

# Memory Management

---

In a multiprogramming system:

- many processes in memory simultaneously, and every process needs memory for:
  - instructions (“code” or “text”),
  - static data (in program), and
  - dynamic data (heap and stack).
- in addition, operating system itself needs memory for instructions and data.

⇒ must share memory between OS and  $k$  processes.

The memory management subsystem handles:

1. Relocation
2. Allocation
3. Protection
4. Sharing
5. Logical Organisation
6. Physical Organisation

# The Address Binding Problem

---

Consider the following simple program:

```
int x, y;  
x = 5;  
y = x + 3;
```

We can imagine this would result in some assembly code which looks something like:

```
str #5, [Rx]           // store 5 into 'x'  
ldr R1, [Rx]           // load value of x from memory  
add R2, R1, #3         // and add 3 to it  
str R2, [Ry]           // and store result in 'y'
```

where the expression '`[ addr ]`' should be read to mean “the contents of the memory at address `addr`”.

Then the address binding problem is:

*what values do we give  $Rx$  and  $Ry$  ?*

This is a problem because we don't know where in memory our program will be loaded when we run it:

- e.g. if loaded at 0x1000, then x and y might be stored at 0x2000, 0x2004, but if loaded at 0x5000, then x and y might be at 0x6000, 0x6004.

# Address Binding and Relocation

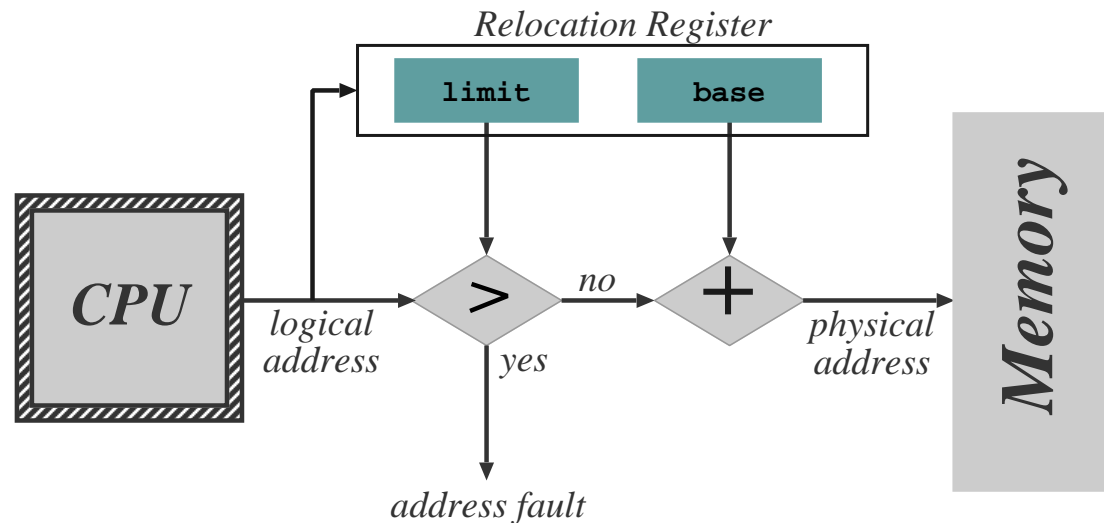
---

To solve the problem, we need to set up some kind of correspondence between “program addresses” and “real addresses”. This can be done:

- at compile time:
  - requires knowledge of absolute addresses; e.g. DOS .com files
- at load time:
  - when program loaded, work out position in memory and update every relevant instruction in code with correct addresses
  - must be done every time program is loaded
  - ok for embedded systems / boot-loaders
- at run-time:
  - get some hardware to automatically translate between program addresses and real addresses.
  - no changes at all required to program itself.
  - most popular and flexible scheme, providing we have the requisite hardware, viz. a memory management unit or MMU.

# Logical vs Physical Addresses

Mapping of logical to physical addresses is done at run-time by Memory Management Unit (MMU), e.g.



1. Relocation register holds the value of the base address owned by the process.
2. Relocation register contents are added to each memory address before it is sent to memory.
3. e.g. DOS on 80x86 — 4 relocation registers, logical address is a tuple  $(s, o)$ .
4. NB: process never sees physical address — simply manipulates logical addresses.
5. OS has privilege to update relocation register.

# Contiguous Allocation

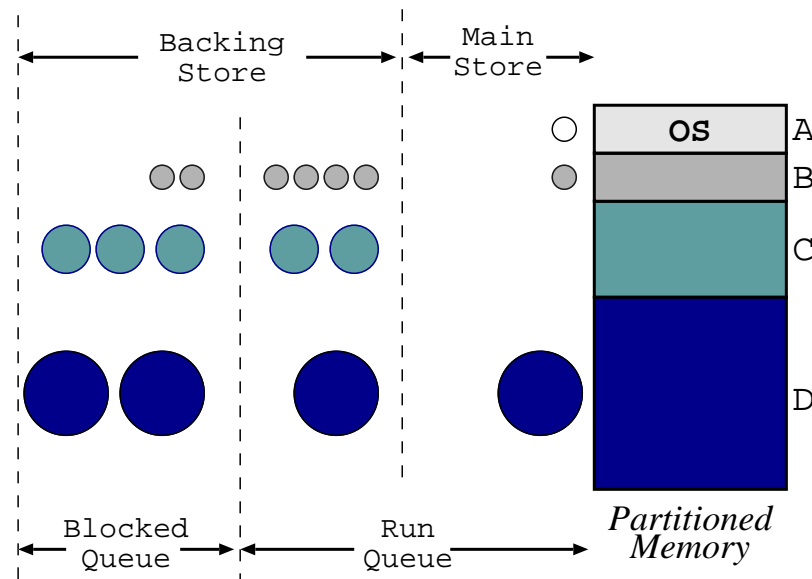
---

Given that we want multiple virtual processors, how can we support this in a single address space?

Where do we put processes in memory?

- OS typically must be in low memory due to location of interrupt vectors
- Easiest way is to statically divide memory into **multiple fixed size partitions**:
  - each partition spans a contiguous range of physical memory
  - bottom partition contains OS, remaining partitions each contain exactly one process.
  - when a process terminates its partition becomes available to new processes.
  - e.g. OS/360 MFT.
- Need to protect OS and user processes from malicious programs:
  - use base and limit registers in MMU
  - update values when a new processes is scheduled
  - NB: solving both relocation and protection problems at the same time!

# Static Multiprogramming



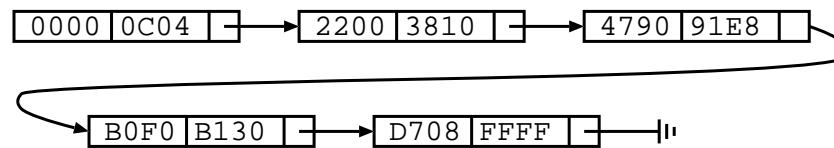
- partition memory when installing OS, and allocate pieces to different job queues.
- associate jobs to a job queue according to size.
- swap job back to disk when:
  - blocked on I/O (assuming I/O is slower than the backing store).
  - time sliced: larger the job, larger the time slice
- run job from another queue while swapping jobs
- e.g. IBM OS/360 MFT, ICL System 4
- **problems:** fragmentation (partition too big), cannot grow (partition too small).

# Dynamic Partitioning

---

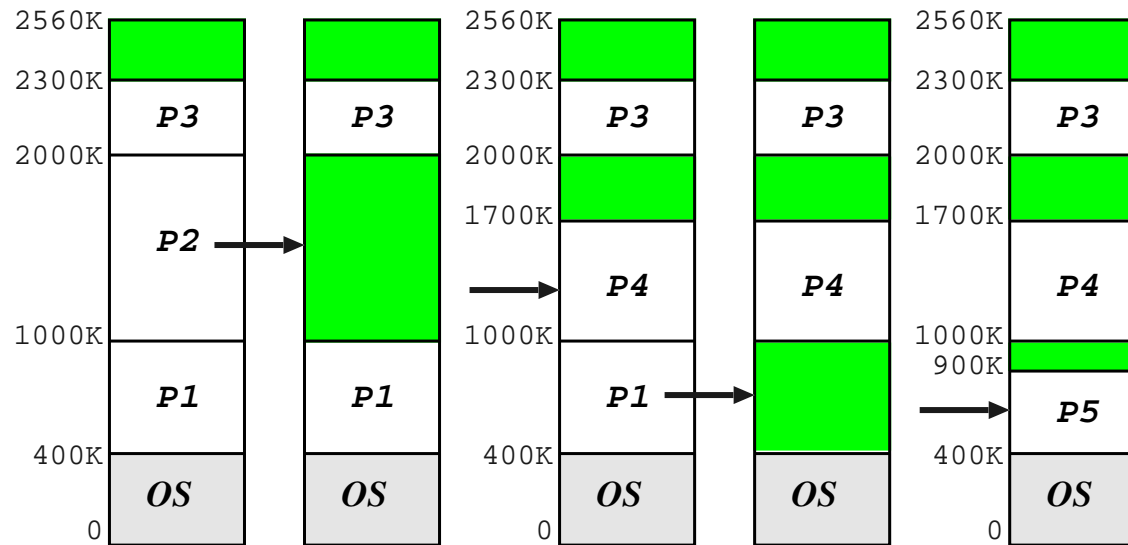
Get more flexibility if allow partition sizes to be dynamically chosen, e.g. OS/360 MVT (“Multiple Variable-sized Tasks”):

- OS keeps track of which areas of memory are available and which are occupied.
- e.g. use one or more *linked lists*:



- When a new process arrives into the system, the OS searches for a hole large enough to fit the process.
- Some algorithms to determine which hole to use for new process:
  - **first fit**: stop searching list as soon as big enough hole is found.
  - **best fit**: search entire list to find “best” fitting hole (i.e. smallest hole which is large enough)
  - **worst fit**: counterintuitively allocate largest hole (again must search entire list).
- When process terminates its memory returns onto the free list, coalescing holes together where appropriate.

# Scheduling Example

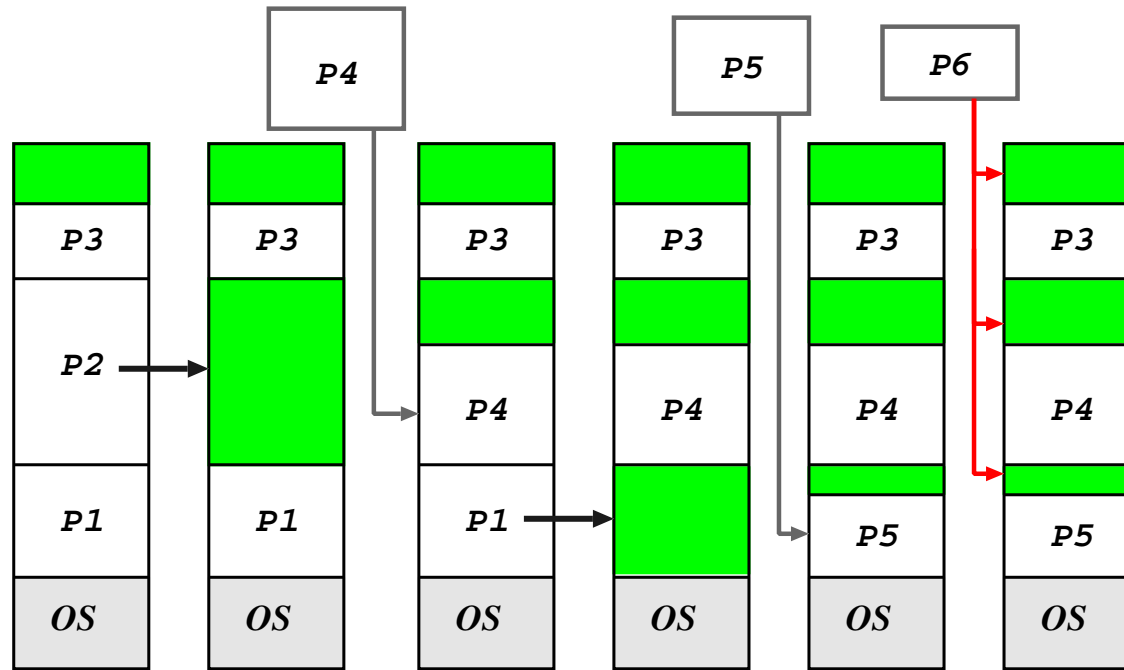


- Consider machine with total of 2560K memory, where OS requires 400K.
- The following jobs are in the queue:

Process	Memory Req'd	Total Execution Time
$P_1$	600K	10
$P_2$	1000K	5
$P_3$	300K	20
$P_4$	700K	8
$P_5$	500K	15



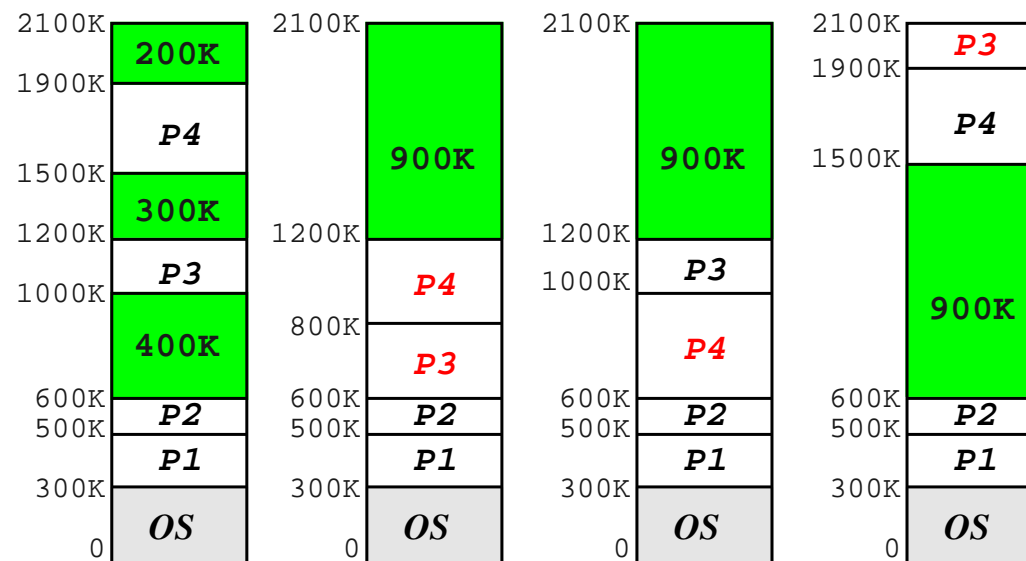
# External Fragmentation



- Dynamic partitioning algorithms suffer from external fragmentation: as processes are loaded they leave little fragments which may not be used.
- External fragmentation exists when the total available memory is sufficient for a request, but is unusable because it is split into many holes.
- Can also have problems with tiny holes

Solution: compact holes periodically.

# Compaction



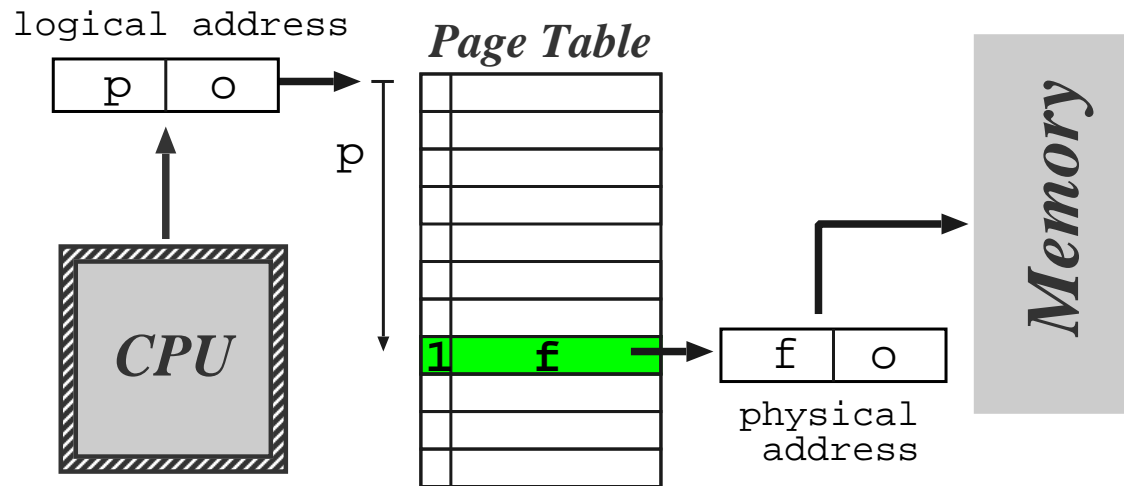
Choosing optimal strategy quite tricky. . .

Note that:

- We require run-time relocation for this to work.
- Can be done more efficiently when process is moved into memory from a swap.
- Some machines used to have hardware support (e.g. CDC Cyber).

Also get fragmentation in *backing store*, but in this case compaction not really a viable option. . .

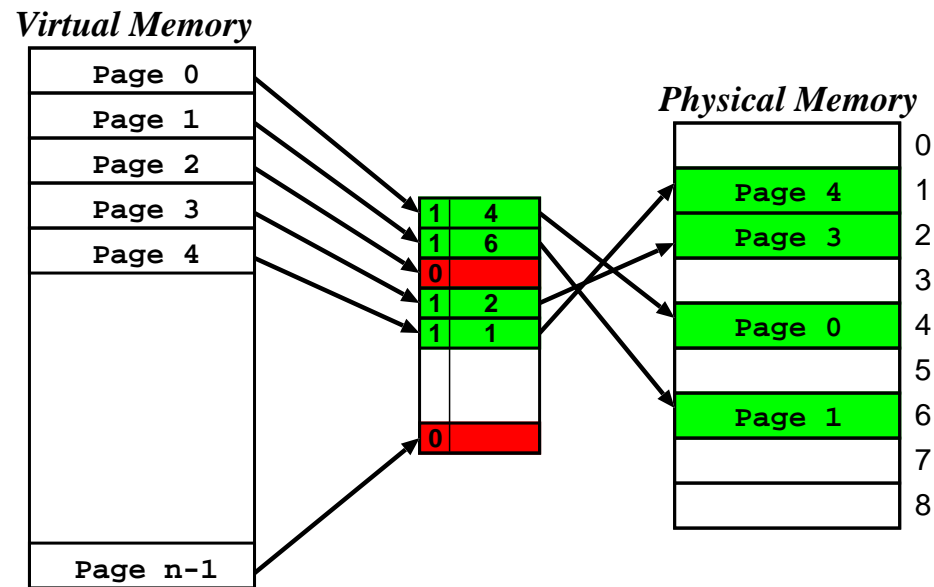
# Paged Virtual Memory



Another solution is to allow a process to exist in non-contiguous memory, i.e.

- divide physical memory into relatively small blocks of fixed size, called **frames**
- divide logical memory into blocks of the same size called **pages**
- (typical page sizes are between 512bytes and 8K)
- each address generated by CPU comprises a page number  $p$  and page offset  $o$ .
- MMU uses  $p$  as an index into a **page table**.
- page table contains associated frame number  $f$
- usually have  $|p| \gg |f| \Rightarrow$  need **valid bit**

# Paging Pros and Cons



- ✓ memory allocation easier.
- ✗ OS must keep page table per process
- ✓ no external fragmentation (in physical memory at least).
- ✗ but get **internal fragmentation**.
- ✓ clear separation between user and system view of memory usage.
- ✗ additional overhead on context switching

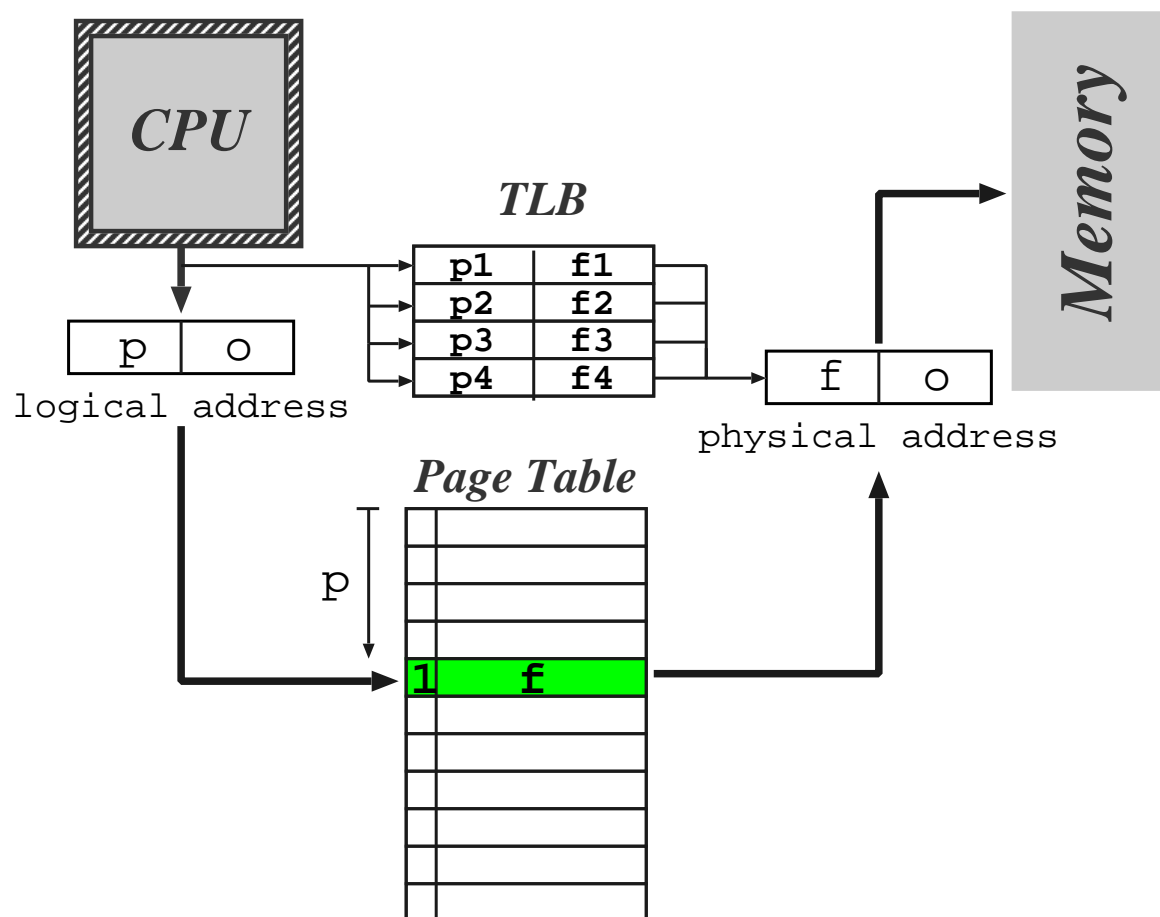
# Structure of the Page Table

---

Different kinds of hardware support can be provided:

- Simplest case: set of dedicated relocation registers
  - one register per page
  - OS loads the registers on context switch
  - fine if the page table is small. . . but what if have large number of pages ?
- Alternatively keep page table in memory
  - only one register needed in MMU (page table base register (PTBR))
  - OS switches this when switching process
- **Problem:** page tables might still be very big.
  - can keep a page table length register (PTLR) to indicate size of page table.
  - or can use more complex structure (see later)
- **Problem:** need to refer to memory *twice* for every ‘actual’ memory reference. . .
  - ⇒ use a **translation lookaside buffer** (TLB)

# TLB Operation



- On memory reference present TLB with logical memory address
- If page table entry for the page is present then get an immediate result
- If not then make memory reference to page tables, and update the TLB

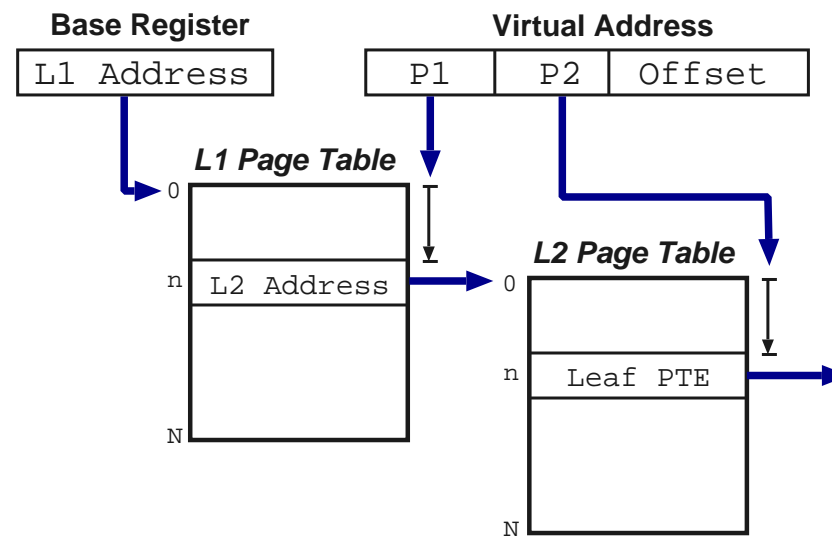
## TLB Issues

---

- Updating TLB tricky if it is full: need to discard something.
  - Context switch may requires TLB flush so that next process doesn't use wrong page table entries.
    - Today many TLBs support **process tags** (sometimes called **address space numbers**) to improve performance.
  - Hit ratio is the percentage of time a page entry is found in TLB
  - e.g. consider TLB search time of  $20ns$ , memory access time of  $100ns$ , and a hit ratio of 80%
- ⇒ assuming one memory reference required for page table lookup, the *effective* memory access time is  $0.8 \times 120 + 0.2 \times 220 = 140ns$ .
- Increase hit ratio to 98% gives effective access time of 122ns — only a 13% improvement.

# Multilevel Page Tables

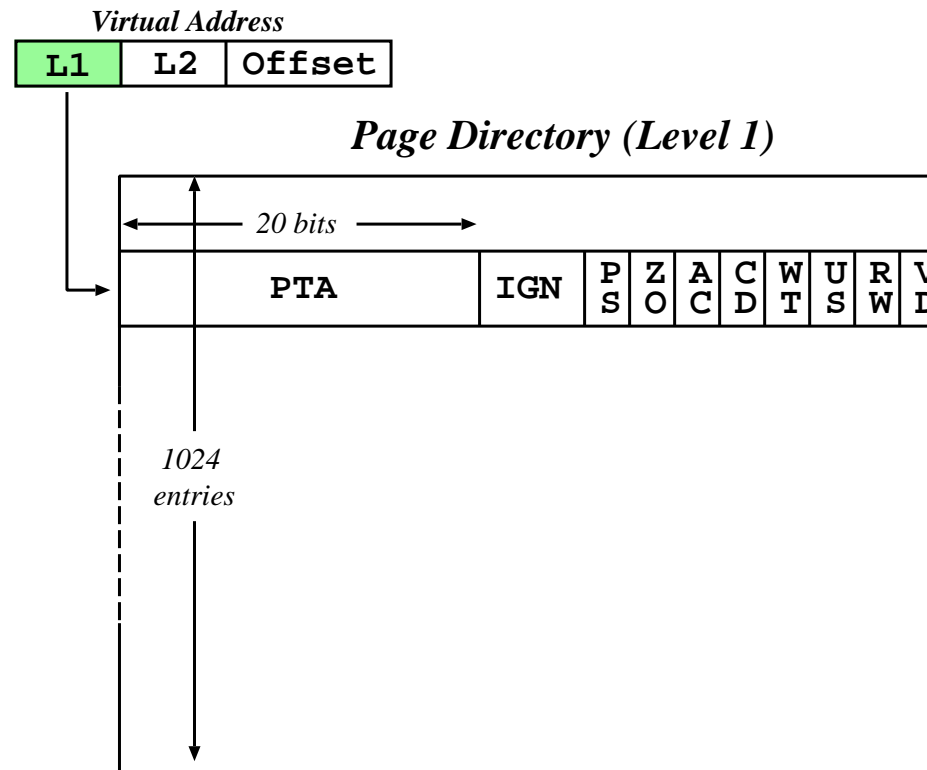
- Most modern systems can support very large ( $2^{32}$ ,  $2^{64}$ ) address spaces.
- Solution – split page table into several sub-parts
- Two level paging – page the page table



- For 64 bit architectures a two-level paging scheme is not sufficient: need further levels (usually 4, or even 5).
- (even some 32 bit machines have  $> 2$  levels, e.g. x86 PAE mode).

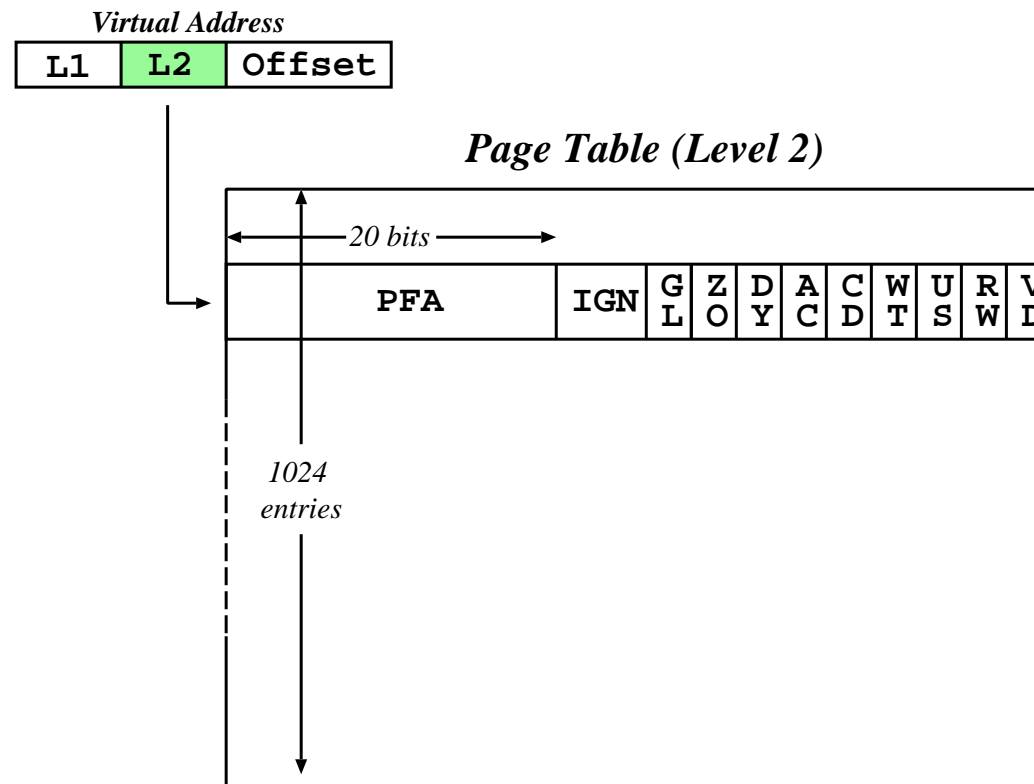


## Example: x86



- Page size 4K (or 4Mb).
- First lookup is in the *page directory*: index using most 10 significant bits.
- Address of page directory stored in internal processor register (cr3).
- Results (normally) in the address of a *page table*.

## Example: x86 (2)



- Use next 10 bits to index into page table.
- Once retrieve page frame address, add in the offset (i.e. the low 12 bits).
- Notice page directory and page tables are exactly one page each themselves.

# Protection Issues

---

- Associate protection bits with each page – kept in page tables (and TLB).
- e.g. one bit for read, one for write, one for execute.
- May also distinguish whether a page may only be accessed when executing in *kernel mode*, e.g. a page-table entry may look like:

Frame Number	K	R	W	X	V
--------------	---	---	---	---	---

- At the same time as address is going through page translation hardware, can check protection bits.
- Attempt to violate protection causes h/w trap to operating system code
- As before, have *valid/invalid* bit determining if the page is mapped into the process address space:
  - if invalid  $\Rightarrow$  trap to OS handler
  - can do lots of interesting things here, particularly with regard to sharing. . .

# Shared Pages

---

Another advantage of paged memory is code/data sharing, for example:

- binaries: editor, compiler etc.
- libraries: shared objects, dlls.

So how does this work?

- Implemented as two logical addresses which map to one physical address.
- If code is *re-entrant* (i.e. stateless, non-self modifying) it can be easily shared between users.
- Otherwise can use [copy-on-write](#) technique:
  - mark page as read-only in all processes.
  - if a process tries to write to page, will trap to OS fault handler.
  - can then allocate new frame, copy data, and create new page table mapping.
- (may use this for lazy data sharing too).

Requires additional book-keeping in OS, but worth it, e.g. over 100MB of shared code on my linux box.

# Virtual Memory

---

- Virtual addressing allows us to introduce the idea of **virtual memory**:
  - already have valid or invalid pages; introduce a new “**non-resident**” designation
  - such pages live on a non-volatile backing store, such as a hard-disk.
  - processes access non-resident memory just as if it were ‘the real thing’.
- Virtual memory (VM) has a number of benefits:
  - **portability**: programs work regardless of how much actual memory present
  - **convenience**: programmer can use e.g. large sparse data structures with impunity
  - **efficiency**: no need to waste (real) memory on code or data which isn’t used.
- VM typically implemented via **demand paging**:
  - programs (executables) reside on disk
  - to execute a process we load pages in *on demand*; i.e. as and when they are referenced.
- Also get *demand segmentation*, but rare.

# Demand Paging Details

---

When loading a new process for execution:

- we create its address space (e.g. page tables, etc), but mark all PTEs as either “invalid” or “non-resident”; and then
- add its process control block (PCB) to the ready-queue.

Then whenever we receive a page fault:

1. check PTE to determine if “invalid” or not
2. if an invalid reference  $\Rightarrow$  kill process;
3. otherwise ‘page in’ the desired page:
  - find a free frame in memory
  - initiate disk I/O to read in the desired page into the new frame
  - when I/O is finished modify the PTE for this page to show that it is now valid
  - restart the process at the faulting instruction

Scheme described above is *pure* demand paging:

- never brings in a page until required  $\Rightarrow$  get lots of page faults and I/O when the process first begins.
- hence many real systems explicitly load some core parts of the process first

# Page Replacement

---

- When paging in from disk, we need a free frame of physical memory to hold the data we're reading in.
- In reality, size of physical memory is limited  $\Rightarrow$ 
  - need to discard unused pages if total demand exceeds physical memory size
  - (alternatively could swap out a whole process to free some frames)
- Modified algorithm: on a page fault we
  1. locate the desired replacement page on disk
  2. to select a free frame for the incoming page:
    - (a) if there is a free frame use it
    - (b) otherwise select a **victim page** to free,
    - (c) write the victim page back to disk, and
    - (d) mark it as invalid in its process page tables
  3. read desired page into freed frame
  4. restart the faulting process
- Can reduce overhead by adding a **dirty bit** to PTEs (can potentially omit step 2c)
- **Question:** how do we choose our victim page?

# Page Replacement Algorithms

---

- First-In First-Out (FIFO)
  - keep a queue of pages, discard from head
  - performance difficult to predict: have no idea whether page replaced will be used again or not
  - discard is independent of page use frequency
  - in general: pretty bad, although very simple.
- Optimal Algorithm (OPT)
  - replace the page which will not be used again for longest period of time
  - can only be done with an oracle, or in hindsight
  - serves as a good comparison for other algorithms
- Least Recently Used (LRU)
  - LRU replaces the page which has not been used for the longest amount of time
  - (i.e. LRU is OPT with -ve time)
  - assumes past is a good predictor of the future
  - **Question:** how do we determine the LRU ordering?



# Implementing LRU

---

- Could try using **counters**
  - give each page table entry a time-of-use field and give CPU a logical clock (e.g. an  $n$ -bit counter)
  - whenever a page is referenced, its PTE is updated to clock value
  - replace page with smallest time value
  - **problem**: requires a search to find minimum value
  - **problem**: adds a write to memory (PTE) on every memory reference
  - **problem**: clock overflow. . .
- Or a **page stack**:
  - maintain a *stack* of pages (a doubly-linked list)
  - update stack on every reference to ensure new (MRU)) page on top
  - discard from bottom of stack
  - **problem**: requires changing 6 pointers per [new] reference
  - possible with h/w support, but slow even then (and extremely slow without it!)
- Neither scheme seems practical on a standard processor  $\Rightarrow$  need another way.

# Approximating LRU (1)

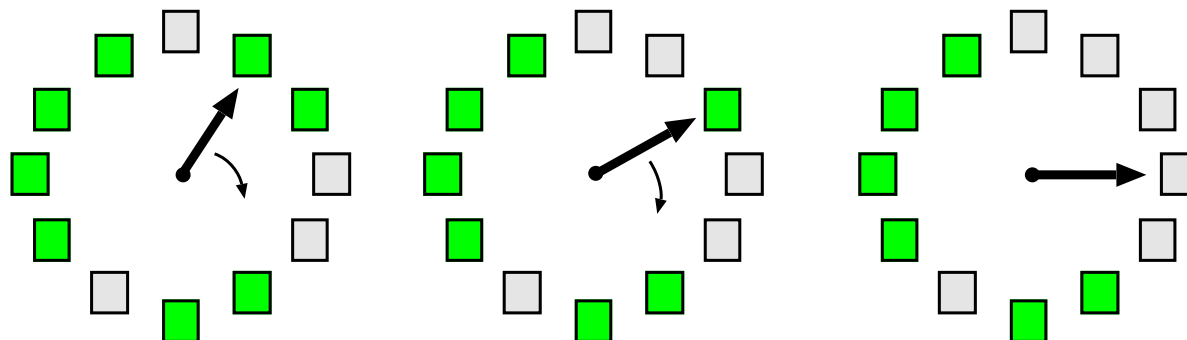
---

- Many systems have a **reference bit** in the PTE which is set by h/w whenever the page is touched
- This allows **not recently used (NRU)** replacement:
  - periodically (e.g. 20ms) clear all reference bits
  - when choosing a victim to replace, prefer pages with clear reference bits
  - if we also have a **modified bit** (or **dirty bit**) in the PTE, we can extend NRU to use that too:

Ref?	Dirty?	Comment
no	no	best type of page to replace
no	yes	next best (requires writeback)
yes	no	probably code in use
yes	yes	bad choice for replacement

- Or can extend by maintaining more history, e.g.
  - for each page, the operating system maintains an 8-bit value, initialized to zero
  - periodically (e.g. every 20ms), shift the reference bit onto most-significant bit of the byte, and clear the reference bit
  - select lowest value page (or one of them) to replace

## Approximating LRU (2)



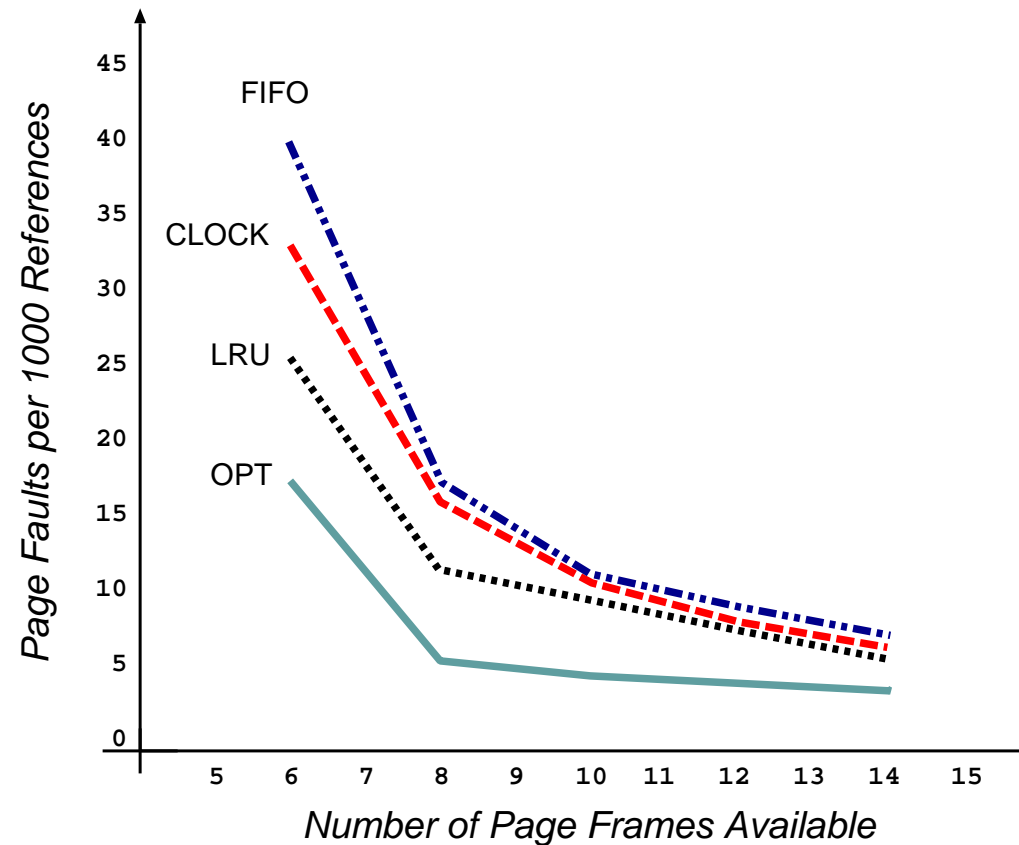
- Popular NRU scheme: **second-chance FIFO**
  - store pages in queue as per FIFO
  - before discarding head, check its reference bit
  - if reference bit is 0, then discard it, otherwise:
    - \* reset reference bit, and add page to tail of queue
    - \* i.e. **give it “a second chance”**
- Often implemented with circular queue and head pointer: then called **clock**.
- If no h/w provided reference bit can emulate:
  - to clear “reference bit”, mark page no access
  - if referenced  $\Rightarrow$  trap, update PTE, and resume
  - to check if referenced, check permissions
  - can use similar scheme to emulate modified bit

# Other Replacement Schemes

---

- **Counting Algorithms:** keep a count of the number of references to each page
  - LFU: replace page with smallest count
  - MFU: replace highest count because low count  $\Rightarrow$  most recently brought in.
- **Page Buffering Algorithms:**
  - keep a min. number of victims in a free pool
  - new page read in before writing out victim.
- **(Pseudo) MRU:**
  - consider access of e.g. large array.
  - page to replace is one application has *just finished with*, i.e. most recently used.
  - e.g. track page faults and look for sequences.
  - discard the  $k^{\text{th}}$  in victim sequence.
- **Application-specific:**
  - stop trying to second guess what's going on.
  - provide hook for app. to suggest replacement.
  - must be careful with denial of service. . .

# Performance Comparison



Graph plots page-fault rate against number of physical frames for a pseudo-local reference string.

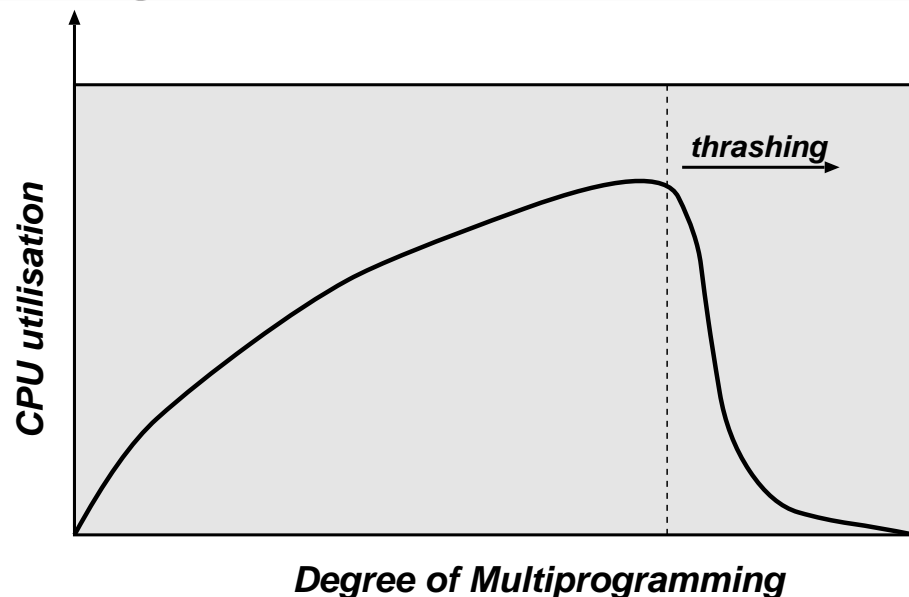
- want to minimise area under curve
- FIFO can exhibit Belady's anomaly (although it doesn't in this case)
- getting frame allocation right has major impact. . .

# Frame Allocation

---

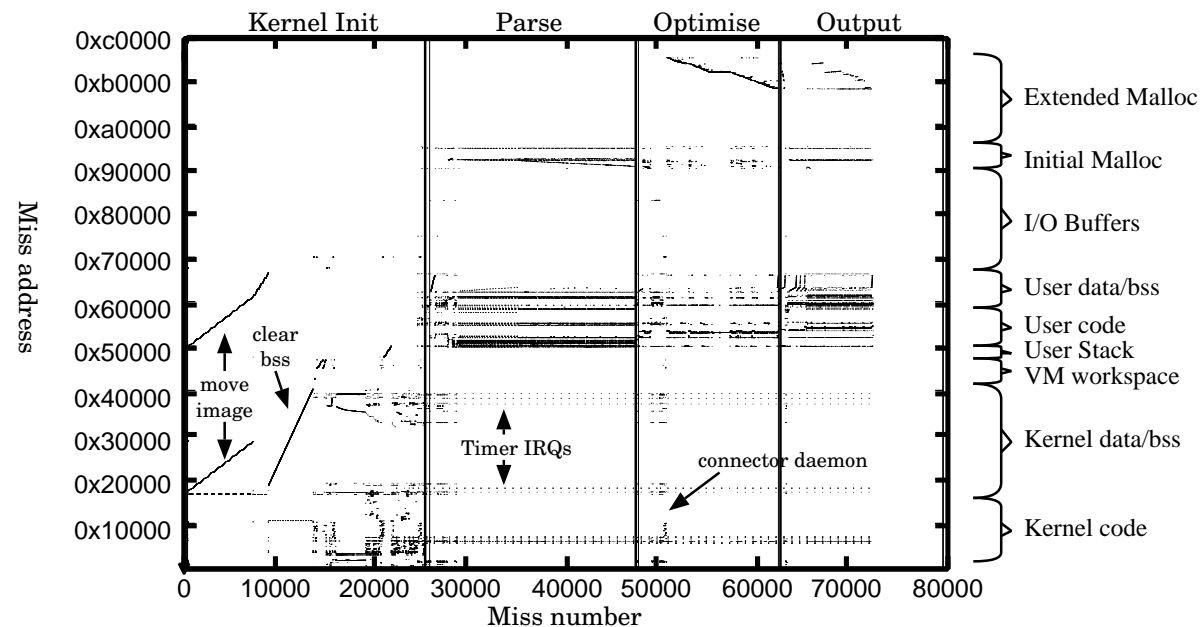
- A certain fraction of physical memory is reserved per-process and for core operating system code and data.
  - Need an *allocation policy* to determine how to distribute the remaining frames.
  - Objectives:
    - Fairness (or proportional fairness)?
      - \* e.g. divide  $m$  frames between  $n$  processes as  $m/n$ , with any remainder staying in the free pool
      - \* e.g. divide frames in proportion to size of process (i.e. number of pages used)
    - Minimize system-wide page-fault rate?  
(e.g. allocate all memory to few processes)
    - Maximize level of multiprogramming?  
(e.g. allocate min memory to many processes)
  - Most page replacement schemes are *global*: all pages considered for replacement.
- ⇒ allocation policy implicitly enforced during page-in:
- allocation succeeds iff policy agrees
  - ‘free frames’ often in use ⇒ steal them!

# The Risk of Thrashing



- As more and more processes enter the system (multi-programming level (MPL) increases), the frames-per-process value can get very small.
- At some point we hit a wall:
  - a process needs more frames, so steals them
  - but the other processes need those pages, so they fault to bring them back in
  - number of runnable processes plunges
- To avoid **thrashing** we must give processes as many frames as they “need”
- If we can't, we need to reduce the MPL: **better page-replacement won't help!**

# Locality of Reference



Locality of reference: in a short time interval, the locations referenced by a process tend to be grouped into a few regions in its address space.

- procedure being executed
- . . . sub-procedures
- . . . data access
- . . . stack variables

**Note:** have locality in both **space** and **time**.



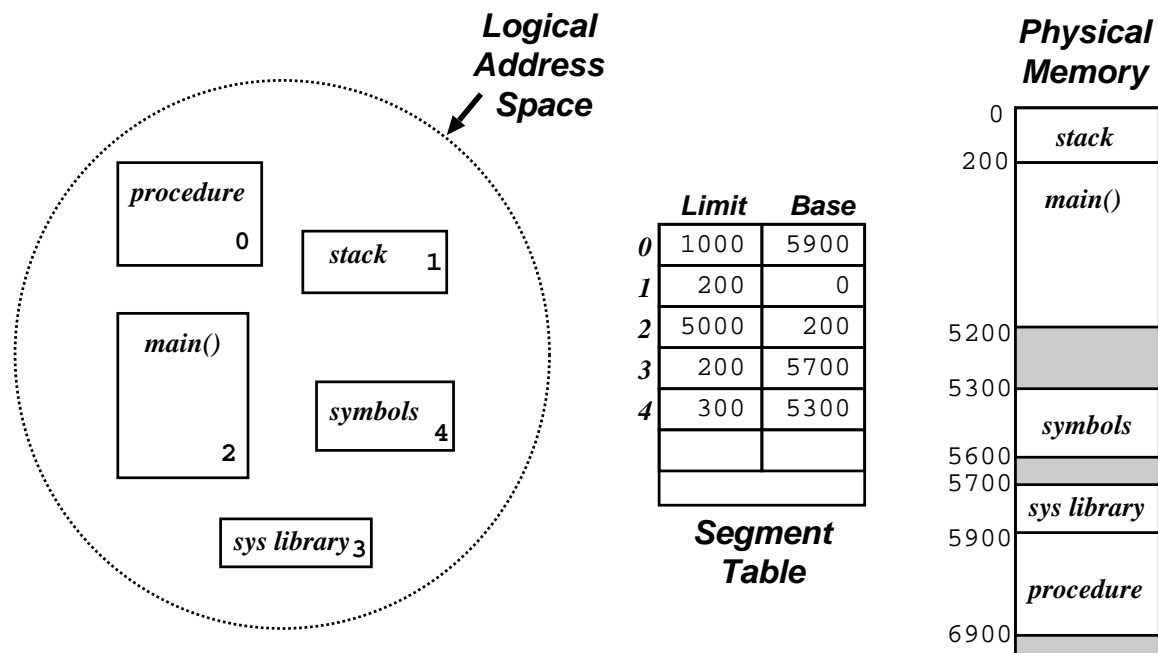
# Avoiding Thrashing

---

We can use the locality of reference principle to help determine how many frames a process needs:

- define the **Working Set** (Denning, 1967)
  - set of pages that a process needs to be resident “the same time” to make any (reasonable) progress
  - varies between processes and during execution
  - assume process moves through *phases*:
    - \* in each phase, get (spatial) locality of reference
    - \* from time to time get *phase shift*
- OS can try to prevent thrashing by ensuring sufficient pages for current phase:
  - sample page reference bits every e.g. 10ms
  - if a page is “in use”, say it’s in the working set
  - sum working set sizes to get total demand  $D$
  - if  $D > m$  we are in danger of thrashing  $\Rightarrow$  suspend a process
- Alternatively use **page fault frequency (PFF)**:
  - monitor per-process page fault rate
  - if too high, allocate more frames to process

# Segmentation



- When programming, a user prefers to view memory as a set of “objects” of various sizes, with no particular ordering
- Segmentation supports this user-view of memory — logical address space is a collection of (typically disjoint) segments.
  - Segments have a name (or a number) and a length.
  - Logical addresses specify segment and offset.
- Contrast with paging where user is unaware of memory structure (one big linear virtual address space, all managed transparently by OS).

# Implementing Segments

---

- Maintain a segment table for each process:

Segment	Access	Base	Size	Others!

- If program has a very large number of segments then the table is kept in memory, pointed to by ST base register STBR
- Also need a ST length register STLR since number of segs used by different programs will differ widely
- The table is part of the process context and hence is changed on each process switch.

Algorithm:

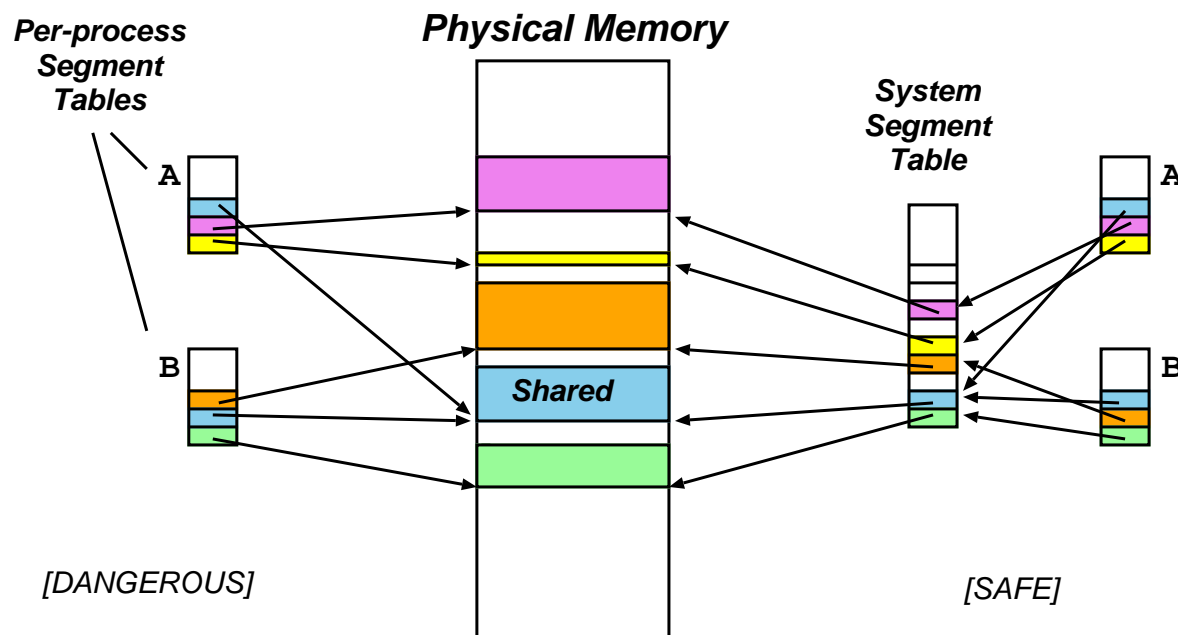
1. Program presents address  $(s, d)$ .  
Check that  $s < \text{STLR}$ . If not, fault
2. Obtain table entry at reference  $s + \text{STBR}$ , a tuple of form  $(b_s, l_s)$
3. If  $0 \leq d < l_s$  then this is a valid address at location  $(b_s, d)$ , else fault

# Sharing and Protection

---

- Big advantage of segmentation is that **protection is per segment**; i.e. corresponds to logical view (and programmer's view)
- Protection bits associated with each ST entry checked in usual way
  - e.g. instruction segments (should be non-self modifying!) can be protected against writes
  - e.g. place each array in own seg  $\Rightarrow$  array limits checked by h/w
- Segmentation also facilitates sharing of code/data
  - each process has its own STBR/STLR
  - sharing enabled when two processes have identical entries
  - for data segments can use copy-on-write as per paged case.
- Several subtle caveats exist with segmentation — e.g. jumps within shared code.

# Sharing Segments



## Sharing segments: dangerously (lhs) and safely (rhs)

- wasteful (and dangerous) to store common information on shared segment in each process segment table
  - want *canonical* version of segment info
- assign each segment a unique System Segment Number (SSN)
- process segment table maps from a Process Segment Number (PSN) to SSN