

## 2.2\_1\_处理机调度的概念和层次

- gxy总结:

需要掌握调度的概念。

理解三种各自的应用场景。

在理解应用场景的基础之上，结合背后的场景对各自的次数进行比较。

**调度的概念：**当有多个任务需要进行处理的时候，由于资源有限，所以有一些任务不能同时处理，需要根据**某种规则**来决定处理任务的**顺序**。

处理机调度：从就绪队列中按照一定的算法选择一个进程并将处理机分配给它运行，以实现进程的并发执行。

- 高级调度作业调度

作业的概念：是一个具体的任务。

进行调度的背景：启动一个程序，相关数据一定要从外存放到内存，但是内存有限，如果内存满了，不能直接放进去，会进行**高级调度**。

**高级调度做的事情：**OS会按照一定的规则从作业后备队列中选择一个作业调入内存，并且创建进程。高级调度对每一个作业，开始的时候调入一次，结束的时候调出一次。

- 低级调度处理机调度/进程调度

处理机调度的背景：有时候要处理很多进程，但是cpu资源有限，所以需要调度，从进程就绪队列中选一个进程分配CPU。

低级调度做的事情：按照某一种策略从**就绪队列**中取出一个进程，将处理机分配给它。

处理机调度是一件会经常发生的事情，频率很高。

- 中级调度内存调度

产生的背景：在内存不够的时候，OS可以将某些进程的数据调出到外存。等到内存空闲或者进程需要运行的时候再重新调入内存。

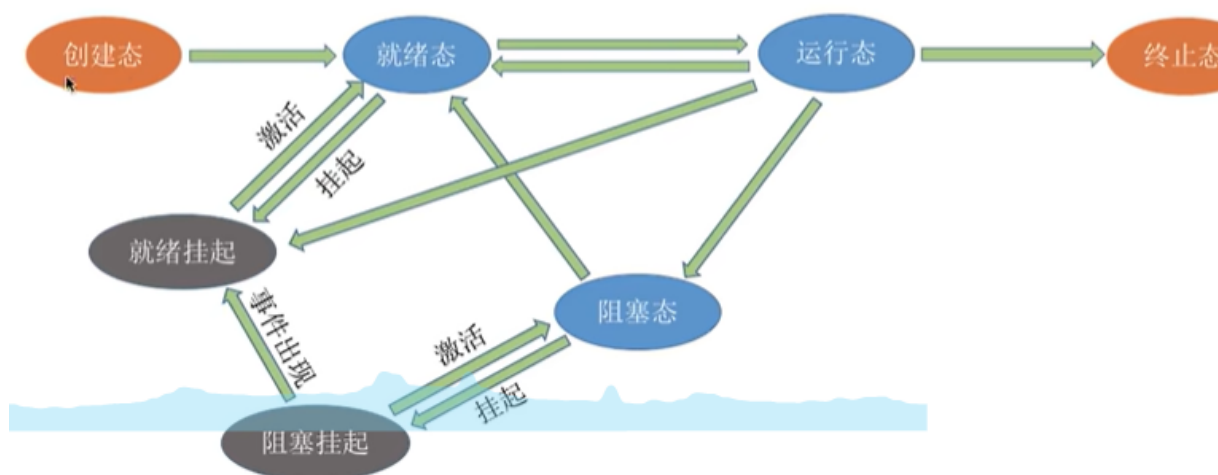
暂时调到外存等待的进程状态时**挂起状态**，被挂起的进程PCB会组织成**挂起队列**。

中级调度做的事情：按照某种策略决定将哪个处于挂起状态的进程重新调入内存。一个进程可能会被多次调出调入内存，中级调度的次数可能比高级调度的次数多。

- 补充：挂起状态和七状态模型。

如果内存不够用，可能把一个处于就绪态的进程**挂起**成为**就绪挂起**的状态，阻塞状态的进程也可以被挂起成为**阻塞挂起**的状态。

暂时调到外存等待的进程状态为挂起状态（挂起态，suspend）  
挂起态又可以进一步细分为就绪挂起、阻塞挂起两种状态  
五状态模型 → 七状态模型



注意“挂起”和“阻塞”的区别，两种状态都是暂时不能获得CPU的服务，但挂起态是将进程映像调到外存去了，而阻塞态下进程映像还在内存中。

解释：

阻塞，就是当前没有得到某一种资源，然后继续在这里等，其他资源还在内存里面。  
挂起，我的进程映像直接被调到了外存。

- 三种调度的对比：



### 三层调度的联系、对比

	要做什么	调度发生在..	发生频率	对进程状态的影响
高级调度 (作业调度)	按照某种规则，从后备队列中选择合适的作业将其调入内存，并为其创建进程	外存→内存 (面向作业)	最低	无→创建态→就绪态
中级调度 (内存调度)	按照某种规则，从挂起队列中选择合适的进程将其数据调回内存	外存→内存 (面向进程)	中等	挂起态→就绪态 (阻塞挂起→阻塞态)
低级调度 (进程调度)	按照某种规则，从就绪队列中选择一个进程为其分配处理机	内存→CPU	最高	就绪态→运行态

- gxy总结:

官方说这一节不是很重要。但时我觉得临界区和广义、狭义上的进程调度都是可以帮助理解的知识点。

其实一句话总结，就是如果进程访问了 **内核程序临界区**，就会导致不能进行进程调度，因为进程调度要访问进程的就绪队列，而当前进程访问了 **内核程序临界区**，就把这里上锁了。

- 进程调度的时机，什么时候需要进程调度？

- 当前运行的进程主动放弃处理机：

- 进程的正常终止

- 进程运行中发生异常而终止

- 进程主动阻塞（比如等待打印机）

- 当前运行的进程被动放弃处理机：

- 进程分到的时间片用完、有更紧急的比如中断的事情需要处理、有更高优先级的进程进入就绪队列。

- 不能进行进程调度的情况：

- 1.处理中断的时候。因为中断比较复杂。

- 2.进程在 **os内核程序临界区** 中。

- 3.原子操作中，原语中。

- 理解：进程在**OS**内核程序临界区中不能进行调度和转换。

- 进程在临界区时可以进行处理机调度。

前置知识：

**临界资源**：只允许同一时刻一个进程访问的资源，各个进程需要 **互斥** 地访问临界资源。

**临界区**：访问临界资源的代码。

因此各个进程也只能互斥的进入临界区。

什么是 **内核程序临界区**？一般用来访问某种 **内核数据结构**，比如进程的就绪队列。

当一个进程处于 **内核程序临界区**，并且需要访问进程的就绪队列的话，在访问之前会把就绪队列锁上，如果这个进程不退出这个**内核程序临界区**，这个就绪队列不允许其他进程访问。就绪队列很明显和 **进程调度** 有关系，所以这个时候会导致不能进行进程调度。从另外一个方面来看：如果一个进程访问了 **内核程序临界区**，特别是进程就绪队列，那么应该尽快完成自己的工作并释放，不然会影响正常计算机工作。

如果进程访问的是普通的临界资源：比如打印机。

访问的时候，会先上锁，在进程结束之前，打印机一直上锁，但是打印机比较慢，这个时候CPU会空闲。所以这个时候应该进行 **进程调度**。

其实最关键的是，访问的普通临界区和进程调度没有关系，所以也不会影响进程调度。

- 进程剥夺的方式:



**非剥夺调度方式**，又称**非抢占方式**。即，只允许进程主动放弃处理机。在运行过程中即便有更紧迫的任务到达，当前进程依然会继续使用处理机，直到该进程终止或主动要求进入阻塞态。



**剥夺调度方式**，又称**抢占方式**。当一个进程正在处理机上执行时，如果有一个更重要或更紧迫的进程需要使用处理机，则立即暂停正在执行的进程，将处理机分配给更重要紧迫的那个进程。

非剥夺：简单、开销小但是不能处理紧急任务，适合早期批处理系统。

剥夺：可以优先处理紧急的进程，也可以实现时间片轮转。适合分时OS。

- 进程的切换与过程。

狭义的进程调度，只是指选择一个进程要开始运行，要分配处理机器。

广义的进程调度，包含之前的一个进程让出来处理机，之后的进程占用处理机。

广义的进程调度包括 **选择一个进程** 和 **进程切换** 两个步骤。

进程切换的时候做了什么：

- 1.对原来的运行进程中的数据进行保存。
- 2.对新的进程的数据进行恢复。

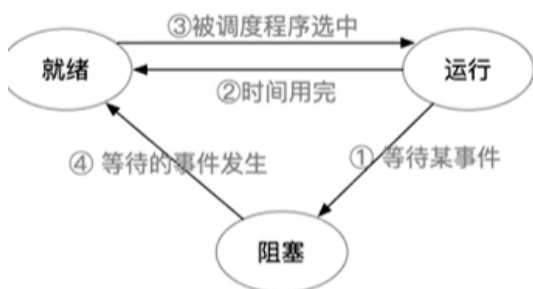
上面所说的数据 是一些运行信息，比如PC,PSW,各种数据寄存器等等。

- 进程的切换不是越多越好:

首先，进程的切换有代价，如果非常频繁的一直切换进程，会使得OS大量时间都用在切换进程上，真正用于执行进程的时间减少。

## 2.2\_3\_调度器和闲逛进程

gxy总结：感觉就是把之前的有一些点又说了一遍。



②、③由调度程序引起，调度程序决定：

让谁运行？——调度算法  
运行多长时间？——时间片大小

调度时机——什么事件会触发“调度程序”？

- 创建新进程
- 进程退出
- 运行进程阻塞
- I/O中断发生（可能唤醒某些阻塞进程）
- 非抢占式调度策略，只有运行进程阻塞或退出才触发调度程序工作
- 抢占式调度策略，每个时钟中断或k个时钟中断会触发调度程序工作

- 不支持内核级线程的OS，调度程序的处理对象是进程。  
支持内核级线程的OS，调度程序的处理对象是线程。
- 闲逛进程：  
CPU如果没有就绪进程，会选择让闲逛进程上CPU。优先级最低，能耗最低。啥也不干。0地址指令，不访问寄存器，同时也会指令周期末尾检查中断。

## 2.2\_4\_调度算法的评价指标

- gxy总结：cpu利用率，系统吞吐量，周转时间，带权周转时间，等待时间，响应时间。  
各自掌握定义即可。

有一个题目[链接](#)

- CPU利用率：  
CPU忙碌的时间占总时间的比率。利用率 =  $\frac{\text{忙碌的时间}}{\text{总时间}}$
- 系统吞吐量：单位时间内完成作业的数量。  
系统吞吐量 =  $\frac{\text{总共完成多少道作业}}{\text{总共花的时间}}$
- 周转时间：

**定义**：周转时间是指从作业被提交给系统开始，到作业完成为止的时间间隔。

它包括四个部分：作业在外存后备队列上等待作业调度（高级调度）的时间、进程在就绪队列上等待进程调度（低级调度）的时间、进程在CPU上执行的时间、进程等待I/O操作完成的时间。后三项在一个作业的整体处理过程中，可能发生多次。

平均周转时间 = 所有作业周转时间和 / 作业数目



- 带权周转时间：

是一个比率吗？之后细说

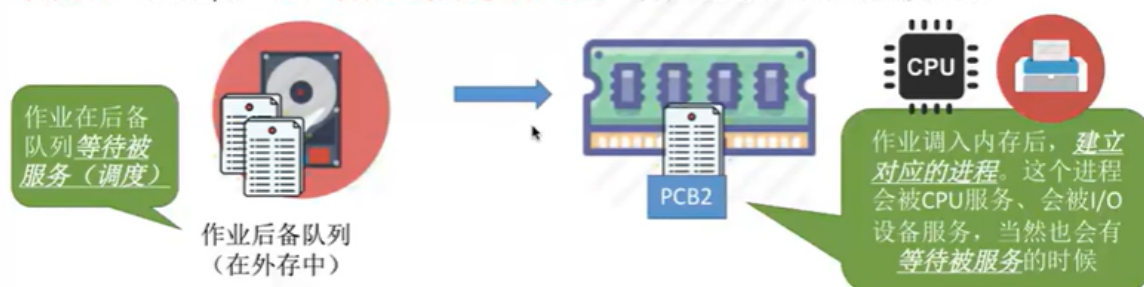
说明：带权周转时间一定 $\geq 1$ ，分子是大的，而且我们希望带权周转时间越小越好，和

$$\text{带权周转时间} = \frac{\text{作业周转时间}}{\text{作业实际运行时间}} = \frac{\text{作业完成时间} - \text{提交时间}}{\text{作业实际运行时间}} \quad (1)$$

$$\text{平均带权周转时间} = \frac{\text{各个作业带权周转时间之和}}{\text{作业个数}} \quad (2)$$

- 等待时间：

**定义：**指进程处于等待处理机状态之和，等待时间越长，用户的满意度越低。



对于进程来说，等待时间就是指进程建立后等待被服务的时间之和，在等待I/O完成的期间其实进程也是在被服务的，所以不计入等待时间。

对于作业来说，不仅要考虑建立进程后的等待时间，还要加上作业在外存后备队列中等待的时间。

一个作业总共需要被CPU服务多久，被I/O设备服务多久一般是确定不变的，因此调度算法其实只会影响作业/进程的等待时间。当然，与前面指标类似，也有“平均等待时间”来评价整体性能。

注意对于进程和作业的等待时间不一样。

- 响应时间：用户从提交命令到首次响应所使用的时间。

## 2.2.5\_调度算法：

- FCFS：

**算法规则：**按照作业/进程到达的先后顺序进行服务。

**处理作业/进程：**作业调度中，考虑的是哪个作业先到达后备队列。进程调度考虑的是哪一个作业先到达就绪队列。

**非抢占式算法。**

**优缺点：**

优点：公平，实现简单。

缺点：排在长作业后面的短作业需要等待很长时间，导致带权周转时间大。

**对长作业有利，对短作业不友好。**

是否会导致饥饿：

饥饿 就是一个进程/作业长时间得不到服务。

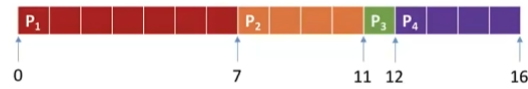
不会导致饥饿。

## ■ FCFS例子:

例题: 各进程到达就绪队列的时间、需要的运行时间如下表所示。使用**先来先服务**调度算法, 计算各进程的等待时间、平均等待时间、周转时间、平均周转时间、带权周转时间、平均带权周转时间。

进程	到达时间	运行时间
P1	0	7
P2	2	4
P3	4	1
P4	5	4

先来先服务调度算法: 按照到达的先后顺序调度, 事实上就是等待时间越久的越优先得到服务。  
因此, 调度顺序为: P1 → P2 → P3 → P4



周转时间 = 完成时间 - 到达时间

P1=7-0=7; P2=11-2=9; P3=12-4=8; P4=16-5=11

带权周转时间 = 周转时间/运行时间

P1=7/7=1; P2=9/4=2.25; P3=8/1=8; P4=11/4=2.75

等待时间 = 周转时间 - 运行时间

P1=7-7=0; P2=9-4=5; P3=8-1=7; P4=11-4=7

平均周转时间 = (7+9+8+11)/4 = 8.75

平均带权周转时间 = (1+2.25+8+2.75)/4 = 3.5

平均等待时间 = (0+5+7+7)/4 = 4.75

注意: 本例中的进程都是纯计算型的进程, 一个进程到达后要么在等待, 要么在运行。如果是又有计算、又有I/O操作的进程, 其等待时间就是周转时间 - 运行时间 - I/O操作的时间

## • SJF/SPF: short job first short process first:

**思想:** 追求最少的平均等待时间, 最少的平均周转时间, 最少的平均带权周转时间。

**规则:** 最短的作业/服务优先得到服务。

**用于作业/进程调度:** 用于进程时, 一般称为短进程优先, short process first.

**抢占否:** 一般默认非抢占式, 但是也有抢占式版本: **最短剩余时间优先算法**。

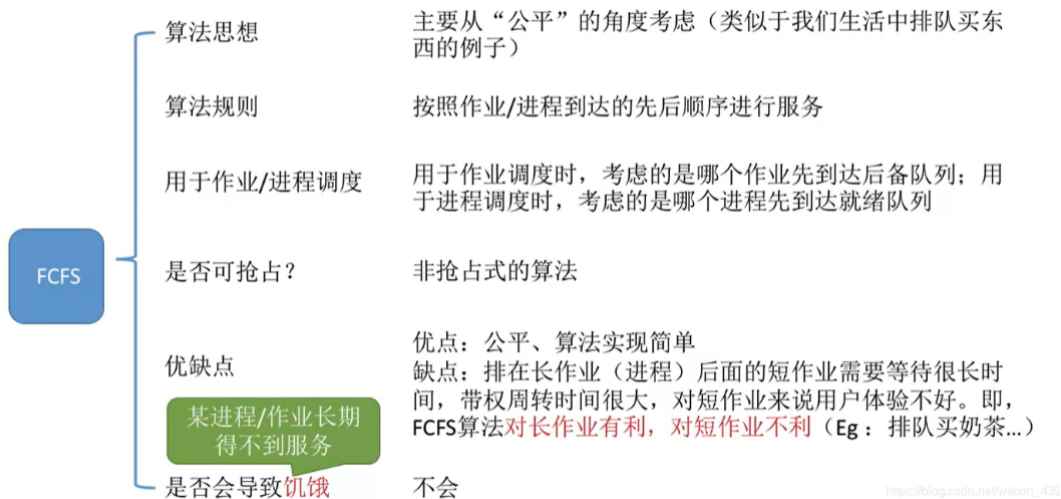
一句话总结:

非抢占式短作业优先, 就是每一次运行的时候, 都处理当前已经到达的运行时间最短的进程。

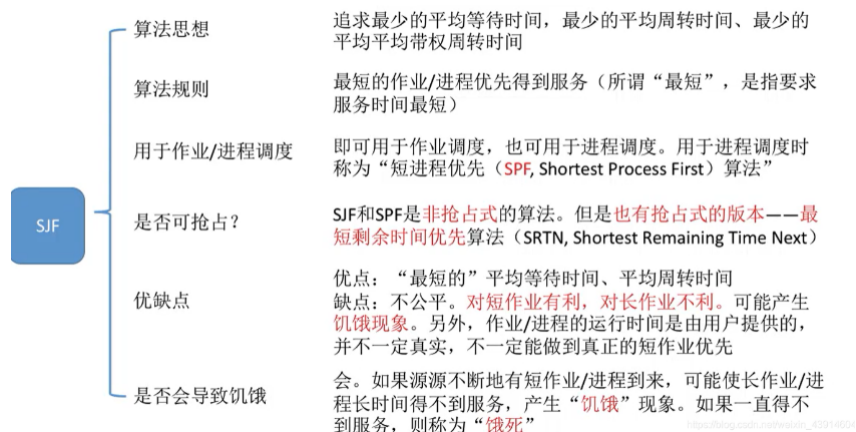
抢占式的短作业优先, 当有进程加入就绪队列的时候就需要调度, 如果新到达的比之前的时间小。就改变运行的进程。

## 调度算法总结:

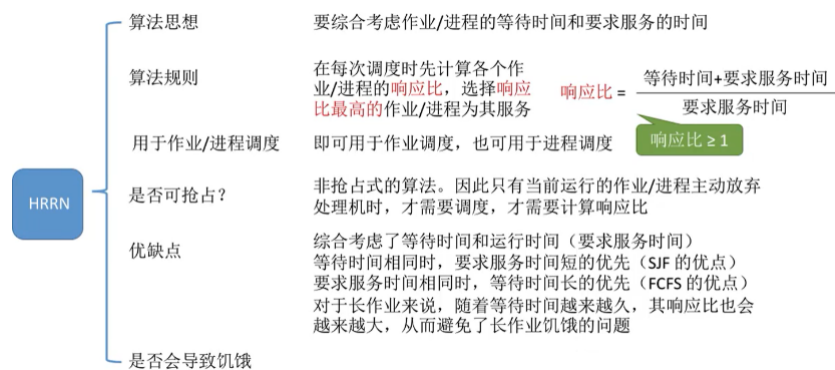
先来先服务。FCFS: (非抢占) 不饿



## 短作业优先。SJF: (可抢占、不抢占) 饥饿



## 高相应比优先：非抢占、不饿



## 时间片轮转：抢占式 不会饥饿



时间片轮转	算法思想	公平地、轮流地为各个进程服务，让每个进程在一定时间间隔内都可以得到响应
	算法规则	按照各进程到达就绪队列的顺序，轮流让各个进程执行一个时间片（如 100ms）。若进程未在一个时间片内执行完，则剥夺处理机，将进程重新放到就绪队列队尾重新排队。
	用于作业/进程调度	用于进程调度（只有作业放入内存建立了相应的进程后，才能被分配处理机时间片）
	是否可抢占？	若进程未能在时间片内运行完，将被强行剥夺处理机使用权，因此时间片轮转调度算法属于抢占式的算法。由时钟装置发出时钟中断来通知CPU时间片已到
	优缺点	优点：公平；响应快，适用于分时操作系统； 缺点：由于高频率的进程切换，因此有一定开销；不区分任务的紧急程度。
	是否会导致饥饿	不会
	补充	时间片太大或太小分别有什么影响？ <a href="https://blog.csdn.net/weixin_43914">https://blog.csdn.net/weixin_43914</a>

## 优先级调度算法：抢占、非抢占都会饥饿。

优先级调度	算法思想	随着计算机的发展，特别是实时操作系统的出现，越来越多的应用场景需要根据任务的紧急程度来决定处理顺序
	算法规则	调度时选择优先级最高的作业/进程
	用于作业/进程调度	既可用于作业调度，也可用于进程调度。甚至，还会用于在之后会学习的I/O调度中
	是否可抢占？	抢占式、非抢占式都有。做题时的区别在于：非抢占式只需在进程主动放弃处理机时进行调度即可，而抢占式还需在就绪队列变化时，检查是否会发生抢占。
	优缺点	优点：用优先级区分紧急程度、重要程度，适用于实时操作系统。可灵活地调整对各种作业/进程的偏好程度。 缺点：若源源不断地有高优先级进程到来，则可能导致饥饿
	是否会导致饥饿	会 <a href="https://blog.csdn.net/weixin_43914">https://blog.csdn.net/weixin_43914</a>

## 多级反馈队列调度算法：抢占式会饥饿。

多级反馈队列	算法思想	对其他调度算法的折中权衡
	算法规则	1. 设置多级就绪队列，各级队列优先级从高到低，时间片从小到大 2. 新进程到达时先进入第1级队列，按FCFS原则排队等待被分配时间片，若用完时间片进程还未结束，则进程进入下一级队列队尾。如果此时已经是在最下级的队列，则重新放回该队列队尾 3. 只有第k级队列为空时，才会为k+1级队头的进程分配时间片
	用于作业/进程调度	用于进程调度
	是否可抢占？	抢占式的算法。在k级队列的进程运行过程中，若更上级的队列（1~k-1级）中进入了一个新进程，则由于新进程处于优先级更高的队列中，因此新进程会抢占处理机，原来运行的进程放回k级队列队尾。
	优缺点	对各类型进程相对公平（FCFS的优点）；每个新到达的进程都可以很快就得到响应（RR的优点）；短进程只用较少的时间就可完成（SPF的优点）；不必实现估计进程的运行时间（避免用户作假）；可灵活地调整对各类进程的偏好程度，比如CPU密集型进程、I/O密集型进程（拓展：可以将因I/O而阻塞的进程重新放回原队列，这样I/O型进程就可以保持较高优先级）
	是否会导致饥饿	会 <a href="https://blog.csdn.net/weixin_43914">https://blog.csdn.net/weixin_43914</a>

对比：



FCFS算法的优点是公平

SJF 算法的优点是能尽快处理完短作业，  
平均等待/周转时间等参数很优秀

时间片轮转调度算法可以让各个进程得到及时的响应

优先级调度算法可以灵活地调整各种进程被服务的机会