

Operating Systems

Memory

Malay Bhattacharyya

Assistant Professor

Machine Intelligence Unit
Indian Statistical Institute, Kolkata

May, 2022



oooooooooooooooooooo

ooooooo

oooooooooooooooooooo

oooooo

ooooo

Basics

Contiguous Memory Allocation

Paging

Swapping

Virtual Memory

Why do we need to manage memory?

The processes, together with the data they access, must be (at least partially) in main memory during execution.

For a better performance, we must we must share memory (i.e., keep many processes in memory). Memory management schemes control the way a shared memory is managed alongside the other tasks.

Selection of a memory-management scheme for a system depends on the system's hardware design.

The association between CPU and memory

The CPU fetches instructions from memory according to the value of the program counter.

Machine instructions can take memory addresses as arguments but not disk addresses. Therefore, any instructions in execution, and any data being used by the instructions, must be in one of these direct-access storage devices.

Note: If the data is not in memory, they must be moved there before the CPU can operate on them.

Looking at the hardware

Main memory and the registers built into each processing core are the only general-purpose storage that the CPU can access directly.

Built-in registers are generally accessible within one cycle of the CPU clock. But completing a memory access may take many cycles of the CPU clock. This is because the main memory is accessed via a transaction on the memory bus.

While accessing main memory, the processor normally needs to *stall*, since it does not have the data required to complete the instruction that it is executing. To manage this efficiently, we can add a relatively faster type of memory between the CPU and main memory, typically on the CPU chip for quick access. Such memories are known as *cache*.

Memory address space

We need to make sure that each process has a separate memory space. Separate per-process memory space protects the processes from each other and is fundamental to having multiple processes loaded in memory for concurrent execution.

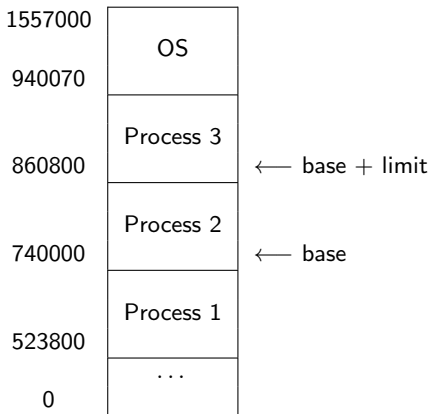
We can provide this protection by using the following two registers.

- ▶ The *base register* holds the smallest legal physical memory address.
- ▶ The *limit register* specifies the size of the range.

For example, if the base register holds 100780 and the limit register is 140200, then the program can legally access all addresses from 100780 through 240199 (inclusive).

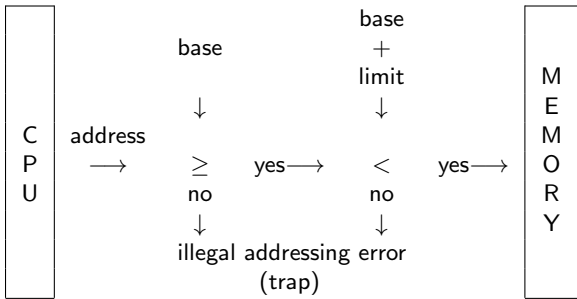
Memory address space

Memory consists of a large array of bytes, each with its own address. A base and a limit register define a logical address space.



How can we protect the memory space?

Protection of memory space can be achieved 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 the OS memory or other users' memory results in an error (known as trap).



Address binding

High-level code:

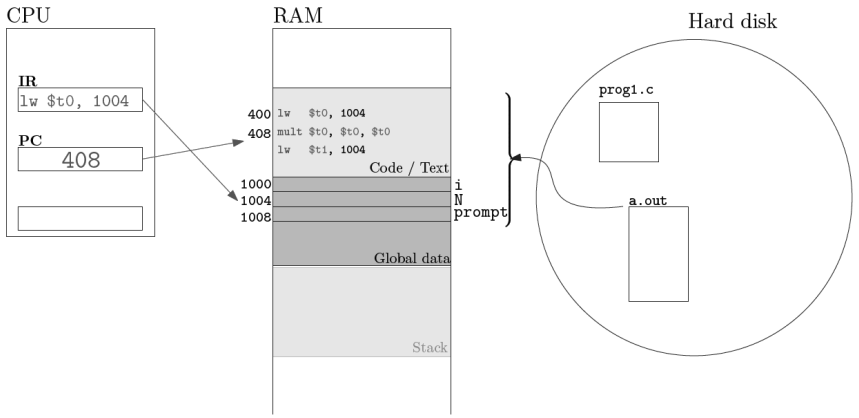
```
char prompt[] = "i";
int N = 20;
int i;
i = N*N + 3*N;
```

⇒

MIPS machine instructions:

```
lw $t0, 1004 # fetch N
mult $t0, $t0, $t0 # N*N
lw $t1, 1004 # fetch N
ori $t2, $zero, 3 # 3
mult $t1, $t1, $t2 # 3*N
add $t2, $t0, $t1 # N*N + 3*N
sw $t2, 1000 # i = N*N + 3*N
```

Address binding



Address binding

Usually a program resides on a disk as a binary executable file. For execution, the program is brought into memory. Most systems allow a user process to reside in any part of the physical memory. Hence, the address space does not necessarily start at 00000.

The binding of instructions and data to memory can happen at any time toward the execution.

1. **Compile time:** The *absolute code* is generated and is bound to actual physical addresses.
2. **Load time:** The *relocatable code* is generated and is bound to relative addresses.
3. **Execution time:** The *absolute code* and *relocatable code* are bound together.

Address binding during compilation

- ▶ Location of program in physical memory must be known at compile time
- ▶ Compiler generates absolute code
 - ▶ Compiler binds names to actual physical addresses
- ▶ Loader copies executable file to appropriate location in memory
- ▶ If starting location changes, program will have to be recompiled

Address binding during loading

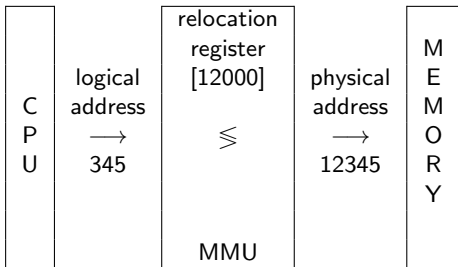
- ▶ Compiler generates relocatable code
 - ▶ Compiler binds names to relative addresses (offsets from starting address)
 - ▶ Compiler also generates relocation table
- ▶ Linker resolves external names and combines object files into one loadable module
- ▶ Loader converts relative addresses to physical addresses
- ▶ No relocation allowed during execution

Address binding during execution

- ▶ Programs/compiled units may need to be relocated during execution
- ▶ CPU generates relative addresses
- ▶ Relative addresses bound to physical addresses at run-time based on location of translated units
- ▶ Suitable hardware support required

Dynamic relocation

The user program never accesses the real physical addresses. The memory management unit (MMU) maps the logical address to a physical address.



Dynamic relocation

Binding addresses at either compile or load time generates identical logical and physical addresses. However, the address binding that happens during execution time results in different logical and physical addresses.

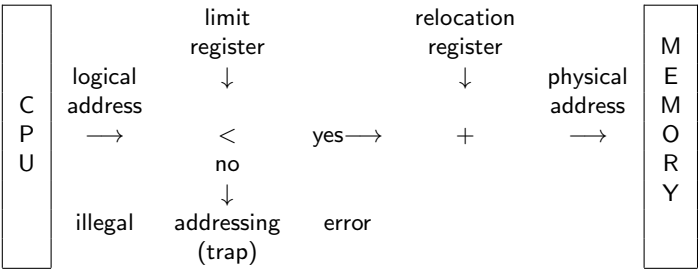
Dynamic loading

All routines are kept on disk in a relocatable load format. With dynamic loading a routine is not loaded until it is called.

When the main program is loaded into main memory during execution, it checks whether the called routine has been loaded. If not, the relocatable linking loader is employed to serve the purpose.

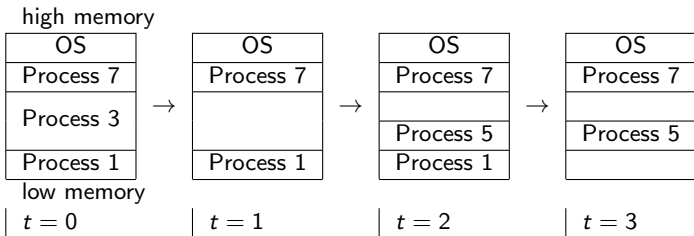
Protection of memory space

The MMU maps the logical address dynamically by adding the value in the relocation register.



Dynamic allocation of memory

Processes are dynamically assigned memory partitions of variable sizes on arrival in memory. As soon as a finishing process leaves memory, it might result into contiguous or non-contiguous holes.



Dynamic allocation of memory

It is a challenging task to allocate a memory partition of size n to a process from the free memory space (i.e., a list of free holes).

The following are some popular strategies to tackle this problem.

- ▶ **First fit:** Allocate the first hole that is large enough by searching either at the beginning of the set of holes or at the location where the previous first-fit search ended.
- ▶ **Best fit:** Allocate the smallest hole that is large enough by searching the entire list, unless the list is ordered by size.
- ▶ **Worst fit:** Allocate the largest hole by searching the entire list, unless it is sorted by size.

Note: Both the first fit and best fit perform better in terms of decreasing time and storage utilization.

Internal and external fragmentation

- ▶ **Internal Fragmentation:** This is caused by the unused memory that is internal to a partition.
- ▶ **External Fragmentation:** This is caused by the freed memory that is non-contiguous.

How does dynamic memory allocation affect fragmentation?

Both the first fit and best fit strategies suffer from *external fragmentation*. As processes arrive and leave memory, they turn the free memory space broken into little pieces.

The 50-percent rule

Even in the optimized version of the first fit strategy, given N allocated memory blocks, another $0.5N$ blocks will be lost due to fragmentation.

Addressing the fragmentation problem

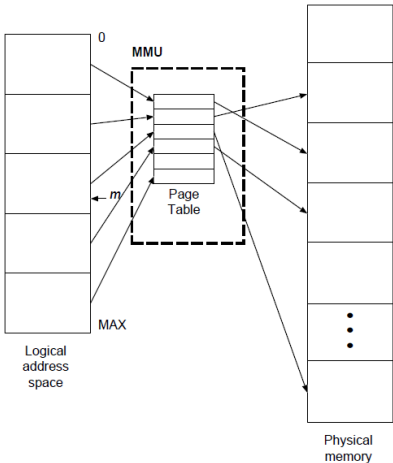
The internal fragmentation can be solved by assigning space to the process via *dynamic partitioning*. It allocates only the amount of space requested by the process.

The following strategies help to solve the problem of external fragmentation.

1. Through *compaction*, we can shuffle the memory contents such that all free memory get placed together in one large block.
2. Through *paging*, we can allow the logical address space of processes to be non-contiguous, thereby allocating physical memory wherever such memory is available.

Basics of paging

When a process with n pages has to be loaded, n free frames have to be found. The kernel keeps track of free frames and the page table translates logical page numbers to physical frame addresses.



Address translation

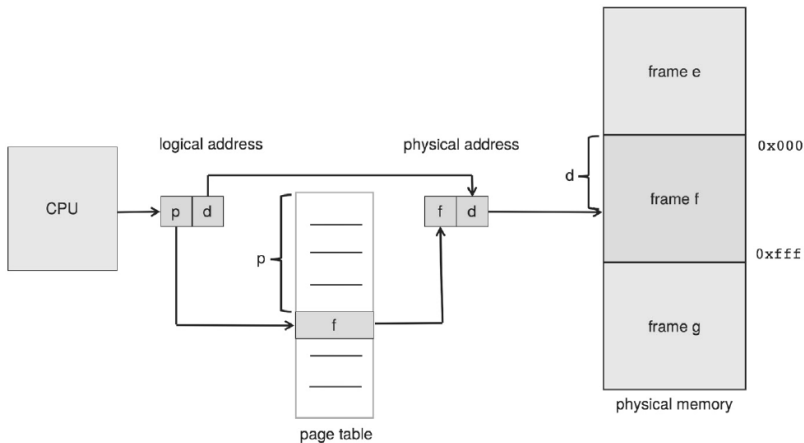
As the page table (PT) accommodates an array of frame addresses, the following properties hold:

- ▶ If the logical address space of process i has k pages (say p_0, p_1, \dots, p_{k-1}), then it must have k valid entries in PT.
- ▶ $PT[i]$ contains the starting physical address of frame in which the process i is stored.
- ▶ If m is the logical address, then the entry $\lfloor m/s \rfloor$ in PT should be consulted, where s denotes the page size.
- ▶ If m is the logical address, then the byte $m \bmod s$ within the frame pointed by $PT[i]$ should be accessed.

Implementation of paging

Every address generated by the CPU is divided into two parts.

1. A page number (p)
2. A page offset (d)



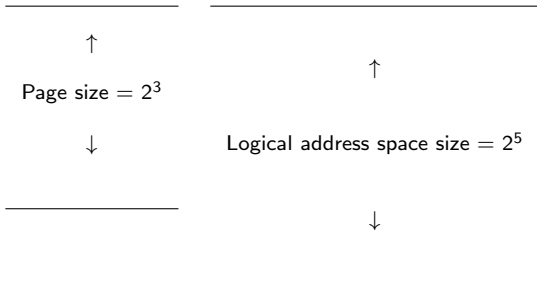
Implementation of paging

Let the size of logical address space is 2^m and the size of page is 2^n , then

1. p denotes the $m - n$ higher order bits of logical address
2. d denotes the n lower order bits of logical address

page number

	p	d
0	00	000
	00	001
	00	010
	00	011
	00	100
	00	101
	00	110
	00	111
1	01	000
...
3	11	111

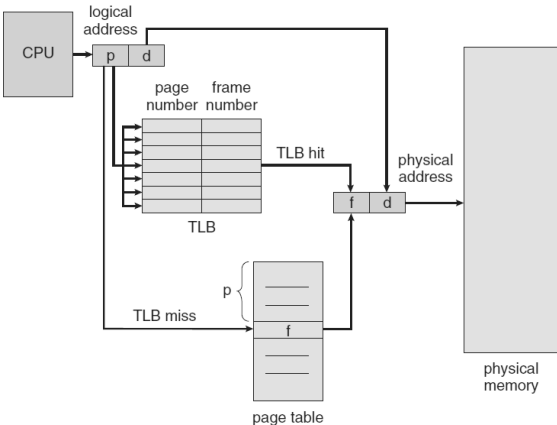


Advantages and disadvantages

- ▶ Memory protection is automatic and process cannot address memory outside its own address space.
- ▶ No external fragmentation happens.
- ▶ Internal fragmentation can happen on half a page per process, on average.
- ▶ Employing small frames causes less fragmentation
- ▶ Employing large frames causes page table overhead but I/O is more efficient.

Paging hardware – With RAM + TLB

Translation look-aside buffers (TLBs) are costly, small, but fast-lookup hardware cache (built using high-speed memory).



Note: TLB has to be flushed on context switch.

Using Translation look-aside buffer (TLB)

The percentage of times a page number of interest is found in the TLB is called the *hit ratio*. A hit ratio of 0.8 (or 80%) signifies that the desired page number can be found in the TLB 80% of the time.

Suppose it takes 10 ns to access memory, then a mapped-memory access takes a 10 ns when the page number is in the TLB. If we fail to find the page number in the TLB then we must first access memory for the page table and frame number (10 ns) and then access the desired byte in memory (10 ns), for a total of 20 ns.

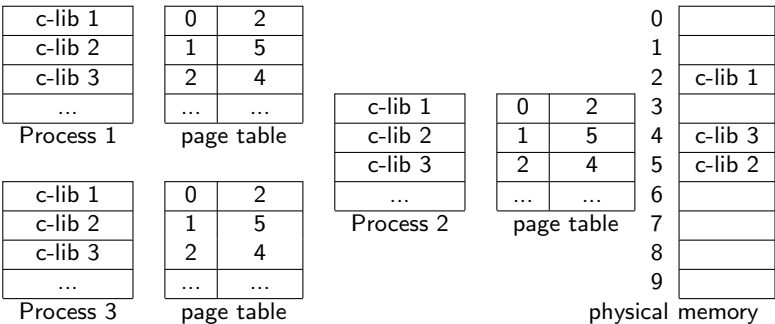
Hence, the *effective memory-access time* turns out to be:

$$0.80 * 10 + 0.20 * 20 = 12 \text{ ns.}$$

Note: We are assuming that a page table lookup takes only one memory access, but it can take more in reality.

Sharing pages

Paging enables the sharing of common code, which is of practical demand in multiprocessor environments. This is primarily used for sharing *reentrant* (read-only) codes associated with heavily used programs.



Here, we consider that “c-lib i” denotes a standard C library.

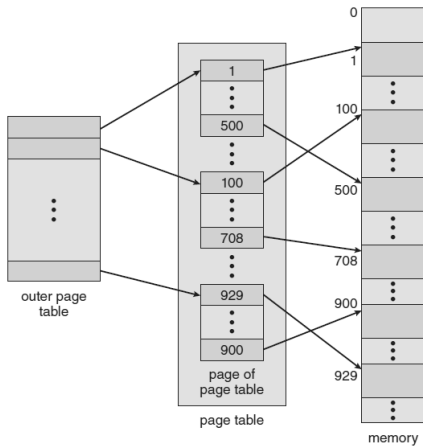
Configuring a page table

The structure of a page table can be configured in different ways as listed below.

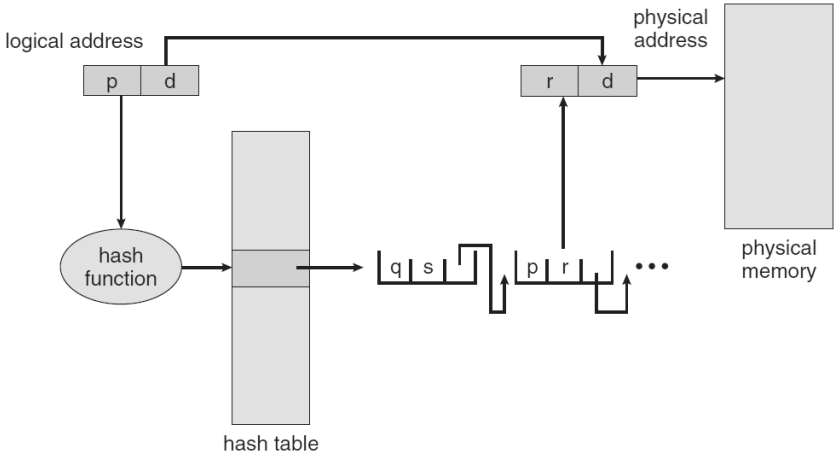
- ▶ Hierarchical page table
- ▶ Hashed page table
- ▶ Inverted page table

Hierarchical page table

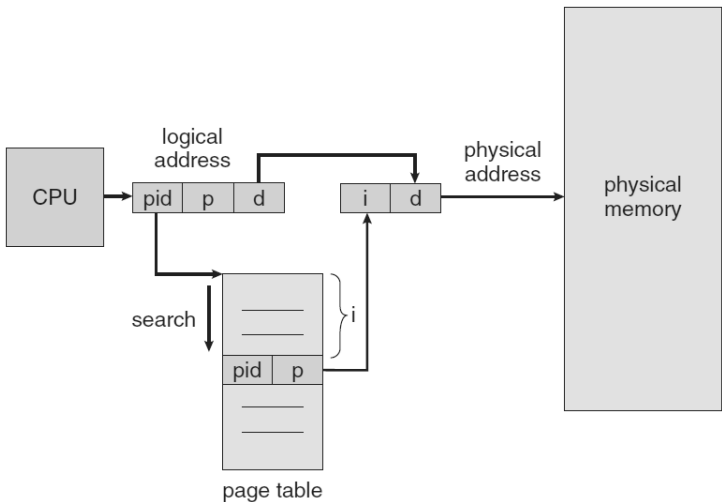
If the page table is large, we might require a secondary (outer) page table to keep track of the pieces of primary (inner) page table. However, many levels of paging may be required for state-of-the-art architectures.



Hashed page table



Inverted page table



Basics of swapping

Motivation:

Consider that a set of processes are resident in memory and occupy all available physical spaces. Now one of the processes forks to create a child.

Swapping allows a process (or a portion of the process) to be temporarily swapped out of memory to a *backing store*, to release memory for another process, and then swapped back in memory for continued execution.

Approaches of swapping

1. Round-robin

- ▶ When a process's quantum expires, it is swapped out, another process is swapped into freed memory.
- ▶ Scheduler allocates next quantum to some other process in memory.

2. Priority-based (roll out, roll in)

- ▶ When higher priority process arrives, lower-priority process is swapped out.
- ▶ When higher priority process finishes, lower priority process can be swapped in.

Note: The time quantum is same for every process.

Performance issues

Problem:

Swapping increases the time of context switch because it involves disk transfer is involved. Consider the example where: Process size = 100K, Transfer rate = 1Mbps.

Hence, swap-out time + swap-in time = $2 * (100K / 1 \text{ Mbps}) = 200\text{ms}$.

Note that this will also include some operational time.

Solution:

The time quantum should be large compared to swap time for a good utilization.

Performance issues

Problem:

If a process is swapped out while waiting for input into buffer in user memory, addresses used by I/O devices may be wrong.

Solution:

The processes with pending I/O should never be swapped.

OR

I/O operations are always done using OS buffers (data can be transferred from OS to user buffer when a process is swapped in).

Performance issues

Problem:

Swapping might cause external fragmentation.

Solution:

Apply compaction as follows.

1. Processes which have to be moved are swapped out
2. Memory is compacted by merging holes
3. Swapped-out processes are swapped in to different memory locations to minimize fragmentation

Background

High-level code:

```

switch(operator){
    case '+':
        // Statement(s)
        break;
    case '-':
        // Statement(s)
        break;
    case '*':
        // Statement(s)
        break;
    case '/':
        // Statement(s)
default:
    // Default statement(s)
}

```

MIPS machine instructions:

???

Background

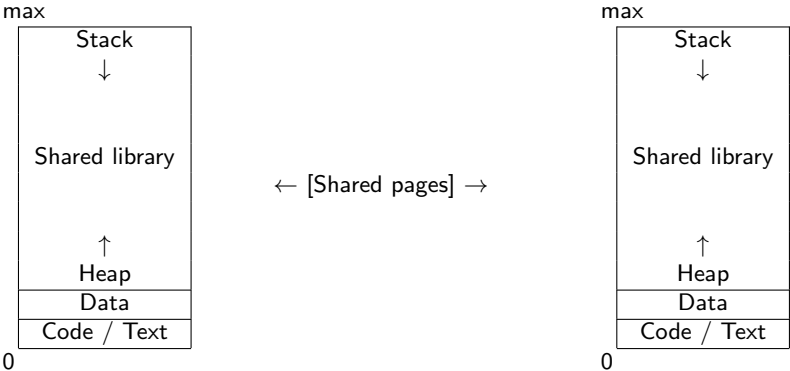
Entire logical address space does not have to be loaded into memory.

- ▶ Some code may be executed rarely (e.g., error handling routines for unusual error conditions, code implementing rarely used features).
- ▶ Arrays/tables may be allocated more memory than required.

Virtual memory is a mechanism to allow execution of processes without requiring the entire process to be in memory.

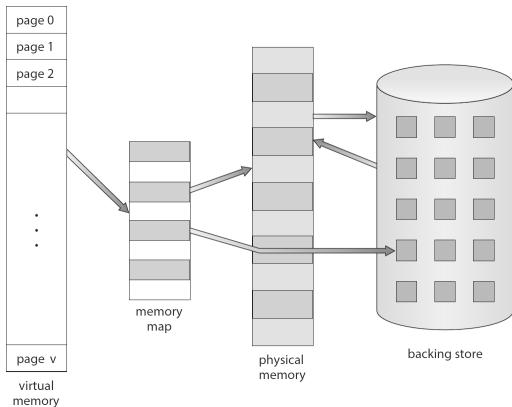
Virtual address space

The heap is allowed to grow upward in memory as it is used for dynamic memory allocation. Similarly, the stack is allowed to grow downward in memory through successive function calls.



Virtual memory

Virtual memory involves the separation of logical memory as perceived by developers from physical memory. This separation allows an extremely large virtual memory to be provided for programmers when only a smaller physical memory is available.



Demand paging

With demand-paged virtual memory, pages are loaded only when they are demanded during program execution. Pages that are never accessed are thus never loaded into physical memory.

