

self study

Pages Tables for large Memories

- Instead of using page-table-in-memory, we can use TLB to speed up virtual address to physical address translation

Multi-level Pages Table

- we have a 32-bit virtual Address that is partitioned into:
 - a 10-bit PT1 field
 - a 10-bit PT2 field
 - and 12-bit offset field
- since offset is 12-bit, pages are 4KB thru there are 2^{20} pages



- Avoid keeping all the pages tables in memory all the time, in particular those that are not needed should not be kept around

What should we avoid with pages table?

- In a two-level page works in this way

A. On the left we have the top-level page table with 1024 entries, corresponding to the 10-bit PT1 field

- When a virtual address is presented to the MMU

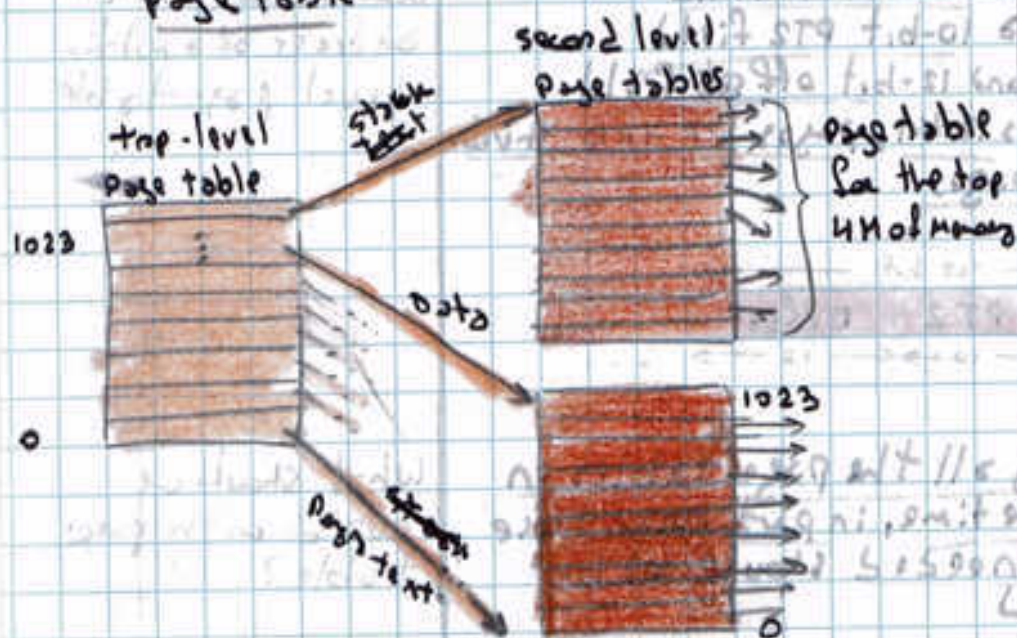
① extract it the PT1 field and uses this value as an index into the top level page table



Multi-level Pages Table (Continue)

- Each of these ~~1024~~ 1024 entries represent 4MB because the entire 4GB virtual address space has been chopped into chunks of 4096 bytes.

B. The top entry located by indexing into top level page table yields the address of the page frame number of a second-level page table



- Entry 0 of the top-level page table points to the page table for the program text.
- Entry 1 points to the page table for the data.
- Entry 1024 points to the page table for the stack.
- The PT2 field is now used as an index into the selected second-level page table to find the page frame number for the page itself.

self study

- Additional level give more flexibility, but it is doubtful that additional complexity is worth it beyond three levels.

Ex: ① 32 bit virtual address: $0x00403004$
 (4,206,596 decimal)
 (12,292 bytes into the data)

② virtual address correspond to

- PT1: 1
- PT2: 2
- Offset: 4

• MMU uses PT1 to index into top-level table and obtain entry 1, which correspond to address 4M to 8M

• Then MMU uses PT2 to index into the second level page table just found and extract entry 3 which correspond to addresses 12292 to 16392 within 4M chunk

→ This entry contains the page frame number of the page containing virtual address $0x00403004$

if page is not in memory, the present/absent bit in the page entry will be zero, causing page fault

if page is in memory, the page frame number taken from the second-level page table is combined with the offset (4) to construct the physical address.

↓
 The address is put on the bus and sent to memory.

What does entry 1 correspond to?

What is PT2 used for?

What happens if page is not in memory?

What happens if page is in memory?

~~Page Tables~~

Multi-level Pages table (continue)

- 32 bit addresses with 2 page table fields
- two-level page tables
- PT to sig for MMU
 - Place it in main memory
- How does MMU know where to find PT?
 - Registers (CR2 on Intel)

Inverted Page table

old level - out of as was as at 1T9 zero UMM.

of ST9 zero UMM. but should take old of page level base

level page alt. missing page entry. level page alt. be valid

entry alt, given is tan; 10 9 1 0

level page alt. given is tan; 10 9 1 0

level page alt. given is tan; 10 9 1 0

self study

Impact of Page Size on Page table

Small page size:

- Advantages:
- less internal fragmentation
 - Better fit for various data structures, code, sections

- Disadvantages:
- Program needs more pages and has larger page table.

Notes:

Bigger

if small page size then small data structure it will

if small page size then small data structure it will

will

will

will

A. Define the virtual address conversion

to swap table, two pages against given address

in virtual address conversion process, convert virtual address into physical address

the virtual address is first converted into physical address, then the physical address is converted into virtual address

the table is maintained in memory, it is called as page table.

Concurrency versus Parallelism

Concurrency:

OS Involvement with Page Table Management

Four times OS deals with with page-tables

1. Process Creation: create a page table

2. Upon Context switch:

A. Load MMU context for a new process
(ie. load CR3 register w/ base address of pagetable)

B. TLB is flushed

3. Page fault time:

A. Determine the virtual address causing fault (read CR2 register, for faulting address)

B. swap target page out, bring needed page in.

4. Process termination time: release page table and other pages

• the goal is that everytime we bring a page to memory, we have to update the pagetable

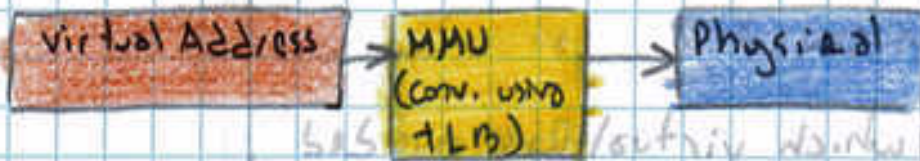
• the table is a mapping between virtual page and the physical page.

copy

self study

• let say in your program you have the following code: $x = x$

you need to access memory of x and the memory address of y so



• Page fault is when accessing virtual page address, there is not a page address.

This is not exactly an error and it is handled by the MMU checking the virtual page address, informing the CPU the virtual address produces the error

→ Page fault is an exception where the CPU tells the OS. Then the OS will handle the CPU and it will bring the page back.

• When a virtual Address page is not located, it may be located in the secondary device = Page fault

When you try to execute or write a code from a page that you don't have permission it will produce a page fault

ex Lets say you do a malloc for a force 100 MB of space. The OS will wait until the program try to access the memory address. Then the OS will incrementally add pages. The first time trying to access it will produce a page fault

Page fault Handling

When the page fault occurs, the OS has to choose a page to evict (remove from memory) to make room for the incoming page.

1. Hardware (MMU) raises a trap to the OS
2. General registers saved
3. OS determines which virtual page needed
4. OS checks validity of virtual address
5. OS selects victim page; if victim dirty (victim dirty refers to whenever you write a piece of memory it becomes dirty) then write to disk
 - this step is normally done out-of-band. i.e. before a page-fault occurs, so that page fault can be handled faster
 - Note: if there is a not good algorithm for this step, the machine will become extremely slow.
6. OS schedules to bring new page in from disk

Dirty bit indicates the page is dirty.
→ if the page is dirty then the context of the page has to be saved
→ if the page is not dirty
→ if you can't have to write the page down and then you have to bring the other page up. This is over-write page.

Page fault handling (continue)

- Victim Dirty: two ways to free space
 1. In case context size over write the page, but it can happen that an application have some data there, so you have to save it. 2. Data that is not sure that you want to save means data here

- Every swap should have twice the amount of memory

7. Pages tables updated to reflect new mappings

8. Fault to instruction based upto whom it began

(we have saved the execution process but now we are losing it here)

9. Fault to process no-schedule

10. registers restored

11. Program continues

• Any algorithm in this has to make sure that the # of pages is minimized.

86 Another way to see this:

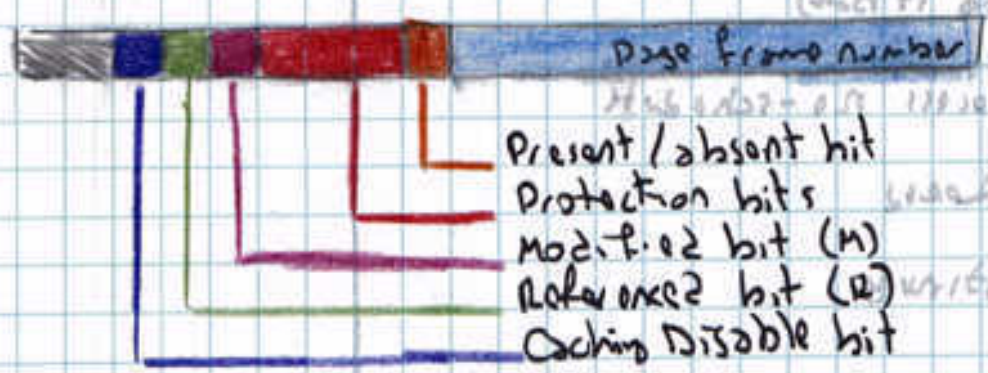
• When a page fault occurs, the OS chooses a page to evict (remove from memory) in order to make space for the incoming page.

→ If the page to be moved has been modified while in memory, ~~it~~ it must be rewritten to the disk to bring the disk copy up to date.

→ If the page has not been changed (i.e. it contains program text), the disk copy is already up to date, so there is no need for rewriting. The page to be read is just overwritten as the page being evicted.

Locating Pages in Memory

Not Recently Used Page Replacement Algorithm



• Two references are used by the OS:

1. Reference bit which is set whenever the page is referenced (read or written)
2. Modified bit which is set whenever is written to (i.e. modified)

• These two bits must be updated on every memory reference ~~and~~ by hardware.

self study

• the Referenced and Modified bits are used as follows:

a. When a process is started up, both bits for all pages are set to 0 by the OS

b. Periodically (by clock interrupt) the referenced bit is cleared in order to distinguish pages that have not been referenced recently from those who haven't have been referenced

c. When page faults occur, the OS inspects all the pages and divides them into 4 categories based on the current values of their referenced and Mod. fied bits:

	Referenced bit	Mod. fied bit
Class 0:	Not Referenced	Not Modified
Class 1:	Not Referenced	Mod. fied
Class 2:	Referenced	Not Modified
Class 3:	Referenced	Modified

i. class 1 occurs when when a class 3 has its reference bit cleared by a clock interrupt

ii. Clock interrupt do not clear the mod. fied bit because this information is needed to know whether the page has to be rewritten to disk or not

iii. clearing reference bit but not mod. fied bit leads to class 1.

38. Not Recently Used Page Replacement Algorithm (continued)

NRU (Not Recently Used): algorithm:

1. removes a page at random from the lowest-numbered non-empty class.
2. The idea is that it is better to remove a modified page that has not been referenced in at least one clock tick than clean page that is in heavy use.
3. ~~NRU~~ NRU is ~~considered~~ considered efficient to implement, and provide optimal performance.

Locating Pages in Memory

- Some times you don't want to kick out pages from the main memory
- Virtual memory an I/O occasionally interact
- Problem: ex.
 - a. Process P1 issues call for read (read system call) for device into buffer (Directly to the device without intervention of the CPU)

b. While waiting for I/O, another process P2 start up

c. P2 has a page fault

d. Buffer for the process P1 may be chosen to be paged out, resulting in DMA (Direct Memory Address) error

Locating Pages in Memory (cont'd)

- Need to specify some pages as "loaded" or "pinned" in memory
 - Those pages are exempted from being victim pages (i.e. in the main memory, you don't want to find it from memory)
- If there are too many processes when you start the computer then the computer will be busy and not responding. This means the computer is swapping a lot of pages in and out of the system.

Page Replacement Algorithms

- Page fault forces choice:
 - Which page must be removed
 - Make room for new memory
- Modified page must first be saved
 - unmodified just overwritten
- Better not to choose an often used page
 - will probably need to be brought back in soon
- Goal of page replacement algorithm is to minimize the number of page faults incurred by the system

90 Page Replacement Algorithms

1. Optimal page Replacement (OPR)
2. Not Recently Used (NRU)
3. First In First Out (FIFO)
4. Second chance
5. Least Recently Used (LRU)
6. Not Frequently Used (NFU)
7. AQS
8. Working set
9. Clock
10. WS Clock

• cpu scheduling Algorithm: to minimize the shortest algorithm first.

1. Optimal Page Replacement Algorithm (OPR)

- A. Replace page that is needed at the farthest point in the future
- i. "Optimal" means that ~~there is~~ no other algorithm can give fewer page faults than optimal page replacement (OPR).
 - ii. Optimal page replacement (OPR) is optimal but impossible
 - iii. It can serve as a useful tool of measurement (baseline) to measure the performance of other algorithms.

1. Optimal Page Replacement Algorithm (OPT) continue

B. Why is it optimal?

- i. OPT says to remove the page with the highest label.
- ii. Pushing page not use farthest away so when there is another page to be pushed away but closed, the lower remove the same. Pushes the page that will fetch soon as far - into the future as possible.
- iii. What if there is another algorithm which does not pick the page needed farthest point in the future and can still yield fewer number of page fault than OPT?

C. Other Algorithms try to approximate OPT

- i. Estimate the pages not needed for a long time by recording the sequence of pages used in the past.

2. Not Recently Used Page Replacement (NRU)

- A. Each page has a Reference bit and a Modified bit which are used to estimate timespan.
 - i. These bits are set when a page is referenced, modified, or close interrupt.

Category	reference bit	mod. bit	status
Class 0	not referenced	not mod. bit	easy to write
Class 1	not referenced	modified	unaccounted cost
Class 2	referenced	not mod. bit	quantity access, locality
Class 3	referenced	mod. bit	used often & dirty

How you can have no reference but not modified? After the Direct Memory Access (DMA) is over the page would be dirty but not referenced.

Not Recently Used Page Replacement (NRU) Algorithm (Cont: NI)

• What happens when someone cursor to page to go to category 2?

Periodically this algorithm will scan, select, and then reset all reference bits back to zero.

i. This means that there is some pages that were reintroduced and some that were moved from category 2 to category 1.

ii. Category 2 means page is dirty but not referenced.

B. NRU periodically scan all memory accesses.

Same (i. If there is low memory, then it would pick a page from the lowest number set of accessed pages

ii. If system is under memory pressure, NRU removes page at random from the lowest numbered non-empty class.

iii. Reset all reference bits

I. This implies that pages moved from $3 \rightarrow 1$ and from $4 \rightarrow 2$.

C. Before the next NRU scan, some pages may be referenced/modified again.

i. Pages modified from $1 \rightarrow \{3, 4\}$ and $2 \rightarrow 4$.

• Maintaining time span requires additional space and process time, so every time this page is accessed or written you have to read the current time span and update it. This is very expensive to maintain so we try to estimate time span.

self study

Not Recently Used Page Replacement Algorithm

- How do you know a page was recently accessed the last time?
reference is equal to access mod. by is write

3. First-In-First-Out Algorithm (FIFO)

A. maintain a linked list of all pages
i. in order they came into memory

B. Page at the beginning of list replaced

C. Disadvantage: Page in memory the longest time may be often used

- Suppose you know that the next the last used page, and an application want to use it. then some way to fix FIFO so second chance can do live.

As the OS maintain a list of all pages currently in memory, with the most recent arrival at the tail and the least arrival at the head.

→ on a page fault, the page head is removed and the new page added to the tail of the list

- FIFO in its pure form is rarely used.

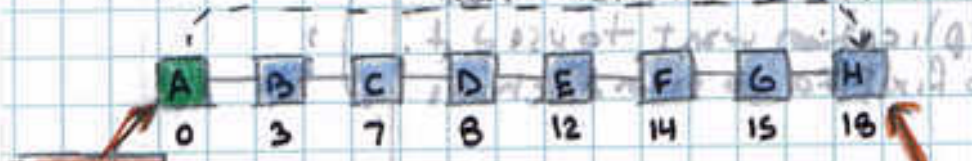
Second Chance Page Replacement Algorithm

• Second chance is a modification to the FIFO algorithm that avoids the problem of throwing out a heavily used page by inspecting the reference bit of the oldest page.

~~Pages are stored in FIFO order~~

A. If reference bit is 0, the page is both old and unused, so it is replaced immediately

B. If reference bit is 1, the bit is cleared, the page is put onto the end of the list of pages, and its load time is updated as though it had arrived in memory. Then the search continues



Page loaded first

Presume A is the oldest page, we look at A and we say: is the reference bit to page A? We check and A will be given a second time so we will move it to the top of the list

Most recently loaded page



A is treated like a newly loaded page

self study second chance page replacement algorithm (cont)

second chance page replacement algorithm (cont)

• When a page-fault occurs:

1. Look at the oldest page in the list
2. If referenced bit of the page is set in its (RTE), then set referenced bit to 0 and move the page to the front of the FIFO list; go back to step 1
3. Select the oldest page with reference bit to 0 as the victim page.
4. If no page with referenced bit (unlikely), then select the oldest page (FIFO)

• If all of them were referenced, then all of them have same priority, so it continues the regular FIFO terminology

• Why is this algorithm bad?

• Linux have two list FIFO so the page would move from inactive list to active list (better).