# **Page Replacement**



**Operating Systems** Baochun Li University of Toronto

### **We have seen that the OS allocates memory frames to programs on demand (i.e., page fault)**

- If no frame is available, then OS needs to evict another page to free a frame
- Which page should be evicted?

### **A page "cache" miss is similar to a TLB miss or a memory cache miss**

- However, a miss may require accessing the disk
- So miss handling can be very expensive
- Disk access times are at least 1000x memory access times

## **When will paging work well?**

- **Paging can only work well if page replacement occurs rarely**
- **Paging schemes depend on the locality of reference Spatial locality**
	- Programs tend to use a small fraction of their memory, or Memory accesses are close to memory accessed recently

#### **Temporal locality**

Programs use same memory over short periods of time, or Memory accessed recently will be accessed again

### **Programs normally have both kinds of locality, and the overall cost of paging is not very high**

### **Why not just evict a random page?**

**If a page evicted is used again in the near future, it needs to be brought back into memory**

**Challenge: How to find a page that is least used to evict?**

Same problem applies to other cache systems (such as memory cache and web cache)

#### <u>i ne opumai Page Replacement Aigorithm</u> **The Optimal Page Replacement Algorithm**

### **The page that is not needed for the longest time in the future should be evicted**

Assuming that the future can be perfectly predicted



### **The Optimal Page Replacement Algorithm**

#### **Problem**

- We cannot accurately predict the future
- So we do not know when a given page will be next needed
- The optimal algorithm is not realizable

### **However it can be used in simulation studies**

- Run the program once
- Generate a log of all page references
- Use the log in the **second run** to simulate optimal algorithm
- Use the "optimal" algorithm as a control case for evaluating other algorithms

%1.+#,)()/0#(51%6#,)(\$&+#7

### **Replace the page that has been in memory for the longest time (oldest page)**



Baochun Li, Department of Electrical and Computer Engineering, University of Toronto

### **Implementation for replacing the oldest page**

Maintain a linked list of all pages in memory

- Keep the list in order of when pages came into memory
- Add the new page to the end of list

### **On a page fault**

Choose page at the front of the list (oldest page)

### **Problem**

The oldest page may be needed again soon Some page may be important throughout execution When it gets old, replacing it may cause immediate page fault

## **FIFO suffers from Bélády's anomaly**

#### **Bélády's anomaly —**

Increasing the number of page frames results in an increase in the number of page faults

This is very bad!

**Need to predict page access pattern in the future**

But we can only learn from the past

## **One idea — pages used in the recent past should not be evicted**

Assumption: pages used recently are likely to be used in the near future

Need a way to track past page references

Requires hardware support!

### **Page Table Bits Revisited**

### **Each page table entry has:**

#### **Referenced** bit

Set by CPU when the page is read or written

Cleared by OS software (never cleared by hardware)

#### **Modified (dirty)** bit

Set by CPU when the page is written

Cleared by OS software (never cleared by hardware)

### **TLB may have the most recent copy of them**

Hardware/OS must synchronize it with page table entry bits

### **Can we use page table bits to estimate the page access pattern in the past?**

Baochun Li, Department of Electrical and Computer Engineering, University of Toronto

## **FIFO, but give a second chance to referenced pages Maintain a linked list of all pages in memory**

New pages are added to the end of list

### **On a page fault**

Look at the first page in the list (oldest page)

If its referenced bit is 0, select it for replacement

Else, **clear** referenced bit, move page to end as if it is a **new** page, repeat

### **If every page was referenced, then second chance reverts back to FIFO**

#### **Clock Algorithm !" A ! A ! A ! A ! A ! A**

## **An implementation of second chance Maintain a circular list of pages in memory On a page fault**

The "hand of the clock" sweeps over circular list The "hand of the clock" sweeps over

Looks for a page that does not have the referenced bit set (instead of moving pages in FIFO list) Looks for a page that does not have the referenced bit set (ins



Baochun Li, Department of Electrical and Computer Engineering, University of Toronto

- **Also called Not Recently Used (NRU)**
- **Replace the page that is not recently used**
- **Initially, all pages have** 
	- Referenced bit  $= 0$
	- Dirty bit  $= 0$

### **Periodically (e.g., every timer interrupt) clear the referenced bit of all pages**

Then, the referenced bit indicates that a page was recently accessed

#### **Enhanced Second Chance Not Recently Used (NRU)**

### **On a page fault, pages are in 4 classes**



**Choose a random page from the lowest non-empty Class to evict** 

Baochun Li, Department of Electrical and Computer Engineering, University of Toronto

#### ! #\$%&'()&\*&)+\$,'\$-./\$+01+\$%&2314&5\$+0&\$216&\$+01+\$ **Least Recently Used (LRU)**

### A refinement of NRU that replaces the page that has not been used for the longest period **of time**

Needs to track how recently a page was used



### **Keep a stack of all pages**

### **On each memory reference**

Move corresponding page to the top of the stack Best implemented as a doubly linked list

## **On a page fault**

Choose page at the bottom of the stack (least recently used)

#### **!"#\$%&'()&)\*+,+-.\*\$/ 0'+-.\*\$1 LRU Implementation — Option 1**

#### **Move referenced page to the top of the stack**



Baochun Li, Department of Electrical and Computer Engineering, University of Toronto

### **Problems**

- Requires moving list elements on each memory access
- Each memory access becomes several accesses!
- Not implementable with hardware

**MMU (hardware) maintains a counter that is incremented for every memory reference** 

### **Every time a page table entry is used**

MMU writes the value of the counter to the page table entry This timestamp value is the "time of last use"

### **On a page fault**

OS software looks through the page table, and identifies the entry with the oldest timestamp

#### **Problem**

Updating of the timestamps must be done for every memory reference

#### **Additional-Reference-Bits Algorithm**

#### **Maintain an 8-bit counter for each page in software**

initially zero

#### **At each timer interrupt (or any regular interval) —**

The referenced bit is shifted into the high-order bit of the counter

- The other bits are shifted to the right
- The low-order bit is discarded
- The referenced bit is then cleared

### **On a page fault, evict the page with the lowest counter**

arbitrarily evict one if more than one candidate

## **The Additional-Reference-Bits algorithm**



Baochun Li, Department of Electrical and Computer Engineering, University of Toronto

### **Problems with Additional-Reference-Bits**

## **Granularity of record-keeping is limited by the frequency of the timer interrupts**

Records only one bit per interval

Lost the ability to distinguish references early in the clock interval from those occurring later

### **Counters have a finite number of bits**

Limits its past time horizon

All we can do is to pick one of them at random

### **Comparison of Page Replacement Algorithms**



## **Working Set Model**

#### **Locality of reference revisited**

- Spatial locality: processes tend to use small a fraction of their memory
- Temporal locality: processes tend to use the same memory over short periods of time

#### **Working Set**

- The set of pages a process needs currently
- If working set is in memory, no page faults occur

#### **What if you can't get the working set into memory?**

- Thrashing not enough frames to accommodate the WS Page faults occur every few instructions
- The user-level program makes little progress



### **Demand Paging vs. Prefetching**

- **Demand paging: loading pages on demand initially when a process starts to run**
- **Prefetching: load the working set of the process before letting it run, minimizes page faults**
- **But how does the OS determine the working set?**

#### **How Big Is the Working Set?** \$ 4,000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.00<br>\$ 5.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.000 \$6.0 " 9&5,&()\*&:\$/%52;&,\*(&(5.\*&52(\*/6'+

#### **Look at the last K memory references**  <u> Look at the last is memory references</u>

Alternatively, look back over last time T, the working set time interval

As K (or T) gets larger, more pages are needed

#### **Goal: Design a page replacement algorithm that keeps this working set in memory**  $\blacksquare$



### **Approximating the working set model**

- **Approximate with interval timer + a reference bit Example: T = 10,000 time units**
- **Timer interrupts after every 5000 time units**
- **Keep in page table: 2 bits for each page**
- **On a timer interrupt, copy and set the value of the Referenced bit to 0**
- **If one of the bits is 1 => the page is in the WS**

## **The Working Set Algorithm**

### **Hardware sets the R bit when a page is referenced**

The virtual time is the time that the process has run on the CPU

#### **On a timer interrupt —**

If R is 1, it is cleared, and the current virtual time is written to "Time of Last Use"

## **The Working Set Algorithm: On a Page Fault**

#### **On a page fault —**

If R is 1, the current virtual time is written to "Time of Last Use"

If R is 0, and the age (current virtual time  $-$  time of last use)  $>$ T, evict

If R is 0, and age  $\epsilon$  = T, record the page with the greatest age

#### **If the entire page table has been scanned**

If one or more pages has  $R = 0$ , evict the one with the greatest age

Otherwise, all pages have been referenced, evict one at random, preferably a clean page (dirty  $= 0$ )

### **Local vs. Global Replacement**

**Say a process gets a page fault and a page needs to be replaced** 

**Which process's page should be replaced?**

#### **Policy 1: Local page replacement**

Choose a page of the same process

#### **Policy 2: Global page replacement**

Choose a page of any process

#### **Some algorithms can be used with either policy**

e.g., LRU can be used with both local or global replacement But not the Working Set algorithms

#### **Local vs. Global Replacement** Example: Process A has a page fault

#### **Example: Process A has a page fault**



A0  $A<sub>1</sub>$ **Replace**   $A2$ **earliest**  A<sub>3</sub> **local page** A4 A<sub>5</sub> **Replace**  B<sub>0</sub> **earliest**  B<sub>1</sub> **global page**B<sub>2</sub>  $\bigcirc$ **B4 B5** B<sub>6</sub>  $C<sub>1</sub>$ C<sub>2</sub> C<sub>3</sub>

Baochun Li, Department of Electrical and Computer Engineering, University of Toronto

### **Problem With Local Page Replacement**

- **Suppose there are 10 processes and 5,000 frames in memory**
- **Should each process get 500 frames?**

### **No**

Small processes do not need all those pages

Large processes may benefit from more frames

#### **Idea**

Look at the needs of each process and give each an adequate number of frames to prevent thrashing But how?

### **Page Fault Frequency**

#### **Page fault frequency declines as a process is assigned more pages** Page fault frequency declines as a program is age ruurt rieguerieg u



Number of page frames assigned

## **Page Fault Frequency**

- **The page fault frequency provides an estimate of the working set needs of a process**
- **Goal: Allocate frames so that the page fault frequency is roughly equal for all processes**
- **How should the page fault frequency be measured?**

#### **For each process**

On each fault, increment a counter **f**

 $f = f + 1$ 

Every second, update Fault Frequency (**ff**, in faults/second) via **aging**

ff =  $(1 - a)$  \* ff +  $a$  \* f, f = 0  $(0 < a < 1$ , when  $a \rightarrow 1$ , history is ignored)

#### **This global page allocation algorithm can then be combined with a local page replacement algorithm**

**It is expensive to run replacement algorithm on each page fault**

### **Instead, OS can use a paging daemon to maintain a pool of free frames**

Runs replacement algorithm when pool reaches low watermark

Writes out dirty pages and frees them

Frees enough pages until pool reaches high watermark

### **Frames in pool still hold previous contents**

Can be rescued if page is referenced before reallocation

#### **Demand paging**

#### **Working-set minimum and working-set maximum (usually 50 and 345 pages)**

If page fault occurs for a process below its working-set maximum: allocates a page from a list of free pages

If at the working-set maximum, select a page for replacement using a local page replacement policy (variation of the Clock algorithm)

#### **If the amount of free pages falls below a threshold —**

Evaluate the number of pages allocated to a process

If it is more than the working-set minimum, evict pages until it reaches the minimum

### **What we've covered so far**

### **Three Easy Pieces: Chapter 22 (Beyond Physical Memory: Policies)**