

的末尾。

Thread Scheduling 可以 FCFS/SJF

Process synchronization

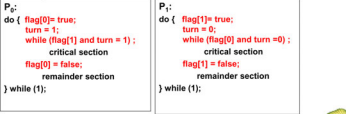
The Critical-Section Problem

N 个进程都计算同样的共享数据，每个进程都有一个临界区，其中共享数据被访问。问题：需要保证只有一个进程进入临界区。临界区问题的解决必须满足三个要求：**互斥 (mutual exclusion)**、**空闲让进 (progress)**、**有限等待 (bounded wait)**。让权等待不是必须的。

Peterson 算法

只用于两个进程的情况，并且假设 load 和 store 是原子操作，是一种软件解决方法。

The two processes share two variables:
• int turn;
• boolean flag[2];
The variable turn indicates whose turn it is to enter the critical section. The flag array is used to indicate if a process is ready to enter the critical section. flag[i] = true implies that process P_i is ready!



Meets all three requirements, solves the critical-section problem for two processes!

现代 OS 中不适用(若编译器交换两条指令则奇)

硬件同步方法

单处理器：在 CS 禁止中断。
多处理器：**Memory barriers**(an instruction forcing any change in memory to be propagated (made visible) to all other processors)

Hardware instruction original

```
boolean test_and_set (boolean *target) {
    boolean rv = *target;
    *target = true;

    return rv;
}

while(1) {
    key = TRUE;
    while (key == TRUE)
        Swap(&lock,&key);
    critical section
    lock = false;
    remainder section
}
```

Test_and_set 存在 busy waiting。会忙等待的信号

量是 spin lock

抽象出 2 个硬件实现的不可以被中断的原子操作：赋值和交换，然后来解决临界区，testandset 的共享变量是 lock 初始 false, swap 也是一样，但是多了局部变量 key 不是共享的。

硬件方法优点：进程数随意，简单，支持多个临界区；缺点：无法让权等待，可能饥饿，可能死锁。会引起 busy waiting

semaphores

atomic 操作 wait(P,不用 waiting queue 会 busy waiting), signal(V,不会 busywaiting)

count 信号量，值域不受限；

binary 信号量=互斥锁 mutex，只能是 0/1(可以被>2 个线程使用,不会 block,会 busywaiting)。为消除 wait(P)时的 busy waiting，对信号量增加 block(run->wait)和 wakeup(wait->ready)来避免忙等 wait(semaphore *S) { S->value--; if (S->value < 0) { add this process to S->list; bl ock(); } signal(semaphore *S) { S->value++; if (S->value <= 0) { remove a process P fro m S->list; wakeup(P); } 具有忙等的信号量值非负，但是这种实现可以为负，负数绝对值代表等待该信号量的进程数，0 代表无资源可用。Mutex of 2 concurrent processes has value 0: 1 process entered CS, 0 process blocked Count S of 3 concurrent processes has value 1: 0 process entered CS, 0 process blocked 0: 1 process entered CS, 0 process blocked

-1: 1 process entered CS, 1 process blocked
-2: 1 process entered CS, 2 process blocked
Wait and signal 成对出现，互斥操作就在同一进程中出现，同步操作在不同进程。连续的 wait 顺序是需要注意的，但连续 signal 无所谓。同步 wait 和互斥 wait 相邻时，要先同步 wait。
优点：简单、表达能力强；缺点：不够安全，使用不当会死锁
优先级倒置 (priority inversion)：当优先级较低的进程持有较高优先级进程所需的锁定时调度问题。解决方法：**priority inheritance**: temporary assign the highest priority of waiting process to the process holding the lock.

Monitor 管程

一次只有 1 个进程能在管程内活动.内部可定义 condition 类型的变量以提供同步机制

Bounded-Buffer Problem 有限缓冲区 生产者-消费者问题

很多相互合作进程
设置 N 个缓冲项；信号量 mutex 初始化为 1，用来保证访问 buffer 的互斥要求；full 初始化为 0，表示满缓冲项的个数；empty 初始化为 N 表示空缓冲项的个数(这里不能只用一个 count 表示 buffer 空个数，否则就无法实现锁)生产者：do {…produce an item…wait(empty); wait(mutex); …add to buffer …signal(mutex); signal(full); while (1);
消费者：do {wait(full); wait(mutex); …remove an item from buffer…signal(mutex); signal(empty); …consume the item …} while (1);
Readers-Writers Problem 数据库第一读写问题：允许多个读者同时读，但是只有一个写者，也就是没有读者会因为写者在等待而等待其他读者的完成，写者可能饿死。第二读写问题：写者就绪后，写者就立即开始写操作，也就是说写者等待时，不允许新读者进行操作，读者可能饿死。

共享数据有访问的数据、reader_count_lock 初始 1,保证更新 readcount 时互斥；write_lock 初始 1,为读写公用，供写者作为互斥信号量，被第一个进入 CS 和最后一个离开 CS 读者使用，其他读者不适用；Readcount 初始 0，跟踪多少进程正在读。
写进程：do {wait (write_lock) ; …writing… signal (write_lock) ;} while (1);
读进程：do { wait(reader_count_lock); readcount++; if (readcount == 1) wait(write_lock); signal(reader_count_lock); …reading … wait(reader_count_lock); readcount–; if (readcount == 0) signal(write_lock); signal(reader_count_lock); while (1);

Dining-Philosophers Problem 哲学家进餐

问题描述：N 个哲学家坐在圆桌，每个哲学家和邻居共享一根筷子；哲学家吃饭要用身边的两只筷子一起吃；邻居不允许同时吃饭；哲学家只会思考或者吃饭。
共享数据：一碗米饭；共享变量 chopstick[5]初始为 1 表示筷子未被持有；lock 初始化为 1 表示某个哲学家是否可以同时拿起左右两根筷子。解决方案之一：只允许在 CS 内同时拿起两根筷子；轮询每人是否能够拿起两根筷子，如果能则拿起，如果不能则需要等待那些筷子放下

do {wait(lock);wait(chopstick[i]);wait(chopstick [(i+1) % 5]); signal(lock);…eat…signal(chop stick[i]);signal(chopstick[(i+1) % 5]); … thin k …} while (1);
防止死锁的可能解决：最多只允许 4 个哲学家坐在桌上 do {…thinking…

P(count);/*初始 count=4，是否超过四人准备进餐*/P(chopstick [i]); P(chopstick [(i+1) % 5]);…eat..

V(chopstick [i]); V(chopstick [(i+1) % 5]); V(c ount);
}while(1);
使用非对称的解决方法：奇数先拿左手，偶数先拿右手 do {…thinking…
if(%2 == 1) {P(chopstick [i]);//判断左边的筷子是否可用 P(chopstick [(i+1) % 5]);//判断右边筷子是否可用};else {P(chopstick [(i+1) % 5]);//判断右边的筷子是否可用 P(chopstick [i]);//判断左边的筷子是否可用}…eat..

V(chopstick [i]); V(chopstick [(i+1) % 5]);
}while(1);

Deadlock 死锁

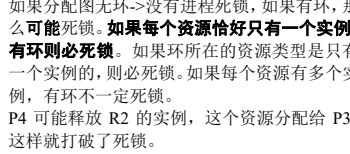
死锁指多个进程因竞争共享资源而造成的一种僵局，若若无外力作用，这些进程都将永远不能再向前推进。进程按以下顺序使用资源：申请 使用 释放。条件和释放为系统调用。

四个必要条件：

Mutual exclusion：至少有一个资源处于共享模式，一次只能有一个进程使用该资源；hold and wait：一个进程必须至少占有一个资源并且等待另一个资源，且该资源被其他进程占有；No preemption：资源不能被抢占，只能在进程使用完成后释放；circular wait：进程间循环等待资源，A 等 B 占的 B 等 C 占的 C 等 A 占的。

资源分配图，由点 V 和边 E 组成，V 被分为两部分：系统活动进程的集合 系统所有资源类型的集合。进程 Pi 到资源 Rj 的有向边记为 Pi->Rj，表示进程 Pi 已经申请了资源类型 Rj 的一个实例，叫请求边；资源 Rj 到进程 Pi 的有向边表示资源类型 Rj 的一个实例

已经分配给了进程 Pi，叫做分配边。进程用圆表示，资源类型用方表示，资源实体是方内部的点。
如果分配图无环->没有进程死锁，如果有环，那么可能死锁。如果每个资源恰好只有一个实例，有环则必死锁。如果环所在的资源类型是只有一个实例的，则必死锁。如果每个资源有多个实例，有环不一定死锁。
P4 可能释放 R2 的实例，这个资源分配给 P3，这样就打破了死锁。



死锁处理

保证系统不进入死锁：prevention avoidance；允许进入死锁但是需要恢复：detection recovery。U L W 三个系统都忽略问题假设没有死锁，是鸵鸟方法。

死锁预防 Prevention

通过限制请求的方式来预防死锁。
Mutual exclusion: 非共享资源必须互斥，例如一台打印机不能被多个进程共享，而共享资源不需要互斥，也不导致死锁，类似只读文件。
Hold and wait: 必须保证：一个进程申请一个资源时不能占有其他资源。进程在执行前就要申请并分配资源，是资源静态预分配的方法；缺点：低资源利用率、可能饥饿。
No preemption: 如果一个进程占有资源并且申请了另一个不能立即分配的资源，那么它现已分配的资源都可以被抢占，也就是被隐式释放了。抢占资源分配到进程所等待的资源列表中。进程需要获取到原有的资源和申请的新资源后才能运行。

Circular wait: 给资源设置显式序号，请求必须按照资源序号递增的方式进行，通过资源的有序申请破坏了循环等待条件。然而有时候资源的“顺序”是模糊的（from, to）

死锁避免

前面的方法虽然避免了死锁，但是降低了吞吐率，我们可以通过获取一些额外的事先信息从而避免死锁 prior information。
最简单和最有用的模型要求每个进程声明它可能需要的每种类型的资源的最大数量。死锁避免算法动态检查资源分配状态，确保永远不会出现循环等待。资源分配状态由可用和已分配资源的数量以及进程的最大需求定义。

安全状态：对于所有进程，如果存在一个安全序列，那么系统就处于安全状态。对于进程序列 P1,P2,…,Pn，如果对于每一个 Pi,Pi 仍然可以申请的资源数小于当前可用的资源加上所有进程 Pj(j<i)所占有的资源，那么这一序列是安全序列。这种情况下，进程 Pi 的资源即使不能立即可用，那么 Pi 可等待知道所有 Pj 释放其资源，当它们完成时 Pi 就可以运行，Pi 运行结束后，Pi+1 就可以获得到所需的资源，如此进行。
安全状态->没有死锁；不安全状态->可能有死锁；避免->保证系统永远不进入不安全状态。

资源分配图，single instance 死锁避免算法：

Single instance: 每种资源只有一个引入一种新边 claim edge 需求边，Pi->Rj 表示进程 Pi 在未来可能请求资源 Rj，用虚线表示。当进程真正请求资源时，用请求边覆盖掉需求边。当资源被分配给进程后，用 assignment edge 分配边来覆盖掉请求边，当资源被释放后，分配边恢复为需求边。系统必须事先说明需求边。算法：假设进程 Pi 申请资源 Rj。只有在需求边 Pi->Rj 变成成分配边 Rj->Pi 而不会导致资源分配图形成环时，才允许申请。

用该算法循环检测，如果没有环存在，那么资源分配会使系统继续安全状态，否则就会不安全，Pi 就要等待。

Banker, 多实体资源类型避免算法：

每个进程实现说明最大需求；进程请求资源时可能会等待；进程拿到资源后必须在有限时间内释放它们。

数据结构：
N 进程数，m 资源类型的种类数；Available: 长度为 m 的向量，表示每种资源的现有实例数量，available[j]=k 表示 j 型资源还有 k 个；Max: n*m 的矩阵，定义每个进程的最大需求，max[i][j]=k 表示进程 Pi 最多可以申请 k 个 Rj 型资源；Allocation: n*m 的矩阵，表示每个进程所分配的各种资源类型的实例数，allocation[i][j]=k 表示已经为 Pi 分配了 k 个 Rj 型实例；Need: n*m 矩阵，表示每个进程还需要剩余的資源，need[i][j]=k 表示进程 Pi 还可能继续申请 k 个 Rj 型的实例。Need=max-allocation。

安全状态检测算法：

1. 设 work 和 finish 分别是长度为 m 和 n 的向量，初始化：work=available, finish[i]=false;
2. 寻找 i 满足 finish[i]=false 且 need[i]<=work, 如果 i 不存在跳到第四步；
3.work=work+allocation[i],finish[i]=true, 返回第二步；
4.如果所有的 finish 都是 true，那么系统处于安全状态。

算法需要 m*n*n 的操作数量级确定系统状态

资源请求算法：

Request[i]为 Pi 的请求向量，如果 request[i][j]=k 那么进程 Pi 需要资源类型 Rj 的数量为 k。当进程 Pi 请求资源时，动作如下：
1.如果 request[i]<=need[i]跳到第二步，否则出错，因为进程 Pi 已经超过了其最大需求。
2.如果 request[i]<=available 跳到第三步，否则 Pi 必须等待，因为没有可用资源

3.假定系统可以分配给进程 Pi 请求的资源，进行下面的操作：Avail=avail-request[i]; allocation[i]=allocation[i]+request[i];need[i]=need[i]-request[i]; 如果产生的资源分配状态是安全的，那么交易完成且进程 Pi 可以分配到资源，如果新状态不安全，那么进程 Pi 必须等待 Request[i]并且恢复到原有的资源分配状态。

死锁检测

允许系统进入死锁状态的话，那么系统就需要提供检测算法和恢复算法。

等待图，单实体资源类型检测算法：

等待图是资源分配图的变形，节点都是进程，Pi->Pj 表示 Pi 在等待 Pj 释放 Pi 所需的资源。当且仅当等待图中有一个环，系统死锁，检测环的算法需要 n*n, n 为点数。
多实体资源类型检测算法：

数据结构：Available, allocation 是一样的，request: n*m 的矩阵，表示当前各进程的資源请求状况，request[i][j]=k 表示 Pi 正在请求 k 个资源 Rj。算法：1.设 work 和 finish 分别是长度为 m 和 n 的向量，初始化：work=available, 如果 allocation[i]非 0, finish[i]=false 否则初始化为 true;

2.寻找 i 满足 finish[i]=false 且 request[i]<=work, 如果 i 不存在跳到第四步；
3.work=work+allocation[i],finish[i]=true, 返回第二步；

4.如果某个 finish 是 false，那么系统处于死锁状态，且对应下标的进程 Pi 死锁。

算法需要 m*n*n 的操作数量级确定系统状态

死锁检测算法的应用

检测算法的调用时机及频率取决于：死锁发生频率以及思索发生时受影响的进程数。如果经常发生死锁，那么就要经常调用检测。如果在不确定的时间调用检测算法，资源图可能有很多环，通常不能确定哪些造成了死锁

死锁恢复

检测到死锁后的措施：通知管理员 系统自己恢复。打破死锁的两种方法：抢占资源 进程终止。

进程终止

两种方法来恢复死锁：终止所有死锁进程 一次终止一个进程直到不死锁。许多因素都影响终止进程的选择：优先级 进程已经计算了多久，还要多久完成 进程使用了哪些类型的资源 进程还需要多少资源 多少进程需要被终止 进程是交互的还是批处理的

抢占资源

抢占资源需要处理三个问题：
选择一个牺牲品 victim: 要代价最小化
回收：退回安全状态，但是很难，一般需要完全终止进程重新执行
饥饿：保证资源不会总是从同一个进程中被抢占。常见方法是代价因素加上回收次数。

Main Memory 主存

层次存储中主存 cache 寄存器为 volatic 易失的。逻辑地址/虚地址/相对地址：由 CPU 生成，首地址为 0，逻辑地址无法在内存中读取信息。物理地址/实地址/绝对地址：内存中存储单元的地址，可以直接寻址。物理地址中的逻辑地址空间是通过一对基址寄存器和界限地址寄存器控制的 base and limit register。若 base reg 为 300040, limit reg 为 120900, 那么程序的合法访问从 300040 到 420910(含)的所有地址。

Address binding3 种情况：

编译时间：如果编译时就知道进程在内存中的地址，那么就可以生成 absolute code。

装载时间：编译时不知道在哪，那么编译器生成

可重定位代码 relocatable code。

执行时间：如果进程在执行时可以移动到另一个内存段，需要硬件支持也就是 base and limit 目前绝大多数都是采用这种。

Memory-Management Unit (MMU)

就是将虚拟地址映射到物理地址的硬件设备。在 MMU 中，base 寄存器叫 relocation register，用户进程送到内存时，都要加上重定位寄存器的值。PA=relocation reg+LA。用户程序只能处理 LA，永远看不到真正的 PA。

Dynamic Loading

进程大小会收到物理内存大小的限制，为了更好的空间使用率，采用动态加载。一个子程序只有在调用时才被加载，所有子程序都以 relocatable 形式存在 disk，需要的时候装入 memory。OS 不需要特别支持，是程序设计做的事。当需要大量的代码来处理一些不常发生的事时很有用，如错误处理。

Swapping (交换技术)

进程可以暂时从内存中交换到 backing store 上，当需要再次执行时再调回。需要 dynamic relocate 支持

backing store: 是快速硬盘，而可以容纳所有用户的所有内存映像，并为这些内存映像提供直接访问，如 Linux 交换区 windows 的交换文件 pagefile.sys

Roll out roll in: 如果有一个更高优先级的进程需要服务，内存交换出低优先级的进程以便装入和执行高优先级进程，高执行完后低再交换回内存继续执行。

交换时间中的主要部分是转移时间 transfer time。总转移时间与所交换的内存大小成正比。系统维护一个就绪的可立即运行的进程队列，并在磁盘上有内存映像。

Contiguous Allocation (连续分配)

内存通常分为两个区域：一个驻留 resident 操作系统，一个用于用户进程，由于中断向一般位于低内存，所以 OS 也放在低内存。重量位寄存器用于保护各个用户进程以及 OS 的代码和数据不被修改。Base 是 PA 的最小值；limit 包含了 LA 的范围，每个 LA 不能超过 Limit。MMU 地址映射是动态的。

Multiple-partition allocation: 分区式管理将内存划分为多个连续区域叫做分区，每个分区放一个进程。有固定分区和动态分区两种。

动态分区：

动态划分内存，在程序装入内存时切出一个连续的区域 hole 分配给进程，分区大小恰好符合需要。操作系统需要维护一个表，记录哪些内存可用哪些已用。从一组可用的 hole 选择一个空闲 hole 的常用算法 first best worst-fit 三种。分别是分配第一个是足够的/分配最小的是足够的/分配最大的。First 和 best 在时间和空间利用率都比 worst 好。还有一个 next-fit 是每次都从上次查找结束的位置开始找，找到第一个足够大的。

碎片 fragmentation

first 和 best 都存在 external frags。外碎片指所有的总可用内存可以满足请求，但是并不连续。外碎片可通过 compaction 拼接 defragmentation 减少。重定向是动态并且在执行时间完成可以进行紧凑操作 重新排列内存将来将碎片拼成一个大块，但是是拼接的开销很大。

Internal frags 是进程内部无法使用的内存，这是由于零头和块大小造成的，比如块大小 8B，进程有 9B，那么不得不给他 16B 的内存，就出现了 7B 的内碎片。

分页存储管理

分页允许进程的 PA 空间非连续；将物理内存分

的技术。

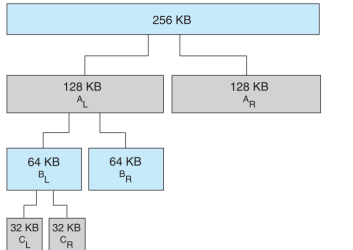
开始的文件访问按照普通按需请求调度，会出现页错误。这样，一页大小的部分文件从文件系统中读入物理页，以后的文件访问就可以按照通常的内存访问来处理，这样就可以用内存操作文件，而非 read write 等系统调用，简化了文件访问和使用的。多个进程可以允许将同一文件映射到各自的虚存中，达到数据共享的目的。

Allocating Kernel Memory 内核内存分配

与对待用户内存不同：内核内存从空闲内存池中获取，两个原因：1.内核需要为不同大小的数据结构分配内存。2.一些内核内存需要连续。

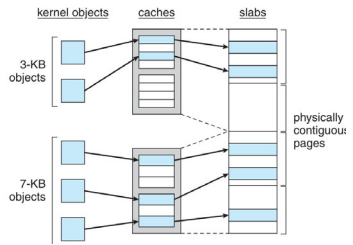
Buddy system

从物理上连续的大小固定的段上进行分配。内存分配按 2 得幂的大小来分配：请求大小必须是 2 的幂；如果不是，那么调整到一个更大的 2 得幂；当需要比可用的更小的分配时，当前块分成两个下一个较低幂的段。继续这一过程直到适当大小的块可用。优点是可以通过合并快速形成更大的段。缺点是由于调整到下一个 2 的幂容易产生 **internal 碎片**。



Slab 分配

为了解决 Buddy 碎片损失的问题，slab 是由一个或多个物理上连续的页组成的。Cache 包含一个或者多个 slab。每个内核数据结构都有一个 cache。每个 cache 都含有内核数据结构的对象实例。当 cache 被创建时，起初包括若干标记为空闲的对象。对象的数量和slab大小有关，12KB的 slab(包含三个连续的页)可以存储 6 个 2KB 的对象。当需要内核数据结构的对象时，可以直接从 cache 上取，并将该对象标记为使用 used。Slab 首先从部分空闲的 slab 中分配，如果没有则从全空的 slab 进行分配。如果没有，从物理连续页上分配新的 slab，把他赋给一个 cache，再从 slab 分配空间。Slab 优点：没有碎片引起的内存浪费；内存请求可以快速满足。



预调页 prepaging

为了减少冷启动时大量 page fault。同时将所有需要的页一起调入内存，但是如果预调页没有被用到，那么 IO 就被浪费了。假设 s 页被预调到内存，其中 a 部分被用到了。问题在于节省的 s*a 个页错误的成本是大于还是小于其他 s*(1-a)不必要的预调页开销。如果 a 接近于 0，调页失败，a 接近 1，调页成功。

页大小

页大小的考虑因素：碎片↓、页大小↓、IO 开销↑、局部性↓、page fault 数量↑、TLB 大小及效率↑

TLB 范围 TLB reach：通过 TLB 可以访问到的内存量。TLB Reach=TLB size * Page Size

理想情况下，每个进程的 WS 应该位于 TLB 中，否则就会有不通过 TLB 调页导致的大量 IO 增大页大小来缓解 TLB 压力；但可能会导致不需要大页表的进程带来的内碎片

提供多种页大小的支持：那么 TLB 无法硬件化，性能降低。

IO 互锁

IO 互锁指页面必定有时被所在内存中。必须锁住用于从设备复制文件的页，以便通过页面置换驱逐。

File System Interface 文件系统接口

文件是存储某种介质上的（如磁盘、光盘、SSD 等）并具 有文件名的一组相关信息的集合

File attributes

name identifier 唯一标识该文件:File metadata)

类型 位置 大小 保护 时间日期 用户标识

所有的文件信息都保存在目录结构中，而目录结构保存在 disk。

文件操作

文件是 ADT 抽象数据类型，其操作有：创建 读写 文件内重定位 删除 截短 truncate Open(Fi) 在硬盘上寻找目录结构并且移动到内存中 Close(Fi)将内存中的目录结构移动到磁盘中。

Open

每个打开文件都有以下信息：文件指针：跟踪上次读写位置作为当前文件位置指针

每次访问文件前都要 open 文件.create(首次 open)会在 disk 上创建 FCB,之后 open 从 disk 读 FCB

文件打开计数器 file-open count：跟踪文件打开和关闭的数量，在最后关闭时，计数器为 0，系统可以移除该条目。

文件磁盘位置 disk location of file：用于定位磁盘上文件位置的信息

访问权限：访问模式信息

锁机制：mandatory lock:根据目前锁请求与占有情况拒绝 access; advisory lock:进程查看锁情况来决定访问策略

文件内部结构 File Structure

None 字 节节的序列 流文件结构

Simple record structure 记录文件结构：lines, fixed length, variable length

Complex Structures：formatted document, relocatable load file

可以通过插入适当的控制字符，用第一种方法模拟最后两个

这些模式由 OS 和程序所决定。

访问方法

Sequential access 顺序访问

文件信息按顺序，一个记录接着一个记录处理。访问模式能够够用，编辑器和编译器用这种方式。读操作读取文件下一文件部分，并自动前移文件指针，跟踪 IO 位置。写操作向文件尾部增加内容，相应文件指针指向新文件结尾。顺序访问基于文件的磁带模型，也适用于随机访问设备。可以重新设置指针到开始位置或者向前向后跳过记录。No read after last write

Direct access 直接访问

文件由固定长度的逻辑记录组成，允许程序按任意顺序进行快速读写，直接访问是基于文件

的磁盘模型。文件可作块或记录的编号序列。读写顺序没有限制。可以立即访问大量信息，DB 常用。往往用指向 blocks 的 index 实现。文件操作必须经过修改从而能将块号作为参数，有读 n 操作，而不是读下一个；写 n 操作：定位到 n：要实现读 n 只需要定位 n 再读下一个即可。注意 n 是相对块号，相对于文件开始的索引号。

sequential access	implementation for direct access
reset	cp = 0;
read_next	read cp ; cp = cp + 1;
write_next	write cp ; cp = cp + 1;

Indexed block access 索引顺序访问访问

目录结构

目录是包含所有文件信息节点的集合。目录结构和文件在磁盘上。

磁盘结构

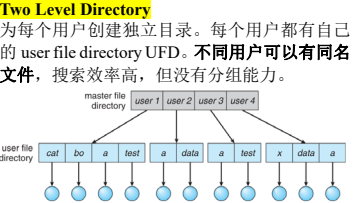
磁盘可以装多种文件系统，分区或片 minidisk slice。

目录操作

搜索文件 创建文件 删除文件 遍历 list 目录重命名文件 遍历 traverse 文件系统

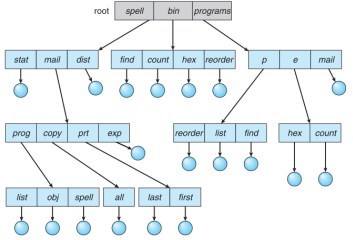
单级目录

所有文件包含在同一目录中，一个文件系统提供给所有用户。由于所有文件在同一级，不能有重名，此外存在着分组问题



树形目录(avoid name collision)

将二级目录拓展即可。搜索高效 有分组能力。



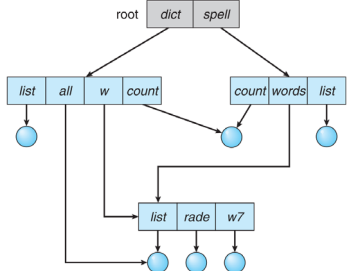
无环图目录 Acyclic graph

树形结构禁止共享文件和目录。无环图允许目录含有共享子目录和文件。

实现文件盒目录共享，UNIX 采用创建一个叫做链接的新目录条目。链接实际上是另一个文件的指针。连接通过使用路径名定位真正文件。注意，无环图目录倒置一个文件可以有多个绝对路径名。不同文件名可能表示同一文件，出现了别名问题。存在 dangling pointer：删除一个文件后指向该文件的其他链接成为 dangling。对于采用符号链接实现共享的系统，删除链接并不影响源文件，如果文件本身被删除，链接也被删除。

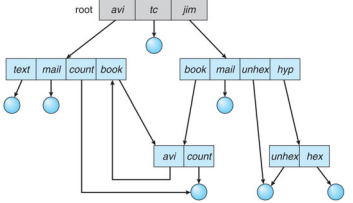
删除的另一方法是保留文件知道删除其全部引用，所以为文件引入了计数，删除一次链接或者条目就计数-1，到 0 时完全删除文件，UNIX 的

硬链接采用这种方法，在 inode 中保留一个 reference count。通过禁止对目录的多重引用，可以维护无环图结构。



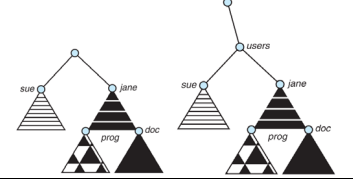
普通图目录 General graph

采用这种目录必须确保没有环，仅允许向文件链接，不允许目录，或者允许环，但使用 garbage collection 回收空间。每次加入链接都要执行环检测算法。



文件系统挂载 mount

FS 在访问前必挂载。没挂载的 FS 会在 mount point 挂载。左图是未安装的卷，右图的用户为挂载点。mount /dev/dsk /users(mount point)



Linux disk partition	
装置	装置在Linux内的文件名
IDE硬盘机	/dev/hd[a-d]
SCSI/SATA/USB硬盘机	/dev/sd[a-p]
USB快闪碟	/dev/sd[a-p](与SATA不同)
软盘驱动器	/dev/fd[0-1]
打印机	251: /dev/lp[0-2] USB: /dev/usb/lp[0-15]
鼠标	PS2: /dev/psaux USB: /dev/usb/mouse[0-15]
当前CDROM/DVDROM	/dev/cdrom
当前的鼠标	/dev/mouse
一块磁盘硬盘分区表 64B，一块分区占 16B，主分区最多 4 个(扩展分区也是一个主分区)，扩展分区最多 1 个，扩展分区上可有多逻辑分区(必须建立在扩展分区上)。主分区范围是 1-4，逻辑分区从 5 开始。	
文件共享 file sharing	
多用户系统的文件共享很有用。文件共享需要通过一定的保护机制实现：在分布式系统，文件通过网络访问；网络文件系统 NFS 是常见的分布式文件共享方法。NFS 是 UNIX 文件共享协议 CIFS 是 WIN 的协议。	
保护 Protection	
访问类型：读 写 执行 追加 append 删除 列表清单 list	
访问控制列表 access-control list ACL	

3 种用户类型：拥有者 owner access 组 group access 其他 public access。UNIX 一个类型有 rwx 3 个权限，所以一个文件需要 3*3=9 位说明文件访问权限

File System Implementation

文件系统：是操作系统中以文件方式管理计算机软件资源的软件和被管理的文件和数据结构（如目录和索引表等）的集合。文件系统储存在二级存储中，磁盘。

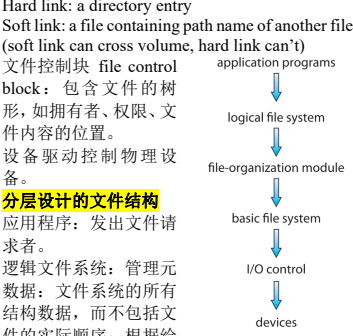
write() & fsync(): write()为在未来某个时间将数据写回到二级存储中，具体时间由文件系统根据性能决定；fsync()强制将 dirty data 写回 disk。

Hard link: a directory entry

Soft link: a file containing path name of another file. (soft link can cross volume, hard link can't)

文件控制块 file control block：包含文件的树形，如拥有者、权限、文件内容的位置。

设备驱动控制物理设备。



分层设计的文件结构

应用程序：发出文件请求者。

逻辑文件系统：管理元数据：文件系统的所有结构数据，而不包括文件的实际顺序；根据给定符号文件名来管理目录结构；逻辑文件系统通过 FCB 来维护文件结构。

文件组织模块知道文件逻辑块到物理块的映射，包括空闲空间管理器。做 translation。

基本文件系统：向合适的设备驱动程序发送一致命令就 可对磁盘上的物理块进行读写 IO 控制：由设备驱动程序和中断处理程序组成，实现内存与磁盘之间的信息转移

文件系统实现

On-Disk FS structure

磁盘上，文件系统可能包括：如何启动所存储的 OS，总块数，空闲块数目和位置，目录结构及各个具体文件。

磁盘结构包括：每个卷的 boot control block 包括从该卷引导操作系统所需信息；每个卷的 volume control block 包括卷的详细信息；目录结构来组织文件；每个文件的 FCB(不含 file name)

file permissions
file dates (create, access, write)
file owner, group, ACL
file size
file data blocks or pointers to file data blocks

In-Memory FS structure

In-memory partition table 分区表 in-memory directory structure 目录结构 system-wide open-file table 系统打开文件表 per-process open-file table 进程打开文件表

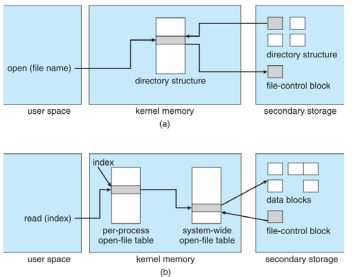


Figure 12.3 In-memory file-system structures. (a) File open. (b) File read.

System-Wide Open-File Table：记录所有被加载到内存中的 FCB inode：Per-Process Open-File Table：指向 System-Wide Open-File Table 中的项（包含当前在文件中位置、文件访问模式等信息）。

虚拟文件系统 VFS

VFS 提供面向对象的方法实现文件系统。允许将相同的系统调用接口 API 用于不同类型的文件系统。（Write syscall -> vfs_write）

目录实现

线性列表 linear list：使用储存文件名和数据块指针的线性表。哈希表：线性表与哈希指向，哈希表根据文件名得到一个值返回一个指向线性表中元素的指针。

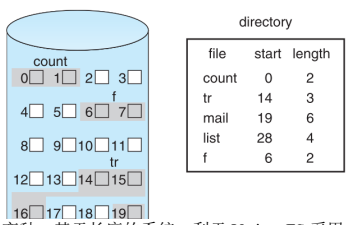
分配方法 Allocation Method

常见的主要磁盘空间分配方法：连续、链接和索引。

连续分配 Contiguous Allocation

每个文件在磁盘上占有一组连续的块。优点：访问很容易，只需要起始块位置和块长度就可以读取。支持寄存器 random access。但是浪费空间，存在动态存储分配问题。First 和 best 表现差不多，第一时间快很多。存在外碎片问题，此外文件大小不可增长。

逻辑到物理的映射：LA/512 分为两部分商 Q 和余数 R。Block to be accessed = Q + starting address Displacement into block = R。LA 是存取文件逻辑地址，512 是块大小



变种：基于长度的系统。利于 Veritas FS 采用。解决了文件大小无法增长的问题，增加了另一个叫做 extent 的连续空间给空间不够的文件，然后与原文件块之间有个指针。一个文件可以有多个 extent。

Linked Allocation 链接分配

解决了连续分配的所有问题。每个文件都是磁盘块的链表。访问起来只需要一个起始地址。没有空间管理问题，不会浪费空间，但是不支持 random access。

地址映射：LA/(512-1)得到商 Q 和余 R，Block to be accessed is the Qth block in the linked chain of blocks representing the file. Displacement into block = R + 1。因为每个索引块的末尾节点是用来链接下一个索引块的，不链数据块，所以要 512-1

性检查。若调用号大于或等于NR_syscalls,

系统调用处理程序终止。(sys_call_table)
4.若系统调用号无效,函数就把-ENOSYS 值存放在栈中 eax 寄存器所在的单元,再跳到 ret from_sys_call()
5.根据 eax 中所包含的系统调用号调用对应的特定服务例程
添加文件系统:
文件类型:
普通文件 (文件名不超过 255)
目录文件
字符设备文件和块设备文件:
fd0 (for floppy drive 0)
hda (for harddisk a)
lp0 (for line printer 0)
tty(for teletype terminal)
管道(FIFO)文件 链接文件 socket 文件

文件系统分三大类: 基于磁盘的文件系统, 如 ext2/ext3/ext4、VFAT、NTFS(win10)等。网络文件系统, 如 NFS 等。特殊文件系统, 如 proc 文件系统、devfs、sysfs (/sys) 等。
Linux 以 ext2/ext3 做为基本的文件系统, 所以它的虚拟文件系统 VFS 中也设置了 inode 结构。物理文件系统的 inode 在外存中并且是长期存在的, VFS 的 inode 对象在内存中, 它仅在需要时才建立, 不再需要时撤消。物理文件系统的 inode 是静态的, 而 VFS 的 inode 是一种动态结构。
dd: 用指定大小的块拷贝一个文件, 并在拷贝的同时进行指定的转换
命令语法: dd [选项]
if = 输入文件 (或设备名称)
of = 输出文件 (或设备名称)
bs = bytes 同时设置读/写缓冲区的字节数 (等于设置 ibs 和 obs)
count=blocks 只拷贝输入的 blocks 块
conv = ucase 把字母由小写转换为大写
conv = lcase 把字母由大写转换为小写。
例: dd if=/dev/zero of=myfs bs=1M count=1 /dev/zero: 零设备 “0”
/dev/loop: loopback device (回环设备、或虚拟设备) 指是用文件来模拟块设备

习题

Banker algorithm

Suppose a system had 12 resources, 3 processes P0, P1 and P2.

	max. current. need	
P0	10	5
P1	4	2
P2	9	3

Currently there are 2 resources available. This system is in an unsafe state as process P1 could complete, thereby freeing a total of four resources. But we cannot guarantee that processes P0 and P2 can complete. However, it is possible that a process may release resources before requesting any further. For example, process P2 could release a resource, thereby increasing the total number of resources to five. This allows process P0 to complete, which would free a total of nine resources, thereby allowing process P2 to complete as well.

A system has 3 concurrent processes, each of which requires 4 items of resource R. What is the minimum number of resource R in order to avoid the deadlock. Answer: 3+3+(3+1)=10
The system design the structure File Control Block (FCB) to manage the files. Commonly, File control block is created on disk when the open system call is invoked.
Which kind of swap space is fastest: raw partition

2. 文件 F 由 200 记录组成, 记录从 1 开始编号, 用户打开文件后, 欲将其中的一条记录输入文件 F 中, 作为其第 30 条记录, 请回答下列问题, 并说明理由。

(1) 若文件系统为顺序存取方式, 每个存储块存放一条记录, 文件 F 的存储区域前后均有足够空闲的存储空间, 则完成上述操作最少要访问多少存储块? F 的文件控制区中会有哪些信息?

(2) 若文件系统为链接分配方式, 每个存储块存放的一条记录和一个链接指针, 则完成上述操作最少要访问多少存储块? 若每个存储块大小为 1KB, 其中 4 个字节存放指针, 则该系统支持文件的最大长度是多少?

【答案】(1) 因为最少访问, 所以选择将前 29 块都移一个存储块单元, 然后将要写入的记录放入到指定的第 30 条的位置上。由于每块都移一位而需存储块数增加, 再访问目标存储块将数据写入, 所以最少需要访问 29+1=30 块存储块。
F 的文件区的文件长度加 1, 起始地址加 1。
(2) 采用链接方式则需顺序访问前 29 块存储块, 然后将新记录的存储块插入链中即可, 把新的块存入磁盘要 1 次访存, 然后修改第 29 块的链接指针磁盘要又一次访存, 一共就是 29+1+1=31 次。
4 个字节的地址的地址范围为 2³²。
所以此系统能支持文件的最大长度为 2³²*(1KB-4)=4096GB

Q: 文件系统有一个 20MB 大文件和一个 20KB 小文件, 当分别采用连续、链接、链接索引、二级索引和 LINUX 分配方案时, 每块大小为 4KB, 每块地址用 4B 表示。问:

(1) 各文件系统管理的最大文件是多少?
(2) 每种方案对大、小两文件各需要多少专用块来记录文件的物理地址(说明各块的用途)?
(3) 如需要读大文件前面第 5.5KB 的信息和后面第 (16M+5.5KB) 的信息, 则每个方案各需要多少次盘 I/O 操作?

A: (1)
连续分配: 理论上是不受限制, 可大到整个磁盘文件区。
隐式链接: 由于块地址为 4B(4*8=32bit), 所以能表示的最多块数为 2³²=4G, 而每个盘中存放文件大小为 4096-4=4092B。链接分配可管理的最大文件为: 4G×4092B=16368GB
链接索引: 由于块的地址为 4B, 所以最多的链接索引块数为 2³²=4G, 而每个索引块有 1023 个文件块地址的指针, 盘块大小为 4KB。假设最多有 n 个索引块, 则 1023×n+n=2³², 算出 n=2²², 链接索引分配可管理的最大文件为: 4M*1023*4KB=16368GB
二级索引: 由于盘块大小为 4KB, 每个地址用 4B 表示, 一个盘块可存 1K 个索引表目。二级索引可管理的最大文件容量为 4KB×1K×1K=4GB。

LINUX 混合分配: LINUX 的直接地址指针有 12 个, 还有一个一级索引, 一个二级索引, 一个三级索引。因此可管理的最大文件为 48KB+4MB+4GB+4TB。
(2)
连续分配: 对大小两个文件都只需在文件控制块 FCB 中设二项, 一是首块物理块号, 另一是文件总块数, 不需专用块来记录文件的物理地址。
隐式链接: 对大小两个文件都只需在文件控制块 FCB 中设二项, 一是首块物理块号, 另一是末块物理块号; 同时在文件的每个物理块中设置存放下一个块号的指针。
一级索引: 对 20KB 小文件只有 5 个物理块大小, 所以只需一块专用物理块来作索引块, 用来保存文件的各个物理块地址。对于 20MB 大文件共有 5K 个物理块, 由于链接索引的每个索引块只能保存 (1K-1) 个文件物理块地址 (另有一个表目存放下一个索引块指针), 所以它需要 6 块专用物理块来作链接索引块, 用于保存文件各个的物理地址。
二级索引: 对大小文件都固定要用二级索引, 对 20KB 小文件, 用一个物理块作一级索引, 用另一块作二级索引, 共用二块专用物理块作索引块, 对于 20MB 大文件, 用一块作一级索引, 用 5 块作二级索引, 共用六块专用物理块作索引块。
LINUX 的混合分配: 对 20KB 小文件只需在文件控制块 FCB 的 i_addr[15]中使用前 5 个表目存放文件的物理块号, 不需专用物理块。对 20MB 大文件, FCB 的 i_addr[15]中使用前 12 个表目存放大文件前 12 块物理块块号 (48K),

用一级索引一块保存大文件接着的 1K 块块号 (4M), 剩下还有不到 16M, 还要用二级索引存大文件以后的块号, 二级索引使用一级索引引 1 块, 二级索引引 4 块 (因为 4KB×1K×4=16M)。总共也需要 6 块专用物理块来存放文件物理地址。

(3)
连续分配: 1. 计算信息在文件中相对块数, 前面信息相对逻辑块号为 5.5K / 4K=1 (从 0 开始编号), 后面信息相对逻辑块号为 (16M+5.5K) /4K=4097。再计算物理块号=文件首块号+相对逻辑块号, 最后化一次盘 I/O 操作读出该块信息。

链接分配: 为读大文件前面 5.5KB 的信息, 只需先读一次文件头块得到信息所在块的块号, 再读一次第 1 号逻辑块得到所需信息, 共 2 次。而读大文件 16MB+5.5KB 处的信息, 逻辑块号为 (16M+5.5K) /4092=4107, 要先把该信息所在块前面块顺序读出, 共化费 4107 次盘 I/O 操作, 才能得到信息所在块的块号, 最后化一次 I/O 操作读出该块信息。所以总共需要 4108 次盘 I/O 才能读取 (16MB+5.5KB) 处信息。
链接索引: 为读大文件前面 5.5KB 处的信息, 只需先读 1 次第一个索引块得到信息所在块的块号, 再读一次第 1 号逻辑块得到所需信息, 共化费 2 次盘 I/O 操作。为读大文件后面 16MB+5.5KB 处的信息, (16MB+5.5KB)/(4KB×1023)=4, 需要先化 5 次盘 I/O 操作依次读出各索引块, 才能得到信息所在块的块号, 再化一次盘 I/O 操作读出该块信息。共化费 6 次盘 I/O 操作。

二级索引: 为读大文件前面和后面信息的操作相同, 首先进行一次盘 I/O 读第一级索引块, 然后根据它的相对逻辑块号计算应该读第二级索引的那块, 第一级索引块表目号=相对逻辑块号 / 1K, 对文件前面信息 1 / 1K=0, 对文件后面信息 4097 / 1K=4, 第二次根据第一级索引块的相应表目内容又化一次盘 I/O 读第二级索引块, 得到信息所在块块号, 再化一次盘 I/O 读出信息所在盘块, 这样读取大文件前面或后面处信息都只需要 3 次盘 I/O 操作。

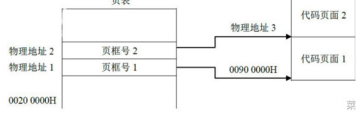
LINUX 混合分配: 为读大文件前面 5.5KB 处信息, 先根据它的相对逻辑块号, 在内存文件控制块 FCB 的 i_addr 第二个表目中读取信息所在块块号, 而只化费一次盘 I/O 操作即可读出该块信息。为读大文件后在 (16MB+5.5KB) 信息, 先根据它的相对逻辑块号判断要读的信息是在二级索引管理范围内, 先根据 i_addr 内容化一次盘 I/O 操作读出第一级索引块, 再计算信息所在块的索引块号在第一级索引块的表目号为 (4097-12-1024) / 1024=2, 根据第一级索引块第 3 个表目内容再化费一次盘 I/O 操作, 读出第二级索引块, 就可以得到信息所在块块号, 最后化一次盘 I/O 读出信息所在盘块, 这样总共需要 3 次盘 I/O 操作才能读出文件后面的信息。
inverted page table
Assume we have a simple, demand paging environment, with no segmentation. Processes P1 and P2 both have logical memory addresses in the range 0...99, inclusive. **page size is 5.** hash table portion of the structure uses a hash function which calculates **index** by adding numerical portion of **process identifier and logical page number**. hash table and IPT shown are below.

Hash Table			Inverted Page Table		
Index	Process ID / logical page	Frame ID	Index	Logical address (minimized for sim)	
...		
3		0			
4	P2 2	1	1	P2:RW	
5	P1 3	4	2	P1:RW	
...		
18	P1 17	2	3	P1:R	
19	P1 18	3	4	P1:RW	

(1) Calculate **physical address** for the following logical addresses. Assume first page of a process is page number 0, and that a mechanism exists to find a free frame in memory.
P1:17 Logical page=17/5=3,offset=2, hash table index = 3 + 1(pid)=4,P1 本来 hash table 里对应 frameID=1, 但是 **IPT 里面 frameID=1 映射到的是进程 P2 而非 P1**, 所以 P1 在 hash table 里面映射后移 1,即 hash index=4+1==> **follow chain to entry 5**. From hash table, physic frame=4. Address = (frame)*(frame size) + offset = 4*5 + 2 = 22
P1:92 Logical page=92/5=18,offset=2. Hash table index = 19. From lookup logical frame = 3. Address = 17
(2) When process P2 attempts to read logical address 97, what happens? This causes page fault. Frame 0 is allocated, and the entry in IPT is updated with logical address information for P2 frame 19. Next the hash table entry at index 21 (19 + 2) is mapped from P2:19 to frame 0. After this processing is complete, the P2 resumes execution as before.

进程同步与死锁
4 processes share a buffer space named B: R1 reads a number from the keyboard and saves it into B, so that the saved number is only consumed by W1, which prints the number to the screen. R2 reads a character from a mouse and also saves character into B, so that only W2 can print character to a printer. Please write the synchronization code for the four processes so that no race condition may arise among them.
Semophore S=1(B 为空, 初始为 1), S1=0(R1 做完, W1 可以做), S2=0(R2 做完, W2 可以做)

计数信号量 S
若干进程对 S 进行 28 次 P 操作和 18 次 V 操作后, 信号量 S 的值为 0(此时无资源可用), 然后又对信号量 S 进行了 3 次 V 操作,此时有 0 个进程等待在信号量 S 的队列中
FAT 计算
FAT 的 entrySize=2B, 则 FAT 最大 size=2¹⁶*2 B(2B=16bit->2¹⁶)
Waiting time/turnaround time cal.
Turn around time=waiting time+burst time
Address computing in page table(VPN→PPN)
PA,VA 均为 32bit,pageSize=4KB,pageTableEntrySize=4B.代码段起始 VA=0x8000,长度=8KB,被装载到从 PA=0x900000 开始的**连续**主存空间。页表从主存 0x200000 开始的物理地址处连续存放



如图所示,代码页面 2 的起始 PA=0x900000+8K。代码页面 1 对应 VA1=0x8000->VPN1=8->物理

地址 1=0x200000+VPN1*pageTableEntrySize=0x200020,同时 PA1=0x900000->PPN1=0x900

10个10上10次上下文交换=>10*1=10*0.1 ms => 10.1-10.1 ms
1个CPU算一次上下文交换=>0.1 ms
10.1-10.1 ms
according to question given at every 1ms i/o operation done, and it will be completed in 10 ms. so 10 i/o operation take
time = 10*1+1*10 = 11ms...[here 0.1 switching overhead]
cpu task, for which we can consider it runs for 10ms bcoz i/o operation will be completed in 10ms...[given as all processes are long running tasks]
so if case 1... time quantum is 1ms then cpu task takes = 10*1+1*0.1[switching after every 1ms quantum time] = 11ms
so cpu utilization = useful work / total work
(10*1+10*1)/(1+1+1+20/22)=90.90% **最佳的情况**
case 2 time quantum =10ms so cpu takes = 10*1+1*1+10.1
cpu utilization = 20/1+10.1+20/21.1 = 94.7%

```
pid = fork();
if (pid == 0)
{
    /* childprocess */
    fork();
    thread_create(. . .);
}
fork();
```

a. How many unique processes are created? _6_ (包括第一次运行该程序的进程)
第一次运行该程序的进程 P0 fork() 创建一个 sub process P1, 没有创建 thread P1 在 pid=0 的代码块里面创建 1 个 sub process P2, P1P2 一共创建 2 个 thread POP1P2 在外面的 fork 各自创建一个 sub process 分别是 P3,P4,P5
b. How many unique threads are created? _2_ (没有主线程)