
Operativsystem

Minneshantering II
(kap 4 i boken)

Sidutbytesalgoritmer

- Då ett sidfel (page fault) uppstår, måste operativsystemet byta ut innehållet i en sidram mot det innehåll som den som förorsakade sidfelet behöver
- Då innehållet i en sidram byts ut
 - om förändrat (skrivet i), måste förändringarna sparas någonstans (= på hårddiskiva)
 - om ej förändrat, kan man helt enkelt "radera" det gamla innehållet (eftersom innehållet redan finns sparat någon annanstans)

Sidutbytesalgoritmer II

- Det nya innehållet kan tas från flere ställen
 - om programkod: direkt från filen som innehåller det körbara programmet
 - om data: från swap-området på hårddisken
 - om ej tidigare använt, en ny "tom" sida
- Problemställning !
 - vilken sidram skall tas i användning då ett sidfel uppstår ???
 - Val av sidutbytesalgoritm

Utbytesalgoritmer generellt

- Problemet med vilken sidram innehållet skall bytas i är generellt
 - Cache-minnen: Vilken cache-rad skall bytas
 - Disk cache: Vilket / vilka block skall bytas?
 - Web server och minnes-cache: Vilken fil skall bytas?

Sidutbytesalgoritmer

- Den optimala
- Icke nyligen använd sida (not recently used, NRU)
- Först in, först ut (First In, First out, FIFO)
- Andra möjligheten (Second Chance)
- "Klock"-algoritmen
- Minst nyligen använda sidan (Least Recently Used, LRU)
- Inte ofta använd (Not Frequently Used)
- Working set

Den optimala algoritmen

- Då innehållet i en sidram skall bytas, numrera sidramarna enligt hur länge det tar innan sidramen nästa gång accesseras
- Byt innehållet i den ram som längst i framtiden kommer att accesseras!!
- Problemet: Denna algoritm är omöjlig att implementera

Not recently used

- Varje gång man läser från en sida, sätts en bit R (read), varje gång man skriver sätts en bit M (modified)
- Med jämna mellanrum nollställer operativsystemet R-biten (t.ex. vid klock-avbrott)
- Vid ett sidfel, kan klassificera sidramarna enligt
 - Klass 0: ej refererad, ej modifierad
 - Klass 1: ej refererad, modifierad
 - Klass 2: refererad, ej modifierad
 - Klass 3: refererad, modifierad
- Byt ut innehållet i en slumpmässigt vald sidram av lägsta möjliga klass

First-In, First-Out

- Algoritm: Sidram väljs i löpande ordning
- Används sällan, ger dåliga resultat, eftersom sidor som används mycket byts ut lika ofta som sådana som använts t.ex. bara en gång

Second chance

- Som FIFO, men använder sig också av R-biten, dvs. sida som refererats nyligen byts ej ut, utan man ser på nästa i listan
 - dvs. sidorna får en ”andra möjlighet”
- Klock-algoritmen är samma som FIFO/Second chance, det är bara fråga om hur algoritmen är implementerad

Least recently used (LRU)

- Introducera en räknare C som inkrementeras för varje instruktion som CPU:n utför
- Då en sida refereras, skriv värdet på räknaren C till ett fält i sidtabellen
- Vid sidfel, välj den sidram som har det lägsta räknarvärdet sparat för att byta ut
- Problem: Hårdvara implementerar sällan denna funktionalitet

Mjukvaru-LRU

- Vid varje klock-avbrott, inkrementera en räknare C
- Vid samma klock-avbrott, kopiera räknarens värde till sidtabellen för sådana sidor som har R-biten satt
- Vid sidfel, byt ut innehållet i den sidram med lägsta motsvarande räknare

Not Frequently Used (NFU)

- Vid klockavbrott, om R-biten satt, inkrementera ett räknare-fält i sidtabellen
- Vid sidfel, ersätt sida med lägsta värdet på räknaren
- Problem: Detta system "glömmer" aldrig, en sida med högt räknarvärde förblir i minne, även om den inte använts på länge

NFU med åldrande

- Innan R-biten läggs till räknaren, högerskiftas räknaren
- Dessutom, R-biten adderas till den mest signifikanta biten i räknaren
- Motsvarar nu ganska långt LRU

Working set algoritmer

- Tidigare presenterades "Demand paging" sidor laddas efter behov
 - Problem: Då en sida laddas, tar det tid p.g.a. I/O -> proc. blockeras. Då sidan laddats, fortsätter processen och behöver med stor sannolikhet en ny sida i närheten av den förra -> ett nytt sidfel -> "trashing" (kontinuerliga sidfel)
 - Borde ha laddat flere sida i samband med samma I/O!
- Working set = den mängd av sidor i minnet som processen för tillfället använder sig av
- Man strävar efter att processens "working set" skall finnas i sidramar då processen skeduleras ("prepaging")

Working set

- Working set kan definieras som $w(k, t)$ vid en tidpunkt t den mängd av sidor som använts vid de senaste k minnesreferenserna
- En sidbytesalgoritm är härmed: Använd en sidram vars sida inte finns i working set
- Det är dock svårt att använda sig av k (de senaste referenserna), istället använder man sig ofta av approximationen de sidor som refererats under den senaste tiden T

Working set algorithm

- Upprätthåll en räknare C för nuvarande tid
- Vid klockavbrott, kopiera värdet på C till ett fält (*current_time*) i sidtabellen
- Vid sidfel, beräkna "åldern" på sidorna ($C - \text{current_time}$), byt ut en sida med ålder större än T

Sidstorlek ?

- Små sidor -> stora sidtabeller (ryms t.ex. ej i TLB)
- Stora sidor -> intern fragmentering (eftersom i medeltal 50 % av den sista sidan är oanvänd)
- Overhead kan beräknas
$$overhead = se/p + p/2$$
där s är medelprocesstorleken, e är radstorleken i sidtabellen, p är sidstorleken

Design av paging systems

- Globala eller lokala sidbytesalgoritmer
 - Globala: Sidramarna gemensamt för alla processer
 - Dynamiskt avgöra hur många sidramar till en process
 - Lokala: Varje process blir tilldelad en fixerad mängd sidramar
 - Minnesbehovet ökar: trashing
 - Minnesbehovet minskar: slöseri med minne

Cleaning policy

- Vi vill ha sidor lediga
- Paging daemon: undersöker periodiskt vilka sidor som kan skrivas till minnet

Committed Memory

- Committed memory = totalt antal sidor som måste ha en representation i antingen RAM eller swap-utrymme på skiva
 - Vad händer om vi skriver till för mycket minne?

NT (4.0) Minneshantering

- i386: 2 GB user space / 2 GB kernel
 - genom boot switch /3GB även 3GB/1GB
- Varje process har en minimum samt maximim ”working set”, samt en nuvarande storlek på working set
- Sidhämtning: ”Clustered demand paging”, hämtar ett kluster av sidor kring den sida som begärs (0-7 sidor)
- Sidbyte: Klock-algoritmen
 - NT 5.0 – om R-bit satt, nollställ räknare, annars inkrementera räknare, byt ut sidram med största värde på räknare
 - Multiprocessor-maskiner: FIFO

Processlista NT

Image Name	User Name	CPU	Mem Usage	Peak Mem Usage	Page Fa...	VM Size	Paged Pool	NP Pool
AcroRd32.exe	jbjorkqv	00	35 856 K	46 204 K	485 006	31 284 K	83 K	8 K
SshClient.exe	jbjorkqv	00	6 508 K	6 516 K	1 691	2 032 K	38 K	4 K
SshClient.exe	jbjorkqv	00	1 720 K	6 340 K	3 392	1 920 K	37 K	3 K
POWERPNT.EXE	jbjorkqv	00	8 404 K	50 544 K	263 459	70 712 K	91 K	15 K
OPSCAN.EXE	jbjorkqv	00	780 K	5 676 K	1 828	3 388 K	45 K	5 K
rundll32.exe	jbjorkqv	00	840 K	3 408 K	1 411	2 216 K	33 K	2 K
QCWLICON.EXE	jbjorkqv	00	2 892 K	11 216 K	387 140	7 892 K	44 K	8 K
nslookup.exe	jbjorkqv	00	128 K	3 052 K	796	1 352 K	32 K	2 K
SshClient.exe	jbjorkqv	00	1 812 K	6 404 K	4 074	3 144 K	39 K	4 K
QCTRAY.EXE	jbjorkqv	00	10 564 K	15 720 K	33 924	8 168 K	77 K	42 K
ibmprc.exe	jbjorkqv	00	576 K	1 644 K	602	412 K	16 K	1 K
cmd.exe	jbjorkqv	00	76 K	2 536 K	971	2 012 K	31 K	2 K
ibmmessages.exe	jbjorkqv	00	2 968 K	8 028 K	42 087	4 676 K	43 K	5 K
iexplore.exe	jbjorkqv	00	43 872 K	66 596 K	269 840	61 080 K	109 K	22 K
Tp5crex.exe	jbjorkqv	00	308 K	1 976 K	640	680 K	26 K	1 K
CCAPP.EXE	jbjorkqv	00	6 412 K	16 124 K	27 195	6 908 K	56 K	11 K
TRAYAP~1.EXE	jbjorkqv	00	652 K	3 992 K	1 712	1 672 K	34 K	5 K
TPONSCR.exe	jbjorkqv	00	600 K	2 124 K	869	704 K	27 K	1 K
taskmgr.exe	thinrkqv	02	2 840 K	5 008 K	2 033	1 508 K	34 K	3 K

Show processes from all users

End Process

Processes: 64 CPU Usage: 2% Commit Charge: 496M / 1225M

NT Sidramar; status

Valid (or Active)

The page is part of at least one process' working set.

Transition

The page is not on any page list, and the Memory Manager is in the process of reading the page in from or writing the page out to a file (e.g., a paging file).

Standby

The Memory Manager has recently removed the page from a working set, and any copy of the page on disk is consistent with the in-memory copy.

Modified

The Memory Manager has recently removed the page from a working set, but the process modified the page, and a copy of the file on disk (e.g., in a paging file) does not have the new content.

Modified No-write

The Memory Manager recently removed the page from a working set, but something, typically a file system driver, marked the page so that it didn't automatically write back to a file, even if a process modified the page. A file system indicated that the Memory Manager need not write the data to disk yet.

Free

The page is not part of any working set, and the process that referenced it removed any other references to it (e.g., it was deallocated). The zero-page thread has not cleared the page (i.e., filled it with zeroes).

Zeroed

The Memory Manager freed the page, and the zero-page thread cleared it.

Bad

The MMU detected the page as faulty, and the page is now off-limits.

NT Working set

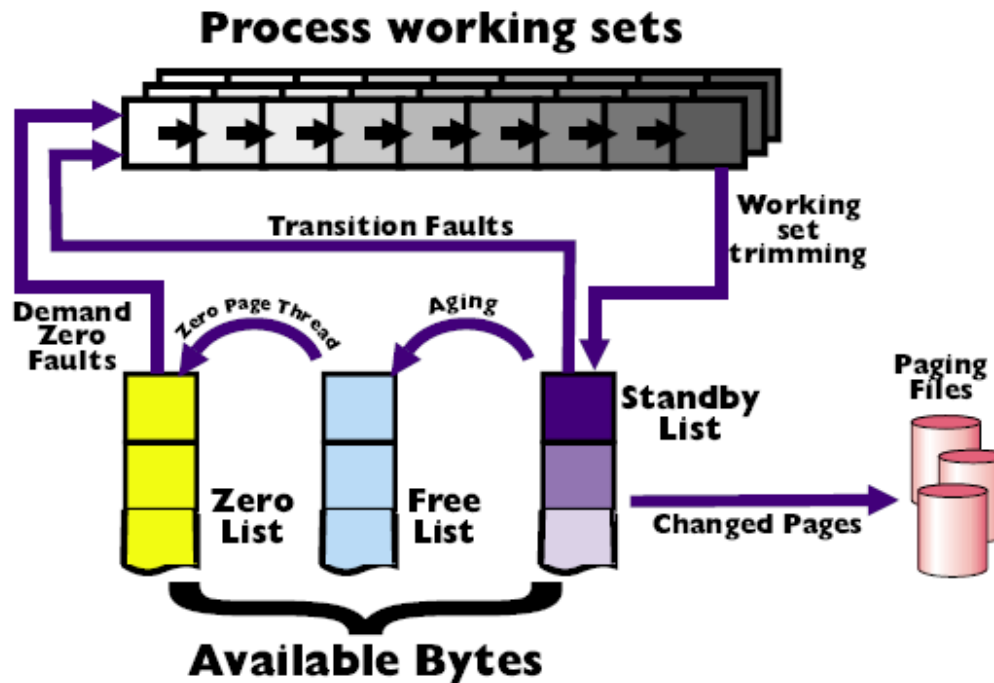
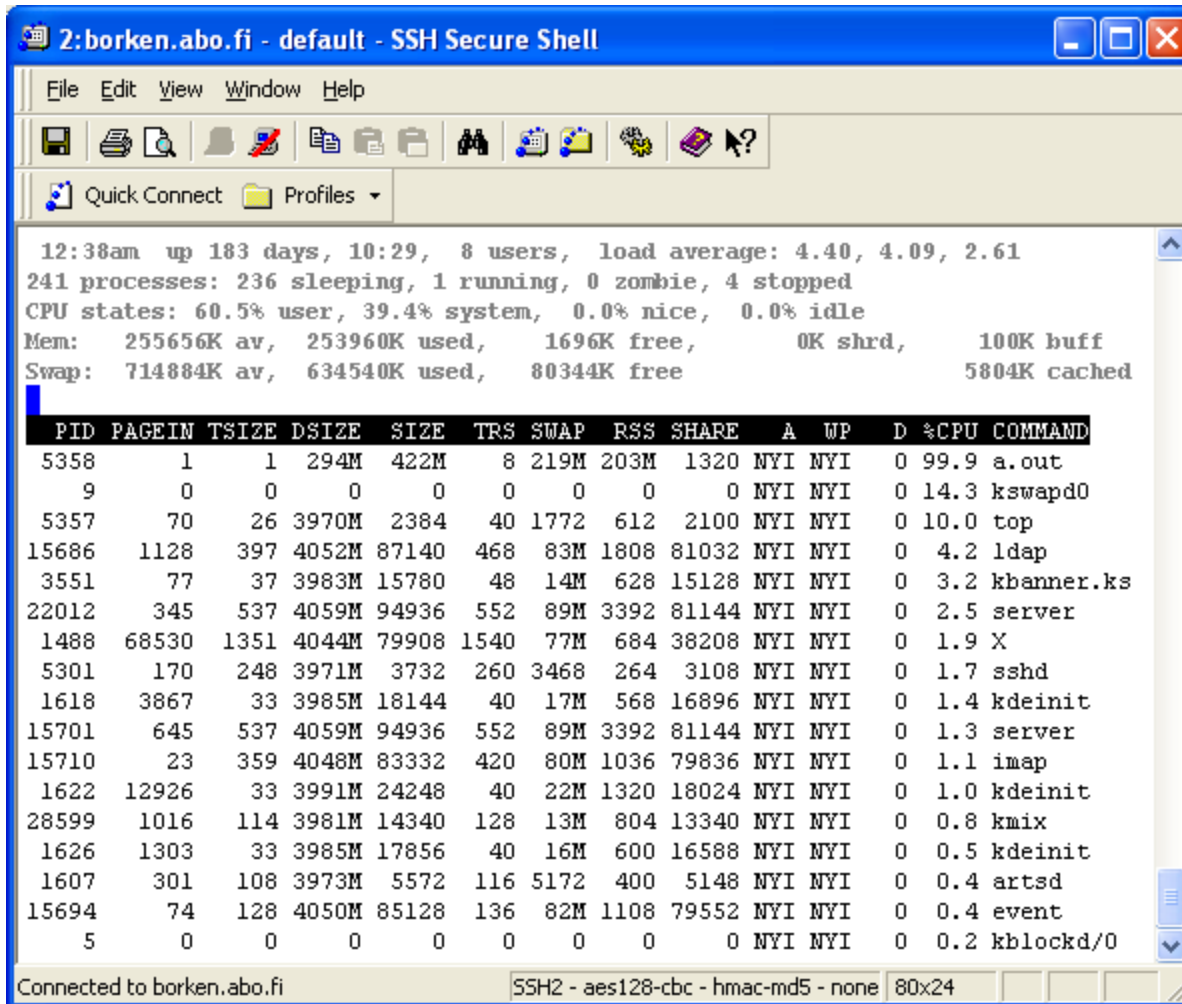


FIGURE 8. PROCESS WORKING SET TRIMMING IN WINDOWS NT. PERIODICALLY, ALL PAGES ABOVE THE PROCESS MINIMUM WORKING SET ARE STOLEN AND PLACED ON THE STANDBY LIST. TRIMMED PAGES AGE AND MOVE TO THE FREE LIST AND, ULTIMATELY, TO THE ZERO LIST. MEANWHILE, TRIMMED PAGES WHICH ARE STILL ACTIVE ARE "SOFT-FAULTED" BACK INTO THE PROCESS WORKING SET WHEN REFERENCED. MODIFIED PAGES MUST BE COPIED TO A PAGING FILE BEFORE MOVING TO THE FREE LIST. A MODIFIED PAGE WRITER THREAD OF THE KERNEL WRITES PAGES TO DISK IN BULK WHENEVER A DESIGNATED THRESHOLD OF CHANGED PAGES ACCUMULATES.

Linux minneshantering

- 3 GB User space / 1 GB kernel space
- Upprätthåller resident set size (RSS), mängden av sidor en process har allokerade i sidramar
- Sidhämtning: demand paging (ingen prepaging)
- Sidbyte: Soft-LRU, räknare som uppdateras vid accesser, högerskiftas då ej accesserad

Linux minneshantering



2:borken.abo.fi - default - SSH Secure Shell

File Edit View Window Help

Quick Connect Profiles

12:38am up 183 days, 10:29, 8 users, load average: 4.40, 4.09, 2.61
241 processes: 236 sleeping, 1 running, 0 zombie, 4 stopped
CPU states: 60.5% user, 39.4% system, 0.0% nice, 0.0% idle
Mem: 255656K av, 253960K used, 1696K free, 0K shrd, 100K buff
Swap: 714884K av, 634540K used, 80344K free 5804K cached

PID	PAGEIN	TSIZE	DSIZE	SIZE	TRS	SWAP	RSS	SHARE	A	WP	D	%CPU	COMMAND
5358	1	1	294M	422M	8	219M	203M	1320	NYI	NYI	0	99.9	a.out
9	0	0	0	0	0	0	0	0	NYI	NYI	0	14.3	kswapd0
5357	70	26	3970M	2384	40	1772	612	2100	NYI	NYI	0	10.0	top
15686	1128	397	4052M	87140	468	83M	1808	81032	NYI	NYI	0	4.2	ldap
3551	77	37	3983M	15780	48	14M	628	15128	NYI	NYI	0	3.2	kbanner.ks
22012	345	537	4059M	94936	552	89M	3392	81144	NYI	NYI	0	2.5	server
1488	68530	1351	4044M	79908	1540	77M	684	38208	NYI	NYI	0	1.9	X
5301	170	248	3971M	3732	260	3468	264	3108	NYI	NYI	0	1.7	sshd
1618	3867	33	3985M	18144	40	17M	568	16896	NYI	NYI	0	1.4	kdeinit
15701	645	537	4059M	94936	552	89M	3392	81144	NYI	NYI	0	1.3	server
15710	23	359	4048M	83332	420	80M	1036	79836	NYI	NYI	0	1.1	imap
1622	12926	33	3991M	24248	40	22M	1320	18024	NYI	NYI	0	1.0	kdeinit
28599	1016	114	3981M	14340	128	13M	804	13340	NYI	NYI	0	0.8	kmix
1626	1303	33	3985M	17856	40	16M	600	16588	NYI	NYI	0	0.5	kdeinit
1607	301	108	3973M	5572	116	5172	400	5148	NYI	NYI	0	0.4	artsd
15694	74	128	4050M	85128	136	82M	1108	79552	NYI	NYI	0	0.4	event
5	0	0	0	0	0	0	0	0	NYI	NYI	0	0.2	kblockd/0

Connected to borken.abo.fi SSH2 - aes128-cbc - hmac-md5 - none 80x24

Linux sidramar

- Kärnan håller reda på tillståndet för varje sidram genom en räkka av deskriptorer av typen *struct page*
 - Räckan med sidramsdeskriptorer kallas *mem_map*
 - Håller reda på statusen för varje sidram (*== 0 is free, > 0 is used*)
 - Flaggor för *dirty, locked, referenced, etc.*
- Kärnan allokerar och frigör sidramar genom *__get_free_pages(gfp_mask, order)* och *free_pages(addr, order)*

Linux sidramar

```
struct page {
    unsigned long flags;                /* atomic flags, some possibly
                                        updated asynchronously */
    atomic_t count;                    /* Usage count, see below. */
    struct list_head list;             /* ->mapping has some page lists. */
    struct address_space *mapping;     /* The inode (or ...) we belong to. */
    unsigned long index;               /* Our offset within mapping. */
    struct list_head lru;              /* Pageout list, eg. active_list;
                                        protected by zone->lru_lock !! */

    union {
        struct pte_chain *chain; /* Reverse pte mapping pointer.
                                    * protected by PG_chainlock */
        pte_addr_t direct;
    } pte;
    unsigned long private;             /* mapping-private opaque data */
#ifdef WANT_PAGE_VIRTUAL
    void *virtual;                     /* Kernel virtual address (NULL if
                                        not kmapped, ie. highmem) */
#endif /* WANT_PAGE_VIRTUAL */
};
```

Linux sidramar, statusflaggor

PG_locked - I/O på gång

PG_error - I/O error när sidan flyttades

PG_referenced - Nyligen accesserad för I/O

PG_uptodate - Satt efter I/O

PG_dirty - Sidan har modifierats

PG_lru - Finns i aktiv / inaktiv lista

PG_active - i aktiva listan

PG_slab - intern minnesallokering

PG_reserved - kernelen använder sidan

Linux minneshantering

```
#include <malloc.h>
main() {
    long i,j;
    void *p;
    for (i=0; i<4000; i++) {
        p = malloc(1<<20);
        if (p) ((int *)p)[0] = 12;
        printf("allokerat %lu bytes, p: %p\n", i*(1<<20), p);
        if (p) {
            for (j=0; j<1<<20; j++) {
                ((char *)p)[j] = 3;
            }
        }
    }
    while (1);
}
```

Linux sidramar

- In theory, paging eliminates the need for contiguous memory allocation, but...
 - Some operations like DMA ignores paging circuitry and accesses the address bus directly while transferring data
- As an aside, some DMA can only write into certain addresses
 - Contiguous page frame allocation leaves kernel paging tables unchanged, preserving TLB and reducing effective access time
- As a result, Linux implements a mechanism for allocating contiguous page frames

Allokering av löpande sidramar

Buddy system algorithm

- All page frames are grouped into 10 lists of blocks that contain groups of *1, 2, 4, 8, 16, 32, 64, 128, 256, and 512* contiguous page frames, respectively
 - The address of the first page frame of a block is a multiple of the group size, for example, a 16 frame block is a multiple of 16×212
- The algorithm for allocating, for example, a block of 128 contiguous page frames
 - First checks for a free block in the *128* list
 - If no free block, it then looks in the *256* list for a free block
 - If it finds a block, the kernel allocates 128 of the 256 page frames and puts the remaining 128 into the *128* list
 - If no block it looks at the next larger list, allocating it and dividing the block similarly
 - If no block can be allocated an error is reported

Repetition

- Minneshantering

- Se till att processer har tillgång till minne;

- Virtuellt minne:

- Varje process får en virtuell adressrymd (=möjlighet att läsa / spara data)
 - OS ser till att det tillgängliga fysiska minnet används ändamålsenligt (bl.a. delas mellan processer)
 - Upprätthåller per process en minneskarta över tillgängligt minne
 - Upprätthåller datastrukturer över hur det fysiska minnet används
 - Ser till att data vid behov laddas / flyttas till och från den fysiska minnet