

Chapter 8 Main Memory

8.1 Background

Memory consists of a large array of bytes, each with its own address. The CPU fetches instructions from memory according to the value of the program counter.

The memory unit sees only a stream of memory addresses; it does not know how they are generated (by the instruction counter, indexing, indirection, literal addresses, and so on) or what they are for (instructions or data). Accordingly, we can ignore how a program generates a memory address. We are interested only in the sequence of memory addresses generated by the running program.

8.1.1 Basic Hardware.

Main memory and the registers built into the processor itself are the only general-purpose storage that the CPU can access directly. Therefore, any instructions in execution, and any data being used by the instructions, must be in one of these direct-access storage devices. If the data are not in memory, they must be moved there before the CPU can operate on them.

For proper system operation we must protect the operating system from access by user processes, and processes from accessing other processes' memory. Since the operating system shouldn't intervene between the CPU and its memory accesses (because of the resulting performance penalty), this protection is implemented at hardware level.

First we need to make sure that each process has a separate memory space. We can provide this protection by using two registers: a *base register*, which holds the smallest legal physical memory address of the process; and the *limit register*, which specifies the size of the range of memory available for the process (so that all memory addresses n that satisfy $\text{base} \leq n < \text{base} + \text{limit}$, are available for a given process with such base and limit registers).

Protection of memory space is accomplished by having the CPU hardware compare every address generated in user mode with the registers. Any attempt by a program executing in user mode to access operating-system memory or other users' memory results in a trap to the operating system, which treats the attempt as a fatal error. The base and limit registers can be loaded only by the operating system, which uses a special privileged instruction. Since privileged instructions can be executed only in kernel mode, and since only the operating system executes in kernel mode, only the operating system can load the base and limit registers.

8.1.2 Address Binding.

A program residing on disk as a binary executable file, to be executed, must be brought into memory and placed within a process. Depending on the memory management in use, the process may be moved between disk and memory during its execution. The processes on the disk that are waiting to be brought into memory for execution form the *input queue*.

Most systems allow a user process to reside in any part of the physical memory. Thus, although the address space of the computer may start at 0x00000, the first address of the user process need not be 0x00000.

A user program goes through several steps before being executed. Addresses may be represented in different ways during these steps.

Addresses in the source program are generally symbolic. A compiler typically binds these symbolic addresses to relocatable addresses (such as "14 bytes from the beginning of this module"). The linkage editor or loader in turn *binds* the relocatable addresses to absolute addresses (such as 0x74014). Each binding is a mapping from one address space to another.

Binding of instructions and data to memory addresses can be done at any step along the way:

- **Compile time.** If you know at compile time where the process will reside in memory, then *absolute code* can be generated. If at some time later, the starting location changes, it is necessary to recompile the code.
- **Load time.** If it is not known at compile time where the process will reside in memory, then the compiler must generate *relocatable code*. In this case, final binding is delayed until load time (of the program into memory).
- **Execution time.** If the process can be moved during its execution from one memory segment to another, then binding must be delayed until run time. Most general-purpose operating systems use this method.

8.1.3 Logical and Physical Address Space.

An address generated by the CPU is referred to as a *logical address*, whereas an address seen by the memory unit (the address loaded into the memory-address register of the memory) is commonly referred to as a *physical address*. The compile-time and load-time address-binding methods generate identical logical and physical addresses. The execution-time binding scheme results in differing logical and physical addresses. In this case, we usually refer to the logical address as a *virtual address*. The set of all logical addresses generated by a program is a *logical address space*. The set of all physical addresses corresponding to these logical addresses is a *physical address space*. The run-time mapping from virtual to physical addresses is done by a hardware device called the *memory-management unit* (MMU). The user program deals with logical addresses and never deals with real physical addresses. The memory-mapping hardware converts logical addresses into physical addresses. The user program generates only virtual addresses and thinks that process runs in locations in the virtual address space, however these are mapped to physical addresses before they are used.

8.1.4 Dynamic Loading

So far, it has been necessary for the entire program and all data of a process to be in physical memory for the process to execute. The size of a process has thus been limited to the size of physical memory. To obtain better memory-space utilization, we can use dynamic loading: a routine is not loaded until it is called.

8.1.5 Dynamic Linking and Shared Libraries

Dynamically linked libraries are system libraries that are linked to user programs when the programs are run. Dynamic linking, is similar to dynamic loading. Here, though, linking, rather than loading, is postponed until execution time.

This feature is usually used with system libraries, such as language subroutine libraries. Without this facility, each program on a system must include a copy of its language library in the executable image. This requirement wastes both disk space and main memory.

With dynamic linking, a *stub* is included in the image for each library- routine reference. The stub is a small piece of code that indicates how to locate the appropriate memory-resident library routine or how to load the library if the routine is not already present. When the stub is executed, it checks to see whether the needed routine is already in memory. If it is not, the program loads the routine into memory.

8.2 Swapping

A process must be in memory to be executed. A process, however, can be *swapped* temporarily out of memory to a backing store and then brought back into memory for continued execution. Swapping makes possible for the total physical address space of all processes to exceed the real physical memory of the system, thus increasing the degree of multiprogramming in a system.

8.2.1 Standard Swapping

It involves moving processes between main memory and a backing store. The backing store is commonly a fast disk, which must be large enough to accommodate copies of all memory images for all users, and it must provide direct access to these memory images. The system maintains a *ready queue* of all processes whose memory images are on the backing store and ready to run. Whenever the CPU scheduler decides to execute a process, it calls the dispatcher. The dispatcher checks to see whether the next process in the queue is in memory. If it is not, and if there is no free memory region, the dispatcher swaps out a process currently in memory and swaps in the desired process. It then reloads registers and transfers control to the selected process. Context switch time is given by: $\Delta t = 2 \cdot \text{size}_{\text{process}} / r$, where r is transfer rate of disk and 2 comes from swap in swap out.

Clearly, it would be useful to know exactly how much memory a user process is using, not simply how much it might be using. Then we would need to swap only what is actually used, reducing swap time. For this method to be effective, the user must keep the system informed of any changes in memory requirements. Thus, a process with dynamic memory requirements will need to issue system calls (some `request_memory()` and `release_memory()`) to inform the operating system of its changing memory needs.

If we want to swap a process, we must be sure that it is completely idle. Of particular concern is any pending I/O.

Standard swapping is not used in modern operating systems. It requires too much swapping time and provides too little execution time. Modified versions of swapping, however, are found on many systems.

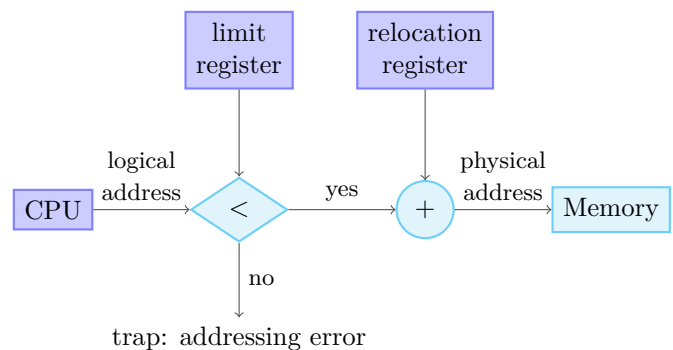
8.3 Contiguous Memory Allocation

The main memory must accommodate both the OS and various user processes. Contiguous memory allocation is an early solution to allocate memory in a somewhat efficient manner. In *contiguous memory allocation*, each process is contained in a single section of memory that is contiguous to the section containing the next process.

8.3.1 Memory Protection

The goal of memory protection is to prevent a process from accessing a region of memory it does not own. We can provide a simple solution with a system having a *relocation register* (which contains smallest memory address reachable by the process) and a *limit register* (range of logical addresses of the process).

When the CPU scheduler selects a process for execution, the dispatcher loads the relocation and limit registers with the correct values as part of the context switch. Because every address generated by a CPU is checked against these registers, we can protect both the operating system and the other users' programs and data from being modified by this running process.



The relocation-register scheme provides an effective way to allow the operating system's size to change dynamically. Dynamically loaded operating system code is called *transient* os code.

8.3.2 Memory Allocation

One of the simplest methods for allocating memory is to divide memory into several fixed-sized *partitions*. Each partition may contain exactly one process. Thus, the degree of multiprogramming is bound by the number of partitions. In this *multiple-partition* method, when a partition is free, a process is selected from the input queue and is loaded into the free partition. When process terminates, the partition becomes available for another process.

In the *variable-partition* scheme, the operating system keeps a table indicating which parts of memory are available and which are occupied. Initially, all memory is available for user processes and is considered one large block of available memory, a hole. Eventually, as you will see, memory contains a set of holes of various sizes.

the memory blocks available comprise a *set* of holes of various sizes scattered throughout memory. When a process arrives and needs memory, the system searches the set for a hole that is large enough for this process. If the hole is too large, it is split into two parts. One part is allocated to the arriving process; the other is returned to the set of holes. When a process terminates, it releases its block of memory,

which is then placed back in the set of holes. If the new hole is adjacent to other holes, these adjacent holes are merged to form one larger hole.

This procedure is a particular instance of the general *dynamic storage-allocation* problem, which concerns how to satisfy a request of size n from a list of free holes. There are many solutions to this problem:

- **First fit.** Allocate the first hole that is big enough.
- **Best fit.** Allocate the smallest hole that is big enough.
- **Worst fit.** Allocate the largest hole.

(Simulations have shown that both first fit and best fit are better than worst fit in terms of decreasing time and storage utilization)

8.3.3 Fragmentation

External Fragmentation. Both the first-fit and best-fit strategies for memory allocation suffer from *external fragmentation*. As processes are loaded and removed from memory, the free memory space is broken into little pieces. External fragmentation exists when there is enough total memory space to satisfy a request but the available spaces are not contiguous.

One solution to the problem of external fragmentation is compaction. The goal is to shuffle the memory contents so as to place all free memory together in one large block. Compaction is not always possible, however. If relocation is static and is done at assembly or load time, compaction cannot be done.

Another possible solution to the external-fragmentation problem is to permit the logical address space of the processes to be noncontiguous, thus allowing a process to be allocated physical memory wherever such memory is available. Two complementary techniques achieve this solution: segmentation (Section 8.4) and paging (Section 8.5). These techniques can also be combined.

Fragmentation is a general problem in computing that can occur wherever we must manage blocks of data.

Internal Fragmentation. Memory fragmentation can be internal as well as external. Consider a multiple-partition allocation scheme with a hole of 18,464 bytes. Suppose that the next process requests 18,462 bytes. If we allocate exactly the requested block, we are left with a hole of 2 bytes. The overhead to keep track of this hole will be substantially larger than the hole itself. The general approach to avoiding this problem is to break the physical memory into fixed-sized blocks and allocate memory in units based on block size. With this approach, the memory allocated to a process may be slightly larger than the requested memory. The difference between these two numbers is *internal fragmentation*—unused memory that is internal to a partition.

8.4 Segmentation

As a programmer, we think of memory as a collection of variable-sized segments with no necessary order among the segments. When writing a program, a programmer thinks of it as a main program with a set of methods, procedures, or functions, which may include various data structures; without

caring what addresses in memory these elements occupy. Segments vary in length, and the length of each is intrinsically defined by its purpose in the program. Elements within a segment are identified by their offset from the beginning of the segment (the first statement of the program, the seventh stack frame entry in the stack, the fifth instruction of the `Sqrt()`, ...).

Segmentation is a memory-management scheme that supports this programmer view of memory. A logical address space is a collection of segments.

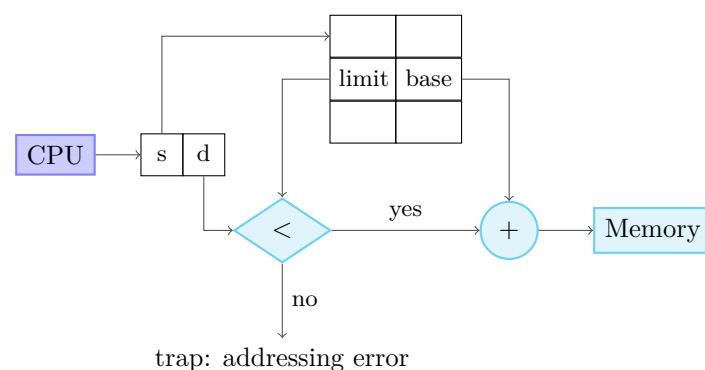
Each segment has a name and a length. The addresses specify both the segment name and the offset within the segment. The programmer therefore specifies each address by two quantities: a segment name and an offset. For simplicity of implementation, segments are numbered and are referred to by a segment number, rather than by a segment name. Thus, a logical address consists of a two tuple:

$$(\text{segment-number}, \text{offset}) \equiv (s, d).$$

Normally, when a program is compiled, the compiler automatically constructs segments reflecting the input program.

8.4.1 Segmentation Hardware

We must define an implementation to map the two-dimensional user defined addresses (s, d) into the one-dimensional physical memory:



This mapping is effected by a *segment table*. Each entry in the segment table has a *segment base* and a *segment limit*. The segment base contains the starting physical address where the segment resides in memory, and the segment limit specifies the length of the segment.

8.5 Paging

Paging is another memory-management scheme that offers that the physical address space of a process to be noncontiguous. However, paging avoids external fragmentation and the need for compaction (unlike segmentation). It also saves the considerable problem of fitting memory chunks of varying sizes onto the backing store.

8.5.1 Basic Method

The basic method for implementing paging involves breaking physical memory into fixed-sized blocks called *frames* and breaking logical memory into blocks of the same size called *pages*. When a process is to be executed, its pages are loaded into any available memory frames from their source (a file system or the backing store). The backing store is divided

into fixed-sized blocks that are the same size as the memory frames or clusters of multiple frames.

In this scheme, every address generated by the CPU is divided into two parts: a *page number* (p) and a *page offset* (d). The page number is used as an index into a *page table*. The page table contains the base address of each page in physical memory. This base address is combined with the page offset to define the physical memory address that is sent to the memory unit. The page size (like the frame size) is defined by the hardware. Taking the size of page to be a power of 2, makes the translation of logical address into page number and page offset quite simple: if size of logical address space is 2^m bytes, and page size is 2^n bytes, then the high-order $m - n$ bits of a logical address designate the page number, and the lower n bits designate the page offset.

When we use a paging scheme, we have no external fragmentation: any free frame can be allocated to a process that needs it. However, we may have some internal fragmentation. Notice that frames are allocated as units. If the memory requirements of a process do not happen to coincide with page boundaries. If process size is independent of page size, we expect internal fragmentation to average one-half page per process. This consideration suggests that small page sizes are desirable. However, overhead is involved in each page-table entry, and this overhead is reduced as the size of the pages increases. Also, disk I/O is more efficient when the amount data being transferred is larger. Generally, page sizes have grown in time.

On a 32-bit CPU, each page-table entry can be up to 4 bytes (32 bits) long. A 32-bit entry can point to one of 2^{32} physical page frames. If frame size is, say, 4KB (2^{12}), then a system with 4-byte entries can address up to 2^{44} bytes (16TB) of physical memory. Therefore, paging lets us use physical memory that is larger than what can be addressed by the CPU's address pointer length.

When a process arrives in the system to be executed, its size, expressed in pages, is examined. Each page of the process needs one frame. Thus, if the process requires n pages, at least n frames must be available in memory. If n frames are available, they are allocated to this arriving process. The first page of the process is loaded into one of the allocated frames, and the frame number is put in the page table for this process. The next page is loaded into another frame, its frame number is put into the page table, and so on.

An important aspect of paging is the clear separation between the programmer's view of memory and the actual physical memory. The programmer views memory as one single space, containing only this one program. In fact, the user program is scattered throughout physical memory, which also holds other programs.

Notice that the user process by definition is unable to access memory it does not own. It has no way of addressing memory outside of its page table, and the table includes only those pages that the process owns.

The operating system keeps track of allocated and free frames in physical memory. This is done by keeping a data structure called a *frame table*. The frame table has one entry for each physical page frame, indicating whether the latter is free or allocated and, if it is allocated, to which page of which process(ess).

The operating system also maintains a copy of the page

table for each process, just as it maintains a copy of the instruction counter and register contents. This copy is used to translate logical addresses to physical addresses whenever the operating system must map a logical address to a physical address manually. It is also used by the CPU dispatcher to define the hardware page table when a process is to be allocated the CPU. Paging therefore increases the context-switch time.

8.5.2 Hardware Support

The hardware implementation of the page table can be done in several ways. The modern solution is to use a special, small, fast-lookup hardware cache called a *translation look-aside buffer* (TLB). The TLB is associative, high-speed memory. Each entry in the TLB is a (key, value) pair. When the associative memory is presented with an item (a potential key), it is compared with all keys simultaneously. If the item is found, the corresponding value field is returned. The search is fast: a TLB lookup is part of the instruction pipeline, essentially adding no performance penalty.

The TLB contains only a few of the page-table entries. When a logical address is generated by the CPU, its page number is presented to the TLB. If the page number is found, its frame number (physical memory address) is immediately available and is used to access memory. As just mentioned, these steps are executed as part of the instruction pipeline within the CPU, adding no performance penalty compared with a system that does not implement paging. If the page number is not in the TLB (known as a *TLB miss*), a memory reference to the page table must be made. Depending on the CPU, this may be done automatically in hardware or via an interrupt to the operating system. When the frame number is obtained, we can use it to access memory (Figure 8.14). In addition, we add the page number and frame number to the TLB, so it will be found quickly on the next reference.

Some TLBs store *address-space identifiers* (ASIDs) in each TLB entry. An ASID uniquely identifies each process and is used to provide address-space protection for that process. When the TLB attempts to resolve virtual page numbers, it ensures that the ASID for the currently running process matches the ASID associated with the virtual page.

If the TLB does not support separate ASIDs, then every time a new page table is selected (for instance, with each context switch), the TLB must be flushed (or erased) to ensure that the next executing process does not use the wrong translation information. Otherwise, the TLB could include old entries that contain valid virtual addresses but have incorrect or invalid physical addresses left over from the previous process. The percentage of times that the page number of interest is found in the TLB is called the *hit ratio*. Modern computers reach a 99% hit ratio. To find the *effective memory-access time*, we calculate the expected time it takes to access memory, if r is hit ratio, and take time Δt to access memory; then on a TLB hit, it takes $2\Delta t$ seconds to access some memory (the time to access the page table + the actual memory access). Thus the effective memory access is:

$$\text{effective access time} = r \cdot \Delta t + (1 - r)2\Delta t.$$

8.5.4 Shared Pages

An advantage of paging is the possibility of sharing common code. This consideration is particularly important in a time-sharing environment. However, code can only be shared if it is *reentrant code* (code that never changes during execution).

Thus, two or more processes can execute the same code at the same time. Each process has its own copy of registers and data storage to hold the data for the process's execution. The data for two different processes will, of course, be different.

8.6 Structure of the Page Table

8.6.1 Hierarchical Paging

Suppose a system with a 32-bit logical address space. If the page size in such a system is 4KB (2^{12}), then a page table may consist of $2^{32}/2^{12} \approx 1$ million entries. Assuming that each entry consists of 4 bytes, it means that each process may need up to 4MB of physical address space for the page table alone. Clearly, one would not want to allocate the page table contiguously in main memory. One simple solution to this problem is to divide the page table into smaller pieces.

In *hierarchical paging*, we page the page table, and each page number is further divided into a 10-bit page number and a 10-bit page offset. Thus, a logical address is as $(p, d) = ((p_1, p_2), d)$. Where d is our 12-bit page offset, p_1 is an index into the outer page table, and p_2 is the displacement within the page of the inner page table.

This solution is unfeasible for a 64-bit logical address space system, since we would end with an outer page index p_1 of 42 bytes. This index could be further divided into a 2nd outer page and an outer page offset. This 2nd outer page index would still be too big (32-bits), which may need further paging... Thus leading to highly nested hierarchical pages, which is inappropriate.

8.6.2 Hashed Tables

A common approach for handling address spaces larger than 32 bits is to use a *hashed page table*, with the hash value being the virtual page number.