

# Virtual Memory

①

Previous: Required entire process to be in memory  
 not always needed: locality of reference

## Loc of Ref:

Programs spend most time in small piece of code

Knuth: 90% of time in 10% of code

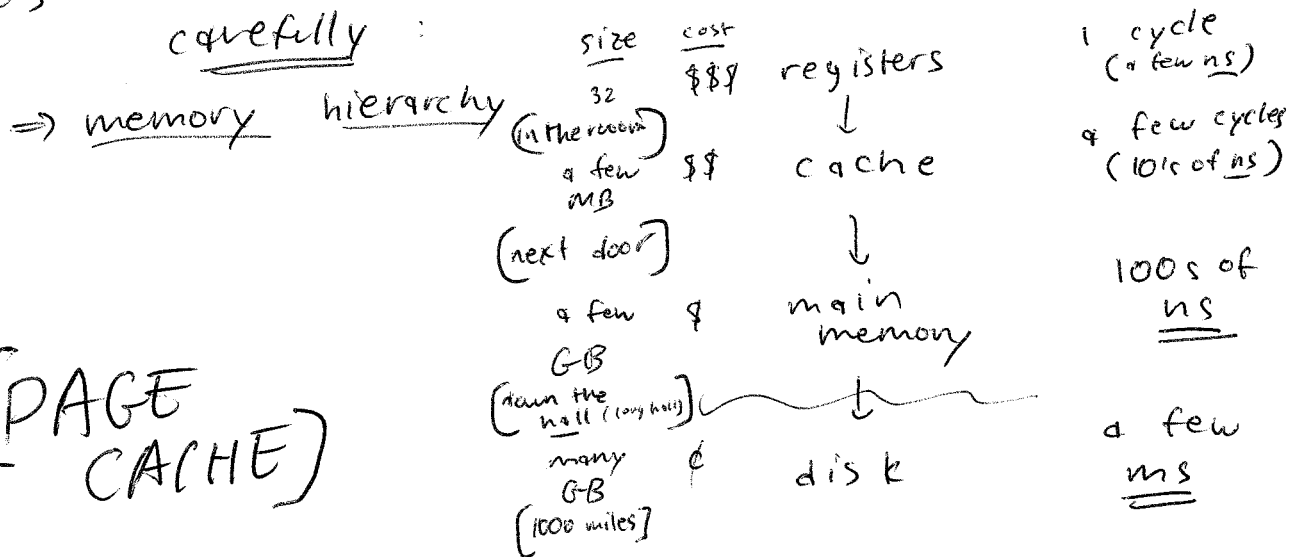
Don't need entire ~~AS~~ in mem @ once

Idea: keep unused pages on ~~store~~ disk!

Process can run even if all pages not loaded!

OS: must manage which pages in mem, disk

carefully:

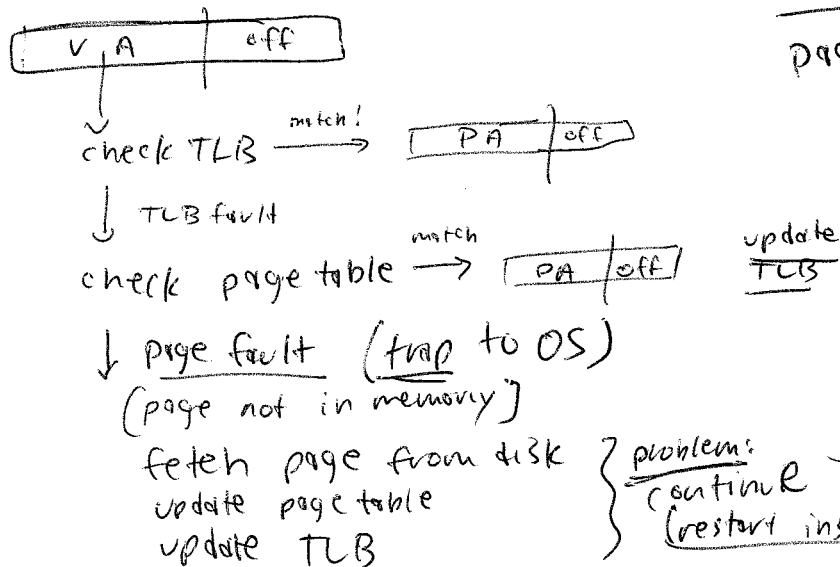


(PAGE CACHE)

## Implication: Page Access

needed:  $H/W + S/W$

page table: present bit



problem:  
 must kick out old page  
 which one?

problem:  
 (continue  
 (restart instruction))

# Continuing a Process: Mechanisms

2

Can be difficult  
Depends on ISA

simple case:

load  $\langle \text{addr} \rangle, R_2 \rightarrow$  page fault on addr  
↓  
just restart: OK

copy  $\langle \text{addr}_1 \rangle, \langle \text{addr}_2 \rangle, R_i$  size  
( $R_{dst}$ ), ( $R_{src}$ ),  $R_{size}$   
only OK if copy is not destructive

[move  $+(SP), R_2$ ] auto increment

H/W support helps

## OS Decisions: Policies

Page Selection: when to bring page from disk to memory

Replacement: which page(s) to move to disk?

## Selection:

Demand Paging: load page only on fault  
start up w/ none in memory, wait until referenced

(No OS choice) Request paging: user-specified/control  
really hard. overlays

Prepaging: Load page before referenced  
hard to know when this is a good idea

Clustering: when bringing in 1 page, bring in many  
(e.g. code pages)

# Page Replacement

③

Optimal: Throw out page that won't be accessed for longest time in future  
 → best if access pattern is known (not practical) ideal

frequency based

LFU:

least frequently used when bad?

MFU:

used little, must be needed soon!

simple

Random: Pick any @ random  
 easy to implement, may work OK

FIFO: simple: throw out oldest page  
 Fair per page

recency based

LRU: Throw out page that hasn't been used in longest time  
 (past predicts future?)  
 w/ locality, LRU approximates OPT

E.g. / Ref string A B C A B D A D B C B

FIFO      OPT      LRU  
 A B C      A B C      A B C

A  
 B  
 D replace! D B C  
 A replace! D A C  
 D replace! D A B  
 B D A B  
 C C A B  
 B fault -

(F) A B D

(F) A B D <sup>LRU</sup>

fault A B C or C B D

fault C B D

String

1	2	3	4	1	2	5	1	2	3	4	5
A	B	C	D	A	B	E	A	B	C	D	E

(4)

FIFO : 3

	1	2	3	}	9
	<del>1</del>	<del>2</del>	<del>3</del>		
	<del>4</del>	<del>1</del>	<del>2</del>		
	5	3	4		

FIFO : 4

	<del>1</del>	<del>2</del>	<del>3</del>	<del>4</del>
	<del>5</del>	<del>1</del>	2	3
	4	5		

$4 + 4 + 2 = 10$

>> Belady's Anomaly <<

# Implementing LRU

TLB miss, page fault  
TLB miss, page hit

(5)

TLB hit, page ~~hit~~ ~~fault~~ hit  
TLB ~~hit~~, page ~~hit~~ ~~fault~~ fault

## Software Perfect LRU:

Keep ordered list of pages

on memory reference, move page to front of list

Replacement: remove page from back of list

## H/W Perfect LRU:

Register / page

Memory ref: put system clock in register

clock overflow

Replacement: scan, find oldest

=> Impractical to support efficiently

=> approximate: old page, not oldest

# Clock Algorithm "second chance"

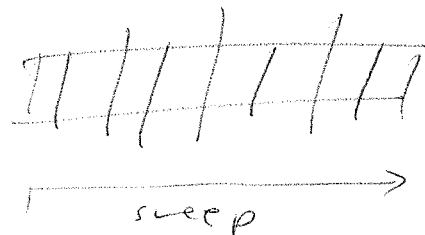
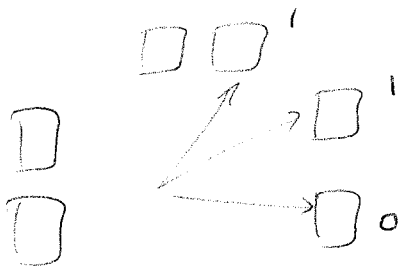
H/W Keep use / reference bit per page

On memory reference, set bit

S/W

Replacement: Look for page w/ use bit cleared

Scan: traverse all pages, clearing bits over time



> what if too fast?

too slow?

## Extension : Dirty Bit

(6)

Add another bit: dirty

page is dirty if it's been modified

when replaced, must be written to disk

⇒ Prefer kicking out ~~#~~ clean : why?

Problems can occur when:

clean pages ref'd frequently (code)

⇒ might get thrown out instead  
of lesser used dirty page

## Extension : Dirty Bit

(6)

Add another bit: dirty

page is dirty if it's been modified

when replaced, must be written to disk

⇒ Prefer kicking out ~~the~~ clean : why?

Problems can occur when:

clean pages ref'd frequently (code)

⇒ might get thrown out instead  
of lesser used dirty page

# How to Allocate Across Processes?

①

2 ways

Global Replacement

Per-process replacement

## Global

All pages lumped in single pool

(all compete for page frames) ⇒ Fault: replace any page

+ flexible, no limits

- one process can dominate

## Per-Process (Local)

Each process has separate pool  
fixed # of pages, frac of phys mem

Page fault: replace your page

+ isolation

- inefficient

⇒ also, could do per user

## Hybrid (VAX/VMS)

Per-process limits  
when removed from per-process list

⇒ Global list

only when removed from here ⇒ disk

miss in per-process list checks global list



# Over committing Memory

Processes: actively accessing  $N+1$  pages  
only  $N$  page frames

Cycle / ~~...~~

(Ref page not in memory  
replace page w/ new one)

## ⇒ Thrashing

System just moving pages back and forth  
(many) each { load instr } ⇒ disk!  
{ instruction fetches }

Illusion of virtual memory breaks  
(no longer an "infinite" amt of fast mem)

## ⇒ E.g.

$H$ : % of hits

(10-  
~100ns)

$C_{mem}$ : in-mem page access

$C_{fault}$ : page not in mem

~25 ms

(3-18 ms)

Overall cost =

$$H \cdot C_{mem} + (1-H) \cdot C_{fault}$$

if ~~just~~  $H = 97\%$

$$\Rightarrow O_{Cost} = 750 \text{ } \underline{\underline{\mu s}} \left( \frac{7500 \times}{\text{in-mem}} \right)$$

⇒ even small miss rate is a problem

System : does not know it has  
too much work

Page replacement:

view is too narrow

(only which page to replace)

Note say:

Student analogy: too many courses  
drop one!

(or focus on one/few @ a time,  
ignore others

=> allows progress)

OS approach: "Admission Control"

don't let everyone run

only run those who fit in memory

rest are swapped to disk

how to detect?

certain amt of paging going on?

what are potential problems?

starvation: large memory jobs  
get stuck waiting

single mem hog: never runs?

# Different Approach: working sets

(4)

Informally: "working set" is ...

> collection of pages ref'd frequently at a given time by process

> must be resident to avoid thrashing

Locality:

recent use  $\Rightarrow$  future?

More formally:

(unique) pages ref'd by process w/ in last  $\sim$  seconds

$\sim$ : the WS parameter

time

A A B C B B B C D E D E B ...

WS = ABC

WS = BC B

Goals:

Ensure WS / process remains in memory

Pages not in WS: Discard whenever

Long-term:

Process not executed unless

WS resident in main memory

# Balance set

Keep all processes in mind

Two groups:

Active : WS is loaded

Inactive : WS on disk

"Balance set":

sum of WS's across active processes

Long-term:

Policy #1:

if Bset > mem,

move some active → inactive

Policy #2:

when to move processes

back from inactive ⇒ active?

# Implementing WS

⑥

Hard to know real WS

what pages have been accessed  
w/ in last  $\sim$  seconds?

analogy to capacitor

charge on mem ref,  
discharge if not ref'd

if charge is "low", not in WS

Really, leverage use bits

+ idle time per page

idle time: amt of CPU received by  
process since last access

Periodically scan resident pages

if use bit on, clear idle time

if off, add CPU (since last scan)  
to idle time

clear use bit (set ~~to~~ to off)

How often to scan?

too fast?

too slow?

# Unresolved

7

what should  $\frac{1}{i}$  be?

too large?

small?

How to pick active processes?

How to compute is pages are shared?

How much memory is needed?

# Trends

Memory is cheap

⇒ Not under contention as in early machines

⇒ If system pages a lot, something is wrong

Large page sizes

⇒ UltraSPARC ~~8k~~, 1MB

Why? ~~fewer~~ fewer pages to manage

⇒ more internal frag

Larger Virtual AS

64-bit VA

32-bit limits some apps (a few GB)

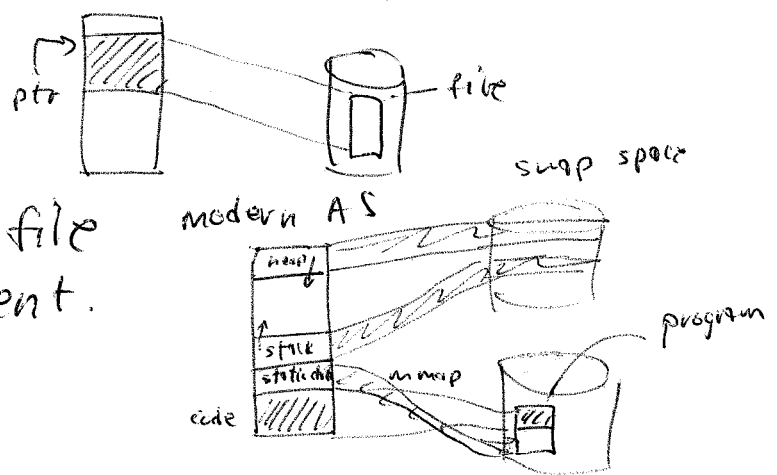
Page tables: must be sparse, or inverted

Integrated VM / File I/O

⇒ madvise()  
CPU: worst case

mmap()

access: fault,  
data read from file  
update: dirty, event.  
flush to disk

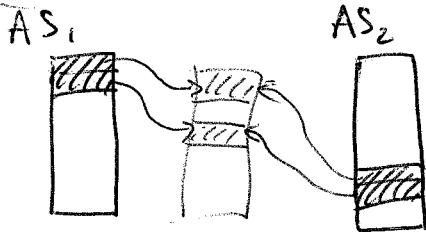


# Advanced Topics

1

## Sharing Memory

`shmctl()` `shmget()`  
`shmat()`



## Multiprocessors

- some memory faster to access than others
- page migration

## Copy-on-write



`fork()`

exact copy



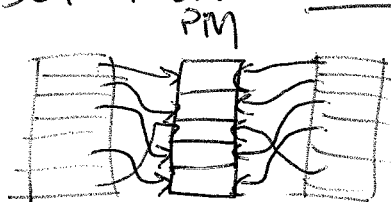
`execv()`

"is"  
(overlay w/  
new  
image)



⇒ wasteful : why go to all the work of that copy

⇒ solution: lazy copy



- mark page as read only (protection)
- when written to, catch fault, ⇒ write make copy



# Advanced Topics

(2)

> Software emulated "use" bit

Q) how to tell if page has been used by process recently?

H/W: use bit

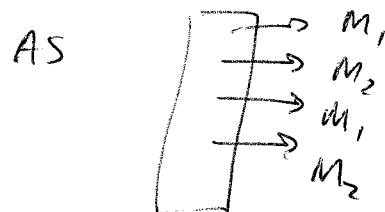
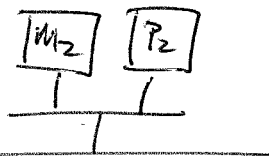
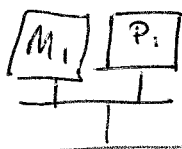
s/w: use protection bits!

mark all as unreadable when accessed, fault, OS

notes that page has been read, written (s/w)  
pages that have not been unmarked can be thrown out!

> Multiprocessors

Some memory is cheaper to access than others



⇒ how to fix?

> page migration (data → process)

> process migration (process → data)