

Page replacement algorithms

OS 2003

32

When a page fault occurs

- OS has to choose a page to evict from memory
- If the page has been modified, the OS has to schedule a disk write of the page
- The page just read overwrites a page in memory (e.g. 4Kbytes)
- Clearly, it's better not to pick a page at random
- Same problem applies to memory caches

OS 2003

33

Optimal page replacement

- At the moment of page fault:
 - Label each page in memory is labeled with the number of instructions that will be executed before that page is first referenced
 - Replace the page with the highest number: i.e. postpone as much as possible the next page fault
- Nice, optimal, but unrealizable
 - The OS can't look into the future to know how long it'll take to reference every page again

OS 2003

34

“Not recently used” algorithm

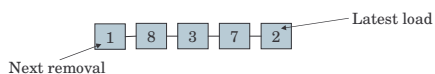
- Use **R**eferenced and **M**odified bits
- **R&M** are in hardware, potentially changed at each reference to memory
 - **R&M** are zero when process is started
- On clock interrupt the **R** bit is cleared
- On page fault, to decide which page to evict:
 - Classify:
 - Class 0 – R=0,M=0
 - Class 1 – R=0,M=1
 - Class 2 – R=1,M=0
 - Class 3 – R=1,M=1
 - Replace a page at random from the lowest class

OS 2003

35

FIFO replacement

- FIFO, first in first out for pages
- Clearly not particularly optimal
- It might end up removing a page that is still referenced since it only looks at the page's age
- Rarely used in pure form...

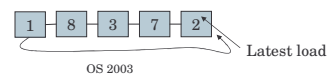


OS 2003

36

“Second chance” algorithm

- Like FIFO but...
- Before throwing out a page checks the **R** bit:
 - If 0 remove it
 - If 1 clear it and move the page to the end of the list (as it were just been loaded)
 - If all pages have R=1, eventually the algorithm degenerates to FIFO (why?)

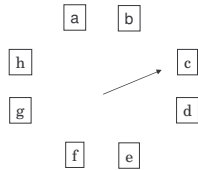


OS 2003

37

Clock page algorithm

- Like “second chance” but...
- ...implemented differently:
 - Check starting from the latest visited page
 - More efficient



OS 2003

38

Least recently used (LRU)

- Pages recently used tend to be used again soon (on average)
- Idea! Get a counter, maybe a 64bit counter
- Store the value of the counter in each entry of the page table
- When is time to remove a page, find the lowest counter value (this is the LRU page)

OS 2003

39

NFU algorithm

- Since LRU is expensive
- NFU: “Not Frequently Used” algorithm
- At each clock interrupt add the R bit to a counter for each page: i.e. count how often a page is referenced
- Remove page with lowest counter value
- Unfortunately, this version tends not to forget anything

OS 2003

40

Aging (NFU + forgetting)

- Take NFU but...
- At each clock interrupt:
 - Right shift the counters (divide by 2)
 - Add the R bit to the left (MSB)
- As for NFU remove pages with lowest counter
- Note: this is different from LRU since the time granularity is a clock tick and not every memory reference!

OS 2003

41

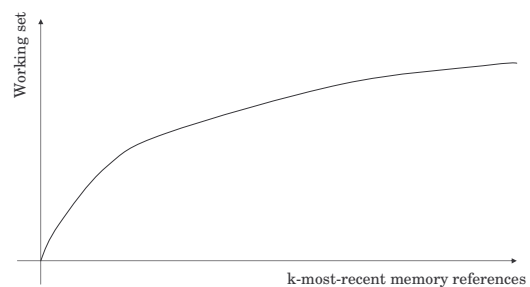
Process' behavior

- Locality of reference: most of the time the last k references are within a finite set of pages $<$ a large address space
- The set of pages a process is currently using is called the *working set* of the process
- Knowing the working set of processes we can do very sophisticated things (e.g. pre-paging)

OS 2003

42

Working set



OS 2003

43

WS based algorithm

- Store time information in the table entries
- At clock interrupt handle R bits as usual (clear them)
- At page fault, scan entries:
 - If R=1 just store current time in the entry
 - If R=0 compute “current-last time page was referenced” and if $> \text{threshold}$ the page can be removed since it's no longer in the working set (not used for *threshold* time)
- **Note:** we're using time rather than actual memory references

OS 2003

44

WSClock algorithm

- Use the circular structure (as seen earlier)
- R=1, page in the WS – don't remove it
- R=0, M=0 no problem (as before)
- M=1, schedule disk write appropriately to procrastinate as long as possible a process switch
 - No write is schedulable (R=1 always), just choose a clean page

OS 2003

45

Summary

Algorithm	Comment
Optimal	Not implementable, useful for benchmarking
NRU (Not recently used)	Very crude
FIFO	Might throw out important pages
Second chance	Big improvement over FIFO
Clock	Realistic (better implementation)
LRU (Least Recently Used)	Excellent but difficult to implement
NFU	Crude approx to LRU
Aging	Efficient in approximating LRU
Working set	Expensive to implement
WSClock	Good and efficient

OS 2003

46

Design issues

OS 2003

47

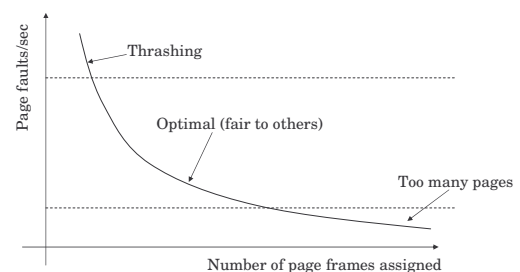
Design issues

- Local vs. global allocation policy
 - When a page fault occurs, whose page should the OS evict?
- Which process should get more or less pages?
 - Monitor the number of page faults for every process (PFF – page fault frequency)
 - For many page replacement algorithms, the more pages the less page faults

OS 2003

48

Page fault behavior



OS 2003

49

Load control

- If the WS of all processes > memory, there's *thrashing*
- E.g. the PFF says a process requires more memory but none require less
- Solution: swapping – swap a process out of memory and re-assign its pages to others

OS 2003

50

Page size

- Page size p , n pages of memory
- Average process size s , in pages s/p
- Each entry in the page table requires e bytes
- On average $p/2$ is lost (fragmentation)
- Internal fragmentation: how much memory is not used within pages
- Wasted memory: $p/2 + se/p$
- Minimizing it yields the optimal page size (under simplifying assumptions)

OS 2003

51

Two memories

- Separate data and program address spaces
- Two independent spaces, two paging systems
- The linker must know about the two address spaces

OS 2003

52

Other issues

- Shared pages, handle shared pages (e.g. program code)
 - Sharing data (e.g. shared memory)
- Cleaning policy
 - Paging algorithms work better if there are a lot of free pages available
 - Pages need to be swapped out to disk
 - Paging daemon (write pages to disk during spare time and evict pages if there are too few)

OS 2003

53

Page fault handling

1. Page fault, the HW traps to the kernel
 1. Perhaps registers are saved (e.g. stack)
2. Save general purpose microprocessor information (registers, PC, PSW, etc.)
3. The OS looks for which page caused the fault (sometimes this information is already somewhere within the MMU)
4. The system checks whether the process has access to the page (otherwise a protection fault is generated, and the process killed)
5. The OS looks for a free page frame, if none is found then the replacement algorithm is run
6. If the selected page is dirty ($M=1$) a disk write is scheduled (suspending the calling process)
7. When the page frame is clean, the OS schedules another transfer to read in the required page from disk
8. When the load is completed, the page table is updated consequently
9. The faulting instruction is backed up, the situation before the fault is restored, the process resumes execution

OS 2003

54

Segmentation

OS 2003

55

Why?

- Many separate address spaces (segments) (e.g. data, stack, code, and many others if needed)
- Each segment is separate (e.g. addresses from 0 to some MAX)
- Segments might have different lengths
- Segment number + address within segment
- Linking is simplified (libraries within different segments can assume addresses starting from 0) – e.g. if a part of the libraries is recompiled the remainder of the code is unaffected
- Shared library (DLL's) implementation is simpler (the sharing is simpler)

OS 2003

56

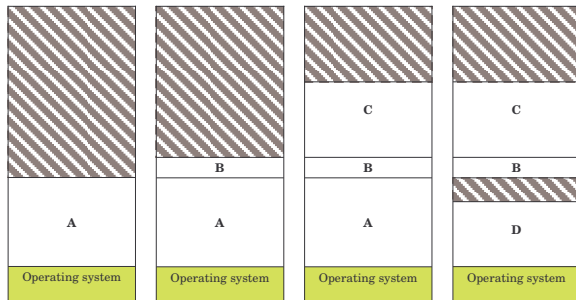
Comparing paging and segmentation

Consideration	Paging	Segmentation
Need the programmer be aware that this technique is being used?	No	Yes
How many linear address spaces are there?	1	Many
Can the total address space exceed the size of physical memory?	Yes	Yes
Can procedures and data be distinguished and separately protected?	No	Yes
Can tables whose size fluctuate be accommodated easily?	No	Yes
Is sharing of procedures between users facilitated?	No	Yes
Why was this technique invented?	To get a large linear address space without having to buy more physical memory	To allow programs and data to be broken up into logically independent address spaces and to aid sharing and protection

OS 2003

57

Pure segmentations



OS 2003

58

Fragmentation

- External fragmentation:
 - Memory fragments not used (we've already seen this)
 - Memory wasted in unused holes

OS 2003

59

Segmentation + paging (Pentium)

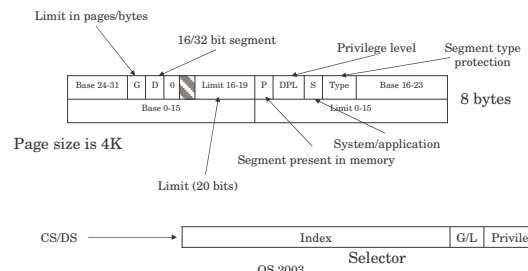
- 16K segments
- 1G 32bit words (DoubleWords)
- Two tables: LDT, GDT – Local (to the process) and global (to the processor) descriptor table
- To work with a segment the machine loads the segment number into a special register (CS, DS, etc.) – CS, DS are 16 bit registers
- The descriptor of the segment (see next slide)

OS 2003

60

The segment descriptor

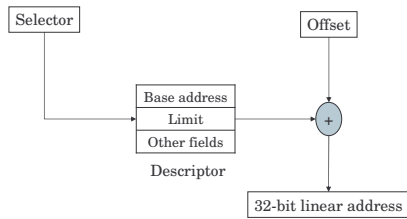
- This is used by the microcode within the Pentium to work with segments



OS 2003

61

Getting the address

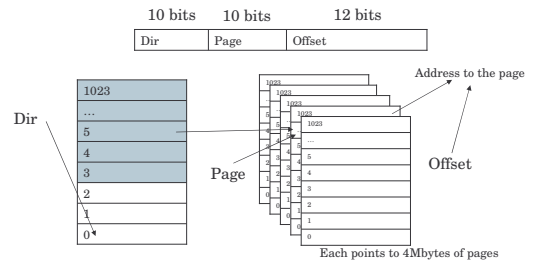


OS 2003

62

Paging on the Pentium

- 2-level page table in memory



OS 2003

63

More on the Pentiums

- TLB, to avoid repeated accesses to memory
- The whole thing can be used with just a single segment to obtain a linear 32bit address space
- Set base and limit appropriately
- Protection (a few bits)

OS 2003

64