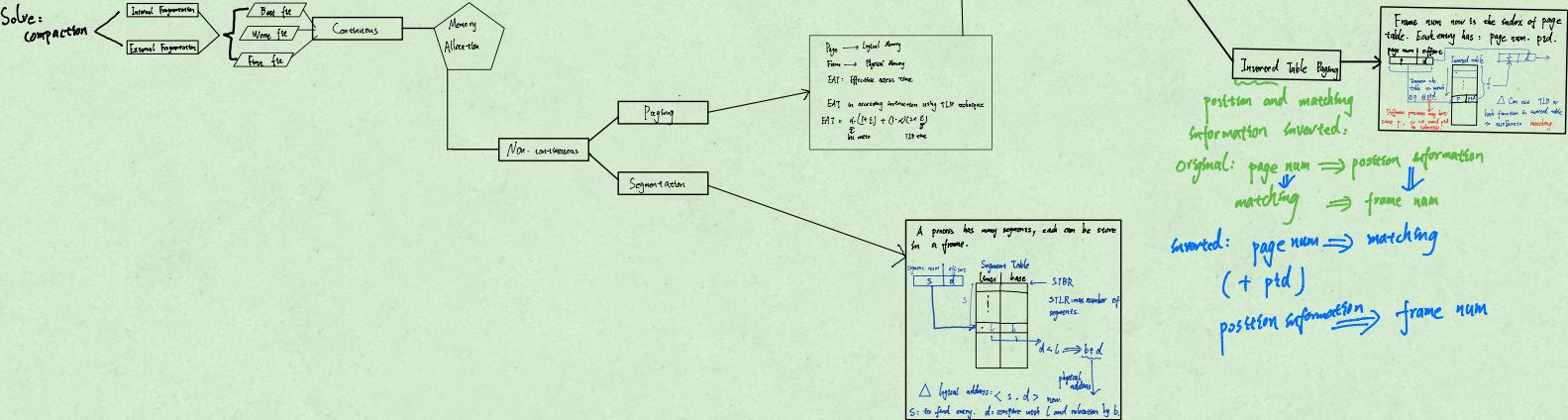
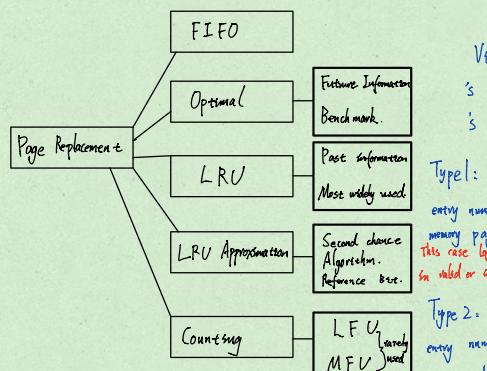


Main Memory:



Virtual Memory

4 elements:
 Logical address. Page Table.
 Main Memory. Virtual Memory
 Page fault handling in page reference:
 ① check if valid or not.
 ② check page table's valid bit.
 ③ find empty frame in main memory.
 ④ Using Replacement strategy to choose victim and replace
 ⑤ update the page table, redo the page fault reference.



Δ In main memory part,
 page table entry num = logical page num
 $=$ main memory frame num.
 So we can always have a valid logical
 page matched with a page table entry.

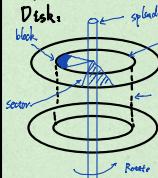
Virtual Memory's page-table
 's difference with main memory
 's page-table:

Type1: Virtual Memory's page-table
 entry num = logical page num = main
 memory page num. \geq main memory
 this case logical page maybe frame num.
 in valid or invalid state in page-table

Type2: Virtual Memory's page-table
 entry num = main
 memory frame num. $<$ virtual memory
 page num = logical page num.
 this case logical page maybe
 in valid or dirty states or
 invalid (\leq 3 in total) state.
 in page-table.

Mass storage system:

Typical: Magnetic Tape. Disk. HDDs. SSDs.



$$\text{Position Time} (\text{Random Access Time}) = \text{Seeking Time} + \text{Rotation Latency}$$

A sequence of disk requests may need to be queued if disk is busy, so we need algorithms to schedule the service of disk.

Scheduling

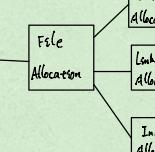
- ① FIFO
- ② SSTF Seeking seek time first
- ③ SCAN
- ④ C-SCAN
- ⑤ LOOK.
- ⑥ C-LOOK.

in the same direction

Algorithms:

File system

File: A collection of relatively data
File activities (same size type) -> long in a directory entry, all the entries are in directory structure.
File operations (read/write open...)
File access & group: user/public/group/working
file organization: single-level, two-level, tree-structure.
directory entry: Name | Size | Type |



External fragmentation problem.
Fast (hard/overhead low) -> Base performance.
Simple
Directory entry: Name | Size | Type |

Slow (many times I/O and disk access).
No external fragmentation.
A file is divided into small blocks in a link list and each block has a pointer to next block.
Directory entry: Name | Size | Block | ... | Size | Block |

File Alloc. Table (FAT)
Stores the file allocation of all file blocks.
Link-List in FAT has FAST (low overhead memory).
Directory entry: Name | ... | Size | Block | ... | Size | Block |

Index table mapping
Index block: R/R Data block.
& Index table block.
(store file information and place pointer to index block, index block's element point to data blocks.
④ Only gets new block when needed)

Convey E effect: share process behind long processes. like if have one CPU bound and many I/O bound but CPU bound share time

① starvation → Aging can solve. ② Need to know the CPU burst of next time
Exponential Aging formula to estimate: $T_{avg} = \frac{1}{N} T_{avg} + (1-\lambda) T_{avg} \frac{(1-\lambda)^{N-1}}{\lambda}$
predicted actual length predicted (convey E)

③ It's optimal. Using SJF can have the min waiting time.

Shortest Remaining-Time Job First (SRPT) → preemptive

① SJF is PS's specific case. ② Starvation and Aging can solve
③ Priority round-robin for each process

① quantum = q , $q > 1$: RR = FCFS
 $q < 1$: Low Latency: $q \gg$ Preempt latency.
② when a process runs out of its slice, if there are n process, it need to wait as much as $(n-1)q$.
③ how-to choose q : 80% of processes CPU burst $\leq q$

① There are many Ready queue, and different scheduling scheme in different queue. (Priority, Round-robin, Preempt, Round-robin)
② Scheduling
Round robin: (May cause starvation, low priority need to wait longer time slice / higher priority queue, larger time slice)

① Process can enter different queue. (In Master level, one process can only be in one queue)
② Much (the Master level), Round-robin scheduling. If don't finish, move to the next level queue.
③ When a process in lower level queue is preempted, then it's done change level and start from a new quantum next turn.

Process Contention Scope: user level threads
→ kernel level thread
① Scheduled by thread library.
② Priority set by programmer.
③ Appear in Many-to-Many / Many-to-one Model

System Contention Scope: kernel level threads → LWP (then to use CPU)

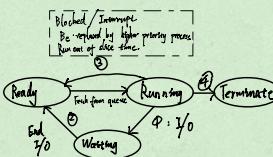
CPU scheduling:

CPU scheduler → preemptive → ① Non-preemptive → ② Dispatcher → Dispatch Latency × Switch Latency

Scheduling Criteria:

- ① CPU Utilization
- ② Throughput
- ③ Turnaround Time { calculate Waiting Time } average
- ④ Waiting Time
- ⑤ Response Time (used in time-sharing)

Process Scheduling



Thread Scheduling

Examples:
Linux
Solaris
Windows XP

They are all similar to each other.

① priority scheduling ② Multi-level queue (Scheduling classes)

At first, ① real-time class
② definite class

③ Relative priority & Global priority.

④ Some processes have fixed priority, others have dynamic priority.

Linux: (1) Two array → Running = higher priority
real-time processes fixed priority Expire = low index
(2) when put into Expire array, = longer quantum
other processes dynamically change priority (nice value) = higher decay factor

(3) when all run array empty done, exchange.

(4) Each time scheduler select: highest class, highest task in the class.

Solaris: (i) six classes. Priority: (Interrupt) Real-time > System > fair-shared > fair > fixed priority > time-shared > interactive

Deadlock:

Deadlock:
what will cause deadlock: when these four conditions are all satisfied:
 ① Mutual Exclusion
 ② Hold and wait
 ③ No preemption
 ④ Circular Waiting.

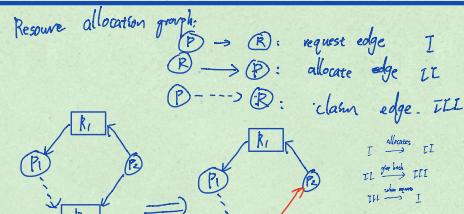
How to solve deadlock problem:
 ① No deadlock.
 ② Possible deadlock:
 Detect and recovery
 Deadlock Prevention
 Deadlock Avoidance

Deadlock prevention:
 break one of four conditions:
 ① Avoid Mutual Exclusion
 ② Avoid Hold and wait process. 把进程
 可抢占 Hold and wait process with lock
 占进程 要等到 (1) do resource get (2) request
 resource yet 才能继续执行.
 ③ No resources Total order. process can
 only request resources 全序 with increasing number

Periodically call Detect algorithm. If has deadlock, need to recovery. Select some process to kill and preempt some resources.
 minimum loss.
 (lowest priority, shortest exec time, largest number of resources, have a long time to complete)

Deadlock avoidance: when allocate resources to a process request, always ensure the system in safe state.
 Process states: U
 Unsafe
 Deadlock
 Resource allocation graph:
 single resource, multiple instances, Banker's Algorithm.

△ Need prior information (maximum demands)



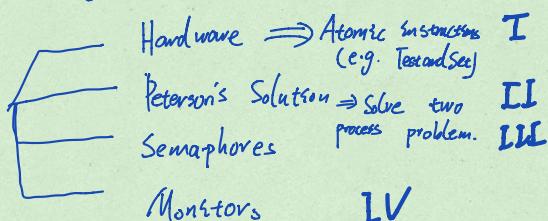
△ request edge turn to allocate edge only when it will not form cycles.

Banker's Algorithm:
 ① Max - Allocation = Need.
 ② Need < Available, 分配给该process.
 ③ Available + = allocation

Each time when we want to allocates resources, assume we have allocated and update the ①, ②, ③ values and call detection Algorithm, find whether its safe.

Process Synchronization:

Synchronization technique



I
II
III
IV

Synchronization Hardware: TestAndSet()
 # lock is initialized to false.
 do {
 waiting[i] = true;
 key = true;
 while (key && !waiting[j]) {
 key = TestAndSet(&lock);
 }
 // do critical section.
 j = (j+1)%size; Assume that we know total
 while (j != i && !waiting[j]) {
 j = (j+1)%size;
 }
 if (j == i) lock = false;
 else waiting[j] = false;
 } while (true)

Requirements
 Mutual Exclusion
 Progress
 Bounded Waiting

Swap(j):
 lock is initialized to false;
 do {
 key = True;
 while (key) {
 Swap(key, &lock);
 }
 // do critical section;
 lock = false;
 } while (TRUE)

均为全局变量
 key为局部变量
 lock为全局变量

Peterson's Solution:

Use flag[i] and turn variables.

△ It's used in two process synchronization.

do {
 turn = j; ⇒ 每turn让给对方, 确保了次序先后.
 flag[i] = true;
 while (turn == j) {
 if (flag[i] == true) break;
 // do critical section
 flag[i] = false; ⇒ First arrived process's while
 is broken by ②, later
 arrived process's while is
 broken by ②.
 } while (True)

i process: do {
 turn = j;
 flag[i] = true;
 while (turn == j && !flag[i]);
 // do critical section
 flag[i] = false;
 } while (true)

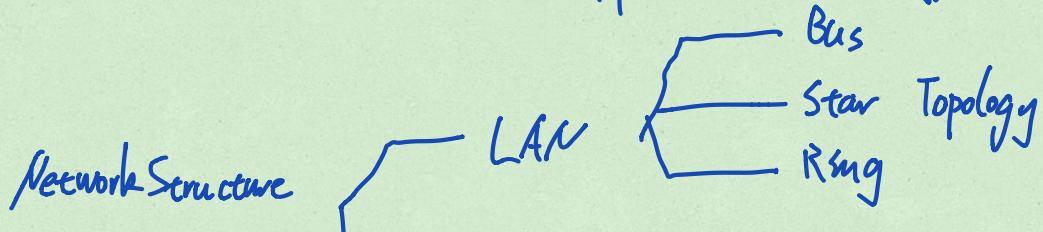
j process: do {
 turn = i;
 flag[j] = true;
 while (turn == i && !flag[j]);
 // do critical section
 flag[j] = false;
 } while (true)

△ Process Synchronization 核心是确定三点:
 ③ 先到的先执行 ①两个 process 不能同时执行

② 后到的在等待 e.g. Test and Set() 中, lock为了满足①,
 的执行完毕后也 key为了满足③, waiting为了满足②
 应该被执行.

Distributed System.

$$\text{Speed up} = S = \frac{T_1}{T_P} \quad \text{Efficiency} = \frac{S}{P}$$



— Network Topology. How to assess:

Installation Cost. Communication Cost. Reliability

fork - v fork.

Create: subprocess
特征=子进程完全copy父进程-1倍,开销很大。

shared space subprocess
应减少不必要的开销

特征:父进程被阻塞直到子进程 exit / exec.

clone. pthread_create.

create LWP.
与 pthread-create 类似,但能访问 kernel

LWP 与 HWP: 一个 process 独自占有资源 \Rightarrow HWP.
与其它 process 共享资源.

create user threads,
由线程库管理,无法访问内核.

创建、同步、销毁.

Lecture 12: Distributed System Basic Knowledge Introduction

Motivation
(分布式系统优点)

Resource Sharing \Rightarrow 可用不同设备特种功能
Computing Speedup \Rightarrow 提速
Communication \Rightarrow 交流信息
Reliability \Rightarrow 更可靠

Metrics

Speed : $\frac{W}{J_p} \Rightarrow$ workload
Speedup : $S = \frac{T_i}{T_p}$
Efficiency : $\frac{\text{Speedup}}{N} \Rightarrow N \text{ processors}$

Types

Network OS : User be aware of Multiplicity

Distributed OS

Not aware of Multiplicity

Data
Computation
Process } + migration

Network structure

LAN Topology

Bus
Ring
Star

Compare Criteria:
① Installation cost

WAN Topology

Tree
Partial
fully connected
Connected

② Communication Cost
③ Reliability

Communication Structure

Naming

Routing

Connection

Contention

Application: ~~Port~~ Port

Transportation \Rightarrow IP: Port

Network \Rightarrow IP

Data \Rightarrow Mac

Fixed
Virtual Circuit } in order

Dynamic \Rightarrow out of order

Circuit switching \Rightarrow fix

Message switching \Rightarrow temporary

Packet switching \Rightarrow divided,

fixed length packet,
out of order

CSMA/CD

Token Passing

Message slot

Communication Protocol

7. Application

6. Presentation

5. Session

4. Transportation

3. Network

2. Data Link

1. Physical

ISO Environment

Network

Environment

Robustness

failure detection

I-am-up

Are-you-up

Site inter

link fail

Message lose

Link fail

通知每个

site, link

不能用 / site

Reconfiguration

A-B (link fail)

B site failed

不能用

Lecture 13: Distributed file system.



Comparative;

① Caching && Remote Service

优势：Caching 在不频繁写入中更有优势，而 Remote Service 在频繁写入有优势

应用: Cache \Rightarrow Desk, large main memory Remote file Access \Rightarrow small memory Deskless

Lower internachine interface: Cache \Rightarrow RAM

KK upper user interface Remote file access: `fscanf`

② Stateful File Service & Stateless File Service.

优势: Stateful: ① Quick

- reduce disk access
- can read ahead the next blocks

② Shorter Request messages: 每个 request 信息只包含 file, #元 position.

Stateless: ① Robustness: 打开关闭时无需建立. 断开 connection identifier

When run into a crash, stateful server lose all states.
While in a stateless server, it doesn't be affected.

