# Paging (cont.)

30/10-2003
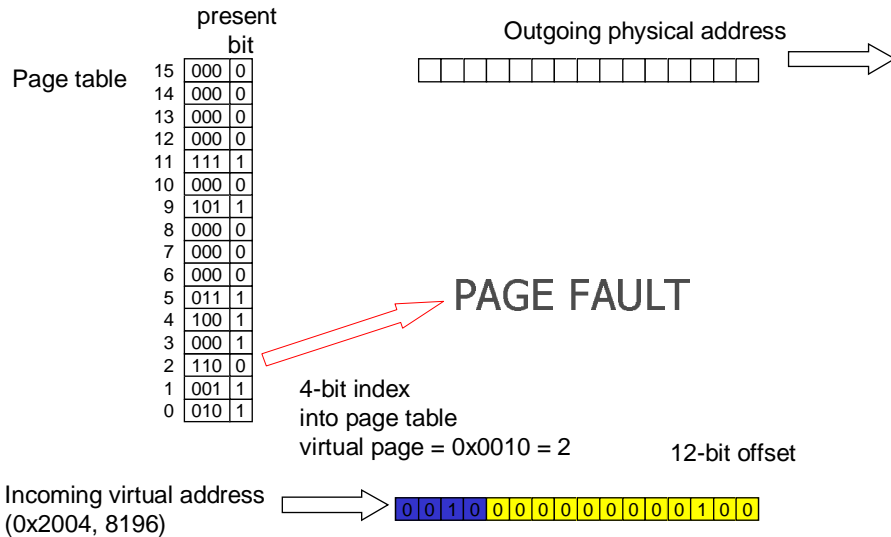
## Pål Halvorsen

(including slides from Andrew Tanenbaum)

---

# Memory Lookup

present
bit

Page table

Outgoing physical address

(0x6004, 24580)

| | | |
|---|---|---|
| 15 | 000 | 0 |
| 14 | 000 | 0 |
| 13 | 000 | 0 |
| 12 | 000 | 0 |
| 11 | 111 | 1 |
| 10 | 000 | 0 |
| 9 | 101 | 1 |
| 8 | 000 | 0 |
| 7 | 000 | 0 |
| 6 | 000 | 0 |
| 5 | 011 | 1 |
| 4 | 100 | 1 |
| 3 | 000 | 1 |
| 2 | 110 | 1 |
| 1 | 001 | 1 |
| 0 | 010 | 1 |

4-bit index
into page table
virtual page = 0x0010 = 2

12-bit offset

Incoming virtual address
(0x2004, 8196)

# Memory Lookup

Page table

present bit

| | | |
|---|---|---|
| 15 | 000 | 0 |
| 14 | 000 | 0 |
| 13 | 000 | 0 |
| 12 | 000 | 0 |
| 11 | 111 | 1 |
| 10 | 000 | 0 |
| 9 | 101 | 1 |
| 8 | 000 | 0 |
| 7 | 000 | 0 |
| 6 | 000 | 0 |
| 5 | 011 | 1 |
| 4 | 100 | 1 |
| 3 | 000 | 1 |
| 2 | 110 | 0 |
| 1 | 001 | 1 |
| 0 | 010 | 1 |

Outgoing physical address

PAGE FAULT

4-bit index
into page table
virtual page = 0x0010 = 2

12-bit offset

Incoming virtual address
(0x2004, 8196)

0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0

---

# Page Fault Handling

1.  Hardware traps to the kernel saving program counter and process state information

2.  Save general registers and other volatile information

3.  OS discover the page fault and tries to determine which virtual page is requested

4.  OS checks if the virtual page is valid and if protection is consistent with access

5.  Select a page to be replaced

6.  Check if selected page frame is "dirty", i.e., updated

7.  When selected page frame is ready, the OS finds the disk address where the needed data is located and schedules a disk operation to bring in into memory

8.  A disk interrupt is executed indicating that the disk I/O operation is finished, the page tables are updated, and the page frame is marked "normal state"

9.  Faulting instruction is backed up and the program counter is reset

10. Faulting process is scheduled, and OS returns to routine that made the trap to the kernel

11. The registers and other volatile information is restored and control is returned to user space to continue execution as no page fault had occured

# Page Replacement Algorithms

- Page fault $\rightarrow$ OS has to select a page for replacement

    - Modified page $\rightarrow$ write back to disk

    - **Not** modified page $\rightarrow$ just overwrite with new data

- How do we decide which page to replace?
    $\rightarrow$ determined by the **page replacement algorithm**
    $\rightarrow$ several algorithms exist:
    - Random
    - Other algorithms take into acount usage, age, etc.
      (e.g., FIFO, not recently used, least recently used, second chance, clock, ...)
    - which is best???

# Optimal

- **Best possible** page replacement algorithm:
    - When a page fault occurs, all pages in memory are labeled with the number of instructions that will be executed before this page will be used again
    - The page with **most** instructions before reuse is replaced

- Easy to describe, but impossible to implement
  (OS cannot look into the future)

- Estimate by logging page usage on previous runs of process

- Useful to evaluate other page replacement algorithm

# Not Recently Used (NRU)

- Two status bits associated with each page:
    - R $\rightarrow$ page referenced (read or written)
    - M $\rightarrow$ page modified (written)

- Pages belong to one of four set of pages according to the status bits:
    - **Class 0**: not referenced, not modified  (R=0, M=0)
    - **Class 1**: not referenced, modified  (R=0, M=1)
    - **Class 2**: referenced, not modified  (R=1, M=0)
    - **Class 3**: referenced, modified  (R=1, M=1)

- NRU removes a page at random from the lowest numbered, non-empty class

- Low overhead

---

# First In First Out (FIFO)

- All pages in memory are maintained in a list sorted by age
- FIFO replaces the oldest page, i.e., the first in the list

**Reference string:** A B C D A E F G H I A J

Now the buffer is full, a page fault results in a replacement

Page most
recently loaded
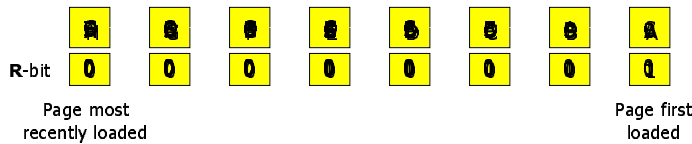
Page first loaded, i.e.,
FIRST REPLACED

- Low overhead
- FIFO is rearly used in its pure form

# Second Chance

- Modification of FIFO
- *R* bit: when a page is referenced again, the R bit is set, and the page will be treated as a newly loaded page

**Reference string:**  A   B   C   Ⓓ Ⓐ Ⓔ Ⓕ Ⓖ Ⓗ Ⓘ

Page I will be inserted, find a page to page out by looking at the first page loaded:
Page A's R-bit = 0, replace, shift chain left, and insert last in the chain (B)
Now the buffer is full, next page causes a replacement
-if R-bit = 0, replace
-if R-bit = 1, clean R-bit, move page last, and finally look at the new first page



**R-bit**

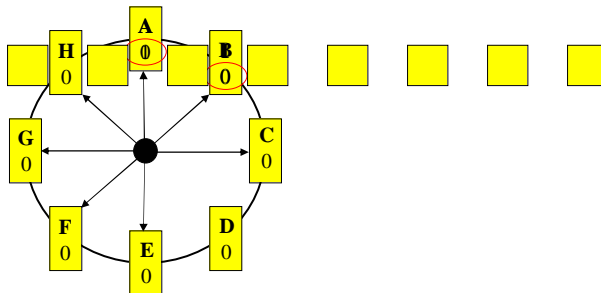| Page most recently loaded | | | | | | | Page first loaded |

- *Second chance* is a reasonable algorithm, but inefficient because it is moving pages around the list

---

# Clock

- More efficient way to implement *Second Chance*
- Circular list in form of a clock
- Pointer to the oldest page:
  - R-bit = 0 $\rightarrow$ replace and advance pointer
  - R-bit = 1 $\rightarrow$ set R-bit to 0, advance pointer until R-bit = 0, replace and advance pointer

**Reference string:**  A   B   C   Ⓓ Ⓐ Ⓔ Ⓕ Ⓖ Ⓗ Ⓘ

# Least Recently Used (LRU)

- Replace the page that has the longest time since last reference

- Based on the observation that
  *pages that are heavily used in the last few instructions will probably be used again in the next few instructions*

- Several ways to implement this algoithm

---

# Least Recently Used (LRU)

- LRU as a linked list:

**Reference string:**  A   B   C   D   A   E   F   G   H   A   C   I

Now the buffer is full. Next page fault will trigger replacement
(most recently used)

Page most
recently used

Page least
recently used

- **Expensive** – maintaining an ordered list of all pages in memory:
  - most recently used at front, least at rear
  - update this list <u>every memory reference</u> !!

# Least Recently Used (LRU)

- LRU by using *aging*:
  - "reference counter" for each page
  - after a clock tick:
    - shift bits in the reference counter to the right (rightmost bit is deleted)
    - add a page's referece bit in front of the reference counter (left)
  - page with *lowest* counter is replaced

| | Clock tick 0<br>1 0 1 0 1 1 | | Clock tick 1<br>1 1 0 1 0 0 | | Clock tick 2<br>1 1 0 1 0 1 | | Clock tick 3<br>1 0 0 0 1 0 | | Clock tick 4<br>0 1 1 0 0 0 |
|---|---|---|---|---|---|---|---|---|---|
| **1** | 00000000 | **1** | 10000000 | **1** | 11000000 | **1** | 11100000 | **1** | 11110000 | **1** | 01111000 |
| **2** | 00000000 | **2** | 00000000 | **2** | 10000000 | **2** | 11000000 | **2** | 01100000 | **2** | 10110000 |
| **3** | 00000000 | **3** | 10000000 | **3** | 01000000 | **3** | 00100000 | **3** | 00010000 | **3** | 10001000 |
| **4** | 00000000 | **4** | 00000000 | **4** | 10000000 | **4** | 11000000 | **4** | 01100000 | **4** | 00110000 |
| **5** | 00000000 | **5** | 10000000 | **5** | 01000000 | **5** | 00100000 | **5** | 10010000 | **5** | 01001000 |
| **6** | 00000000 | **6** | 10000000 | **6** | 01000000 | **6** | 10100000 | **6** | 01010000 | **6** | 00101000 |

---

# Least Recently Used (LRU)

- LRU as a matrix:
  - $N$ pages $\rightarrow$ $N$ x $N$ matrix
  - Page $N$ is referenced $\rightarrow$ row $N$ is set (1)
  - $\rightarrow$ column $N$ is cleared (0)
  - Replace page with *lowest row* value

"Page frame" string: ①②③④③②①④

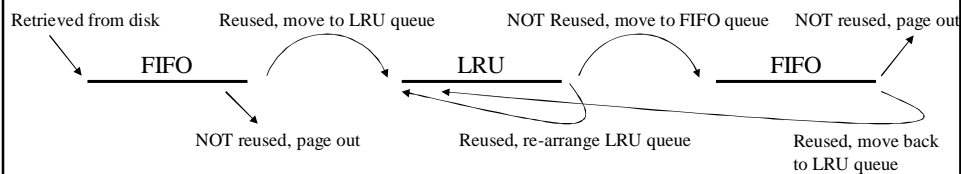| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **1** | 0 | 1 | 1 | 0 |
| **2** | 0 | 0 | 1 | 0 |
| **3** | 0 | 0 | 0 | 0 |
| **4** | 1 | 1 | 1 | 0 |

# Counting Algorithms

- LRU by using a reference counter
  - clear the counter when the page is referenced (counter = 0)
  - increase all counters each clock tick
  - replace the page with the *highest* counter

- *Not/Least Frequently Used* (N/LFU)
  - counter initially 0
  - increase the page's counter *only if* it has been referenced during this clock tick
  - replace the page with *lowest* counter

- *Most Frequently Used* (MFU)
  - counter as LFU
  - replace the page with the *highest* counter
    (assuming low counters mean new, fresh pages)

---

# LRU-K & 2Q

- **LRU-K:** bases page replacement in the *last K references* on a page [O'Neil et al. 93]

- **2Q:** uses 3 queues to hold much referenced and popular pages in memory [Johnson et al. 94]
  - 2 FIFO queues for seldom referenced pages
  - 1 LRU queue for much referenced pages

Retrieved from disk     Reused, move to LRU queue     NOT Reused, move to FIFO queue     NOT reused, page out

FIFO       LRU       FIFO

NOT reused, page out     Reused, re-arrange LRU queue     Reused, move back to LRU queue
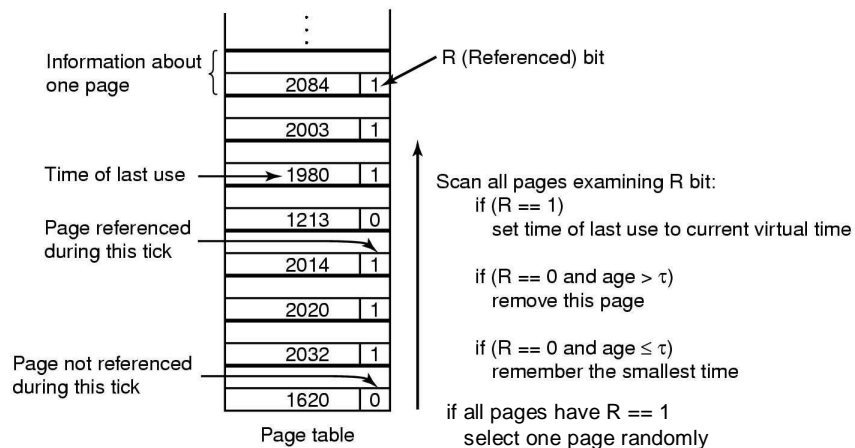
# Working Set Model

- Working set:
  set of pages which a process is currently using


- Working set model:
  paging system tries to keep track of each process' working set and makes sure that these pages is in memory before letting the process run
  $\rightarrow$ reduces page fault rate (prepaging)


- Defining the working set:
  - set of pages used in the last $k$ memory references (must count backwards)
  - approximation is to use all references used in the last XX instructions


# Working Set Page Replacement Algorithm

$\tau$ - time period to calculate the WS over
**age** - virtual time - last reference time

| 2204 | Current virtual time

Information about one page
R (Referenced) bit

| 2084 | 1 |
| 2003 | 1 |

Time of last use → | 1980 | 1 |

Page referenced during this tick
| 1213 | 0 |
| 2014 | 1 |
| 2020 | 1 |

Page not referenced during this tick
| 2032 | 1 |
| 1620 | 0 |

Page table

Scan all pages examining R bit:
    if (R == 1)
        set time of last use to current virtual time

    if (R == 0 and age > $\tau$)
        remove this page

    if (R == 0 and age $\leq \tau$)
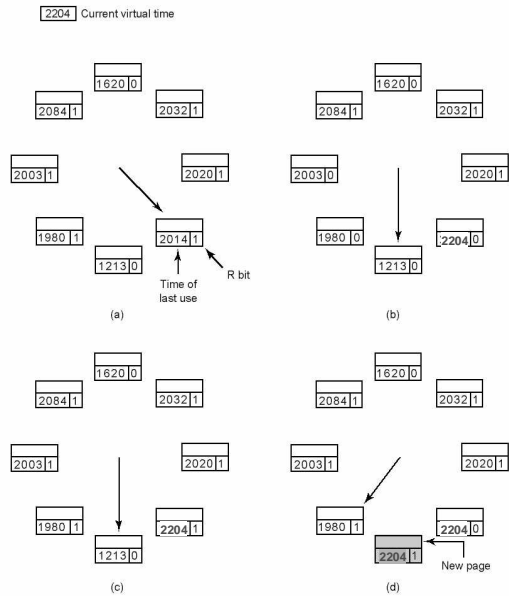        remember the smallest time

if all pages have R == 1
    select one page randomly

- Expensive - must search the whole page table

## WSClock Page Replacement Algorithm

- Organize each page table entry as a clock
- As with clock - the page pointed to is examined first
  - R = 1:
    clear bit, set virtual time, continue **(b)**
  - R = 0: **(c)**
    - age < τ : continue to next
    - age > τ :
      - if page clean, replace **(d)**
      - othervice, write to disk and continue to next
- If all pointer comes back to start
  - writes are scheduled to clean pages (find first)
  - no scheduled writes (all in WS), several option
    - remove first clean
    - remove oldest
    - ...



---

## Belady's Anomaly

- **Question:** *the more page frames, the fewer page faults?*

- Belady's anomaly gives a counter example using the FIFO replacement algorithm and buffers of 3 and 4 page frames:



| Refernce string: | 0 | 1 | 2 | 3 | 0 | 1 | 4 | 0 | 1 | 2 | 3 | 4 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Youngest page: | 0 | 1 | 2 | 3 | 0 | 1 | 4 | 4 | 4 | 2 | 3 | 3 | |
| | | 0 | 1 | 2 | 3 | 0 | 1 | 1 | 1 | 4 | 2 | 2 | ⇒ 9 page faults |
| Oldest page: | | | 0 | 1 | 2 | 3 | 0 | 0 | 0 | 1 | 4 | 4 | |
| | P | P | P | P | P | P | P | | | P | P | | |

| Youngest page: | 0 | 1 | 2 | 3 | 3 | 3 | 4 | 0 | 1 | 2 | 3 | 4 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 2 | 2 | 3 | 4 | 0 | 1 | 2 | 3 | ⇒ 10 page faults |
| | | | 0 | 1 | 1 | 1 | 2 | 3 | 4 | 0 | 1 | 2 | |
| Oldest page: | | | | 0 | 0 | 0 | 1 | 2 | 3 | 4 | 0 | 1 | |
| | P | P | P | P | | | P | P | P | P | P | P | |

# Stack Algorithms

- Observation:
  Every process generates a sequence of memory references as it runs where each memory reference corresponds to a virtual page

- Reference string: ordered list of page numbers (process' memory accesses)

- A paging system can be characrized by 3 items:
  1. Reference string of the executing process
  2. Page replacement algorithm
  3. Number of page frames available in memory
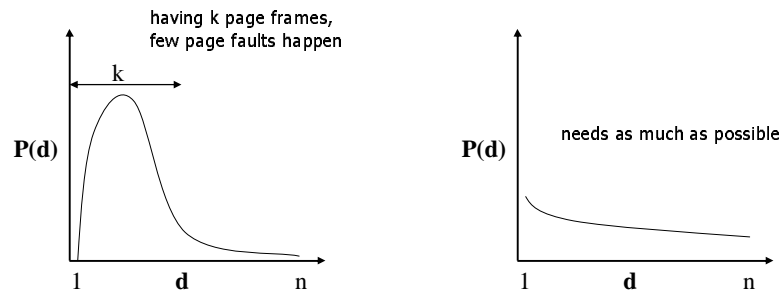
---

# Stack Algorithms

| Reference string: | 0 | 2 | 1 | 3 | 5 | 4 | 6 | 3 | 7 | 4 | 7 | 3 | 3 | 5 | 5 | 3 | 1 | 1 | 1 | 7 | 2 | 3 | 4 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 2 | 1 | 3 | 5 | 4 | 6 | 3 | 7 | 4 | 7 | 3 | 3 | 5 | 5 | 3 | 1 | 1 | 1 | 7 | 2 | 3 | 4 | 1 |
| | | 0 | 2 | 1 | 3 | 5 | 4 | 6 | 3 | 7 | 4 | 7 | 7 | 3 | 3 | 5 | 3 | 3 | 3 | 1 | 7 | 2 | 3 | 4 |
| | | | 0 | 2 | 1 | 3 | 5 | 4 | 6 | 3 | 3 | 4 | 4 | 7 | 7 | 7 | 5 | 5 | 5 | 3 | 1 | 7 | 2 | 3 |
| | | | | 0 | 2 | 1 | 3 | 5 | 4 | 6 | 6 | 6 | 6 | 4 | 4 | 4 | 7 | 7 | 7 | 5 | 3 | 1 | 7 | 2 |
| | | | | | 0 | 2 | 1 | 1 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 4 | 4 | 4 | 4 | 5 | 5 | 1 | 7 |
| | | | | | | 0 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 6 | 6 | 6 | 6 | 4 | 4 | 5 | 5 |
| | | | | | | | 0 | 0 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 6 | 6 | 6 | 6 |
| | | | | | | | | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Page faults: | P | P | P | P | P | P | P | | P | | | | | P | | | P | | | | P | | P | P |
| Distance string: | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ | 4 | ∞ | 4 | 2 | 3 | 1 | 5 | 1 | 2 | 6 | 1 | 1 | 4 | 7 | 4 | 6 | 5 |

**Distance string:** a page reference is denoted by the distance from the top of the stack, i.e., the number of references since the page was referenced last

# Distance String Properties

- The statistical properties of the distance string may show expexted performance of the algorithm (probability density function)

having k page frames,
few page faults happen

k

$P(d)$

1    $d$    n

needs as much as possible

$P(d)$

1    $d$    n

---

# Predicting Page Fault Rates

- $F_m = \Sigma C_k + C_\infty$ , $m + 1 \le k \le n$
  $F_m -$ # page faults occuring with a given distance string and $m$ page frames
  $C_k -$ # occurrences of $k$, i.e., # times we have a distance $k$ since last reference

- Example - computation of the page fault rate:

Distance string: ∞ ∞ ∞ ∞ ∞ ∞ ∞ 4 ∞ 4 2 3 1 5 1 2 6 1 1 4 7 4 6 5

$C_1 = 4$ ← # times 1 occurs in the distance string    $F_1 = 20$ ← $C_2 + C_3 + C_4 + C_5 + C_6 + C_7 + C_\infty$
$C_2 = 2$ ← # times 2 occurs in the distance string    $F_2 = 18$ ← $C_3 + C_4 + C_5 + C_6 + C_7 + C_\infty$
$C_3 = 1$ ← # times 3 occurs in the distance string    $F_3 = 17$ ← $C_4 + C_5 + C_6 + C_7 + C_\infty$
$C_4 = 4$    $F_4 = 13$ ← $C_5 + C_6 + C_7 + C_\infty$
$C_5 = 2$    $F_5 = 11$ ← $C_6 + C_7 + C_\infty$
$C_6 = 2$    $F_6 = 9$ ← $C_7 + C_\infty$
$C_7 = 1$    $F_7 = 8$ ← $C_\infty$
$C_\infty = 8$    $F_8 = 8$ ← $C_\infty$

# Locality

- Reference locality:

  - Time:
    pages that are referenced in the last few instructions will probably be referenced again in the next few instructions

  - Space:
    pages that are located close to the page being referenced will probably also be referenced

# Demand Paging Versus Prepaging

- Demand paging:
  pages are loaded on demand, i.e., *after* a process needs it
    - Should be used if *we have no knowledge* about future references
    - Each page is loaded separatly from disk, i.e., results in many disk accesses

- Prepaging:
  prefetching data in advance, i.e., before use
    - Should be used if *we have knowledge* about future references
    - # page faults is reduced, i.e., page in memory when needed by a process
    - # disk accesses can be reduced by loading several pages in *one* I/O-operation

# Allocation Policies

- How should memory be allocated among the competing runnable processes?

- Equal allocation:
  - all processes get the same amount of pages

- Proportional allocation:
  - amount of pages is depending on process size

# Allocation Policies

- Local page replacement:
  consider only pages of own process when replacing a page
  - corresponds to equal allocation
  - can cause thrashing
  - multiple, identical pages in memory

- Global page replacemet:
  consider all pages in memory when replacing a page
  - corresponds to proportional allocation
  - better performance in general
  - monitoring of working set size and aging bits
  - data sharing

# Allocation Policies
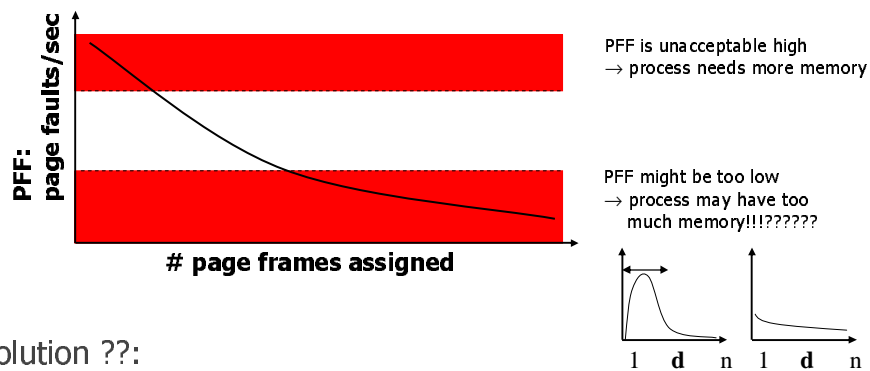
- Example: *local* versus *global* replacement
  insert page **A5** using age replacement

| | Age | | | Age | | | Age |
|---|---|---|---|---|---|---|---|
| A1 | 10 | | A1 | 10 | | A1 | 10 |
| A2 | 7 | **Local replacement**: | A2 | 7 | | A2 | 7 |
| A3 | 4 | Replace the oldest | A5 | 13 | | A3 | 4 |
| A4 | 11 | of A's pages | A4 | 11 | | A4 | 11 |
| B1 | 6 | | B1 | 6 | | B1 | 6 |
| B2 | 12 | | B2 | 12 | | B2 | 12 |
| B3 | 1 | **Global replacement**: | B3 | 1 | | A5 | 13 |
| B4 | 3 | Replace the oldest page | B4 | 3 | | B4 | 3 |
| C1 | 8 | in memory | C1 | 8 | | C1 | 8 |
| C2 | 2 | | C2 | 2 | | C2 | 2 |
| C3 | 9 | | C3 | 9 | | C3 | 9 |
| C4 | 5 | | C4 | 5 | | C4 | 5 |

Original configuration     *Local* replacement     *Global* replacement

---

# Allocation Policies

- Page fault frequency (PFF):
  Usually, more page frames → fewer page faults

**PFF: page faults/sec** (y-axis)

**# page frames assigned** (x-axis)

PFF is unacceptable high
→ process needs more memory

PFF might be too low
→ process may have too
  much memory!!!??????

1   **d**   n    1   **d**   n

Solution ??:
Reduce number of processes competing for memory
- reassign a page frame
- swap one or more to disk, divide up pages they held
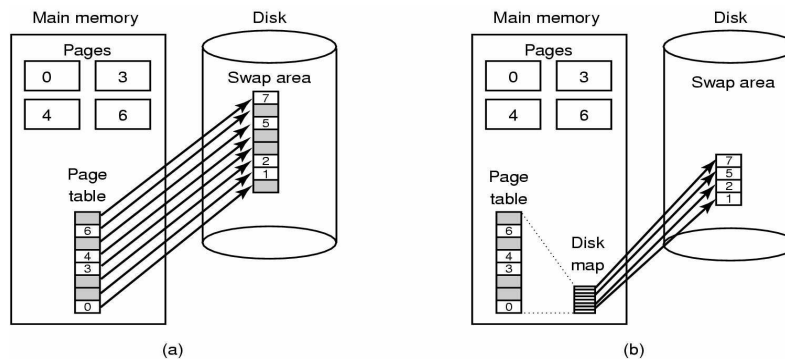- reconsider degree of multiprogramming

# Page Size

- Determining the optimum page size requires balancing several ***competing*** factors:
    - Data segment size ≠ n x page size → internal fragmentation (small size)
    - Keep in memory only data that is (currently) used (small size)

    - Disk operations (large size)
    - Page table size: access/load time and space requirements (large size)
    - Page replacement algorithm: operations per page (large size)

- Usual page sizes is 4 KB – 8 KB, but up to 64 KB is suggested for systems supporting "new" applications  managing high data rate data streams like video and audio

# Locking & Sharing

- Locking pages in memory:
    - I/O and context switches
    - Much used pages
    - …

- Shared pages
  users running the same program at the same time, e.g., editor or compiler
    - Problem 1: not all pages are shareable
    - Problem 2: process swapping or termination
    - …

# Backing Store

- **Backing store (disk management):**
  Where on disk shall the pages be put when paged out?

  a) Special, in-advance allocated swap area
     (problem: growing processes)

  b) Allocate disk space when needed
     (problem: must hold all disk addresses in memory,
     time to allocate a new disk block)
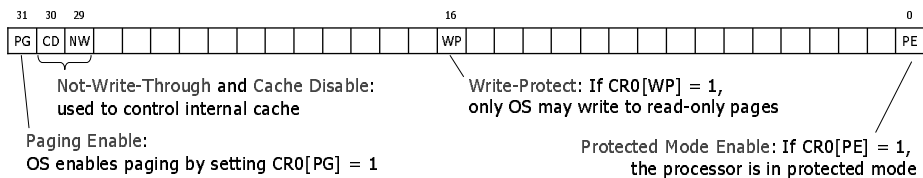


---

# Paging Daemons

- **Paging daemons:**
  Background process which sleeps most of the time, but is for example
  awakened periodically or when the CPU is idle

  - Taking care that enough free page frames are available by writing
    back modified pages before they are reused

  - Prepaging

# Paging on Pentium

- In protected mode, the currently executing process have a 4 GB address space ($2^{32}$) – viewed as 1 M 4 KB pages

    - The 4 GB address space is divided into 1 K page groups (1 level – page directory)

    - Each page group has 1 K 4 KB pages (2 level – page table)

- Mass storage space is also divided into 4 KB blocks of information

- Uses control registers for paging information

---

# Control Registers used for Paging on Pentium

- Control register 0 (CR0):

| 31 | 30 | 29 | | | | | | | | | | | | | 16 | | | | | | | | | | | | | | | 0 |
|----|----|----|---|---|---|---|---|---|---|---|---|---|---|---|----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PG | CD | NW | | | | | | | | | | | | | WP | | | | | | | | | | | | | | | PE |

Not-Write-Through and Cache Disable: used to control internal cache

Paging Enable: OS enables paging by setting CR0[PG] = 1

Write-Protect: If CR0[WP] = 1, only OS may write to read-only pages

Protected Mode Enable: If CR0[PE] = 1, the processor is in protected mode

- Control register 1 (CR1) – does not exist, returns only zero

- Control register 2 (CR2)
    - only used if CR0[PG]=1 & CR0[PE]=1

| 31 | | 0 |
|----|--------------------------|---|
| | Page Fault Linear Address | |

# Control Registers used for Paging on Pentium

- Control register 3 (CR3) – page directory base address:
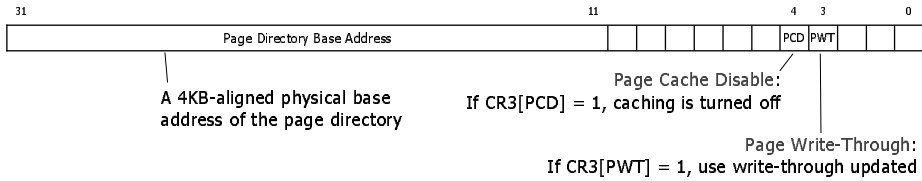    - only used if CR0[PG]=1 & CR0[PE]=1

| 31 | | 11 | | 4 | 3 | | 0 |
|---|---|---|---|---|---|---|---|
| Page Directory Base Address | | | | PCD | PWT | | |

A 4KB-aligned physical base
address of the page directory

Page Cache Disable:
If CR3[PCD] = 1, caching is turned off

Page Write-Through:
If CR3[PWT] = 1, use write-through updated

- Control register 4 (CR4):

| 31 | | 4 | | 0 |
|---|---|---|---|---|
| | | PSE | | |

Page Size Extension: If CR4[PSE] = 1,
the OS designer may designate some pages as 4 **MB**

---

# Pentium Memory Lookup

Incoming virtual address
(0x1402038, 20979768)

| 31 | | | | | | | 22 | 21 | | | | | 12 | 11 | | | | | | | | | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 |

Page directory:

| 31 | 12 | | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| PT base address | | ... | PS | | A | | | U | W | P |

physical base
address of the
page table

page
size

accessed

present

allowed to write

user access allowed

# Pentium Memory Lookup

Incoming virtual address
(0x1402038, 20979768)

| 31 | | | | | | | | | | 22 | 21 | | | | | | | | | | 12 | 11 | | | | | | | | | | | | | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | | |

Index to page directory
(0x6, 6)

| 31 | 12 | | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0...01010101111 | ... | | | | | | | | | 1 |
| 0...01111111000 | ... | | | | | | | | | 0 |
| 0...01110000111 | ... | | | | | | | | | 0 |
| 0...00001010101 | ... | | | | | | | | | 1 |
| 0...01111000101 | ... | | | | | | | | | 0 |
| 0...00000000100 | ... | | | | | | | | | 0 |
| ...... | | | | | | | | | | |

CR3:

| Page Directory Base Address |
|---|

**Page table PF:**
1. Save pointer to instruction
2. Move linear address to CR2
3. Generate a PF exception − jump to handler
4. Programmer reads CR2 address
5. Upper 10 CR2 bits identify needed PT
6. Page directory entry is really a mass storage address
7. Allocate a new page − write back if dirty
8. Read page from storage device
9. Insert new PT base address into page directory entry
10. Return and restore faulting instruction
11. Resume operation reading the same page directory entry again − now P = 1

---

# Pentium Memory Lookup

Incoming virtual address
(0x1402038, 20979768)

| 31 | | | | | | | | | | 22 | 21 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | | | |

Index to page directory
(0x6, 6)

| 31 | 12 | | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0...01010101111 | ... | | | | | | | | | 1 |
| 0...01111111000 | ... | | | | | | | | | 0 |
| 0...01110000111 | ... | | | | | | | | | 0 |
| 0...00001010101 | ... | | | | | | | | | 1 |
| 0...01111000101 | ... | | | | | | | | | 0 |
| 0...00000000100 | ... | | | | | | | | | 1 |
| ...... | | | | | | | | | | |

CR3:

| Page Directory Base Address |
|---|

**Page frame PF:**
1. Save pointer to instruction
2. Move linear address to CR2
3. Generate a PF exception − jump to handler
4. Programmer reads CR2 address
5. **Upper 10 CR2 bits identify needed PT**
6. **Use middle 10 CR2 bit to determine entry in PT − holds a mass storage address**
7. Allocate a new page − write back if dirty
8. Read page from storage device
9. Insert new page frame base address into page table entry
10. Return and restore faulting instruction
11. Resume operation reading the same page directory entry and page table entry again − both now P = 1

Page table:

| 31 | 12 | | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0...01010101111 | ... | | | | | | | | | 1 |
| 0...01010100000 | ... | | | | | | | | | 0 |
| 0...01100110011 | ... | | | | | | | | | 1 |
| 0...00010000100 | ... | | | | | | | | | 1 |
| ...... | | | | | | | | | | |

# Pentium Memory Lookup

Incoming virtual address
(0x1402038, 20979768)

| 31 | | | | | | | | 22 | 21 | | | | | | | | | 12 | 11 | | | | | | | | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 |

Index to page directory
(0x6, 6)

Index to page table
(0x2, 2)

Page:

Page offsett
(0x38, 56)

requested data

| 31 | 12 | | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0...01010101111 | ... | | | | | | | | | 1 |
| 0...01111111000 | ... | | | | | | | | | 0 |
| 0...01110000111 | ... | | | | | | | | | 0 |
| 0...00001010101 | ... | | | | | | | | | 1 |
| 0...01111000101 | ... | | | | | | | | | 0 |
| 0...00000000100 | ... | | | | | | | | | **1** |
| ...... | | | | | | | | | | |

CR3:

Page Directory Base Address

| 31 | 12 | | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0...01010101111 | ... | | | | | | | | | 1 |
| 0...01010100000 | | | | | | | | | | **1** |
| 0...01100110011 | | | | | | | | | | 1 |
| 0...00010000100 | | | | | | | | | | 1 |
| ...... | | | | | | | | | | |

---

# Page Fault Causes

- Page directory entry's P-bit = 0:
  page group's directory (page table) not in memory

- Page table entry's P-bit = 0:
  requested page not in memory

- Attempt to write to a read-only page

- Insufficient page-level privilege to access page table or frame

- One of the reserved bits are set in the page directory or table entry