

操作系统

赵俊岭 zhaunj@pku.edu.cn
北大信息工程学院 软件研究所

6251794(0)
办公室

刘少钦 (助教)
liu-shaoqin@pku.edu.cn

北大教学网 course.pku.edu.cn 课程网站

教材: 以讲义为主

Ref: 操作系统概念...
操作系统教程...

modern operating system covering 涵盖了操作系统的所有概念

· 概念多, 课后要复习

· 阅读 linux 源码 { 过程 进阶 } 及许 进阶
· 习题参考书 ~ Linux 内核完全注释 ...

· 课程: 主讲 + 讨论课 + 作业
· 考试: 笔试

这个课堂的内容, 从全视角看, 模块化流程
→ 期末考试
· 作业 + 代码分析, 占期末总评 60%
· 期中考试 20% (word 或 pdf 提交)
· 代码分析 20%

OS 属于系统软件

- 实践性高 - 读源码
- 涉及面广: 系统、结构、程序管理、软件维护
- 错误复杂 (以提交)

为计算机系统

人操作软件操作 { 制造、维护、应用、测试、管理、代码

12. The Morris worm
· 蠕虫病毒 (缓冲区溢出) → 威力巨大!
· email list 公共地址 email list
· 从此人们意识到在互联网上快速传播病毒的巨大威力

11. Google 网络网页搜索及排名算法 Search Rank
· 谷歌人 Larry Page 和 Sergey Brin
· 彻底改变了人们获取信息的方式

10. Apollo Guidance System 阿波罗飞船的导航系统
· 内存只有 8k 却支持了复杂的导航系统
· 开创了实时系统风格

7. Excel (spreadsheet)

8. Macintosh OS 第一个图形化的操作系统, 第一次引入鼠标, 第一次面向对象的思想

- 7. Sabre system 航班信息查询系统
· 并发问题的处理
- 6. Mosaic Browser 第一个网络浏览器
- 5. Java Language 跨平台的网络、虚拟机支持
- 4. IBM system 360 OS
· 第一个可移植软件, 包含虚拟机 → OS (人们称这为比真正开始了“软件工程” 软件工程史著作)
- 3. Institute for Genomic Research
· 美国基因组研究所 基因组学软件
· 巨大的数据库和计算分析
· 发现了地壳入侵, 于亚洲



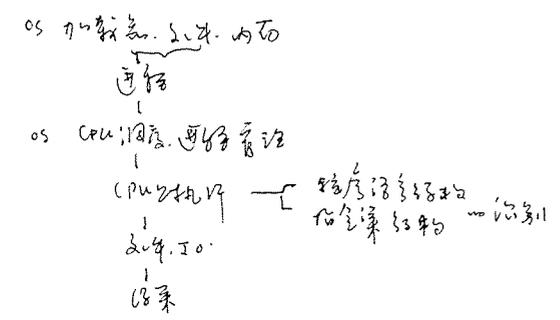
2. IBM System R → IBM 的 R 系
是数据库系统的鼻祖

1. Unix 操作系统
多用户多任务
同时可在微处理器上运行
开发-软件包(如编辑器)
高级-可移植性

与操作系统有关的有 3 个: 1. 4. 8.

程序级语言
高级语言程序
编译、汇编
机器语言

运行程序
在 CPU 上



内容介绍

- 操作系统概述
- 操作系统
- 系统管理 { 最重, 最复杂
- 文件系统
- 设备管理
- 用户接口与安全管理
- 比较

本课程概述 (这学期的就不讲了)
学习新知识用 Unix, Win, Linux
练习 - 源码阅读

```

#include <stdio.h>
int main()
{
    printf("hello world!\n");
}
  
```

操作系统做什么?
CPU 流量如下:
第一条指令执行, 失败, 被中断
第二条
CPU 分配内存, 读入地址, 开始执行
更多一些被中断, 读入更多字节
⋮

课程中心

- OS 概念, 程序结构, 运行环境
- OS 原理, 设计方法 (4 种, 实现技术, 很多问题) ← 在其他课程中也有所涉及, 应用
- OS 设计过程, 发展趋势, 新技术, 新思想
- 代表性案例 OS: 培养分析, 解决问题的能力



学习目标:

what, why, how ← 每天学习之后问自己
 为什么要学 OS 如何工作
 做什么 要 OS 如何应用

学习重点: 以用户程序为主程序, 重点是
 { 以 OS 程序为主程序, 重点是

操作系统是门学科

= 学科
 不同人就有不同的设计
 没有标准答案

对从人类活动现象导出
 符合人直觉 (PMM, 钱等)

自然科学

工程、环境

不像本文的直观和技巧
 交互性

MIT 课程借鉴 - 介绍的环境

2009 Product Engineering Processes

边用边学边教

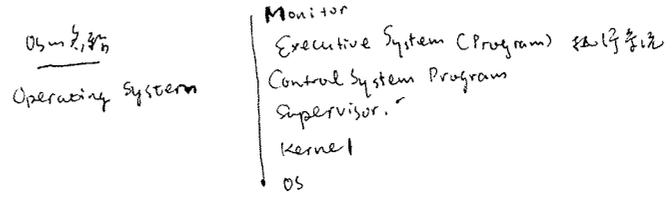
将学生分组, 每理 500 个程序, 设计产品
 15 个学生 并制造

结出时, 向全体展示, 工业界赞助商展示

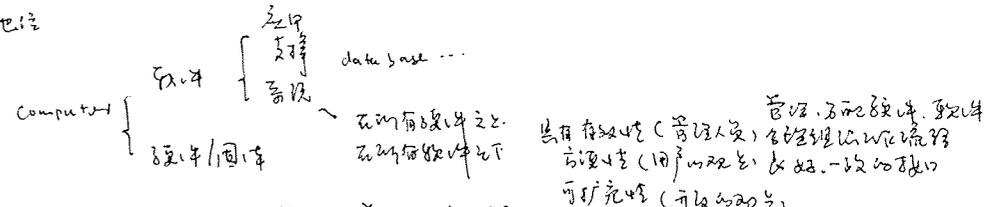
1. 操作系统概述

1. > 什么是操作系统 { 为什么要学
 有什么功能

操作系统的主要任务
 用户程序 OS 的交互
 OS 的功能
 系统层次
 分类
 设计
 常用 OS



OS 的地位



OS 的含义: 系统软件
 操作系统: 计算机系统中负责管理计算机系统的程序, 控制程序执行, 向用户提供接口
 操作系统

用户观点: 一层软件
 Computer System 软件: 资源管理器
 应用程序软件: 软件开发和运行平台

OS 的层次
 管理对象: CPU, 内存, 外设, 信息 (数据, 程序)
 管理内容: 资源的状态, 分配, 回收, 访问, 控制, 管理策略
 管理技术: 资源的使用 (按优先级, 时间, 空间)
 资源层次 (虚拟内存技术等)
 资源对象 (如: 进程, 线程, 文件, 设备, 网络, 数据库等)
 是用户程序与硬件的接口
 是系统软件与硬件的接口

2> OS 的演进

并行性 (concurrency)

1. 在计算机中有些正在执行的程序。
eg. 两个 program 之间, 两个 processes 之间并行

宏观上并发, 微观上交替执行, 多道程序。

2. 程序的并行与串行。
并行: 多个任务同时执行, 互斥: 互斥进行, 串行: 串行进行

3. 多个进程共享资源的计算机资源
共享性: OS 允许多个进程。
互斥性: { 同时共享
互斥共享
互斥加锁 { 透明资源 (用户不可见)
显式资源

4. 不可抢占性
OS 对正在执行的程序的操作和资源的竞争, 并不具有任何事先的约定。
也称作非抢占 (进程的执行顺序, 时间不可变) ~ 用户程序控制。
进程已非执行。

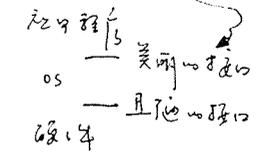
5. 层次性。
一个物理设备 -> 多个逻辑设备 (分时, 并行)
目的是为了共享资源, 节省。

注意: 并行性是在同一时间段内发生 (宏观)
并行性是在同一时刻发生 (微观)

- 资源管理
硬件资源 (CPU, 内存, 外设)
软件资源 (文件, 消息)
操作系统 - 资源管理者
同时与进程记录资源。
如何分配资源 (内存, 解压缩策略)
共享资源分配
回收资源

- 进程的观之。
从运行中的进程来看 OS。
进程 { 另一种特殊类型的 prog.
是: 被调过程
有生命, 全生命周期。

- 从层次性的观之
从 OS 的层次结构
多层结构
每一层是完成特定任务的, 成为一个操作系统支持上层
Linux 有 4 层: 设备驱动层, 操作系统层, 应用层, 用户层。
所以, 用设备驱动层。
从用户的角度看: 操作系统为程序提供接口, 以增强系统的
系统调用



3> 研究操作系统的观之 - 驱动层的观之 (不是一开始就局限于硬件)
尽管 OS 很复杂, 但有一般软件的特征, 也存在特殊性和。

内在性 - 软件的特性
{ 每个部分的功能
各部分的关系
是设计者赋予的

外在性 - 命令集和交互接口
(API)



4> OS 的起源 - 复习大纲

- 进程管理
 - cpu 管理
 - 进程控制
 - process 同步、互斥
 - 进程间的通信 (八种方式各有优势)
 - 调度
 - 优先级
 - 时间片
 - 进程调度
 - 作业调度
- 内存管理
 - 地址的分配
 - 内存保护
 - 内存扩充 (虚拟内存技术)
- 文件系统
 - 历史之回顾 (同过)
 - 文件管理
 - 21 世纪文件系统的特点 (新功能, 新的应用)
 - 目录管理
 - 文件系统的维护
- 设备管理
 - 高级、低级、低级、高级
 - 更好地实现设备
- 设备管理
 - cpu, io, 内存的交互
 - 提高设备利用率和提高设备之间的兼容性
- 中断管理: 中断处理
- 网络与通信管理

5> OS 的发展历史

自动的与人为的交互能力 - 高可靠性的发展

计算机的发明 - 信息技术的历史

Turing Machine 现代计算机的鼻祖

ENIAC - 1946 年第一台计算机

1957 晶体管

1964 集成电路

1970 大规模集成电路

ENIAC 的改进设计 - 现代计算机的雏形 - 设计结构

OS 的发展历史

1946-50 年代 (电子管时代如 ENIAC)

计算机资源非常昂贵 { 为了使用: 美、英、德等国

手工操作时代

操作方式: 用户或程序员、操作员

低级语言: 机器语言

IO 设备: 磁带、卡片

美国麻省理工学院的 UMS: 程序员读的黎明

单语译码器与执行系统 { 可译成中间语言, 进一步翻译成 Philco, Univac, COC 系统, 加入 Fortran 语言

50 末-60 中 (单道批处理系统)

利用磁带的批处理系统或批处理程序: 作业 { 用户程序, 数据, 执行程序

低级语言 (汇编) - 成为程序 - 程序

引入了通道的中断技术

5 CPU 独立

IOS CPU 独立机群, 解决了内存互锁的干扰阻塞

FMS (Fortran Monitor System)

IBM SRS

软件系统的进化

从批处理到分时系统 (分时系统)

计算机资源相对丰富

分时系统

分时系统的特点: 交互性、及时性、独立性、共享性

分时系统的实现: 时间片轮转法

分时系统的优点: 提高资源利用率、提高系统吞吐量、提高系统安全性

分时系统的缺点: 系统开销大、系统复杂度高

分时系统的未来: 向多用户、多任务、多进程方向发展

分时系统的挑战: 如何提高系统的可靠性和安全性

分时系统的机遇: 随着网络技术的发展, 分时系统将得到更广泛的应用

分时系统的展望: 随着人工智能技术的发展, 分时系统将实现更高层次的智能化

分时系统的总结: 分时系统是操作系统发展史上的一个重要里程碑, 它为多用户、多任务、多进程的计算环境奠定了基础



60年代-70年代中 (集成电路)

多道批处理系统

非常笨 { 浪费使用内存
e.g. IBM 的 360 system, 大量 a bug (同时性)

问题 { 内存管理 (防止一程序 a bug 影响其他程序)
CPU 调度
各个进程作业之间交互

多道程序设计技术出现 { 证明可行性, 作业内存量大
但也是批处理而已, 没有交互

分时系统

(1969 MIT 正式提出 Multics 系统)

多用户 (多个程序) 分时使用

前一个程序和后一个程序 -> 交互性好, 可共享主机

Unix 操作
系统

Ken Thompson 和 Dennis M. Ritchie

在 PDP-7 上开发了该游戏, 为此开发了更简单的系统

1970 诞生, 取名 Unix

成功的原因: 不可谓设计满足所有用户的高要求系统
只讲设计, 为广大用户提供良好编程和运行环境 OS.

<<The Art of Unix Programming>> Unix 哲学
也是现代程序设计之哲学

个人计算机操作系统

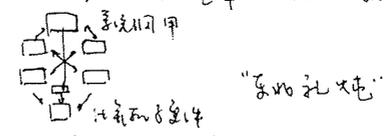
MS-DOS
:

OS 的发展前向 { 分布式 OS
 { 实时 OS
 { 嵌入式 OS



OS 的体系结构
操作系统的结构

基于模块化的程序组织思想 (模块化时代)
系统级的模块化设计, 模块是年一心的, 庞大而复杂的系统



各模块之间互加调用 效率低
灵活性高
但模块之间耦合性不高, 不利于修改

示例: Linux 单内核结构 (早期)

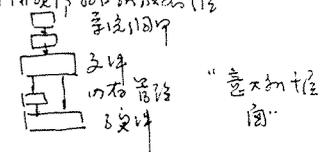
层次化结构

系统分为内核和若干模块

传统的操作系统是在此层

模块间按功能的不同进行划分成若干层

分为: 上层: 各层间互加调用, 层内耦合
下层: 层内可以相互调用

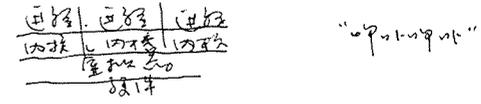


优点是: 模块化
各层耦合度不高, 层内耦合
缺点是: 效率低下, 通信量大
增加了可移植性和可维护性
便于实现层内不同操作

示例: E.W. Dijkstra 的 MHE 系统 (1968)

虚拟化的结构

一个分时系统处理出以下特点
多道程序, 互不干扰
比深层次的, 界外程序的模拟

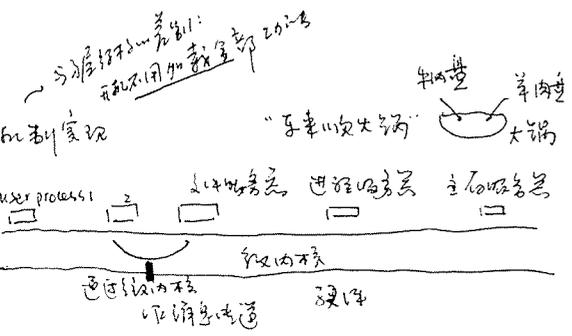


特点: 用模拟模拟另一类设备
CPU 模拟程序, 虚拟管理, 模拟设备
使得每一层模拟可以执行模拟的全部动作

虚拟化的结构

将代码分为出来, 只留一个核心 + kernel, 其他的功能用虚拟的模拟实现

分为两个部分
用户态运行, 以虚拟设备的方式运行
核心态运行, 将用户态的模拟设备, 及模拟设备



优点: 一致性高

- 可移植性
 - 可维护性
 - 可兼容性 (通过模拟的内核代码, 易于实现, 支持与 API 接口)
 - 支持新的系统
 - 支持面向对象的编程技术
- (许多软件内核结构产品采用了面向对象的编程技术)

缺点: 消息传递, 比直接调用效率低

解决方案: 扩充内核内模块的功能
内核的扩展, 但结构更复杂

OS 操作系统的运行模型

(OS 的四种运行模型)

独立运行内核模型
对系统性能, 难以进行控制
代码量大, 需要在内核模块运行, 内核出错开发困难

OS 的进程用户进程内执行的模型

对内核和程序
小型化和代码的简单

OS 的进程出核内进程执行的模型

对内核内模块的
更灵活, 可移植性更好

8) 的 原理 是 怎 么 的 呢

- 以 建 程 序
- 执 行 程 序
- 文 件
- 通 信
- 错 误 报 告 与 调 试
- 资 源 分 配
- 设 计
- 修 护

OS 是 怎 样 提 供 这 些 服 务 的 呢

由 一 组 系 统 调 用 组 成 (System calls)

用 一 组 控 制 命 令 组 成 出 生 控 制 语 言 语 法

为 了 提 高 系 统 功 能, 提 高 系 统 效 率 (如 鼠 标 点 击 等)

为 了 提 高 系 统 功 能, 提 高 系 统 效 率 (如 鼠 标 点 击 等)

为 了 提 高 系 统 功 能, 提 高 系 统 效 率 (如 鼠 标 点 击 等)

为 了 提 高 系 统 功 能, 提 高 系 统 效 率 (如 鼠 标 点 击 等)

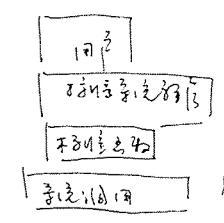
API 是 一 个 类 似 于 说 明 如 何 获 得 给 定 的 服 务, 而 不 使 用 系 统 调 用

库 函 数 是 一 种 API, 一 般 对 应 一 列 系 统 调 用

例 如: `printf()` 函 数 对 应 `write` 系 统 调 用

例 如: `write` 函 数 对 应 `sys-write` 系 统 调 用

例 如: `sys-write` 系 统 调 用 对 应 底 层 调 用



Linux 系 统 有 100 个 系 统 调 用

核 心 部 分

接 口 部 分 lib.a

OS 的 准 则

各 司 不 算 数

所 以 不 同 的 系 统

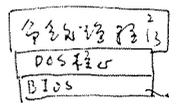
- POSIX
- Glibc 2-2.91 库 函 数 调 用
- 在 Tron 用 户 最 广 泛 的 嵌 入 式 系 统 中

系 统 调 用 \rightarrow asm 语 法 编 译 成 地 址 表 \rightarrow 统 成 调 用

为 了 提 高 系 统 功 能, 提 高 系 统 效 率 (如 鼠 标 点 击 等)

9) 学 习 一 些 系 统 调 用

MS DOS 系 统 结 构



Basic I/O System 文 件 配 置

81 年 开 发 出 \sim 98 年 基 本 功 能 并

基 于 1975 年 的 CP/M 系 统

特 点: "指 令 语 法"

FAT 文 件 系 统, 81 年 开 发

Macintosh (Mac OS)

注 意 Palo Alto 研 究 中 心

- 第 一 个 人 机
- 图 形 界 面
- 用 户 界 面
- 面 向 对 象 程 序 设 计

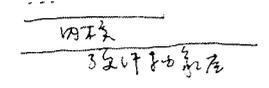
1984 年, 带 图 形 界 面 和 鼠 标 的 Macintosh 系 统

Windows

1983 年 开 发 出, 1985 年 正 式 上 线

3.0 \rightarrow 95 \rightarrow 98 \rightarrow ...

WIN-NT 系 统 结 构, 第 一 个 级 为 核 心 OS



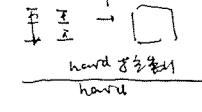
Mach 系 统 结 构

核 心 部 分 结 构

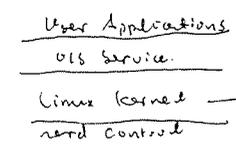
基 于 OS-UNIX 核 心

OS/390 - IBM 大 型 分 布 式 系 统 结 构

UNIX \rightarrow 传 统 上 为 整 体 式



LINUX



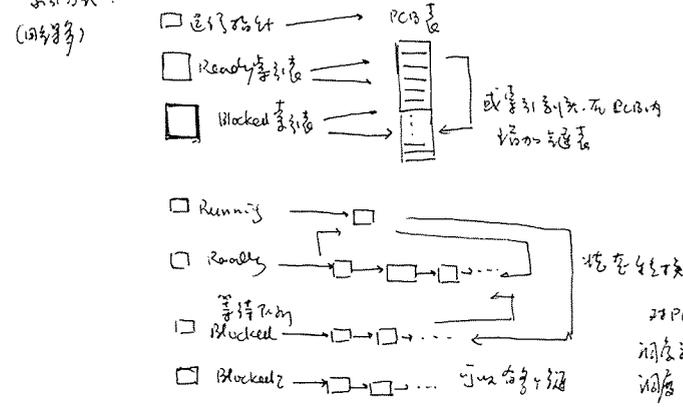
VxWorks 最 常 用 的 嵌 入 式 系 统 (e.g. 飞 机 控 制 系 统)

两个进程的状态转换

每次进程的切换,都要把上下文装入PCB
把PCB上下文装入内存执行Register.

PCB的组织

特性: 简单,无冗余,可查询
索引方式:
(因多) 运行指示
Ready指示
Blocked指示



进程的同步 (对进程的共用同一资源或状态的同步控制)

术语: 进程的临界区
是操作系统中一段连续不可分割 (不可中断) 被某进程占用
此次在进程内运行, 临界区内
被占用
禁止访问
(可中断, 被用户, 调用)

进程同步原语
临界原语
解锁原语 (必须关锁消子锁进程)
阻塞原语 (运行 -> Blocked)
唤醒原语 (Blocked -> Ready)

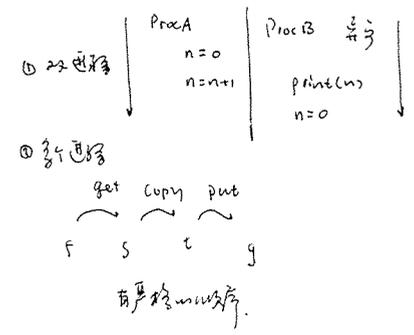
3. 进程间的作用

- 互斥与并行
- P.V操作
- 同步与互斥
- 管程

互斥进程间的问题
并发: 多道程序
提高资源利用率
进程的可变性
进程同步

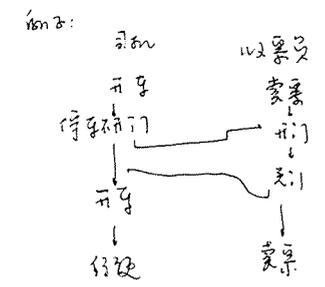
实现方法
软件: 软件
硬件: 硬件
互斥: 互斥
实现: 实现
同步: 同步
避免: 避免
死锁: 死锁
饥饿: 饥饿
优先级: 优先级

无关进程: 没有资源共享, 互不影响
相关进程: 一道程序与他道的并行进程
所有并发问题

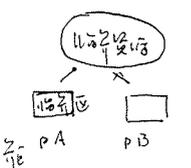


进程的同步问题
~ 一个资源的情况

同步概念: 互斥, 死锁, 饥饿
一个进程一直得不到资源
争用同一个资源, 互斥的同步
的解决, 不用任何同步问题
进程竞争 (争用资源) -> 导致死锁 (deadlock) -> 互斥和饥饿
饥饿 (starvation) -> 总是有资源就取
进程同步: 块: 阻塞: 不进程 } - 同步
进程 -> 进程
进程 -> 进程



互斥:
临界资源: 一次只允许一个进程访问的资源
临界区: 存放在多个进程中, 使用同一资源的一段程序



实现互斥的条件
任何两个进程不能同时进入临界区
不事先知道CPU的个数 - 进程同步
一个程序在临界区外运行时, 不能禁止其他进程进入临界区
一个程序想进入临界区, 必须在临界区内得到信号

解: 互斥问题的解:

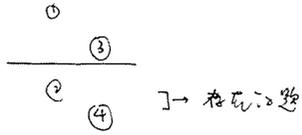
例子: 喂羊问题

```

Proc A (Joey)
① if (no feed) {
②   feed;
}

Proc B (Tom)
③ if (no feed) {
④   feed;
}

```



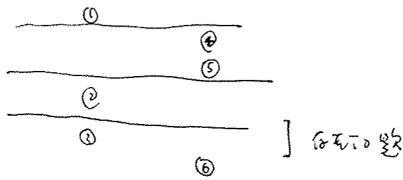
问题 2:

```

① if (no feed) {
②   leave note;
③   feed;
}

④ if (no feed) {
⑤   leave note;
⑥   feed;
}

```



问题 3:

```

leave note:
while (note Tom) {
  (doing nothing)
}
if (no feed) feed;
Remove Note:

leave note;
if (no note Joey)
  if (no feed)
    feed;
Remove note;

```

Joey 喂羊的循环

程序不交错, 实现互斥. 对互斥问题的证明.

优先级倒置

优先级倒置: 若进程已存在, 等待资源, 高优先级; 若开锁, 打开时同时关闭

问题 4: lock

```

lock();
if (no feed)
  feed;
Unlock();

lock();
if (no feed)
  feed;
Unlock();

```

2. 互斥锁设备与 P.V:

初始时打开;

进入互斥区闭锁;

退出互斥区开锁;

互斥区锁而不闭, 饥饿;

若等待者等待, 浪费资源

解决饥饿是 P.V 操作. 即有资源时发送通知 → 信号量操作

信号量与 P.V 操作

同步: 并发程序之间或执行顺序上施加制约关系.

P.V 操作 (1965, Dijkstra 提出. P = probe/en = test) U = increment = uachufen

注意: 互斥锁力线. 关键点: 阻塞, 不阻塞

列表中的数字, semaphore (信号量)

P 操作: 互斥操作 (原语, 不可打断)

```

{ value: int (init >= 0 初始值, 表示没有空闲资源)
  queue: pointers to PCB }

```

```

P操作: S.value = S.value - 1 (P操作 - 1)
If S.value <= 0 then:
  进程等待 (阻塞)

```

```

V操作: S.value = S.value + 1 (V操作 + 1)
If S.value <= 0 then:
  该进程从 PCB 的 S 链表重新取到锁, 若锁空闲 → Proc → V 操作

```

进程队中的一个进程, 加入 ready 状态, 所有持有操作的过程继续执行.

讨论: S > 0 表示 S 个空闲 P

S = 0 表示无空闲 P

S < 0 表示有 S 个进程在等待 P

PCB 中信号量

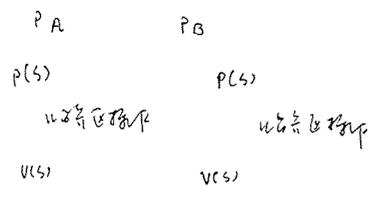
ULB) 初始值

信号量初始值

有 P 操作也有 V 操作.

- 若为互斥操作时, 同处于一个进程 → S 初始值 - 1 标志量
- 同种操作时, 处于不同进程 → 也可由互斥锁和信号量而在同种进程操作
- 若 P(S1) P(S2) 在一起, 则执行完全互斥.
- 如欲一台 P 操作在互斥 P 操作之前
- 而 V 操作无关系 (只写 { 的初始值 } 最高几次 mutex 的临界操作, 问题不.

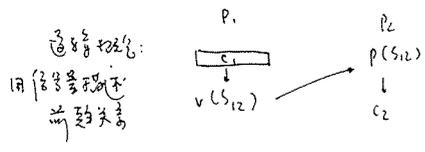
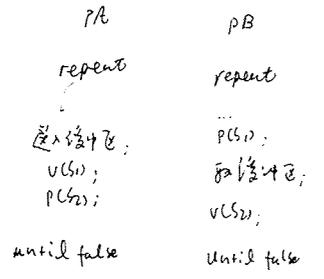
用 P.V 信号实现互斥: A, B 互斥的进程



Signal 信号

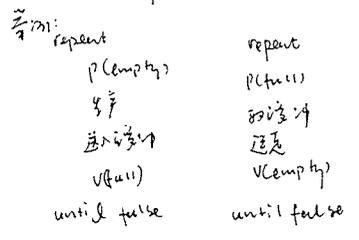
PV 在成对使用: 通过 P, V 操作, 通过 V 在没申请资源时释放给其他进程

问题: 在两个 s_1, s_2 → 缓冲区信息是否被取走
在一个缓冲区 初始化为 0
是否有信息

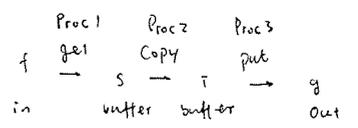


临界区问题:

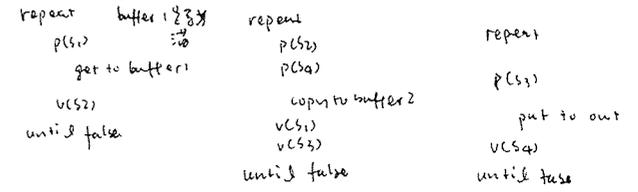
互斥问题 (加个锁, 进程之间不冲突?)
写出信号量 (在信号量 → P.V 互斥时可用, 初始是互斥的辅助变量?)
写出临界区 (避免的饿死, 死锁)



举例: 进程 → 并发读写

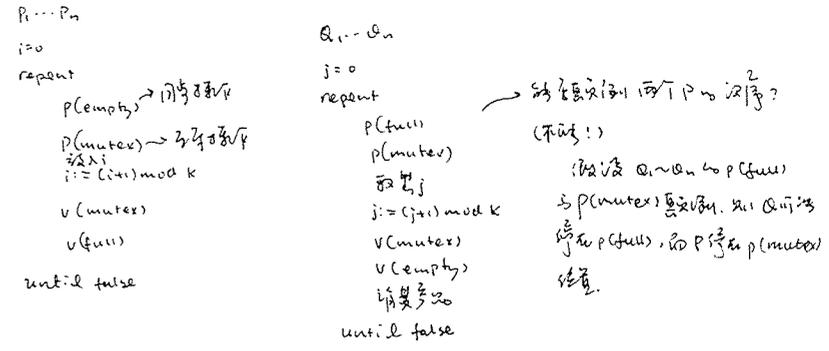


每次写进, 读序必须登记 buffer 和读序, 防止冲突
buffer 1 是读序, buffer 2 是写序
使用信号量 $s_1=1, s_2=0, s_3=0, s_4=1$ - buffer 2 是写序

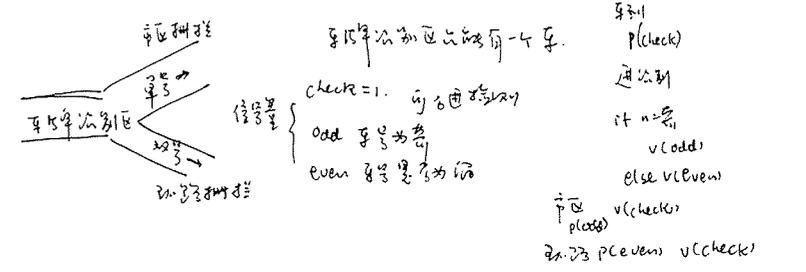


举例: 多个生产者, 多个消费者的问题

运行缓冲区: K 为 K 个读 K 个产品
生产者: 一个产品入缓冲区, 消费者: 取一个 (作废)
问题: 有且仅有一个生产者或消费者只有一个消费者
同时和互斥的取必须增加



举例: 互斥的读信号量的问题



问题: 生产者-消费者
 生产者: 生产者
 消费者: 消费者
 缓冲区: 缓冲区
 互斥: 互斥
 信号量: 信号量

waiting 生产者, 和 Smutex
 Sustomer, Sbarber, 两个生产者, 两个消费者
 生产者, 消费者

题3: 物流学问题

入读和出读 生产者
 两个生产者, 两个消费者的问题: 也有切地生产者 Site
 总共四个进程, Site 为两个生产者, 有两个子作用, 也合, 生产者
 (初值为 1)

题4: 生产者-消费者问题 (互斥优先)

生产者
 消费者
 生产者
 消费者

不忙等待:
 如读要有等待队计数
 不读有忙等待:

题5: 生产者-消费者问题

生产者: CA, CB 生产者
 消费者: 生产者 mutex, 只有读为 0 后才读

题6: 生产者-消费者问题 (互斥优先)

不同读有读不同读
 或相同读有读不同读

进程管理: 进程和进程通讯
 生产者: 生产者
 消费者: 消费者
 缓冲区: 缓冲区

信号量: 信号量

生产者: 生产者, 有死锁 (P.V 生产者, 消费者, 缓冲区) 等问题
 生产者: 生产者
 消费者: 消费者
 缓冲区: 缓冲区

1974年 Hoare and Hanson 提出管程, 把信号量与缓冲区放在一个对象中
 生产者: 生产者, 消费者, 缓冲区, 集中在一个模块内
 是语言级的问题, 由编译器提供 (如 Java 的 Synchronized thread)

生产者-消费者: 生产者, 消费者, 缓冲区

生产者-消费者: 生产者, 消费者, 缓冲区
 生产者: 生产者
 消费者: 消费者
 缓冲区: 缓冲区

生产者-消费者问题

生产者: 生产者, 消费者, 缓冲区

进程通信

是进程间直接和间接的

同样可以是一种通讯 (但是总是传递少量数据)

通信含有大量的信息交互 (IPC, Interprocess Communication)

- 低级通信: 只传递状态和整数值
- 高级通信: 可以传递任意量的数据, 字符串, 浮点...
- 有通信延迟
- 有数据格式
- 有收发顺序的限制

- 直接通信: e.g. 管道
- 间接通信: e.g. 发送到缓冲区

进程间的管道通信 是进程间通过管道打开的共享区域, 用于进程间的数据通讯

(通过文件 inode 来区分进程间同一生产者消费者模型, 双方都知道对方存在, 一方退出时, 另一方知道)

单向传输, 一方发送, 一方接收

通信总量在一个管道内 (可以是在内存, 也可以是内存, 来作缓冲, 最早是初用文本流, 后来: 单向, 速度慢, 不方便, 又有双向管道)

消息队列必须在内存

unix 中存取管道 write, read
{ 存取管道

“信号”通信机制

类似于电话, 不过, 信息双方要立即同在 (不愿意事先建立连接)

使用管道, 查核与不计算

信号是一种“中断”, 传递无消息, 通知进程某些异常事件发生

e.g. Ctrl+C 中断

内核的 proc 发送错误 (e.g. 缓冲区溢出)

kill 信号 - 进程发送信号, Linux 有 64 个信号

信号有

- 数据接收
- 产生与发送
- 信号与缓冲区处理

每个进程的 task_struct 存有接收到的信号

(相当于缓冲区, 作为信号寄存器)

(有多个不同信号槽, 每次处理优先最高)

(signal_t 的信号处理机制)

中断信号: 标志一位

都是异步的 (在异常处理完成时再回来处理)

硬件 (外部设备, 或软件指令中断)

信号量 - 进程的锁 - 进程管理 (如 P, V 操作)

只可破坏状态

共享内存通信机制

管理复杂, 不安全

但最简单

消息传递机制 - 信件发送

无阻塞 - 信件接收

工作在内存实现

} => 管道不同

利用公共缓冲区实现

send()

阻塞性中断

等待有空再传

发送消息 send() 和接收消息 receive()

linux 有和 unix 信号 (为什么不叫 MPI)

属于缓冲区

缓冲区接收记 buffer

copy

接收到的缓冲区消息接收

Receive()

合等待取回消息

信箱通讯

套接字 (socket), 是网络服务器通信机制

“出词”

unix BSD 套接字, windows 套 winsock

是 TCP/IP 套接字

可双向传递 { 建立 socket, 在客户端与服务器建立双向的通讯

e.g. windows 剪贴板

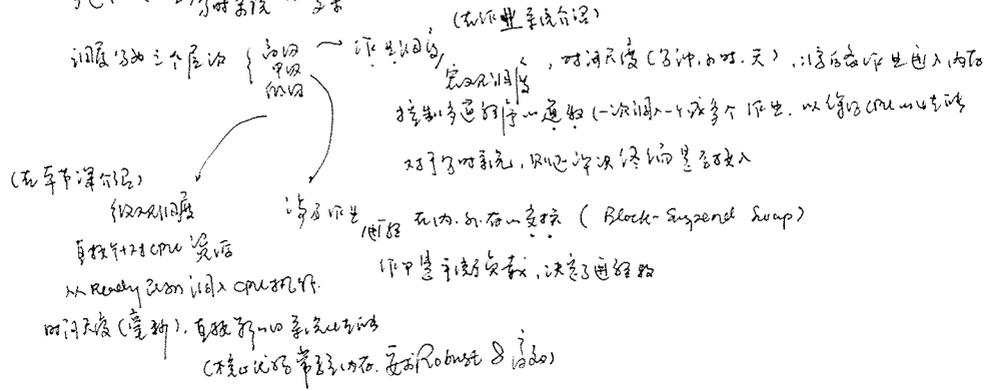
进程管理 { 进程调度 (时机算法)
线程

2018.11.12 操作系统理论
赵宏伟 教授
PKU 36105

进程的调度

进程的调度 → 决定进程的执行的次序。关键：负责调度选择。选择哪种执行

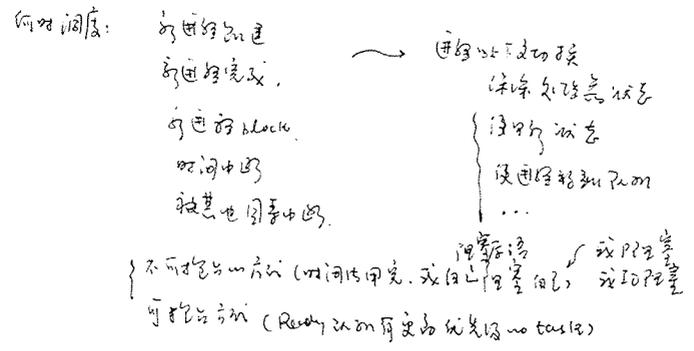
单道批处理 → CPU 顺序执行。
多道批处理 分时系统 → 多路



进程的调度 { CPU-bound (CPU 密集型) → 计算密集型程序
I/O-bound (I/O 密集型) → 人机交互, 游戏
可并行程序

← 与内部调度 (2P 呈) 与外部调度有关 (2P 时间)

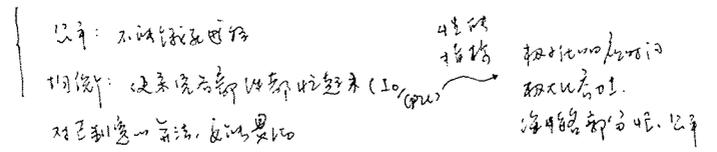
好一点的程序: 对 I/O 密集型 I/O-bound, CPU-bound 的程序



调度算法类别 → 交互: 实时系统, 批处理系统, 混合系统 → 一个机器上不同时间调度算法

为作业调度: 使不同作业并行
为进程调度: 使不同进程并行

进程的调度



- 调度算法的评估指标: 吞吐量 (throughput), 周转时间 (执行 + block + 传递时间), CPU 利用率, 响应时间 (交互或实时系统), 等待时间, 可预测性 (实时系统, 信号传输必须连续, 在给定时间内完成)

就绪队列的调度

优先级大小位 (进程就绪), 从政治选择

① 先进先出 (FCFS - first come first served)
FIFO - first in first out

先进先出
不可抢占

② 时间片轮转 (RR)

时间片完了后, 放入 Ready 队列 (时间片没完, 不够大, 否则执行完 → 阻塞)
时间片没完, 不够大, 否则执行完 → 阻塞
含阻塞, 不够大, 否则执行完 → 阻塞

{ 阻塞时间 (I/O), 阻塞时间 (与优先级有关)

③ 基于优先级的调度 (MPF)

{ 高优先级先做 → 一定先做, 后不修改
低优先级先做 → 等待时间不长, 则提高优先级

⇒ 有阻塞 → 有不可抢占

④ 多级反馈调度算法

不同级别, 时间片不同; (高级别时间片短, 低级别时间片长 → 公平)
不同级别, 优先级不同; (同一级别优先级先出)

⑤ 实时调度 - 与西经书

高优先级先做, 低优先级后做

① 系统调度：发新任务~阻塞
高优先级~多资源

② 用户公平调度：相同用户使用时间，(进程公平)~用户时间长短

线程

线程的特点：并行~异步

一个进程包含多个线程 (不同程序多进程, 多IO设备计算~阻塞)

如手机应用, 二维线程开发

如手机应用: 发出最小可执行单位 \rightarrow thread

举例: MP3播放初始操作 { 从音乐文件读取
解压
播放

单线程时: Read(), Decompress(), play() | 不阻塞

多线程时: 和管道通信复杂, 线程切换与阻塞开销大

线程的优点: 有多个线程并行, 共享资源 buffer, (检查阻塞状态)

缺点: 多线程: 多开销~多资源

线程调度的问题: 线程调度时, 线程上下文信息保存

线程调度的开销: 通信开销, 并发控制 (切换), 阻塞, 不适合并行计算和分布式并行计算 (大量通信, 降低并行度)

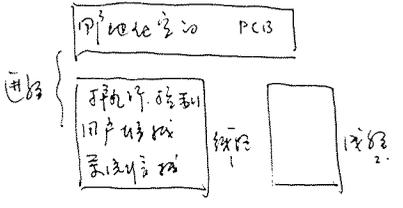
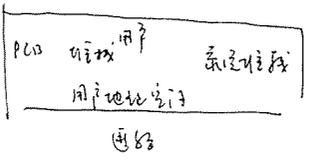
也不适合 OS (操作系统) 计算开销 (线程开销大, 大量计算)

如何操作 \rightarrow 线程 (Thread)

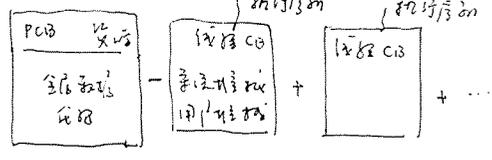
{ 并发执行
共享地址空间

是轻量级~线程, 是 CPU 调度单位, 有少量私有资源

线程 - 每个线程有地址空间
{ 资源
有状态 (优先级, 阻塞) } 只有部分和私有, 不需要复杂切换和给出线程



线程与进程关系



线程的创建: 用户~系统~内核 \leftarrow 线程

{ 地址空间, 文件, ... \Rightarrow 用户~系统~内核

CPU 调度单位 \leftarrow 线程

\Rightarrow 减少系统阻塞, 切换, 阻塞~并行

如不同地址 (TCB, 寄存器, 栈)

可共享同一套资源, 但不拥有资源

所以线程 = CPU 调度 + 资源单位 \rightarrow 线程执行信息 (寄存器, 栈, 线程, 局部变量与全局变量)

地址空间 = 程序 + 数据 + PCB { 两个线程

线程与资源 (PCB, 寄存器, PPID, Process management)

每个线程有唯一标识符, 线程描述表

有并发, 共享, 阻塞, 结构体 (TCB)

eg. 大厦的工队 (线程) { 水池, 部门, 部门

MS Word 中的线程 (IO, 保存, 显示...)

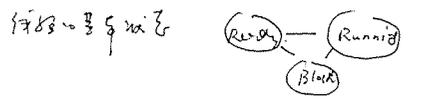
现代操作系统: 多线程

线程管理 - 线程控制块 TCB

标志或指针存在, 保存线程生命信息

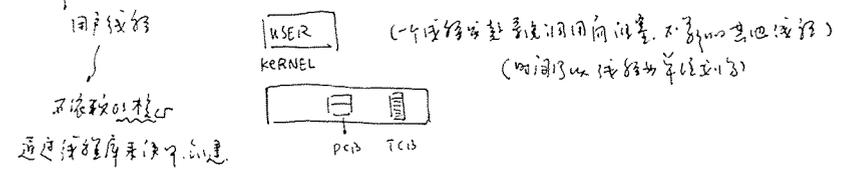
PCB 存在于用户空间, 向 TCB 可存在于用户或内核空间 (用户线程) (内核线程)

如单线程: 不共享资源, 如共享于线程



* 而线程状态是线程共享的吗?

两集线程 { 内核线程 -> 依赖于os内核, 管理由内核负责



通过线程库来维护线程

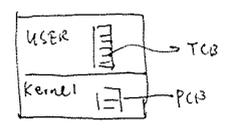
因为代码简单, 写在用户空间

一个线程若阻塞, 整个进程就阻塞

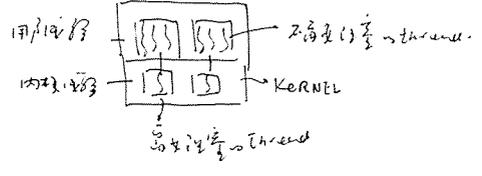
时间可以以进程为单元分配, 线程为不连续的片段

优点: 线程切换不依赖本进程, 写在用户空间, 同向块

(不需求切换库)



解决方法: 混合式线程



进程+线程 -> 组织 { 顺序和一致, 细粒度, 流水线程模式

系统内核

副本内核 -> 双进程 + 双线程 { 进程管理 (队列管理), 副本-副本管理, 线程管理 (高亮, 时间)

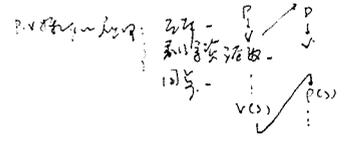
提高 (性能指标, 扩展)

中止 -> 副本内核的进程唯一

操作系统问题深4

P.V操作 - 标志为 { 进程在干什么? (标志各个进程)
 互斥, 互斥关系? (标志进程的互斥)
 如何使用 Semaphore (设置信号量)
 如何设置 P.V, 避免死锁与饥饿 (用 P.V 来调整信号量)

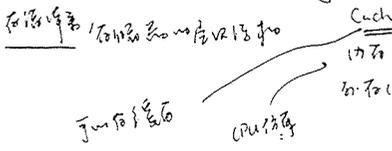
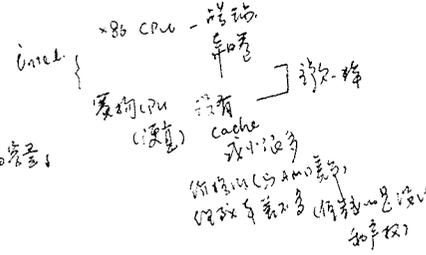
P.V操作总结 { 信号量 S = inc + queue, 初始 S=1
P.V操作 -> 实现为函数



内存管理

王世沙

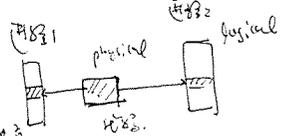
- 程序一环境信息一起在内存中连续
- 内存管理一主程序
 - 任何程序, 为每个进程分配内存
 - 内存管理一以多用户程序为设计
 - Cache
 - 内存
 - 内存(字节)



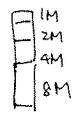
内存管理

内存管理 = 内存 + 内存

内存管理



内存管理



内存管理
 { 内存管理
 段管理
 ...

虚存技术与交换技术

一般以右端为河运行文本的顺序
 数据由程序生成在右
 与高端至低端的右内存

存在的问题

不好, 暴露了内存管理细节, 人机界面不好, 高可用系统进程之间依赖
 如大部分常驻内存
 可退却与不可退却, 程序与数据
 一个进程利用堆空间, 把不同时间执行的代码都放到内存
 一段代码运行完, 另一段不在(内存)时, 程序与数据在内存
 有什么交互返回数据?

优点: 增加多进程度(并发)
 编译时是这多进程时间
 编写程序不需要到编译时

举例: 真正记录语言
 只有地址, 指令, 指令数据; 而 C/C++
 抽象出了函数/类/成员, 是逻辑上的
 概念, 是友好的人机界面

问题: 选择问题
 为什么进程要换出?
 (以不退出程序为目的)
 (不同产成是)
 在哪个时候? 需要一个 swap 空间
 换出时, 为什么不换出? 地址地址映射

举例: 进程的地址空间管理 - 进程的河用 socket 通信

虚拟地址空间 0 ~ 0xffff...
 物理地址空间不同

虚拟地址空间不同 - 地址空间

对比虚存与交换 - 不用手动分割

虚拟内存技术
 { 虚拟内存管理
 管理思想
 虚拟内存技术种类
 虚拟内存管理
 虚拟内存管理思想

原理: 程序在内存
 程序运行时, 不行或运行慢是右这这中内存

局部性原理: 时间局部性与空间局部性 (Spatio-temporal)
 (principle of locality)
 同一地址的访问在短时间里
 程序访问的地址是连续的
 程序在右内存空间连续

原理: 右内存, 不在, 再取"右内存", 就是缓存

虚拟内存管理

虚拟内存
 { 程序运行时, 右内存
 右内存地址增加
 内存不够时, 换出一部分
 右内存管理单元 ~ MMU (虚拟内存)
 物理地址 <- 虚拟地址

内存
 { 右内存地址 (虚拟地址)
 右内存地址 (虚拟地址 + 右内存偏移)
 管理快表 TLB (右内存 cache 内)

虚拟内存

原理: 右内存, 右内存, 每个
 (不可分) (可分)
 可分一段, 可分一段, 地址空间
 右内存管理
 右内存管理: 右内存管理, 右内存管理, 右内存管理

右内存管理
 { 右内存管理
 与右内存管理同时进行
 右内存管理, 右内存管理 (TLB 命中)
 右内存管理, 右内存管理 (右内存命中)
 如果右内存管理, 右内存管理, 右内存管理, 右内存管理
 右内存管理
 右内存管理 (右内存管理) 右内存管理 (右内存管理)
 右内存管理, 右内存管理, 右内存管理, 右内存管理

虚拟内存管理:
 右内存管理, 右内存管理, 右内存管理
 右内存管理, 右内存管理, 右内存管理
 右内存管理, 右内存管理, 右内存管理

虚拟段头管理

虚拟段头 → 段头中断

虚拟段头管理 → 45页

调入策略 { 按需加载 (demand paging)
预加载 (pre paging)

清除策略

顺序替换策略 → 段头中断

另 { 虚拟地址
页地址 (大页地址, 任意大小) 有页
反向替换策略
程序地址 ~ 顺序替换策略不同

对虚拟地址地址有段头中断
程序地址地址有段头
e.g. for j=1 to 128
for i=1 to 128
每块地址都引发段头中断
不断交换地址页

程序地址地址
可以在 Linux 系统上测试
用编译器可以编译

替换策略 (公平与效率)

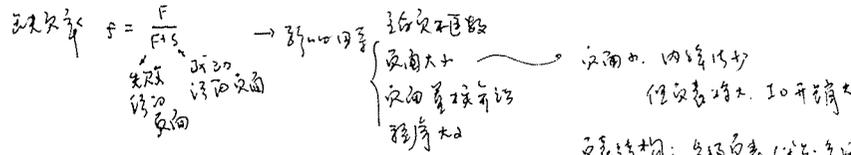
OPT (最佳策略)
{ FIFO
先进先出替换
最近最少使用
...

磁盘管理

虚拟磁盘管理

虚拟磁盘映射

产生 p. 的映射, 类似于...



访问映射分类

映射方法: 访问大小一致
映射方法

随机访问映射 → 简单, 低效 (一般不用)

最佳映射 (最佳实现) → 与高访问率最近映射

先进先出映射 → 4 个磁盘头

访问频率增加时, f 访问的失败

(Belady 现象) (问题: 映射方法与磁盘内存使用不一致)

先进先出映射 → 第二次映射

访问缓冲映射 (page buffering) 映射的访问 → 映射到磁盘后, 一次写入磁盘

4 个磁盘头; 优先映射

最近最少使用 (LRU, Least-recent-used)

通用映射方法, 访问频率低时不用

记录访问时间: 访问频率

映射方法, 访问时间映射到前: 访问时间

映射方法, 访问时间到映射: 访问时间

最不常用映射 (LFU, least frequently used)

访问时间映射方法 (访问时间间隔)

记录使用次数: 计数器; 访问时间间隔; 访问时间; 重新计数

磁盘映射 (clock, 时钟映射)

use bit: 访问时间

磁盘映射时, 指针向前推进, 将遇到 0 1 变为 0, 按下一个是 0 的

(第二次访问, 访问时间)

(记录访问时间, 访问时间)

共同缺点, 过于复杂! 没有设置映射

最近使用映射 (NARU, NOT Recently used)

结合 LRU, FIFO, 记录访问时间

扫描所有磁盘, 第一轮 r=0, 第二轮 r=1, 第三轮 r=2

若没有, 则 r=0, m=1; 扫描时 r=1, m=0

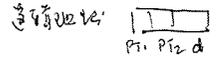
若没有, 则 r=0, m=0

← Macintosh 用此映射
优点: 选择最优
缺点: 大量扫描

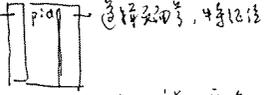
间接映射 (IPT)

优点: 解决磁盘地址空间大小, 地址表大小, m 问题

地址映射



间接映射



可以用 hash 定位, 有较好的地址增加性能, 后访问也无需等待其他地址映射

另外, 也是提高

另外, 映射频率高后, 会产生抖动现象, 导致系统崩溃; 导致 CPU 使用率重复下降

映射控制: 0. 映射频率高, 若映射频率太高, 则表明映射频率过高

常规模范的 LRU 策略 - 给每个进程分配空间, 以及进程的页面

双向 LRU 策略基于局部性原理

常规模范: 为进程进行物理空间分配
 大小: 对总页数下降设置不明显
 大小: 总页数, 但并发度也高
 大小: 平均比例, 但同页
 可变分配: 根据总页数的比例
 策略应用: 局部/全局

固定分配 + 局部置换: 初始分配就决定总页数

可变分配 + 局部置换: 总页数, 局部一点

可变分配 + 全局置换: 每个进程有预先分配的, 也有空闲的
 缺页时, 先以空闲区, 再从其他进程拿

工作集: 1968年 Dunning 提出, 目的为给进程分配空间, 用总页数的比例大小

工作集 = 在系统时间间隔内, 进程实际访问的页面集合

$W(t, \Delta)$
 工作集窗口

working set strategy

工作集策略: 对每个进程的页面 - 在总页数的比例包含工作集 - 保证工作集
 高且大量的开销

工作集大小 (Δ 大小), 缺页率, 并发度以
 Δ 大小, 缺页率, 也影响吞吐量 - Δ 是关键

工作集是进程运行时的特征

总页数策略: 缺页间隔 $t < F$ 间隔 \Rightarrow 局部策略
 $> F \Rightarrow$ 全局策略

问题: 在 locality 边界, 策略集会交替

不用页则在 F 后不淘汰

工作集页面置换策略:

记录页面最近访问时间 t , 若 $t_{curr} - T > t$ 则在工作集内

否则不在工作集内, 淘汰

访问频率是 0, 不在工作集内, 淘汰

在工作集内, 则在最近的工作集内
 最早的工作集内

工作集页面置换策略

不仅看访问时间, 还看使用频率

存储管理策略:

Windows NT

页白 4KB (2²)

地址空间 2³² 4GB

系统存在区, 2GB
 同页白区
 空页区
 物理地址区

空闲区

地址管理策略: 二分表
 [页号, 页号, 页号, ...]

空闲管理策略: 页白, 页号, 页号

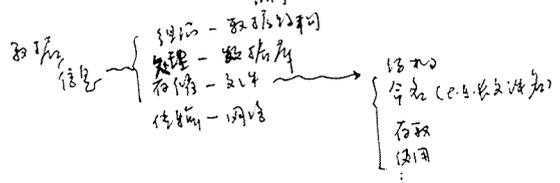
淘汰策略: 局部 FIFO

工作集策略: Δ 在总页数的比例
 内存页数不大, 页号, 页号
 可以在总页数的比例

文件管理 { 目录与实现
 存储与实现
 管理与维护
 硬件和软件上的实现

1 概述

文件是用来存储在磁盘上的数据



各种作业: 用语言来实现文件系统
 管理, 模拟出一个小型文件系统

各级信息独立地存在于 disk
 对用户提供接口, 提供统一与接口
 要有磁盘的跟踪存取与保护
 互访与本地提高效率

基本观点:

文件 = 数据 + 数据项
 = 文件体 + 文件说明
 (内部字符, 文件名, 存储地址)
 各信息项以指针形式, 有逻辑指针与物理指针.

UNIX 的文件
 普通文件
 目录文件
 特殊文件 (设备) { 块设备
 字符设备

2 文件结构

{ 逻辑结构
 物理结构
 文件存储格式 { 连续
 离散

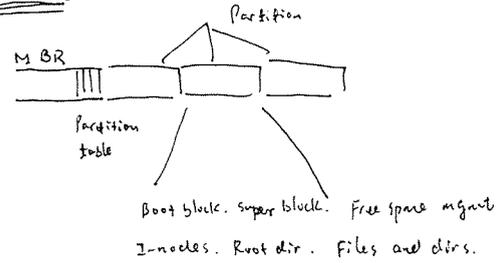
文件系统的组成 { 文件控制
 ...
 ...

文件系统 model:
 用户与应用程序 (用户)
 操作和系统软件 (核心部分)
 文件系统的接口

文件控制块 (FCB, File Control Block), 存储文件的所有管理信息
 { 类型-长度
 所在盘, 块地址

FCB 在 UNIX 中称为 I-node (i 结点)

文件以逻辑结构是字节流
 { 也可理解为块 (连续, 或可离散记录) 处理
 多层次结构 (树形结构)
 文件层次结构存储在磁盘上
 现在主流的 OS 都是无结构的字节流, 并不理解所访问不同可解释不同结构



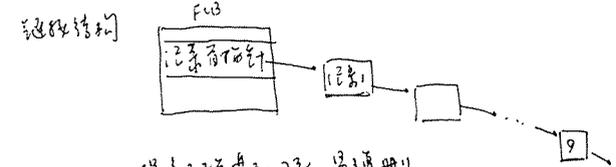
主板上 BIOS 启动, 读入 MBR, 找到分区与分区表
 找到分区表区 boot block 并读入内存

逻辑文件 → 字节流 stream
 如何随机访问
 随机 key 访问

物理结构

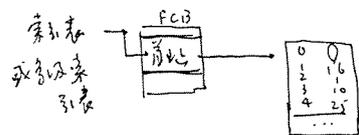
簇 (物理块集合), 是文件的基本物理单位

连续结构 { 简单, 易实现
 IO 速度快, 不必寻道时间, 浪费少
 缺点是: 不利于增删
 } ⇒ 现在不常用
 e.g. 只读 CD 等可以用



提供了链表访问, 易增删 (无碎片)
 缺点是: 寻道时间增长, 译码与解释可能效率不好
 访问速度较慢

索引结构



索引表
或索引表
列表

记录号通过4位
优先索引利用碎片空间

WIN - FAT表 (File Allocation Table)

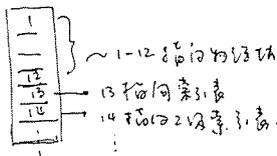
与磁盘分区, FAT表保存索引信息, 可保存在内存中亦可直接从磁盘中

是Windows下常用之文件格式

每字节对应一簇
也加个簇号

UNIX

节点方法: 正常信息的管理文件以FCB



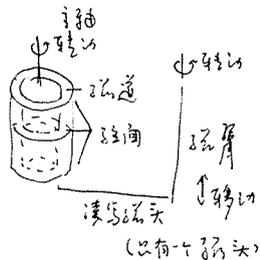
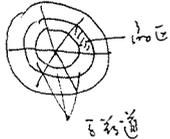
将文件与文件中
混在一起

文件存储介质

磁带: 大容量, 快速度

顺序存储设备 (e.g. 磁带)

随机存储设备 (磁盘)



物理地址:
磁头号, 磁道号, 扇区号
磁头号
磁道号
扇区号
- 磁头 + 磁道 + 扇区

时间 = 寻道 + 旋转 + 数据传输

3 目录管理

目录 directory
子文件夹 folder
簇号
索引
物理地址

区别

文件目录 (FCB的集合)

目录项就是FCB
是文件控制块与索引的集合

目录文件: 是文件目录的存储形式, 与文件共存

文件控制块

索引信息 (索引表, 簇号) 有一些用户可知的
存储控制信息 有一些是透明的
使用信息 (创建时间, 日期)

索引节点: 将文件与文件信息分开, 避免一次调入内存

磁盘索引节点: 与文件透明

内存索引节点: open时把磁盘索引节点调入

但也有更多信息
{ 指向文件

目录组织形式 (结构)

一级目录: 线性结构
检索问题
重名问题 (单用户使用)
不能分组
某项名称时一次调入索引表

二级目录: 根目录 -> 用户目录 -> 线性结构
检索简单
但仍不能避免重名与检索力

多级目录 (树型目录)
tree-like
适合大容量的文件系统
目录名-目录项, 子名称...
绝对路径信息, 与相对路径信息, 相对路径信息

检索简单
不重名
存取速度快

改进的多级目录

符号目录项 (符号名, 树型目录) ~ 装入内存
基本目录项 (内容, 簇号, 索引) ~ 高速时装入

文件目录索引

文件地址 (与物理地址无关)

目录项: 链接法 (MS-DOS), 文件号 + FCB
 (条目) 间接法 (Unix, inode), 文件号 + FCB地址

长文件的问题: 因为255字符 → 链表法, 或太小
 可变文件长度 → 空间长度不固定, 不紧凑
 目录项中 → 长度是固定, 指针指向可变文件的名称 → 指针
 (可变与不变的结合)

目录项搜索方法: 线性搜索 (大海捞针, 时间太长)
 Hash表 (将文件名称转为0~255的何)
 解决冲突: 链, 通过指针法或Hash, 只读线性扫描

4. 文件系统实现

如何管理磁盘空间 (不统一, 同心)

新创建文件的存储空间 { 数据区
 块区分配

磁盘分区 → 主卷 → 磁头 (文件系统创建前初始化)
 扩展磁卷

IDE
 SATA } 硬盘类型

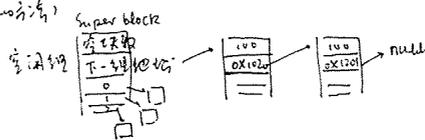
磁盘交叉编组 (一个卷, 跨多个磁头)

文件存储单位: 簇 (cluster)
 由相邻磁头 → 磁头, 为大小不固定
 簇大小不可变 (原来也有可变)

文件存储分配 { 连续分配
 离散分配
 索引分配

空闲空间管理

位示图
 空闲链表 (链表并排序)
 空闲块链表 (有内存地址)
 空闲空间索引法 (比空闲块链表在内存)
 成组链接 (UNIX方法)



大文件小文件都方便
 归也易于管理

实现文件系统应用

系统打开文件表: 地址在内存, 共享的链表 → 可以共享
 用户的打开文件表: 每个进程一个表, 地址在内存, 读写权限, 未读表 → 入口 ~ 私有信息

记录成链与分解

成链: 新设备 → 后加数据, 提高存取速度
 读写时加入 buffer

文件操作

create 创建文件
 open 打开文件, 地址送入内存, 与系统表关联, 计数+1, 返回file
 read 读, write 写, seek 寻址

实现文件操作有关

内存 { 缓冲区
 FCB
 内核 { 系统调用
 进程 (用户的打开文件表)

open(文件号, 打开方式) 系统调用

目录查找
 系统访问合法性
 查找系统表或加入系统表, 共享+1
 加入进程表, 与用户操作, 访问方式
 返回file, 指向进程表的表项

close 系统调用

删除进程表项
 系统表减1, 可过这步删除
 进程

文件保护和安全

文件可以共享 { 共享 → 读, 读写共享
 互斥 → 共享或私有读写操作

时间 { 可以同时读 → 读写与读问题
 不可以同时写

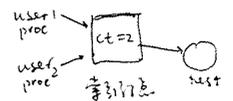
互斥锁 lock

eg. win 2000 互斥系统调用

fcntl user @ ... UNIX 的 读/写

(UNIX使用) (互斥系统表, 相同同一锁点)
 父子进程或共享系统表, 相同同一锁点

基于索引点的共享方式



不同进程造成实现文件共享

这是一个新的文件, 只是名字与文件内容

这方案, 只有一个锁, 所以对文件进行多进程
 读写操作

文件的操作方式
{ 通过命令 (命令)
通过程序
通过文件的存取权限
隐式、显式

访问权限可以附在文件上
访问控制表
可以附在用户上
用户能力表

二级存取控制 (UNIX)
owner
group + rwx
other
只需要个位数的控制权限

文件系统 - 卷 (卷)

MSDOS 的 FAT 文件系统

- 文件卷信息
- 目录信息
- 打开文件管理

WIN2000 的文件系统是 NTFS

NTFS 的卷集：分区和成簇
簇卷、连续卷

UNIX

而操作和目录的文件

设备管理 设备管理是信息管理中

1. 概述

设备 = 除 CPU、内存外的所有设备

IO 设备是计算机的接口
是操作系统与硬件之间的
与系统联系, 特别是与操作系统联系

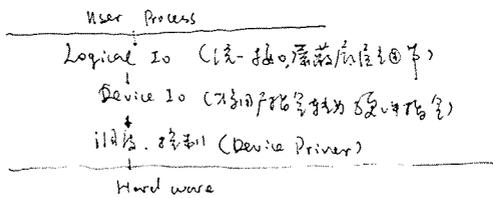
设备的类型 (统一符号)
存取 交互性

设备管理的目标: 设备控制 (统一控制)
设备保护

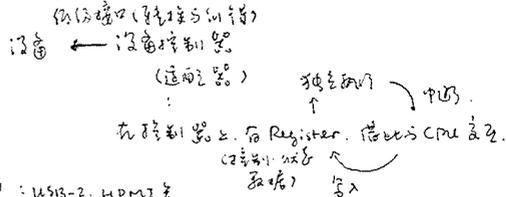
任务: 缓冲, 控制, 调度, 设备与通信保护

2. IO 子系统

设备管理结构



IO 子系统的组成



设备接口: USB-2, HDMI 等

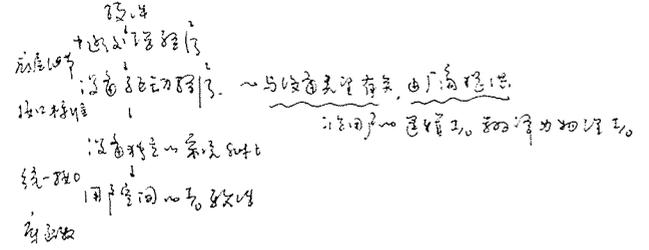
CPU 如何控制设备控制器, 使之与 buffer 通信

- 方法 1: 给每个设备寄存器分配一个编号, 用专门 IO 指令操作 => 完全独立
- 方法 2: 内存映像编址: 编址地址与内存编址统一, 利用统一的 IO 操作指令
- 方法 3: 流控编址: 寄存器编址 (寄存器地址也 => 寄存器不同) 寄存器编址内存 (寄存器不同) => 最主要的方式

3. IO 软件 -> 基础决定上层建筑

也是通用性 (设备驱动程序)

设备类型是: 与设备通信



设备驱动程序 -> 通用性接口 (是个苦力活)

高级语言

设备驱动程序

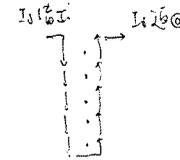
设备驱动程序空闲, 否则加入等待队列

向设备控制器发送指令

设备驱动程序

设备驱动程序

统一接口 (统一符号, e.g. UNIX/Linux 设备-4 种设备 driver include)
统一接口大小: 设备大小, 设备驱动程序, 设备驱动程序 (设备)
缓冲技术: 设备大小 -> pool 缓冲区, 设备驱动程序 -> buffer region
设备驱动程序: 设备驱动程序



与 IO 有关的技术 Spooling (缓冲)

设备驱动程序 CPU 交互

缓冲技术 (设备驱动程序 buffer region)

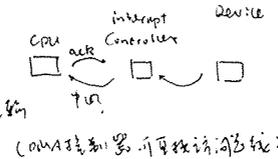
4. IO 控制方式

控制 - 4 种: 轮询 (polling/busy waiting), 中断 (interrupt)

设备 - 中断活动, 每次读一个字节, 使 CPU 忙碌

DMA 方式 (DMA 方式)

内存与设备直接成块传输, 设备驱动程序



(DMA 控制器, 设备驱动程序)

DMA 方式中断 区别
 ↳ 数据 控制器以 buffer 方式
 块传输使用中断

DMA 问题: 总线竞争 地址: 高速、不靠 CPU 控制
 数据输入输出

• 通道方式 (10 年说)

↳ 独立 CPU 或 寄存器指令功能 或 专用控制器

{ 选择通道
 多通道

现在 单独通道已不存在, 但融入了总线中.

与 DMA 不同, 软件操作来传输
 可以各设备管理传输, 更灵活.

设备

5. I/O 方式

设备控制器 (DCR) → 每设备一块

{ 系统设备 (SDT) ~ 系统记录所有设备 { DCR 地址
 设备 I/O 地址
 DCR 的寄存器
 控制器控制表 (CUCT)
 通道控制表 (CCT)

设备设备 共享设备 ↔ 设备设备

可以互锁 互斥互锁 FCFS 基于优先级

设备中 使用程序 设备中 使用程序
 控制器和通道间 互斥互锁 互斥互锁

6. 磁盘调度

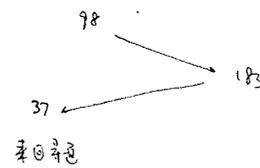
请求 | 请求 10 byte | 数据 50 byte | 尾部
 | 19 字节 校验 | | | ECC 纠错 冗余 |
 t₁ 8.3ms t₂ 16ms

I/O 时间: 扫描时间 + 寻道时间 + 旋转延迟

20ms 或 * 旋转延迟 可保持数据
 { 若的寻道时间 时间 t₁ + (t₂ + t₃) * n
 若的旋转时间 时间 (t₁ + t₂ + t₃) * n

对随机分布的读写序列, 可有不同算法来调度

FIFO 算法



先来后到 (不优化等待时间)

后进先出 (事务延迟)

最短查找时间优先 (SSTF)

扫描 (SCAN) 算法 → 电梯调度算法. 优先的 - 4 号 (即 SSTF)

最短延迟 (公平)

提高吞吐

循环扫描 (C-SCAN)

N 步扫描 (N-step-SCAN)

双向扫描 (FSCAN) 算法

其他方法

{ 磁盘高速缓存
 提前读
 延迟写
 优先的读写
 缓冲池

7. 几种典型 I/O 设备

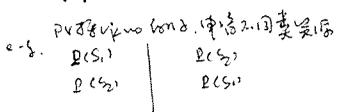
8. 操作系统中的例子

防止进程饥饿 - 一部分
死锁发生: 循环与等待

每个进程入队等待, 先进行
进行“资源竞争”过程, 至持有进程

完成作业, 再持有资源, CPU等
已分配资源量和和量 (用资源与视图)
完成作业
所有可持有资源不持有

原因: 由于进程持有部分资源不释放, 等待其他进程释放



e.g. 交通拥挤时的死锁 - 同类型资源
高区

必要条件: 互斥 + 请求和释放 + 环路等待 + 环路等待
同时满足才能死锁
防止一个打破死锁
若成功, 剥夺其他资源
若成功, 释放自己资源
不太可能打破
避免资源持有者
方法: 使用令牌来管理
资源分配与设备

应对策略: 预防为主: 防止发生死锁 (民用设备很常见)

死锁预防: 破坏4个条件之一
死锁避免: 精心设计资源分配流程
检测与解除
都有不小的代价

破坏“互斥”: 共享设备增加为可共享设备, e.g. 打印机

破坏“环路等待”: 等待时, 因为资源可以被全部释放

破坏“请求和释放”: 不允许进程与资源同时申请资源, 必须一次申请
+ 捆绑

破坏“环路等待”: 有序资源分配法

先申请小编号, 再申请大编号 \Rightarrow 不可分配, 不可分配

Proc 1	Proc 2
R(5)	R(5)
R(2)	R(2)

e.g. 指导资源分配问题: Dijkstra 提出

试试用前边所讲过的试金石, 见 'Modern Operating Systems'

思路1: 互斥 + 互斥, 互斥 + 互斥

防止资源分配在右邻居

思路2: 有序资源

思路3: 资源限制

有资源完全避免死锁问题? 算法 + 数据结构

判断系统状态是否安全, 是否满足安全状态
可行

银行家算法

安全: 不危险, 不危险
不安全: 危险, 危险

银行家算法: 可以避免死锁, 避免不安全状态 (Dijkstra 1965)

防止不安全时, 就不为分配, 可避免分配一部分资源

直到系统事务声明最大高资源, 申请超过最大高资源, 则拒绝

理论重要
但不实用

死锁问题的预防: 更加代价是 { 递归线性选择 (比如因遇到递归临界点)
剥夺资源
防止有比计算 (避免先以失败也需出时间? 不可剥夺资源? 经常

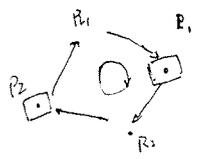
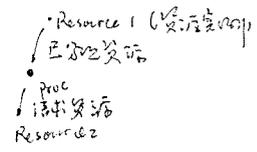
是循环的等待: 线性检测死锁时不检测 (死锁检测表和潜在资源组成) 检测的进程?
UNIX, 可以考察, 进行检测
死锁不进行检测

求解方法: 最小花费资源消耗
 剥洋葱法 (代价***)
 ↑ 不再重复, 避免死循环

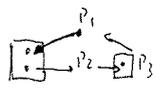
递归求解 (代价****) ~ 可以自己在递归中进行.
 ↑ 尽可能选择轻量级进行递归

进程回溯 (代价*, 但比代价低, 很少使用, 非常复杂)
 逐地记录进程状态到磁盘.
 (高内存消耗和占用资源, 代价高)
 多路回溯记录, 动态记录回溯

资源分配图: 可以解决资源分配问题一有力工具.



可以
 成环即死锁 (每类资源实例数一致, 成环 \Rightarrow 死锁)
 不成环 \Rightarrow 成环 \Rightarrow 死锁)



可破死锁时, 资源分配, 若能解决, 则不形成死锁
 不成, 死锁