

概述

- 计算机的四个组成部分：处理器 内存 输入 / 输出模块 系统总线
- 定义处理器寄存器的两种类别：用户可见寄存器控制和状态寄存器
- 多个中断的处理方式：
 - 正在处理一个中断时，禁止再发生中断
 - 定义中断优先级
 - 内存层次各个元素间的特征：价格、容量和访问时间
- 两种局部性区别：
 - 空间局部性：最近被访问的元素的周围将来可能会被访问；利用更大的缓冲块并且在存储器控制逻辑中加入预处理机制
 - 时间局部性：近被访问的元素将来可能会被再次访问；在高速缓冲存储器中保留最近使用的指令及数据
- ## 操作系统概述
- 系统设计的三个目标：方便 有效 扩展的能力
- 进程的组成：可执行的程序 需要的相关数据 执行上下文
- 地址的区别：
 - 实地址：指的是主存中的地址，实际的主存储器的地址，即物理空间
 - 虚地址：指的是存在于虚拟内存中的地址，可能在磁盘或者主存中
 - 时间片轮流调度技术：所有的进程存放在一个环形队列中并按固定循序依次激活
 - 对称多处理器操作系统设计的关键问题：并发调度 同步 内存管理 可靠性和容错性
 - API 侧重于向上层提供服务，而系统调用则侧重于通过软中断向下层的内核发出一个明确的请求
 - 中断与当前正在运行的进程无关的某些类型的外部事件相关，陷阱与当前正在运行的进程所产生的错误或异常条件相关
- ## 进程
- 创建进程的事件：新的批处理作业 交互登录 提供一项服务而创建 由现有的进程派生
- 进程控制块中的三类信息：进程标识 处理器状态信息 进程控制信息
- 程模型中的状态：新建态 就绪态 运行态 阻塞态 挂起态 退出态
- 进程的几个定义：正在执行的程序 程序实例 由处理器执行的实体 执行指令、当前状态和相关的系统资源表征的活动单元
- 创建一个新进程所执行的步骤：
 - 分配一个唯一的进程标识号
 - 初始化进程控制块
 - 设置正确的连接
 - 创建或扩充其他的数据结构
- 进程映像：程序 + 数据 + 栈 + 属性
- 进程标识信息：存储在 PCB 中的数字标识符，包括：进程 ID、父进程 ID、用户 ID
- 把操作系统作为一个系统进程来实现的优点是：
 - 鼓励模块化操作系统设计原理，使模块间接口最小且最简单
 - 非关键系统功能可简单的用独立的进程来实现
- 在多处理器和多机环境中很有用
- 用户态到内核态切换的三种方式：系统调用 异常 外设中断
- ## 线程
- 线程间的模式切换比进程间的模式切换开销更低的原因：包含的状态信息更少
- 单用户多处理器上使用线程场景：前台和后台操作 异步处理 加速执行 模块化程序结构
- 一个进程中的所有线程共享的资源：地址空间 文件资源 执行特权
- 用户级线程优于内核级线程的优点：
 - 线程切换不需要内核级模式的特权
 - 调用可以是应用程序专用的
 - 可以在任何操作系统中运行
- 用户级线程较之内核级线程的缺点：
 - 许多系统调用都会引起阻塞
 - 不能利用多处理器技术
- 线程的优点：

- 创建和终止的时间少于进程
- 在内线程间切换的时间少于进程
- 提高了不同执行程序间的通信效率
- 基本操作：派生、阻塞、解除阻塞、结束
- 线程同步：同步线程的活动是它们互不干扰且不破坏数据结构，如两个线程向一个链表加入元素，则可能会丢失
- ## 并发
- 程间的三种互相知道的程度：
 - 相互不知道对方的存在 是竞争关系
 - 间接知道对方的存在 合作行为
 - 直接知道对方的存在 合作行为
- 三个控制问题：
 - 互斥 访问一个不可共享的资源(临界资源)
 - 死锁 每个进程都在等待另一个被对方占用的资源
 - 饥饿 一个进程被无限地拒绝访问资源
- 管程的特点：
 - 外部过程不能访问局部数据变量
 - 一个进程通过调用管程的一个过程进入管程
 - 只能有一个进程在管程中执行,任何其他进程都被阻塞
- 处理死锁的三种方法：预防、检测和避免
- 并发的原理：
 - 全局资源的共享充满了危险
 - OS 很难对资源进行最优化分配
 - 定位程序设计错误非常困难
- ## 死锁
- 产生死锁的必要条件：互斥 占有且等待 不可抢占 循环等待(前三者的结果)
- 死锁避免的优点：无须死锁预防的抢占和回滚进程，且与死锁预防相比限制较少。
- 死锁避免的限制：
 - 必须声明实现每个进程请求的最大资源
 - 所讨论的进程必须是无关的，即它们的执行顺序必须没有同步要求的限制
 - 分配的资源数量必须是固定的
 - 在占有资源时，进程不能退出
 - 死锁检测： 不限制资源访问或约束进程行为，只要有可能就会给进程分配其所需要的资源，操作系统周期性的执行一个算法来检测前面的条件
- ## 内存管理
- 内存管理满足的需求： 重定位 保护 共享 逻辑组织 物理组织
- 页表中能找到的典型元素：页号 进程标志符 控制位 链指针
- 固定分区方案简单但有缺点：
 - 限制了系统中活动进程的数量
 - 小作业不能有效地利用分区空间
- 伙伴系统：二分法确定分配区的大小，一定是整个内存空间的二的幂次分之一
- ## 虚拟内存
- 进程执行过程中任何时刻都在内存中的部分称为进程的常驻集，当访问一个不再内存中的逻辑地址，会产生一个中断，操作系统会将此进程置于 阻塞态 ，为此操作系统产生一个 磁盘 I/O 请求。
- 虚拟内存的开销受到系统抖动的影响，会导致处理器大部分的时间都用于交换块而非执行指令。
- ## 调度
- 三种类型的处理器调度： 长程调度 中程调度 短程调度
- 最重要的性能要求是：响应时间
- 周转时间和响应时间的区别：
 - 周转时间：一个进程从提交到完成
 - 响应时间：提交一个请求到开始接收响应
- 简单定义调度策略：
 - FCFS 调度：进程就绪后加入就绪队列，当前正运行的进程停止执行时，选择就绪队列中存在时间最长的进程运行。
 - 轮转调度：周期性地产 生时钟中断，正运行的进程会放置到就绪队列中，基于 FCFS 策略选择下一个
 - 最短进程优先调度：下次选择预计处理时间最短的进程(非抢占策略)
 - 最短剩余时间调度：总是选择预期剩余时间最短的进程

- 最高响应优先调度：当前进程完成或被阻塞时，选择 R 值最小的就绪进程。基于对归一化周转时间的估计。
- 反馈调度：基于抢占原则并使用动态优先级机制，基于每个进程的 执行历史 把它们分配到各个队列中
- ## 文件管理
- 文件组织的选择原则：快速修改 易于修改 节约存储空间 维护简单 可靠性
- 基本文件组织：
 - 堆：是简单的文件组织形式，数据按它们到达的顺序被收集
 - 顺序文件：每条记录都使用固定的格式，所有记录具有相同的长度，并由相同数量、长度固定的块按特定顺序组成。每条记录的第一个字节都是关键字
 - 索引顺序文件：用于支持随机访问的文件索引和溢出文件。索引顺序文件极大的减少了访问单条记录的时间
 - 记录组块：定长组块 变长跨越式组块 变长非跨越式组块

- 1.1 列出并简要定义计算机的四个组成部分。处理器：控制计算机的操作，执行数据处理功能。内存：也叫主存储器，存储数据和程序。输入 / 输出模块：在计算机和外部环境之间移动数据。系统总线：在处理器、内存和输入输出间提供通信的设施。
- 1.2 定义处理器寄存器的两种主要类别。用户可见寄存器：优先使用这些寄存器，可以使机器语言或者汇编语言的程序员减少对主存储器的访问次数。对高级语言而言，由汇编编译器负责决定哪些变量应该分配给主存储器，一些高级语言，如 C 语言，允许程序员建议编译器把哪些变量保存在寄存器中。
- 控制和状态寄存器：用以控制处理器的操作，且主要被具有特权的操作系统例程使用，以控制程序的执行。

- 1.3 一般而言，一条机器指令能指定的四种不同操作是什么？处理器 - 寄存器：数据可以从处理器传送到存储器，或者从存储器传送到处理器。处理器 - I/O：通过处理器和 I/O 模块间的数据传送，数据可以输出到外部设备，或者从外部设备输入数据。数据处理：处理器可以执行很多关于数据的算术操作或者逻辑操作。控制：某些指令可以改变执行顺序。
- 1.4 什么是中断？中断是指计算机运行过程中，出现某些意外情况需中断其正在执行的程序并转入正在运行的程序并转入处理新情况的程序，处理完后又返回原被暂停的程序继续运行。
- 1.5 多个中断的处理方式是什么？处理多中断有两种方法。第一种方法是当正在处理一个中断时，禁止再发生中断。第二种方法是定义中断优先级，允许高优先级的中断打断低优先级的中断处理器的运行。
- 1.6 内存层次各个元素间的特征是什么？存储器的三个重要特性是：价格、容量和访问时间。并且各层次从上到下，每位“价格降低，容量递增，访问时间递增。
- 1.7 什么是高速缓存？高速缓冲存储器是比主存小而快的存储器，用以协助主存跟处理器。作为最近存储地址的缓冲区。
- 1.8 多处理器系统和多核系统的区别是什么？多处理器系统 (Multiprocessor Systems) 是指包含两台或多台功能相近的处理器，处理器之间彼此可以交换数据，所有处理器共享内存、I/O 设备、控制器、及外部设备，整个硬件系统由统一 的操作系统控制，在处理器和程序之间实现作业、任务、数据、数级极其元素等级的全面并行。多内核 (multicore chips) 是指在一枚处理器 (chip) 中集成两个或多个完整的计算引擎 (内核)。

- 1.9 空间局部性和时间局部性的区别是什么？空间局部性是指最近被访问的元素的周围在不久的将来可能会被访问。临时局部性 (即时间局部性)：是指最近被访问的元素在不久的将来可能会被再次访问。1.10 开发空间局部性和时间局部性的策略是什么？空间局部性的开发是利用更大的缓冲块并且在存储器控制逻辑中加入预处理机制。时间局部性的开发是利用在高速缓冲存储器中保留最近使用的指令及数据，并且定义缓冲存储的优先级。
- 第二章 操作系统概述 §1 2.1 操作系统设计的三个目标是什么？方便：操作系统使计算机更易于使用。有效：操作系统允许以更有效的方式使用计算机系统资源。扩展的能力：在构造操作系统时，应该允许在不妨碍服务的前提下有效地开发、测试和引进新的系统功能。
- 2.2 什么是操作系统的内核？操作系统内核是计算机上最低层的软件，提供计算机最核心的功能，比如：进程管理、内存管理、I/O 管理、文件管理、网络管理等。
- 2.3 什么是多道程序设计？两个或两个以上程序在计算机系统中同处于开始到结束之间的状态，这就称为多道程序设计。也就是在计算机内存中同时存放几道相互独立的程序，使它们在管理程序控制之下，相互穿插的运行。多道程序设计运行的特征：多道、宏观上并行、微观上串行。
- 2.4 什么是进程？进程由三部分组成：一段可执行的程序。程序所需要的相关数据 (变量、工作空间、缓冲区等)。程序的执行上下文 (也称进程状态)。
- 2.5 操作系统是怎么实现进程上下文的？执行上下文又称为进程状态，是操作系统用来管理和控制所需的内部数据。这种内部信息和进程是分开的，因为操作系统信息不允许被进程直接访问。上下文包括操作系统管理进程以及处理器正确执行进程所需要的所有信息，包括各种处理器寄存器的内容，如程序计数器 and 数据寄存器。它还包含系统使用的信息，如进程优先级以及进程是否在等待特定 I/O 事件的发生。
- 2.6 列出并简要介绍操作系统的五种典型存储管理职责。进程隔离：操作系统必须保护独立的进程，防止互相干扰各自的存储空间，包括数据和指令。自动分配和管理：程序应该根据需要在存储空间间动态的分配，分配对程序员是透明的。因此，程序员无需关心与存储限制有关的问题。操作系统有效的实现分配问题，可以仅在需要时才给作业分配存储空间。支持模块化程序设计：程序员应该能够定义程序模块，并动态地创建、销毁模块，动态地改变模块大小。
- 保护和访问控制：不论在存储层次中的哪一级，存储器的共享都会产生一个程序访问另一个程序存储空间的潜在可能性。当某个特定的应用程序需要共享时，这是可取的。但在其他时候，他可能会威胁到程序的完整性，甚至威胁到操作系统本身。长期存储：许多应用程序需要在计算机关机后长时间的保存信息。
- 2.7 解释实地址和虚地址的区别。实地址：指的是主存中的区域，实际的主存储器的地址。指向主存空间，即物理空间。虚地址：指的是存在于虚拟内存中的地址，它有时候在磁盘上，有时候在主存中。
- 2.8 描述时间片轮流调度技术。轮流调度是一种固定算法，所有的进程存放在一个环形队列中并按顺序依次激活。因为等待一些事件 (例如：等待一个子进程或者一个 I/O 操作) 的发生而不能被处理的进程将控制权交给调度器。
- 2.9 解释单体内核和微内核的区别。单体内核是一个提供操作系统应该提供的功能的大内核，包括调度、文件系统、网络、设备驱动程序、存储管理等。内核的所有功能成分都能够访问

它的内部数据结构和程序。典型情况下，这个大内核是作为一个进程实现的，所有元素都共享相同的地址空间。微内核是一个小的有特权的操作系统内核，只提供包括进程调度、内存管理和进程间通信等基本功能，要依靠其他进程担任起和操作系统内核联系作用。2.10 什么是多线程？多线程技术是指把执行一个应用程序的进程划分为可以同时运行的多个线程。2.11 列出对称多处理器操作系统设计时要考虑的关键问题。并发 (concurrent) 进程或线程：内核程序可重入，以使多个处理器能同时执行同一段内核代码。当多个处理器执行内核的相同或不同部分时，为避免数据损坏和无效操作，需要妥善管理内核和数据结构。

调度：任何一个处理器都可以执行调度，这既增加了执行调度策略的复杂度，也增加了保证调度相关数据结构不被损坏的复杂度。如果使用的是内核级多线程方式，就存在将同一进程的多个线程同时调度在多个处理器上的可能性。同步：因为可能会存在多个活跃进程须共享地址空间或共享 I/O 资源的情况，因此必须认真考虑如何提供有效的同步机制这一问题。同步用来实现互斥及事件排序。内存管理：多处理器上的内存管理要处理单处理器上内存管理涉及的所有问题。另外，操作系统还要充分利用硬件提供的并行性来实现最优性能。不同处理器上的分页机制必须进行调整，以实现多处理器共享页或段时的数据一致性，执行页面置换。物理页的重用是我们关注的最大问题，即必须保证物理页在重新使用前不能访问到它以前的内容。可靠性 and 容错性：出现处理器故障时，操作系统应能妥善降低故障的影响。调度器和操作系统的其他部分必须能识别出发生故障的处理器，并重新组织管理。第二部分 进程 §1 第三章 进程描述和控制 §1 3.1 什么是指令跟踪？指令跟踪是指为该进程而执行的指令序列。3.2 通常哪些事件会导致创建一个进程？新的批处理作业；交互登录；操作系统因为提供一项服务而创建；由现有的进程派生。3.3 简要定义图 3.6 所示进程模型中的每种状态。运行态：该进程正在执行。就绪态：进程做好了准备，只要有机会就会开始执行。阻塞态：进程在某些事件发生前不能执行，如 I/O 操作完成。新建态：刚刚创建的进程，操作系统还没有把它加入到可执行进程组中。退出态：操作系统从可执行进程中释放出该进程，或者是因为它自身停止了，或者是因为某种原因被取消。3.4 被抢占一个进程是什么意思？处理器为了执行另外的进程而终止当前正在执行的进程，这就叫进程抢占。3.5 什么是交换，某目的是什么？交换是指把主存中某个进程的一部分或者全部内容转移到磁盘。当主存中没有处于就绪态的进程时，操作系统就把一个阻塞的进程换出到磁盘中的挂起队列，从而使另一个进程可以进入主存执行。3.6 为什么图 3.9 (b) 中有两个阻塞态？有两个独立的概念：进程是否在等待一个事件 (阻塞与否) 以及进程是否已经被换出主存 (挂起与否)。为适应这种 2+2 的组合，需要两个阻塞态和两个挂起态。3.7 列出挂起状态进程的 4 个特点。进程可能立即执行。进程可能是或不是正在等待一个事件。如果是，阻塞条件不依赖于挂起条件，阻塞事件的发生不会使进程立即被执行。3.8 解释单体内核和微内核的区别。单体内核是一个提供操作系统应该提供的功能的大内核，包括调度、文件系统、网络、设备驱动程序、存储管理等。内核的所有功能成分都能够访问

除非非代理式地命令系统进行状态转换，否则进程无法从这个状态中转移。3.8 对于哪类实体，操作系统为了管理它而维护其信息？内存、I/O、文件和 进程。3.9 列出进程控制块中的三类信息。进程控制信息：处理器状态信息，进程控制信息。3.10 为什么需要两种模式 (用户模式和内核模式)？用户模式下可以执行的指令和访问的内存区域都受到限制。这是为了防止操作系统受到破坏或者修改。而在内核模式下则没有这些限制，从而使它能够完成其功能。3.11 操作系统创建一个新进程所执行的步骤是什么？给新进程分配一个唯一的进程标识号。给进程分配空间。初始化进程控制块。设置正确的连接。创建或扩充其他的数据结构。3.12 中断和陷阱有什么区别？在程序运行过程中，系统出现了一个必须由 CPU 立即处理的情况，此时，CPU 暂时中止程序的执行转而处理这个新的情况的过程就叫做中断。中断与当前正在运行的进程无关的某些类型的外部事件相关，如完成一次 I/O 操作。陷阱指的是异常或者中断发生时，处理器捕捉到一个执行地址，并且将控制权转移到操作系统中某一个固定地址的机制。陷阱与当前正在运行的进程所产生的错误或异常条件相关，如非法的文件访问。3.13 举出中断的三个例子。时钟中断、I/O 中断、内存失效。3.14 模式切换和进程切换有什么区别？发生模式切换可以不改变当前正处于运行态的进程的状态。进程切换指定另一个进程为运行态，进程切换需要保存更多的状态信息。第四章 线程 §1 4.1 表 3.5 列出了在一个没有线程的操作系统中进程控制块的基本元素。对于多线程系统，这些元素中哪些可能属于线程控制块，哪些可能属于进程控制块？这相对于不同的系统来说通常是不同的，但一般来说，进程是资源的所有者，而每个线程都有它自己的执行状态。关于表 3.5 中的每一项的一些结论如下：进程标识：进程必须被标识，而进程中的每一个线程也必须有自己的 ID。进程控制信息：调度和状态信息主要处于线程级；数据结构在两级都可出现，进程间通信和线程间通信都可得到支持；特权在两级都可以存在；存储管理通常在进程级；资源信息通常也在进程级；处理器状态信息：又名现场信息，一个进程在运行时存放在处理器现场中的各种信息。这些信息通常只与进程有关。进程控制信息，用于管理和调度一个进程。常用的控制信息包括：进程的调度相关信息，如进程状态、等待事件和等待原因、进程优先级、队列指针元等 进程组成信息，如正文段指针、数据段指针 引进程间通信相关信息，如消息队列指针、信号量等互斥和同步机制 进程在缓存存储器内的地址 CPU 资源的占用和使用信息，如时间片余量、进程已占用 CPU 的时间、进程已执行的时间总和、记账信息 进程特权信息，如在内存访问和处理器状态方面的特权 4.2 请列出线程间的模式切换比进程间的模式切换开销更低的原因。包含的状态信息更少。4.3 在进程概念中体现出的两个独立且无关的特点是什么？资源所有权：进程包括存放进程映像的虚拟地址空间；回顾第 3 章的内容可知，进程映像是程序、数据、栈和控制控制块中定义的属性集。进程总具有挂起态，代理可以是进程自己，也可以是父进程或对资源的控制权或所有权，这些资源包括内存、I/O 设备、I/O 设备和文件等。操作系统提供预防进程

