboot [-t [tracefile]] [-lk level] [-lh level] [-s]
```

For example,

```
yboot -t -lk 5 -lh 3
```

The meaning of these switches is as follows:

- **-t:** This turns on “tracing” within the kernel and machine support code, and optionally specifies the name of the file to which the tracing output should be written. To generate traces from your kernel, you may call

```
TracePrintf(int level, char * fmt, args...)
```

where *level* is an integer tracing level, *fmt* is a format specification in the style of `printf()`, and *args* are the arguments needed by *fmt*. You can run your kernel with the “tracing” level set to any integer. If the current tracing level is greater than the *level* argument to `TracePrintf`, the output generated by *fmt* and *args* will be added to the trace. Otherwise, the `TracePrintf` is ignored.

If you just specify **-t** without a *tracefile*, the trace file name will be `TRACE`. You can, if you want, give a different name with **-t foo**.

- **-lk *n*:** Set the tracing level for this run of the kernel. The default tracing level is `-1`, if you enable tracing with the **-t** or **-s** switches. You can specify any *n* level of tracing.

- **-lh n:** Like `-lk`, but this one applies to the tracing level applied to the hardware. The higher the number, the more verbose, complete and incomprehensible the hardware trace.
- **-s:** Send the tracing output to the `stderr` file (this is usually your screen) in addition to sending it to the *tracefile*. This switch enables tracing, even if the `-t` switch is not specified.

## 10.3 Compiling Yalnix user programs

The supplied Makefile contains special options for compiling Yalnix user programs (read the comments in the file for details). In order for these programs to be accessible to your kernel, they also have to be copied to the Yalnix disk drive. The `ydisk` utility is also used to do this. For example, to install a program called “init” on the disk, you would do this:

```
% ydisk cp init
```

Normally, this step is automatically performed in the supplied Makefile. However, if you need to install a program manually, the above command would be used.

As the project proceeds, more and more user programs will be made available. To help with debugging, `ydisk dir` can be used to list all of the programs installed on the yalnix disk. For example:

```
% ydisk dir
start      size name
-----
16236      36360 init
16200      36360 init1
16164      36360 init2
```

The yalnix disk can also be reformatted using `ydisk format`.

## 10.4 Debugging Your Kernel

Due to the restricted execution environment, your kernel only has access to a limited amount of functionality. To produce output to a trace file, the `TracePrintf()` function may be used. Otherwise, all other output must be produced using the raw hardware. One way to produce debugging output is to write a collection of debugging functions that write characters to `MEM_CONSOLE`. In the first part of the project, you will be writing a collection of debugging functions to help you with this.

If you are feeling somewhat masochistic, the `gdb` debugger can be used in a limited sense. However, due to the emulator’s heavy use of signals and special memory mappings, it can

be hard to obtain useful information. If you do choose to use gdb, you may want to enter the following line so that certain signals are ignored:

```
handle SIGILL SIGFPE SIGSEGV SIGUSR1 pass nostop nowrap.
```

As an alternative to gdb, your instructor and TAs have developed a customized set of debugging tools that are specifically designed for the Yalnix project. When your kernel fails, it dumps a special 'ycore' file in the current directory. The program ydb (located in `/yalnix/bin`) can then be used to examine this file and produce a traceback. To use it, simply type 'ydb' after your kernel has crashed. You will get output similar to the following:

```
stonecrusher % ydb
ydb starting.
Reading ycore...
8 stack frames read.
18996 stack bytes read.
Reading symbols from 'yalnix'
Reading symbols from '/yalnix/etc/hardware'
Reading symbols from 'yboot'
Kernel stack trace:
#0 0x01000cac in TrapClock(ffbdd858,411f6ec,1000bf4,10080a8,40f4fe8,1d)
#1 0x0402835c in RealSignalHandler(e,0,ffbddd40,ffbdd898,ffbdd898,6)
#2 0x0402883c in SignalHandler(4,ffbddf78,ffbddcc0,0,0,0)
#3 0x040f51c0 in sigacthandler(4,ffbddf78,ffbddcc0,0,0,40f5400)
#4 0x04026b6c in _InitExceptions(411e0ac,410384c,4103844,3f8,2,40f5400)
#5 0x04023b2c in main(1,ffbefc8c,ffbefc94,80000012,4,0)
#6 0x000115f0 in main(1,ffbefc8c,ffbefc94,22000,0,0)
#7 0x00010e40 in _start(0,0,0,0,0,0)
stonecrusher %
```

## 11 Grading

This project will be graded primarily on correctness and efficiency. However, significant lack of documentation or elegance within your code will adversely affect your grade. Similarly, use of an “elegant” programming feature that results in unreasonable code-bloat will also result in a point deduction. You should use reasonable variable names and include comments in particularly tricky parts of the code, but you need not document code that does what it looks like. Also, when choosing to implement any particular feature, you should focus on solutions that are simple to understand and easy to debug over those that are excessively abstract and overly complicated.

Here is the approximate grading scale for the project:

- 60%. Correctness. Your kernel must be able to run a series of tests.
- 20%. On-time handin. This is for completing the subprojects on time.

- 10%. Programming style and efficiency.
- 10%. Participation (instructor discretion).

## 12 Approximate schedule

Please refer to the following dates when planning your time:

- Monday, January 7. Project handed out.
- Monday, January 14. You need to have formed a project group.
- Friday, January 18. Booting to yalnix.
- Friday, February 1. Booting to init, I/O device drivers, simple processes.
- Friday, February 15. Memory management, virtual memory, context switching, processes.
- Friday, March 1. Semaphores, pipes.
- Friday, March 8. Everything.

Dates are subject to change at any notice. Graduating seniors must complete the project before the last day to submit grades.

## 13 Credits

The Yalnix emulation environment was originally developed by Dave Johnson at Carnegie Mellon University. I have made extensive modifications for my own evil purposes...