

WEEK 5

## Virtual memory algorithms and examples

Stallings, Chapter 8

1

## Design decisions

- Virtual memory or not?
- Paging or segmentation?
- What algorithms?

CONSIDERING ONLY PAGING:

- What page size?

2

## Page size

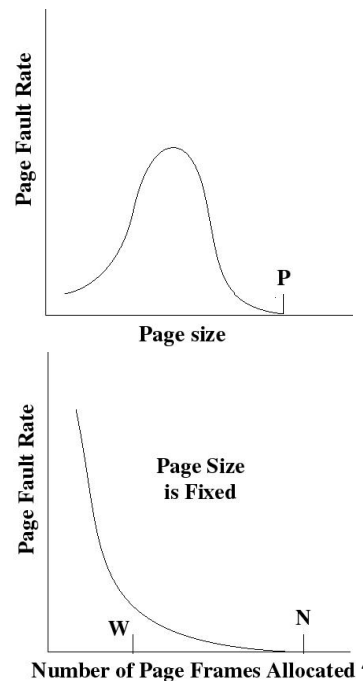
- Reduce internal fragmentation → small
- Reduce page table size → large
- Proportional (1x, 2x, ...) to disk block size
- Optimal value is different for different programs
- Increase TLB hit ratio → large
- Different page size for different applications?
- How does MMU know the page size?

Tbl 8.2 [Stal05]

3

## Page fault rate

- More small pages to the same memory space
- References from large pages more probable to go to a page not yet in memory
- References from small pages often to other pages -> often used pages selected to the memory
- Page size closer to process size
- Too small working set size -> larger page fault rate
- Large working set size -> nearly all pages in the memory



## Algorithms for memory management

- Fetch policy
  - Demand paging vs prepaging
- Placement policy
- Replacement policy
  - Resident set management, selecting the page to be replaced
  - Methods: OPT, LRU, FIFO, random ...
- Cleaning policy
  - Demand cleaning vs precleaning, buffering
- Load control
- Working set size

5

## Resident set mgt

- Number of page frames allocated for a process (resident set size)
  - Fixed: how large, identical for all?
  - Variable: based on what? Working set?
- Replacement policy:
  - Global: over all pages (all processes)
  - Local: just the pages of this process
- Size of the resident set:
  - Too small: higher multiprocessing level, more page faults, trashing
  - Too large: smaller CPU utilisation, wasted memory

6

## Resident set management

	Local Replacement	Global Replacement
<b>Fixed Allocation</b>	<ul style="list-style-type: none"> <li>• Number of frames allocated to process is fixed.</li> <li>• Page to be replaced is chosen from among the frames allocated to that process.</li> </ul>	<ul style="list-style-type: none"> <li>• Not possible.</li> </ul>
<b>Variable Allocation</b>	<ul style="list-style-type: none"> <li>• The number of frames allocated to a process may be changed from time to time, to maintain the working set of the process. (esim. PFF)</li> <li>• Page to be replaced is chosen from among the frames allocated to that process.</li> </ul>	<ul style="list-style-type: none"> <li>• Page to be replaced is chosen from all available frames in main memory; this causes the size of the resident set of processes to vary.</li> </ul>

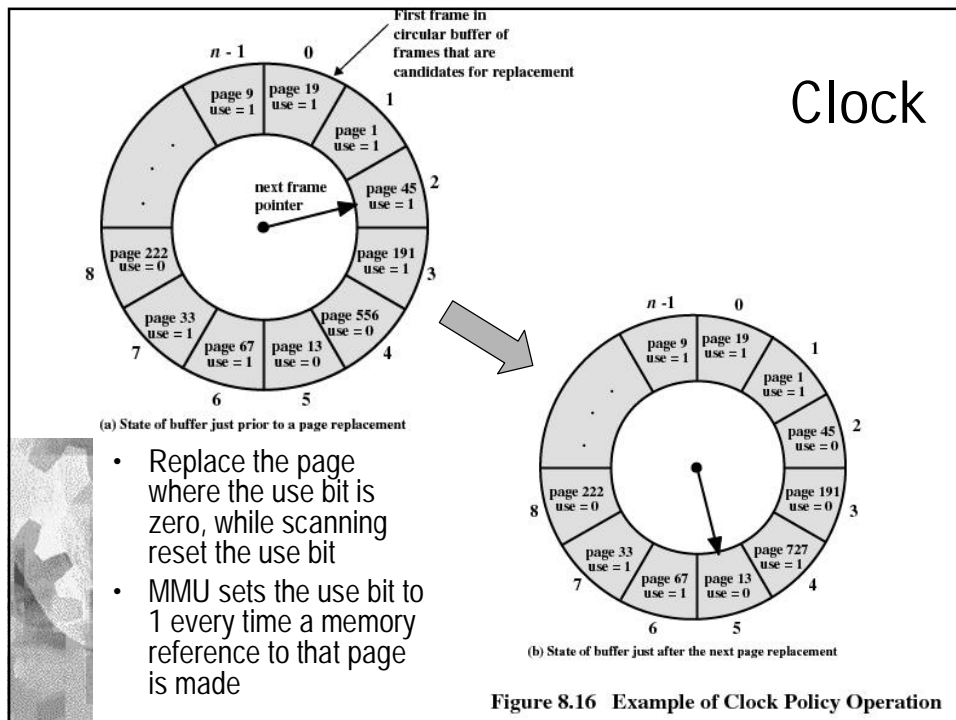
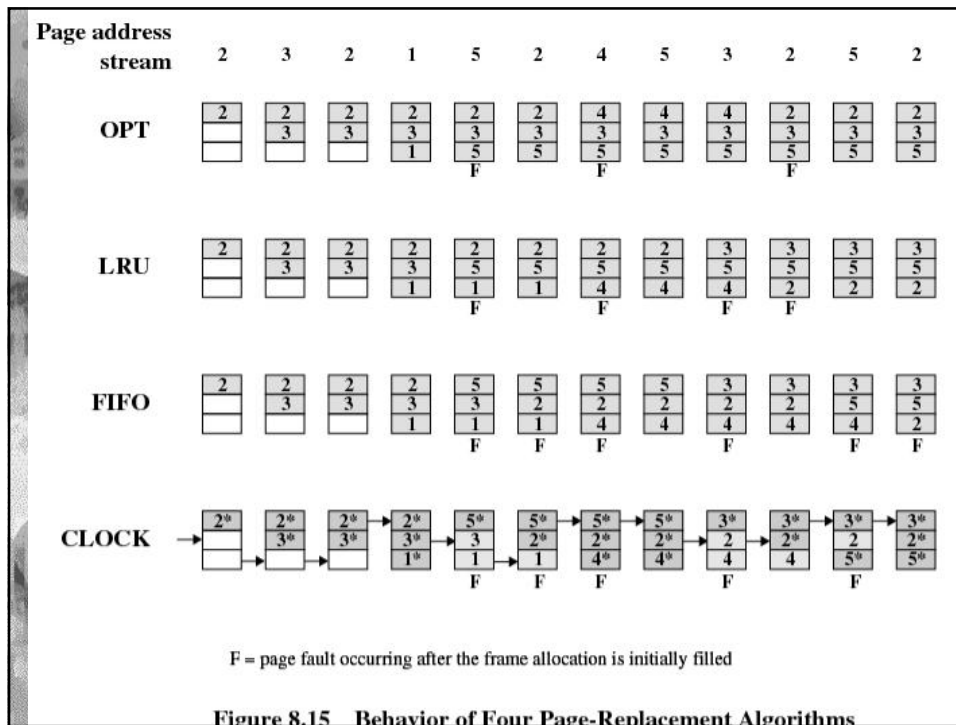
Tbl 8.4 [Stal05]

7

## Replacement policy

- OPT – optimal algorithm
  - Nothing is better than this, unrealistic to be implemented,
  - Clairvoyant, knows the future references and chooses based on them
- LRU, Least recently used
  - High amount of bookkeeping
  - Modification: NRU – Not recently used, easier to implement
    - Only maintain information about recent usage (for example on used bit)
- FIFO, First-in-first-out
  - Easy to implement, but not very useful
  - Modification: second chance – do not replace page that was just recently used
- Clock
- Metrics? Number of page faults or page fault rate

8



## Working set

Sequence of Page References

Window Size,  $\Delta$

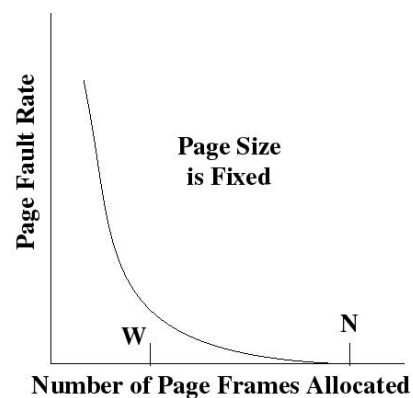
Fig 8.19 [Stal05]

	2	3	4	5
24	24	24	24	24
15	24 15	24 15	24 15	24 15
18	15 18	24 15 18	24 15 18	24 15 18
23	18 23	15 18 23	24 15 18 23	24 15 18 23
24	23 24	18 23 24	•	•
17	24 17	23 24 17	18 23 24 17	15 18 23 24 17
18	17 18	24 17 18	•	18 23 24 17
24	18 24	•	24 17 18	•
18	•	18 24	•	24 17 18
17	18 17	24 18 17	•	•
17	17	18 17	•	•
15	17 15	17 15	18 17 15	24 18 17 15
24	15 24	17 15 24	17 15 24	•
17	24 17	•	•	17 15 24
24	•	24 17	•	•
18	24 18	17 24 18	17 24 18	15 17 24 18

11

## Working set

- Set of pages used during last k references (window size)
- Pages not in the working set are candidates to be released or replaced
- Suitable size for the working set? Estimate based on page fault rate



$$1 \leq |W(t, \Delta)| \leq \min(\Delta, n)$$

12

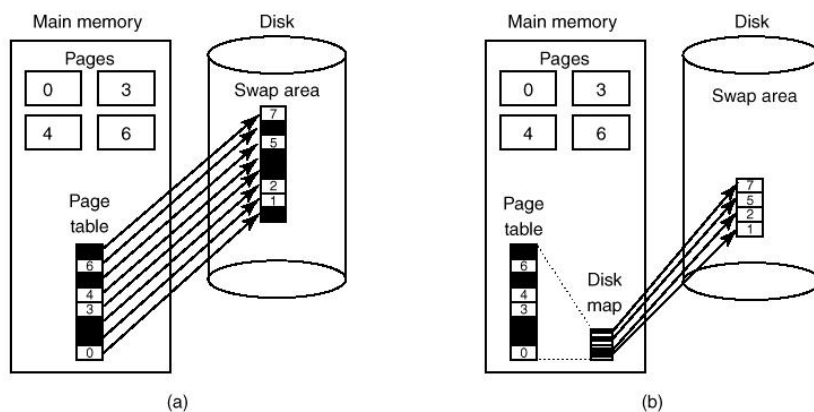
## Dynamically adjusting resident set

- PFF – Page Fault Frequency
- VSWS – Variable-interval Sampled Working Set

13

## Swap area, VM backup store

(Fig 4-33 [Tane01])



Fixed allocation –  
For the whole process  
in advance

Dynamic allocation –  
only for the pages that  
are swapped out

14

# UNIX / Solaris (+4BSD) MUISTINHALLINTA

15

## UNIX/Solaris: data structures

Page table - one for each process

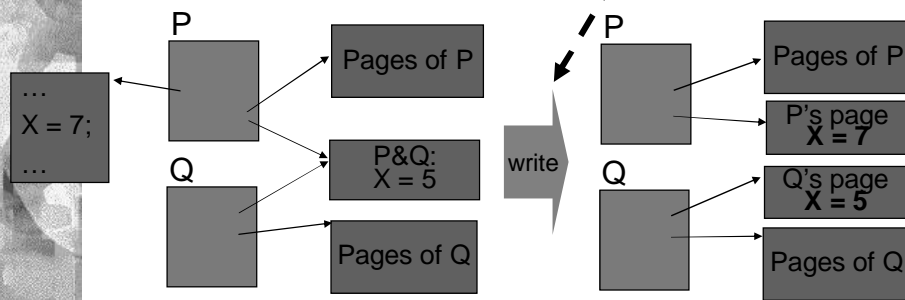
Tbl 8.5 [Stal05]

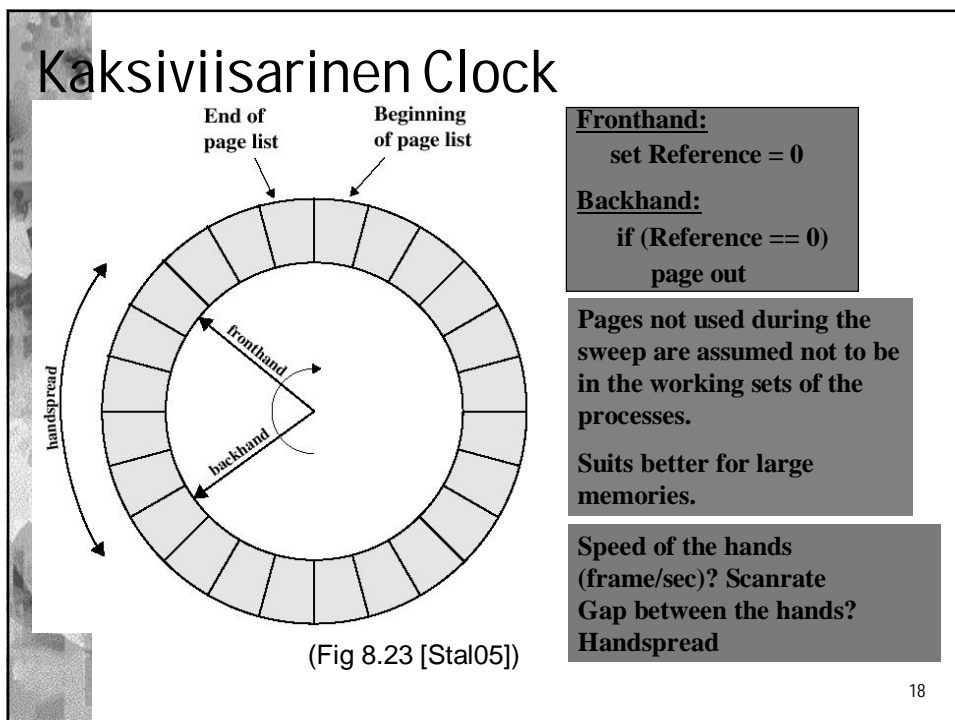
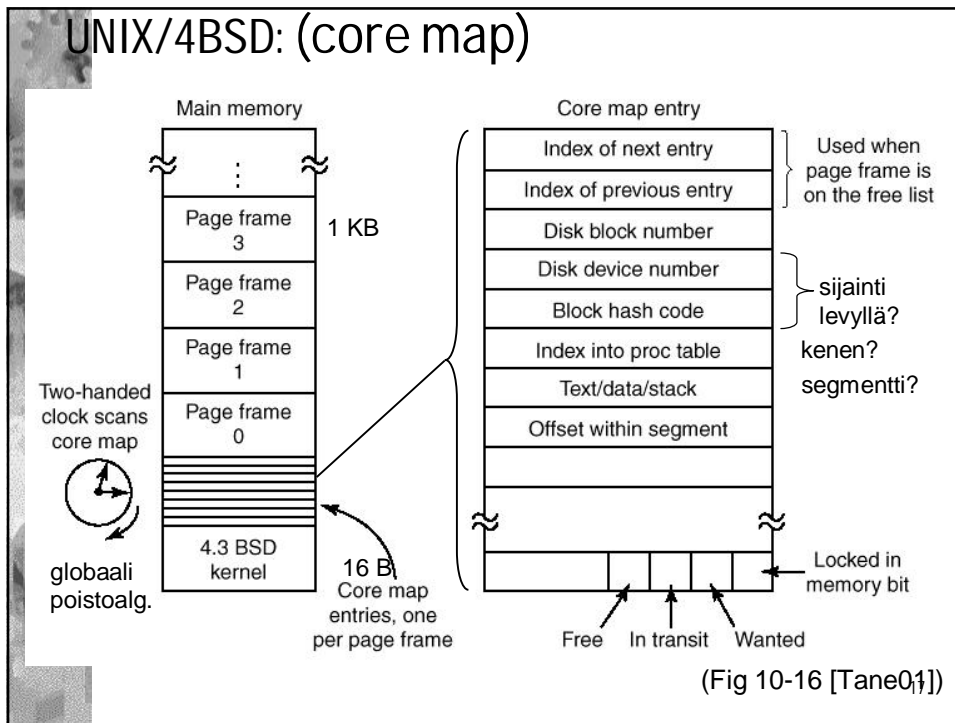
- one entry for each logical page

Location in the main memory?

Page frame number	Age	Copy on write	Modify	Reference	Valid	Protect
-------------------	-----	---------------	--------	-----------	-------	---------

(a) Page table entry





## UNIX/Solaris: Kernel memory allocation

- Based on paging
- Locking pages to memory
- Kernel allocates and deallocates all the time (large number of) small areas and buffers
  - For example: proc, vnode, file descriptor
- Kernel need to allocate and deallocate fragments of pages (not only full pages)
  - Dynamic partitioning for allocations within one page
  - => Lazy Buddy System
  - Delay combining until there are 'too many' fragments of certain size

19

## Linux MUISTINHALLINTA



Ch 8.4 [Stal 05]

20

## Linux: memory management



- 32-bittisissä architectures
  - 3 GB reserved for user processes
  - 1 GB for kernel (including the page tables)
    - at the beginning of the physical memory
  - Kernel mode allowed to make references to the full 4 GB
- Kernel's memory allocation based on slabs
  - Only for the kernel, efficient allocation of small memory areas, avoids allocating full pages

21

## Linux: paging

- Architecture independent

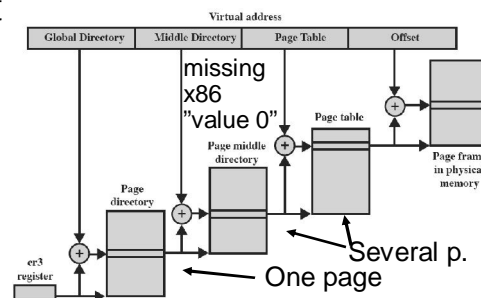
- Alpha: 64b addresses, full support for 3 levels,
  - Page size 8KB, offset 13 bits
- x86: 32b address, only 2-levels,
  - Page size 4KB, offset 12 bits

- 3-level page table

- page directory (PD), 1 page
- page middle directory (PMD), multiple pages
- page table (PT), multiple pages

- Page directory always in memory

- This the page table address given to MMU at process switch
- Other pages can be on disk



(Fig 8.25 [Stal05])



22

## Linux: address translation



- Split address to 4 elements
  - index to page directory (PD-offset)
  - index to page middle directory (PMD-offset)
  - index to page table (PT-offset)
  - offset

```
PMD_base = PD_base + PD[PD_offset]
PT_base  = PMD_base + PMD[PMD_offset]
Page_base = PT_base  + PT[PT_offset]
fyys.os  = Page_base + offset
```

- X86-trick
  - Page middle directory (PMD) reduced to one item, which compiler optimizes away.
  - MMU assumes that the page directory refers directly to the page table (PT)
  - Suitable to all architectures that support only two-level paging

23

## Linux: placement (allocation)

- Reserves a consecutive area (region) for consecutive pages
  - More efficient fetching and storing with disks
  - Optimize the disk utilisation, with the cost of memory allocation
- Buddy System: 1, 2, 4, 8, 16 or 32 page frames
  - Allocates pages in large groups, increases internal fragmentation
  - Implementation: table with list of different sizes unallocated page regions



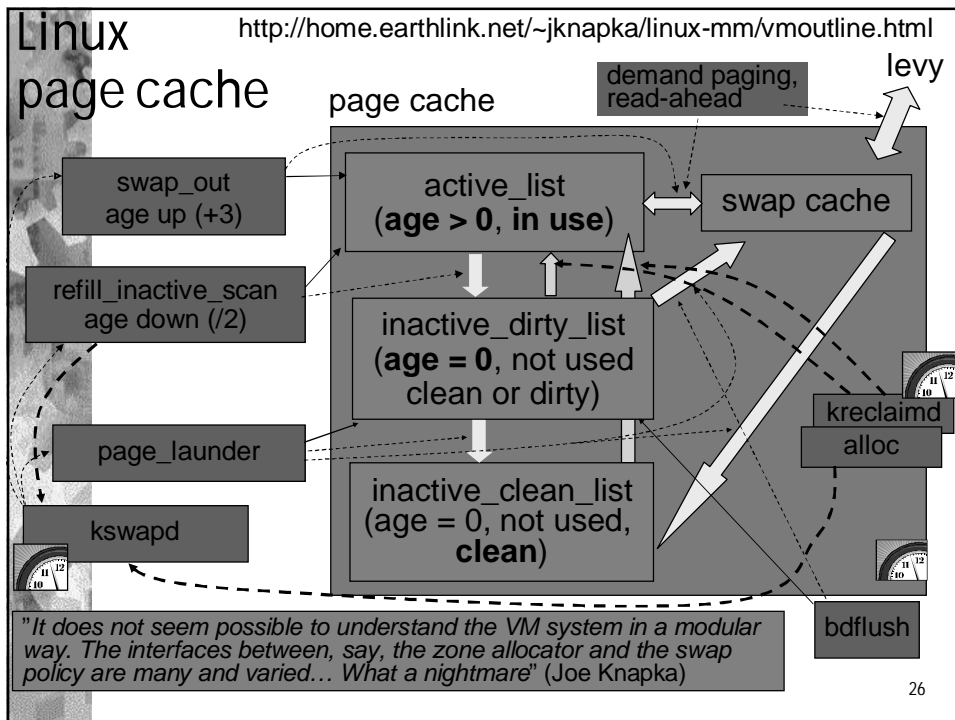
24

## Linux: replacement

- Globaali Clock using M-bit + sort of LRU
  - Counter of dynamically allocated page frames
- 8 bit counter of age vrt. U-bitti
  - MMU increases the counter at each page reference
  - Background process (kswapd) reduces the counter each second
    - if too few page frames free, deallocated unused frames
- Uses pagebuffer
  - To store deallocated pages for a short duration



25



## Linux: kernel memory allocation

- Kernel memory locked to the memory
- Dynamic allocation based on pages
- Dynamically loadable modules fully to the kernel area
  - buddy system on the pages `kmalloc()`
- Additional need for small, short term memory allocations -> use slab
  - One page 'cut' to smaller areas (slabs)
  - buddy system within each page
    - x86: 32, 64, 128, 252, 508, 2040, 4080 tavua (page size 4KB)



27

## Slab allocation

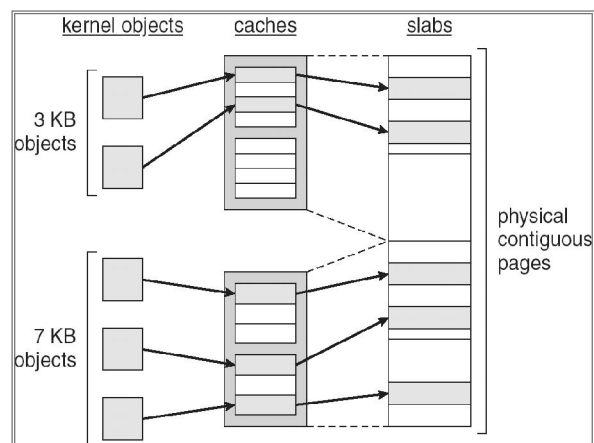



Fig 9.27 – Silberschatz, Galvin & Gagne:  
Operating system concepts with Java

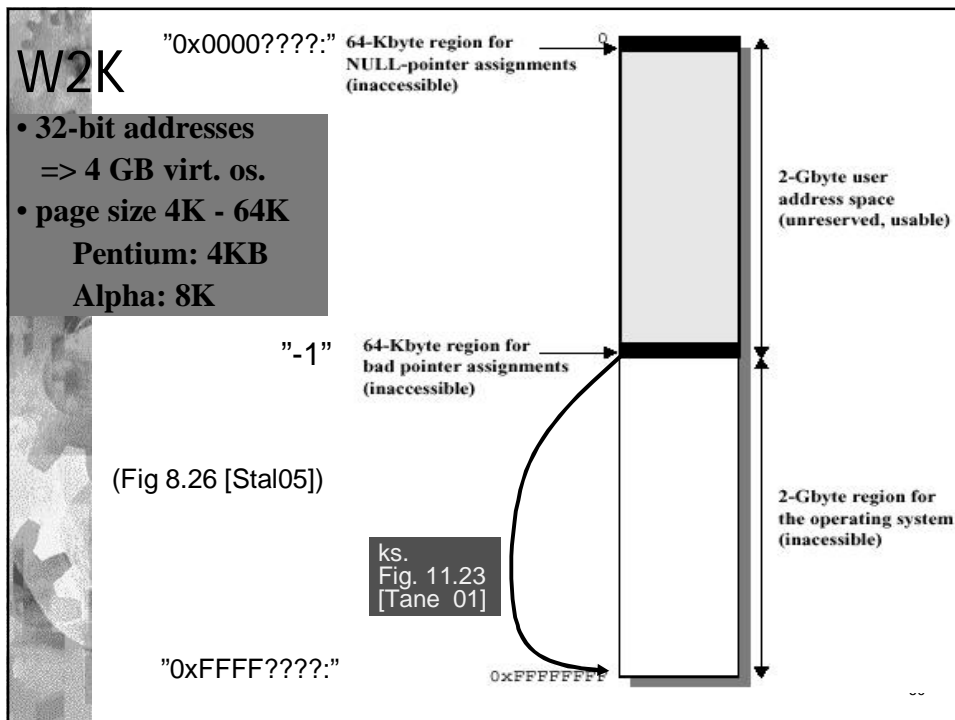
28



# Windows 2000 MUISTINHALLINTA

Ch 8.5 [Stal 05]  
Ch 11.5 [Tane 01]

29



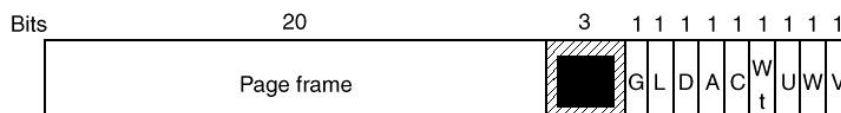
## W2K: Paging

- No segmentation
- Variable resident set size per process
  - Start with a fixed set of frames
  - If large number of free frames, process can get more frames
    - A greedy one is not allowed to allocate last 512 frames
  - If the free memory reduces, the resident set of processes can be decreased
- Demand fetch (no prefetching)
- Disk I/O using larger blocks
  - 1-8 pages at each time

31

## W2K: paging

(Fig 11.26 [Tane01])



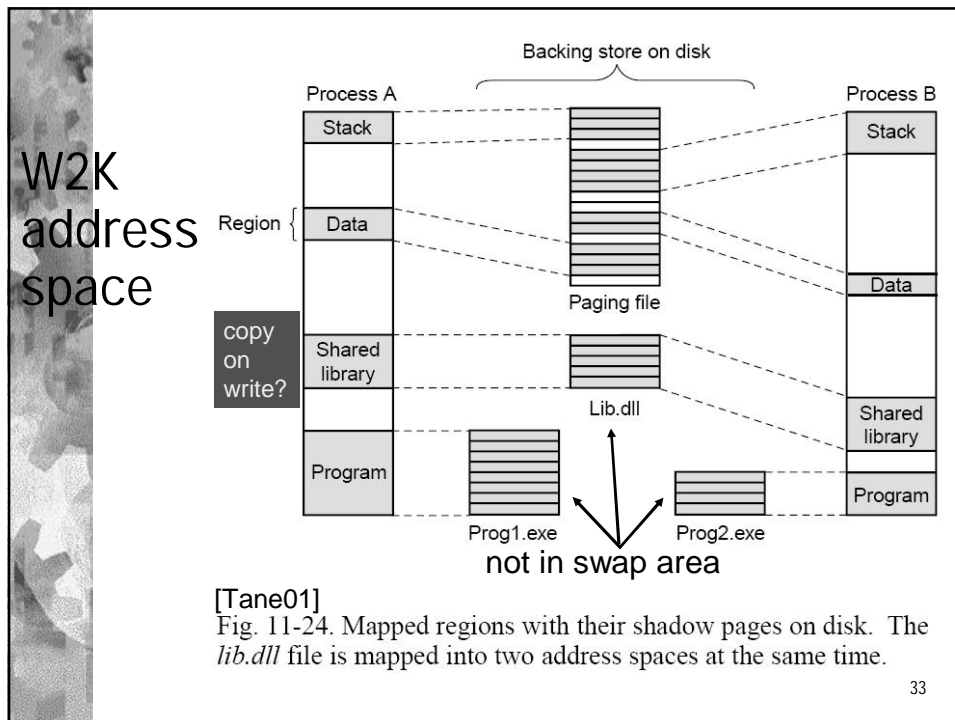
G: Page is global to all processes  
 L: Large (4-MB) page  
 D: Page is dirty  
 A: Page has been accessed  
 C: Caching enabled/disabled

Wt: Write through (no caching)  
 U: Page is accessible in user mode  
 W: Writing to the page permitted  
 V: Valid page table entry

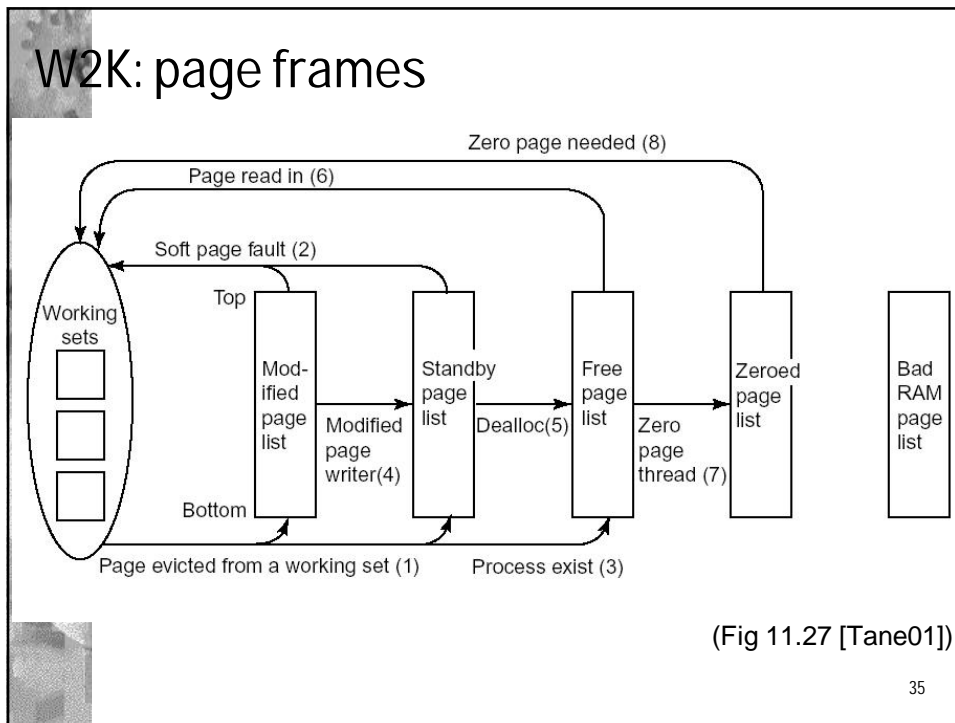
**Pentium**

Fig 11.24 [Tane01]  
 (seur. kalvo)

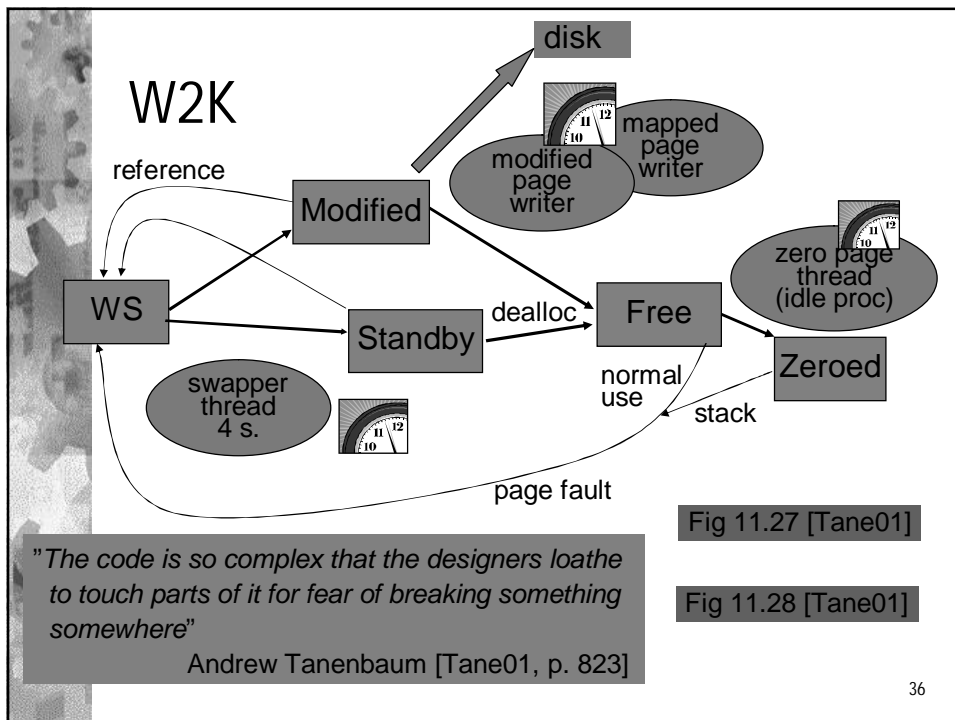
32



- ## W2K: replacement
- Balance set manager -daemon
    - Each second 1 sec check if enough free frames
  - Working set manager
    - Running when needed
    - Dynamic working set size
    - Released page frames of processes
    - Starts with large idle processes
  - Swapper-daemon
    - every 4 sec: search for processes, whose threads (all of them) have been passive for long duration
  - Parts of the OS and buffer pool locked to memory
  - Caching of released pages
- 34



35



36