

操作系统1

1.请你说一下进程与线程的概念，以及为什么要有进程线程，其中有什么区别，他们各自又是怎么同步的

参考回答：

基本概念：

进程是对运行时程序的封装，是系统进行资源调度和分配的基本单位，实现了操作系统的并发；
线程是进程的子任务，是CPU调度和分派的基本单位，用于保证程序的实时性，实现进程内部的并发；

线程是操作系统可识别的最小执行和调度单位。每个线程都独自占用一个虚拟处理器：独自の寄存器组，指令计数器和处理器状态。每个线程完成不同的任务，但是共享同一地址空间（也就是同样的动态内存，映射文件，目标代码等等），打开的文件队列和其他内核资源。

区别：

1.一个线程只能属于一个进程，而一个进程可以有多个线程，但至少有一个线程。线程依赖于进程而存在。

2.进程在执行过程中拥有独立的内存单元，而多个线程共享进程的内存。（资源分配给进程，同一进程的所有线程共享该进程的所有资源。同一进程中的多个线程共享代码段（代码和常量），数据段（全局变量和静态变量），扩展段（堆存储）。但是每个线程拥有自己的栈段，栈段又叫运行时段，用来存放所有局部变量和临时变量。）

3.进程是资源分配的最小单位，线程是CPU调度的最小单位；

4.系统开销：由于在创建或撤消进程时，系统都要为之分配或回收资源，如内存空间、I / o设备等。因此，操作系统所付出的开销将显著地大于在创建或撤消线程时的开销。类似地，在进行进程切换时，涉及到整个当前进程CPU环境的保存以及新被调度运行的进程的CPU环境的设置。而线程切换只须保存和设置少量寄存器的内容，并不涉及存储器管理方面的操作。可见，进程切换的开销也远大于线程切换的开销。

5.通信：由于同一进程中的多个线程具有相同的地址空间，致使它们之间的同步和通信的实现，也变得比较容易。进程间通信IPC，线程间可以直接读写进程数据段（如全局变量）来进行通信——需要进程同步和互斥手段的辅助，以保证数据的一致性。在有的系统中，线程的切换、同步和通信都无须操作系统内核的干预

6.进程编程调试简单可靠性高，但是创建销毁开销大；线程正相反，开销小，切换速度快，但是编程调试相对复杂。

7.进程间不会相互影响；线程一个线程挂掉将导致整个进程挂掉

8.进程适应于多核、多机分布；线程适用于多核

进程间通信的方式：

进程间通信主要包括管道、系统IPC（包括消息队列、信号量、信号、共享内存等）、以及套接字socket。

1.管道：

管道主要包括无名管道和命名管道:管道可用于具有亲缘关系的父子进程间的通信,有名管道除了具有管道所具有的功能外,它还允许无亲缘关系进程间的通信

1.1 普通管道PIPE:

- 1)它是半双工的(即数据只能在一个方向上流动),具有固定的读端和写端
- 2)它只能用于具有亲缘关系的进程之间的通信(也是父子进程或者兄弟进程之间)
- 3)它可以看成是一种特殊的文件,对于它的读写也可以使用普通的read、write等函数。但是它不是普通的文件,并不属于其他任何文件系统,并且只存在于内存中。

1.2 命名管道FIFO:

- 1)FIFO可以在无关的进程之间交换数据
- 2)FIFO有路径名与之相关联,它以一种特殊设备文件形式存在于文件系统中。

2. 系统IPC:

2.1 消息队列

消息队列,是消息的链接表,存放在内核中。一个消息队列由一个标识符(即队列ID)来标记。(消息队列克服了信号传递信息少,管道只能承载无格式字节流以及缓冲区大小受限等特点)具有写权限得进程可以按照一定得规则向消息队列中添加新信息;对消息队列有读权限得进程则可以从消息队列中读取信息;

特点:

- 1)消息队列是面向记录的,其中的消息具有特定的格式以及特定的优先级。
- 2)消息队列独立于发送与接收进程。进程终止时,消息队列及其内容并不会被删除。
- 3)消息队列可以实现消息的随机查询,消息不一定要以先进先出的次序读取,也可以按消息的类型读取。

2.2 信号量semaphore

信号量(semaphore)与已经介绍过的IPC结构不同,它是一个计数器,可以用来控制多个进程对共享资源的访问。信号量用于实现进程间的互斥与同步,而不是用于存储进程间通信数据。

特点:

- 1)信号量用于进程间同步,若要在进程间传递数据需要结合共享内存。
- 2)信号量基于操作系统的PV操作,程序对信号量的操作都是原子操作。
- 3)每次对信号量的PV操作不仅限于对信号量值加1或减1,而且可以加减任意正整数。
- 4)支持信号量组。

2.3 信号signal

信号是一种比较复杂的通信方式,用于通知接收进程某个事件已经发生。

2.4 共享内存(Shared Memory)

它使得多个进程可以访问同一块内存空间,不同进程可以及时看到对方进程中对共享内存中数据得更新。这种方式需要依靠某种同步操作,如互斥锁和信号量等

特点:

- 1)共享内存是最快的一种IPC,因为进程是直接对内存进行存取
- 2)因为多个进程可以同时操作,所以需要进行同步
- 3)信号量+共享内存通常结合在一起使用,信号量用来同步对共享内存的访问

3.套接字SOCKET:

socket也是一种进程间通信机制，与其他通信机制不同的是，它可用于不同主机之间的进程通信。

线程间通信的方式:

临界区: 通过多线程的串行化来访问公共资源或一段代码，速度快，适合控制数据访问；

互斥量Synchronized/Lock: 采用互斥对象机制，只有拥有互斥对象的线程才有访问公共资源的权限。因为互斥对象只有一个，所以可以保证公共资源不会被多个线程同时访问

信号量Semaphore: 为控制具有有限数量的用户资源而设计的，它允许多个线程在同一时刻去访问同一个资源，但一般需要限制同一时刻访问此资源的最大线程数目。

事件(信号), Wait/Notify: 通过通知操作的方式来保持多线程同步，还可以方便的实现多线程优先级的比较操作

2.请你说一说Linux虚拟地址空间

参考回答:

为了防止不同进程同一时刻在物理内存中运行而对物理内存的争夺和践踏，采用了虚拟内存。

虚拟内存技术使得不同进程在运行过程中，它所看到的是自己独自占有了当前系统的4G内存。所有进程共享同一物理内存，每个进程只把自己目前需要的虚拟内存空间映射并存储到物理内存上。事实上，在每个进程创建加载时，内核只是为进程“创建”了虚拟内存的布局，具体就是初始化进程控制表中内存相关的链表，实际上并不立即就把虚拟内存对应位置的程序数据和代码（比如.text .data段）拷贝到物理内存中，只是建立好虚拟内存和磁盘文件之间的映射就好（叫做存储器映射），等到运行到对应的程序时，才会通过缺页异常，来拷贝数据。还有进程运行过程中，要动态分配内存，比如malloc时，也只是分配了虚拟内存，即为这块虚拟内存对应的页表项做相应设置，当进程真正访问到此数据时，才引发缺页异常。

请求分页系统、请求分段系统和请求段页式系统都是针对虚拟内存的，通过请求实现内存与外存的信息置换。

虚拟内存的好处:

1.扩大地址空间;

2.内存保护: 每个进程运行在各自的虚拟内存地址空间，互相不能干扰对方。虚存还对特定的内存地址提供写保护，可以防止代码或数据被恶意篡改。

3.公平内存分配。采用了虚存之后，每个进程都相当于有同样大小的虚存空间。

4.当进程通信时，可采用虚存共享的方式实现。

5.当不同的进程使用同样的代码时，比如库文件中的代码，物理内存中可以只存储一份这样的代码，不同的进程只需要把自己的虚拟内存映射过去就可以了，节省内存

6.虚拟内存很适合在多道程序设计系统中使用，许多程序的片段同时保存在内存中。当一个程序等待它的一部分读入内存时，可以把CPU交给另一个进程使用。在内存中可以保留多个进程，系统并发度提高

7.在程序需要分配连续的内存空间的时候，只需要在虚拟内存空间分配连续空间，而不需要实际物理内存的连续空间，可以利用碎片

虚拟内存的代价:

1.虚存的管理需要建立很多数据结构，这些数据结构要占用额外的内存

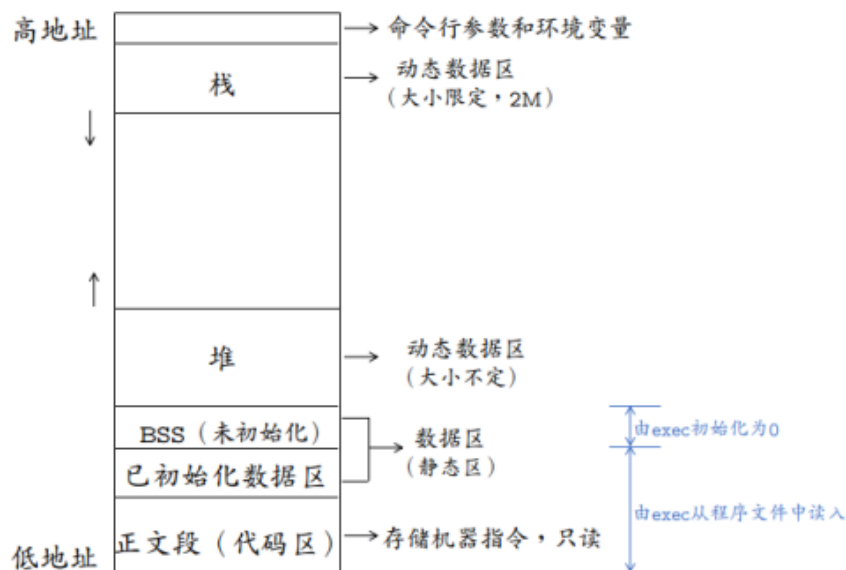
2.虚拟地址到物理地址的转换，增加了指令的执行时间。

3.页面的换入换出需要磁盘I/O，这是很耗时的

4.如果一页中只有一部分数据，会浪费内存。

3.请你说一说操作系统中的程序的内存结构

参考回答：



一个程序本质上都是由BSS段、data段、text段三个组成的。可以看到一个可执行程序在存储（没有调入内存）时分为代码段、数据区和未初始化数据区三部分。

BSS段（未初始化数据区）：通常用来存放程序中未初始化的全局变量和静态变量的一块内存区域。BSS段属于静态分配，程序结束后静态变量资源由系统自动释放。

数据段：存放程序中已初始化的全局变量的一块内存区域。数据段也属于静态内存分配

代码段：存放程序执行代码的一块内存区域。这部分区域的大小在程序运行前就已经确定，并且内存区域属于只读。在代码段中，也有可能包含一些只读的常数变量

text段和data段在编译时已经分配了空间，而BSS段并不占用可执行文件的大小，它是由链接器来获取内存的。

bss段（未进行初始化的数据）的内容并不存放在磁盘上的程序文件中。其原因是内核在程序开始运行前将它们设置为0。需要存放在程序文件中的只有正文段和初始化数据段。

data段（已经初始化的数据）则为数据分配空间，数据保存到目标文件中。

数据段包含经过初始化的全局变量以及它们的值。BSS段的大小从可执行文件中得到，然后链接器得到这个大小的内存块，紧跟在数据段的后面。当这个内存进入程序的地址空间后全部清零。包含数据段和BSS段的整个区段此时通常称为数据区。

可执行程序在运行时又多出两个区域：栈区和堆区。

- 栈区：由编译器自动释放，存放函数的参数值、局部变量等。每当一个函数被调用时，该函数的返回类型和一些调用的信息被存放到栈中。然后这个被调用的函数再为他的自动变量和临时变量在栈上分配空间。每调用一个函数一个新的栈就会被使用。栈区是从高地址位向低地址位增长的，是一块连续的内存区域，最大容量是由系统预先定义好的，申请的栈空间超过这个界限时会提示溢出，用户能从栈中获取的空间较小。
- 堆区：用于动态分配内存，位于BSS和栈中间的地址区域。由程序员申请分配和释放。堆是从低地址位向高地址位增长，采用链式存储结构。频繁的malloc/free造成内存空间的不连续，产

生碎片。当申请堆空间时库函数是按照一定的算法搜索可用的足够大的空间。因此堆的效率比栈要低的多。

4.请你说一说操作系统中的缺页中断

参考回答：

malloc()和mmap()等内存分配函数，在分配时只是建立了进程虚拟地址空间，并没有分配虚拟内存对应的物理内存。当进程访问这些没有建立映射关系的虚拟内存时，处理器自动触发一个缺页异常。

缺页中断：在请求分页系统中，可以通过查询页表中的状态位来确定所要访问的页面是否存在于内存中。每当所要访问的页面不在内存是，会产生一次缺页中断，此时操作系统会根据页表中的外存地址在外存中找到所缺的一页，将其调入内存。

缺页本身是一种中断，与一般的中断一样，需要经过4个处理步骤：

- 1、保护CPU现场**
- 2、分析中断原因**
- 3、转入缺页中断处理程序进行处理**
- 4、恢复CPU现场，继续执行**

但是缺页中断是由于所要访问的页面不存在于内存时，由硬件所产生的一种特殊的中断，因此，与一般的中断存在区别：

- 1、在指令执行期间产生和处理缺页中断信号**
- 2、一条指令在执行期间，可能产生多次缺页中断**
- 3、缺页中断返回是，执行产生中断的一条指令，而一般的中断返回是，执行下一条指令。**

5.请你回答一下fork和vfork的区别

参考回答：

fork的基础知识：

fork:创建一个和当前进程映像一样的进程可以通过fork()系统调用：

```
#include <sys/types.h>
#include <unistd.h>
pid_t fork(void);
```

成功调用fork()会创建一个新的进程，它几乎与调用fork()的进程一模一样，这两个进程都会继续运行。在子进程中，成功的fork()调用会返回0。在父进程中fork()返回子进程的pid。如果出现错误，fork()返回一个负值。

最常见的fork()用法是创建一个新的进程，然后使用exec()载入二进制映像，替换当前进程的映像。这种情况下，派生（fork）了新的进程，而这个子进程会执行一个新的二进制可执行文件的映像。这种“派生加执行”的方式是很常见的。

在早期的Unix系统中，创建进程比较原始。当调用fork时，内核会把所有的内部数据结构复制一份，复制进程的页表项，然后把父进程的地址空间中的内容逐页的复制到子进程的地址空间中。但从内核角度来说，逐页的复制方式是十分耗时的。现代的Unix系统采取了更多的优化，例如Linux，采用了写时复制的方法，而不是对父进程空间进程整体复制。

vfork的基础知识：

在实现写时复制之前，Unix的设计者们就一直很关注在fork后立刻执行exec所造成的地址空间的浪费。BSD的开发者在3.0的BSD系统中引入了vfork()系统调用。

```
#include <sys/types.h>
```

```
#include <unistd.h>
```

```
pid_t vfork(void);
```

除了子进程必须要立刻执行一次对exec的系统调用，或者调用_exit()退出，对vfork()的成功调用所产生的结果和fork()是一样的。vfork()会挂起父进程直到子进程终止或者运行了一个新的可执行文件的映像。通过这样的方式，vfork()避免了地址空间的按页复制。在这个过程中，父进程和子进程共享相同的地址空间和页表项。实际上vfork()只完成了一件事：复制内部的内核数据结构。因此，子进程也就不能修改地址空间中的任何内存。

vfork()是一个历史遗留产物，Linux本不应该实现它。需要注意的是，即使增加了写时复制，vfork()也要比fork()快，因为它没有进行页表项的复制。然而，写时复制的出现减少对于替换fork()争论。实际上，直到2.2.0内核，vfork()只是一个封装过的fork()。因为对vfork()的需求要小于fork()，所以vfork()的这种实现方式是可行的。

补充知识点：写时复制

Linux采用了写时复制的方法，以减少fork时对父进程空间整体复制带来的开销。

写时复制是一种采取了惰性优化方法来避免复制时的系统开销。它的前提很简单：如果有多个进程要读取它们自己的那部门资源的副本，那么复制是不必要的。每个进程只要保存一个指向这个资源的指针就可以了。只要没有进程要去修改自己的“副本”，就存在着这样的幻觉：每个进程好像独占那个资源。从而就避免了复制带来的负担。如果一个进程要修改自己的那份资源“副本”，那么就会复制那份资源，并把复制的那份提供给进程。不过其中的复制对进程来说是透明的。这个进程就可以修改复制后的资源了，同时其他的进程仍然共享那份没有修改过的资源。所以这就是名称的由来：在写入时进行复制。

写时复制的主要好处在于：如果进程从来就不需要修改资源，则不需要进行复制。惰性算法的好处就在于它们尽量推迟代价高昂的操作，直到必要的时刻才会去执行。

在使用虚拟内存的情况下，写时复制（Copy-On-Write）是以页为基础进行的。所以，只要进程不修改它全部的地址空间，那么就不必复制整个地址空间。在fork()调用结束后，父进程和子进程都相信它们有一个自己的地址空间，但实际上它们共享父进程的原始页，接下来这些页又可以被其他的父进程或子进程共享。

写时复制在内核中的实现非常简单。与内核页相关的数据结构可以被标记为只读和写时复制。如果有进程试图修改一个页，就会产生一个缺页中断。内核处理缺页中断的方式就是对该页进行一次透明复制。这时会清除页面的COW属性，表示着它不再被共享。

现代的计算机系统结构中都在内存管理单元（MMU）提供了硬件级别的写时复制支持，所以实现是很容易的。

在调用fork()时，写时复制是有很大优势的。因为大量的fork之后都会跟着执行exec，那么复制整个父进程地址空间中的内容到子进程的地址空间完全是在浪费时间：如果子进程立刻执行一个新的二进制可执行文件的映像，它先前的地址空间就会被交换出去。写时复制可以对这种情况进行优化。

fork和vfork的区别：

1. fork()的子进程拷贝父进程的数据段和代码段；vfork()的子进程与父进程共享数据段

2. fork() 的父子进程的执行次序不确定；vfork() 保证子进程先运行，在调用exec或exit之前与父进程数据是共享的，在它调用exec或exit之后父进程才可能被调度运行。
3. vfork() 保证子进程先运行，在它调用exec或exit之后父进程才可能被调度运行。如果在调用这两个函数之前子进程依赖于父进程的进一步动作，则会导致死锁。
4. 当需要改变共享数据段中变量的值，则拷贝父进程。

6. 请问如何修改文件最大句柄数？

参考回答：

linux默认最大文件句柄数是1024个，在linux服务器文件并发量比较大的情况下，系统会报"too many open files"的错误。故在linux服务器高并发调优时，往往需要预先调优Linux参数，修改Linux最大文件句柄数。

有两种方法：

1. ulimit -n <可以同时打开的文件数>，将当前进程的最大句柄数修改为指定的参数（注：该方法只针对当前进程有效，重新打开一个shell或者重新开启一个进程，参数还是之前的值）

首先用ulimit -a查询Linux相关的参数，如下所示：

```
core file size          (blocks, -c) 0
data seg size           (kbytes, -d) unlimited
scheduling priority     (-e) 0
file size               (blocks, -f) unlimited
pending signals         (-i) 94739
max locked memory       (kbytes, -l) 64
max memory size         (kbytes, -m) unlimited
open files              (-n) 1024
pipe size               (512 bytes, -p) 8
POSIX message queues    (bytes, -q) 819200
real-time priority      (-r) 0
stack size              (kbytes, -s) 8192
cpu time                (seconds, -t) unlimited
max user processes      (-u) 94739
virtual memory          (kbytes, -v) unlimited
file locks              (-x) unlimited
```

其中，open files就是最大文件句柄数，默认是1024个。

修改Linux最大文件句柄数：ulimit -n 2048，将最大句柄数修改为2048个。

2. 对所有进程都有效的方法，修改Linux系统参数

vi /etc/security/limits.conf 添加

```
* soft nofile 65536
* hard nofile 65536
```

将最大句柄数改为65536

修改以后保存，注销当前用户，重新登录，修改后的参数就生效了

7. 请你说一说并发(concurrency)和并行(parallelism)

参考回答：

并发 (concurrency)：指宏观上看起来两个程序在同时运行，比如说在单核cpu上的多任务。但是从微观上看两个程序的指令是交织着运行的，你的指令之间穿插着我的指令，我的指令之间穿插着你的，在单个周期内只运行了一个指令。这种并发并不能提高计算机的性能，只能提高效率。

并行 (parallelism)：指严格物理意义上的同时运行，比如多核cpu，两个程序分别运行在两个核上，两者之间互不影响，单个周期内每个程序都运行了自己的指令，也就是运行了两条指令。这样说来并行的确提高了计算机的效率。所以现在的cpu都是往多核方面发展。

8.请问MySQL的端口号是多少，如何修改这个端口号

参考回答：

查看端口号：

使用命令show global variables like 'port';查看端口号，mysql的默认端口是3306。

(补充：sqlserver默认端口号为：1433；oracle默认端口号为：1521；DB2默认端口号为：5000；PostgreSQL默认端口号为：5432)

修改端口号：

修改端口号：编辑/etc/my.cnf文件，早期版本有可能是my.conf文件名，增加端口参数，并且设定端口，注意该端口未被使用，保存退出。

9.请你说一说操作系统中的页表寻址

参考回答：

页式内存管理，内存分成固定长度的一个个页片。操作系统为每一个进程维护了一个从虚拟地址到物理地址的映射关系的数据结构，叫页表，页表的内容就是该进程的虚拟地址到物理地址的一个映射。页表中的每一项都记录了这个页的基地址。通过页表，由逻辑地址的高位部分先找到逻辑地址对应的页基地址，再由页基地址偏移一定长度就得到最后的物理地址，偏移的长度由逻辑地址的低位部分决定。一般情况下，这个过程都可以由硬件完成，所以效率还是比较高的。页式内存管理的优点就是比较灵活，内存管理以较小的页为单位，方便内存换入换出和扩充地址空间。

Linux最初的两级页表机制：

两级分页机制将32位的虚拟空间分成三段，低十二位表示页内偏移，高20分成两段分别表示两级页表的偏移。

* PGD(Page Global Directory): 最高10位，全局页目录表索引

* PTE(Page Table Entry): 中间10位，页表入口索引

当在进行地址转换时，结合在CR3寄存器中存放的页目录(page directory, PGD)的这一页的物理地址，再加上从虚拟地址中抽出高10位叫做页目录表项(内核也称这为pgd)的部分作为偏移，即定位到可以描述该地址的pgd；从该pgd中可以获取可以描述该地址的页表的物理地址，再加上从虚拟地址中抽取中间10位作为偏移，即定位到可以描述该地址的pte；在这个pte中即可获取该地址对应的页的物理地址，加上从虚拟地址中抽取的最后12位，即形成该页的页内偏移，即可最终完成从虚拟地址到物理地址的转换。从上述过程中，可以看出，对虚拟地址的分级解析过程，实际上就是不断深入页表层次，逐渐定位到最终地址的过程，所以这一过程被叫做page talbe walk。

Linux的三级页表机制：

当X86引入物理地址扩展(Pisycal Addrress Extension, PAE)后，可以支持大于4G的物理内存(36位)，但虚拟地址依然是32位，原先的页表项不适用，它实际多4 bytes被扩充到8 bytes，这意味

着，每一页现在能存放的pte数目从1024变成512了(4k/8)。相应地，页表层级发生了变化，Linux新增加了一个层级，叫做页中间目录(page middle directory, PMD), 变成：

字段	描述	位数
cr3	指向一个PDPT	cr3寄存器存储
PGD	指向PDPT中4个项中的一个	位31~30
PMD	指向页目录中512项中的一个	位29~21
PTE	指向页表中512项中的一个	位20~12
page offset	4KB页中的偏移	位11~0

现在就同时存在2级页表和3级页表，在代码管理上肯定不方便。巧妙的是，Linux采取了一种抽象方法：所有架构全部使用3级页表：即PGD -> PMD -> PTE。那只使用2级页表(如非PAE的X86)怎么办？

办法是针对使用2级页表的架构，把PMD抽象掉，即虚设一个PMD表项。这样在page table walk过程中，PGD本直接指向PTE的，现在不了，指向一个虚拟的PMD，然后再由PMD指向PTE。这种抽象保持了代码结构的统一。

Linux的四级页表机制：

硬件在发展，3级页表很快又捉襟见肘了，原因是64位CPU出现了，比如X86_64，它的硬件是实实在在支持4级页表的。它支持48位的虚拟地址空间¹。如下：

字段	描述	位数
PML4	指向一个PDPT	位47~39
PGD	指向PDPT中4个项中的一个	位38~30
PMD	指向页目录中512项中的一个	位29~21
PTE	指向页表中512项中的一个	位20~12
page offset	4KB页中的偏移	位11~0

Linux内核针对使用原来的3级列表(PGD->PMD->PTE)，做了折衷。即采用一个唯一的，共享的顶级层次，叫PML4。这个PML4没有编码在地址中，这样就能套用原来的3级列表方案了。不过代价就是，由于只有唯一的PML4，寻址空间被局限在(2³⁹)=512G，而本来PML4段有9位，可以支持512个PML4表项的。现在为了使用3级列表方案，只能限制使用一个，512G的空间很快就又不够用了，解决方案呼之欲出。

在2004年10月，当时的X86_64架构代码的维护者Andi Kleen提交了一个叫做4level page tables for Linux的PATCH系列，为Linux内核带来了4级页表的支持。在他的解决方案中，不出意料地，按照X86_64规范，新增了一个PML4的层级，在这种解决方案中，X86_64拥有一个有512条目的PML4，512条目的PGD，512条目的PMD，512条目的PTE。对于仍使用3级目录的架构来说，它们依然拥有一个虚拟的PML4，相关的代码会在编译时被优化掉。这样，就把Linux内核的3级列表扩充为4级列表。这系列PATCH工作得不错，不久被纳入Andrew Morton的-mm树接受测试。不出意外的话，它将在v2.6.11版本中释出。但是，另一个知名开发者Nick Piggin提出了一些看法，他认为Andi的Patch很不错，不过他认为最好还是把PGD作为第一级目录，把新增加的层次放在中间，并给出了他自己的Patch:alternate 4-level page tables patches。Andi更想保持自己的PATCH，他认为Nick不过是玩了改名的游戏，而且他的PATCH经过测试很稳定，快被合并到主线了，不宜再折腾。不过Linus却表达了对Nick Piggin的支持，理由是Nick的做法conceptually least intrusive。毕竟作为

Linux的扛把子，稳定对于Linus来说意义重大。最终，不意外地，最后Nick Piggin的PATCH在v2.6.11版本中被合并入主线。在这种方案中，4级页表分别是：PGD -> PUD -> PMD -> PTE。

10.请你说一说有了进程，为什么还要有线程？

参考回答：

线程产生的原因：

进程可以使多个程序能并发执行，以提高资源的利用率和系统的吞吐量；但是其具有一些缺点：

进程在同一时间只能干一件事

进程在执行的过程中如果阻塞，整个进程就会挂起，即使进程中有些工作不依赖于等待的资源，仍然不会执行。

因此，操作系统引入了比进程粒度更小的线程，作为并发执行的基本单位，从而减少程序在并发执行时所付出的时空开销，提高并发性。和进程相比，线程的优势如下：

从资源上来讲，线程是一种非常"节俭"的多任务操作方式。在linux系统下，启动一个新的进程必须分配给它独立的地址空间，建立众多的数据表来维护它的代码段、堆栈段和数据段，这是一种"昂贵"的多任务工作方式。

从切换效率上来讲，运行于一个进程中的多个线程，它们之间使用相同的地址空间，而且线程间彼此切换所需时间也远远小于进程间切换所需要的时间。据统计，一个进程的开销大约是一个线程开销的30倍左右。

从通信机制上来讲，线程间方便的通信机制。对不同进程来说，它们具有独立的数据空间，要进行数据的传递只能通过进程间通信的方式进行，这种方式不仅费时，而且很不方便。线程则不然，由于同一进程下的线程之间贡献数据空间，所以一个线程的数据可以直接为其他线程所用，这不仅快捷，而且方便。

除以上优点外，多线程程序作为一种多任务、并发的工作方式，还有如下优点：

- 1、使多CPU系统更加有效。操作系统会保证当线程数不大于CPU数目时，不同的线程运行于不同的CPU上。
- 2、改善程序结构。一个既长又复杂的进程可以考虑分为多个线程，成为几个独立或半独立的运行部分，这样的程序才会利于理解和修改。

11.请问单核机器上写多线程程序，是否需要考虑加锁，为什么？

参考回答：

在单核机器上写多线程程序，仍然需要线程锁。因为线程锁通常用来实现线程的同步和通信。在单核机器上的多线程程序，仍然存在线程同步的问题。因为在抢占式操作系统中，通常为每个线程分配一个时间片，当某个线程时间片耗尽时，操作系统会将其挂起，然后运行另一个线程。如果这两个线程共享某些数据，不使用线程锁的前提下，可能会导致共享数据修改引起冲突。

12.请问线程需要保存哪些上下文，SP、PC、EAX这些寄存器是干嘛用的

参考回答：

线程在切换的过程中需要保存当前线程Id、线程状态、堆栈、寄存器状态等信息。其中寄存器主要包括SP PC EAX等寄存器，其主要功能如下：

SP:堆栈指针，指向当前栈的栈顶地址

PC:程序计数器, 存储下一条将要执行的指令

EAX:累加寄存器, 用于加法乘法的缺省寄存器

13.请你说一说线程间的同步方式, 最好说出具体的系统调用

参考回答:

信号量

信号量是一种特殊的变量, 可用于线程同步。它只取自然数值, 并且只支持两种操作:

P(SV): 如果信号量SV大于0, 将它减一; 如果SV值为0, 则挂起该线程。

V(SV): 如果有其他进程因为等待SV而挂起, 则唤醒, 然后将SV+1; 否则直接将SV+1。

其系统调用为:

sem_wait (sem_t *sem): 以原子操作的方式将信号量减1, 如果信号量值为0, 则sem_wait将被阻塞, 直到这个信号量具有非0值。

sem_post (sem_t *sem): 以原子操作将信号量值+1。当信号量大于0时, 其他正在调用sem_wait等待信号量的线程将被唤醒。

互斥量

互斥量又称互斥锁, 主要用于线程互斥, 不能保证按序访问, 可以和条件锁一起实现同步。当进入临界区 时, 需要获得互斥锁并且加锁; 当离开临界区时, 需要对互斥锁解锁, 以唤醒其他等待该互斥锁的线程。其主要的系统调用如下:

pthread_mutex_init: 初始化互斥锁

pthread_mutex_destroy: 销毁互斥锁

pthread_mutex_lock: 以原子操作的方式给一个互斥锁加锁, 如果目标互斥锁已经被上锁, pthread_mutex_lock调用将阻塞, 直到该互斥锁的占有者将其解锁。

pthread_mutex_unlock: 以一个原子操作的方式给一个互斥锁解锁。

条件变量

条件变量, 又称条件锁, 用于在线程之间同步共享数据的值。条件变量提供一种线程间通信机制: 当某个共享数据达到某个值时, 唤醒等待这个共享数据的一个/多个线程。即, 当某个共享变量等于某个值时, 调用 signal/broadcast。此时操作共享变量时需要加锁。其主要的系统调用如下:

pthread_cond_init: 初始化条件变量

pthread_cond_destroy: 销毁条件变量

pthread_cond_signal: 唤醒一个等待目标条件变量的线程。哪个线程被唤醒取决于调度策略和优先级。

pthread_cond_wait: 等待目标条件变量。需要一个加锁的互斥锁确保操作的原子性。该函数中在进入wait状态前首先进行解锁, 然后接收到信号后会再加锁, 保证该线程对共享资源正确访问。

15.请你说一下多线程和多进程的不同

参考回答:

进程是资源分配的最小单位, 而线程是CPU调度的最小单位。多线程之间共享同一个进程的地址空间, 线程间通信简单, 同步复杂, 线程创建、销毁和切换简单, 速度快, 占用内存少, 适用于多核分布式系统, 但是线程间会相互影响, 一个线程意外终止会导致同一个进程的其他线程也终止, 程序可靠性弱。而多进程间拥有各自独立的运行地址空间, 进程间不会相互影响, 程序可靠性强,

但是进程创建、销毁和切换复杂，速度慢，占用内存多，进程间通信复杂，但是同步简单，适用于多核、多机分布。

16.请你说一说进程和线程的区别

参考回答：

- 1) 进程是cpu资源分配的最小单位，线程是cpu调度的最小单位。
- 2) 进程有独立的系统资源，而同一进程内的线程共享进程的大部分系统资源,包括堆、代码段、数据段，每个线程只拥有一些在运行中必不可少的私有属性，比如tcb,线程Id,栈、寄存器。
- 3) 一个进程崩溃，不会对其他进程产生影响；而一个线程崩溃，会让同一进程内的其他线程也死掉。
- 4) 进程在创建、切换和销毁时开销比较大，而线程比较小。进程创建的时候需要分配系统资源，而销毁的时候需要释放系统资源。进程切换需要分两步：切换页目录、刷新TLB以使用新的地址空间；切换内核栈和硬件上下文（寄存器）；而同一进程的线程间逻辑地址空间是一样的，不需要切换页目录、刷新TLB。
- 5) 进程间通信比较复杂，而同一进程的线程由于共享代码段和数据段，所以通信比较容易。

17.游戏服务器应该为每个用户开辟一个线程还是一个进程，为什么？

参考回答：

游戏服务器应该为每个用户开辟一个进程。因为同一进程间的线程会相互影响，一个线程死掉会影响其他线程，从而导致进程崩溃。因此为了保证不同用户之间不会相互影响，应该为每个用户开辟一个进程

18.请你说一说OS缺页置换算法

参考回答：

当访问一个内存中不存在的页，并且内存已满，则需要从内存中调出一个页或将数据送至磁盘对换区，替换一个页，这种现象叫做缺页置换。当前操作系统最常采用的缺页置换算法如下：

先进先出(FIFO)算法：置换最先调入内存的页面，即置换在内存中驻留时间最久的页面。按照进入内存的先后次序排列成队列，从队尾进入，从队首删除。

最近最少使用（LRU）算法：置换最近一段时间以来最长时间未访问过的页面。根据程序局部性原理，刚被访问的页面，可能马上又要被访问；而较长时间内没有被访问的页面，可能最近不会被访问。

当前最常采用的就是LRU算法。

19.请你说一下多进程和多线程的使用场景

参考回答：

多进程模型的优势是CPU

多线程模型主要优势为线程间切换代价较小，因此适用于I/O密集型的工作场景，因此I/O密集型的工作场景经常会由于I/O阻塞导致频繁的切换线程。同时，多线程模型也适用于单机多核分布式场景。

多进程模型，适用于CPU密集型。同时，多进程模型也适用于多机分布式场景中，易于多机扩展。

20.请你说一说死锁发生的条件以及如何解决死锁

参考回答：

死锁是指两个或两个以上进程在执行过程中，因争夺资源而造成的下相互等待的现象。

死锁发生的四个必要条件如下：

互斥条件：进程对所分配到的资源不允许其他进程访问，若其他进程访问该资源，只能等待，直至占有该资源的进程使用完成后释放该资源；

请求和保持条件：进程获得一定的资源后，又对其他资源发出请求，但是该资源可能被其他进程占有，此时请求阻塞，但该进程不会释放自己已经占有的资源

不可剥夺条件：进程已获得的资源，在未使用之前，不可被剥夺，只能在使用后自己释放

环路等待条件：进程发生死锁后，必然存在一个进程-资源之间的环形链

解决死锁的方法即破坏上述四个条件之一，主要方法如下：

资源一次性分配，从而剥夺请求和保持条件

可剥夺资源：即当进程新的资源未得到满足时，释放已占有的资源，从而破坏不可剥夺的条件

资源有序分配法：系统给每类资源赋予一个序号，每个进程按编号递增的请求资源，释放则相反，从而破坏环路等待的条件

21.请问虚拟内存和物理内存怎么对应

参考回答：

1、概念：

物理地址(physical address)

用于内存芯片级的单元寻址，与处理器和CPU连接的地址总线相对应。

虽然可以直接把物理地址理解成插在机器上那根内存本身，把内存看成一个从0字节一直到最大空量逐字节的编号的大数组，然后把这个数组叫做物理地址，但是事实上，这只是一个硬件提供给软件的抽象，内存的寻址方式并不是这样。所以，说它是“与地址总线相对应”，是更贴切一些，不过抛开对物理内存寻址方式的考虑，直接把物理地址与物理的内存——对应，也是可以接受的。也许错误的理解更利于形而上的抽象。

虚拟地址(virtual memory)

这是对整个内存（不要与机器上插那条对上号）的抽象描述。它是相对于物理内存来讲的，可以直接理解成“不真实的”，“假的”内存，例如，一个0x08000000内存地址，它并不对应物理地址上那个大数组中0x08000000 - 1那个地址元素；之所以是这样，是因为现代操作系统都提供了一种内存管理的抽象，即虚拟内存（virtual memory）。进程使用虚拟内存中的地址，由操作系统协助相关硬件，把它“转换”成真正的物理地址。这个“转换”，是所有问题讨论的关键。

有了这样的抽象，一个程序，就可以使用比真实物理地址大得多的地址空间。甚至多个进程可以使用相同的地址。不奇怪，因为转换后的物理地址并非相同的。可以把连接后的程序反编译看一下，发现连接器已经为程序分配了一个地址，例如，要调用某个函数A，代码不是call A，而是call 0x0811111111，也就是说，函数A的地址已经被定下来了。没有这样的“转换”，没有虚拟地址的概念，这样做是根本行不通的。

2、地址转换

第一步：CPU段式管理中——~~逻辑地址转线性地址~~

CPU要利用其段式内存管理单元，先将为一个逻辑地址转换成一个线性地址。

一个逻辑地址由两部份组成，【段标识符：段内偏移量】。

段标识符是由一个16位长的字段组成，称为段选择符。其中前13位是一个索引号。后面3位包含一些硬件细节，如图：



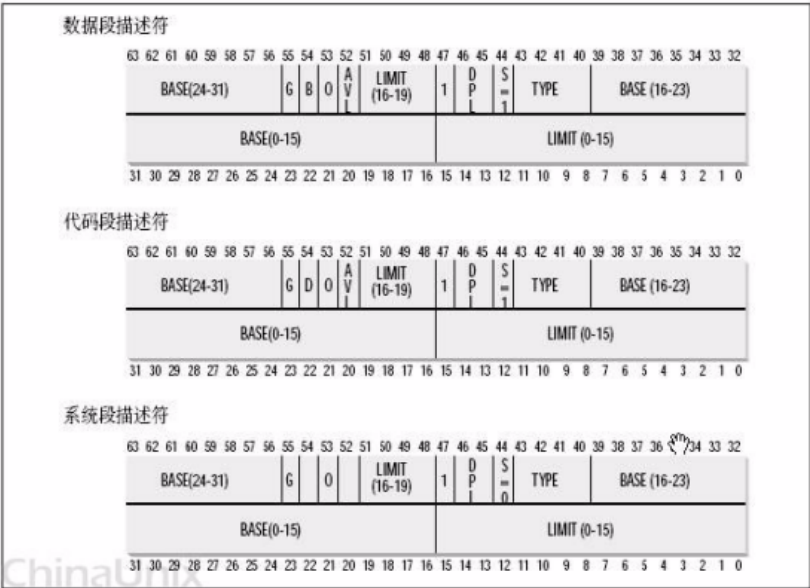
通过段标识符中的索引号从GDT或者LDT找到该段的段描述符，段描述符中的base字段是段的起始地址

段描述符：Base字段，它描述了一个段的开始位置的线性地址。

一些全局的段描述符，就放在“全局段描述符表(GDT)”中，一些局部的，例如每个进程自己的，就放在所谓的“局部段描述符表(LDT)”中。

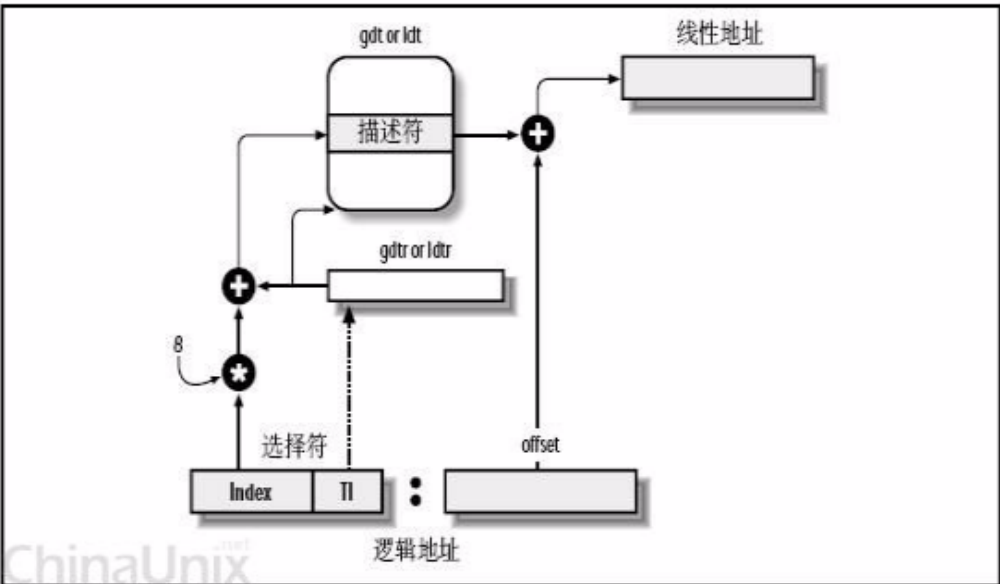
GDT在内存中的地址和大小存放在CPU的gdtr控制寄存器中，而LDT则在ldtr寄存器中。

段起始地址+ 段内偏移量 = 线性地址



首先，给定一个完整的逻辑地址[段选择符：段内偏移地址]，

- 1、看段选择符的TI=0还是1，知道当前要转换是GDT中的段，还是LDT中的段，再根据相应寄存器，得到其地址和大小。我们就有了一个数组了。
- 2、拿出段选择符中前13位，可以在这个数组中，查找到对应的段描述符，这样，它的Base，即基地址就知道了。
- 3、把Base + offset，就是要转换的线性地址了。



第一步：页式管理——线性地址转物理地址

再利用其页式内存管理单元，转换为最终物理地址。

linux假的段式管理

Intel要求两次转换，这样虽说是兼容了，但是却是很冗余，但是这是intel硬件的要求。

其它某些硬件平台，没有二次转换的概念，Linux也需要提供一个高层抽象，来提供一个统一的界面。

所以，Linux的段式管理，事实上只是“哄骗”了一下硬件而已。

按照Intel的本意，全局的用GDT，每个进程自己的用LDT——不过Linux则对所有的进程都使用了相同的段来对指令和数据寻址。即用户数据段，用户代码段，对应的，内核中是内核数据段和内核代码段。

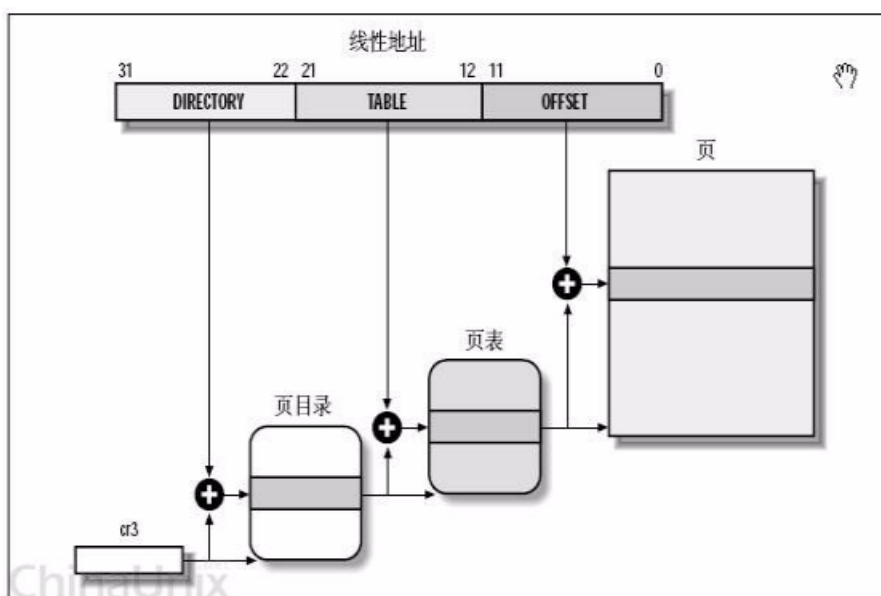
在Linux下，逻辑地址与线性地址总是一致的，即逻辑地址的偏移量字段的值与线性地址的值总是相同的。

linux页式管理

CPU的页式内存管理单元，负责把一个线性地址，最终翻译为一个物理地址。

线性地址被分为以固定长度为单位的组，称为页(page)，例如一个32位的机器，线性地址最大可为4G，可以用4KB为一个页来划分，这页，整个线性地址就被划分为一个total_page[2^{20}]的大数组，共有 2^{20} 个页。

另一类“页”，我们称之为物理页，或者是页框、页帧的。是分页单元把所有的物理内存也划分为固定长度的管理单位，它的长度一般与内存页是一一对应的。



每个进程都有自己的页目录，当进程处于运行态的时候，其页目录地址存放在cr3寄存器中。

每一个32位的线性地址被划分为三部份，【页目录索引(10位): 页表索引(10位): 页内偏移(12位)】

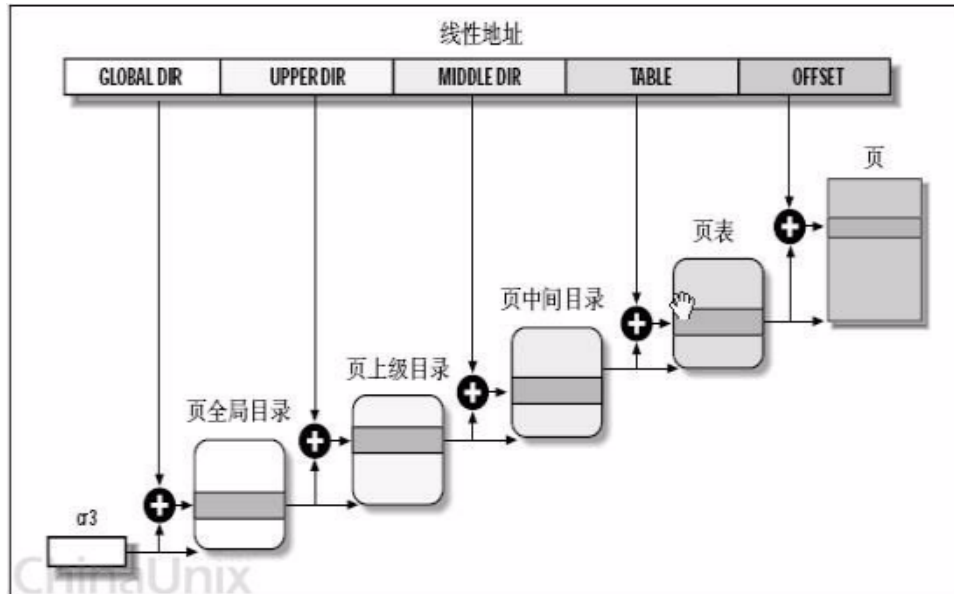
依据以下步骤进行转换：

从cr3中取出进程的页目录地址（操作系统负责在调度进程的时候，把这个地址装入对应寄存器）；

根据线性地址前十位，在数组中，找到对应的索引项，因为引入了二级管理模式，页目录中的项，不再是页的地址，而是一个页表的地址。（又引入了一个数组），页的地址被放到页表中去了。根据线性地址的中间十位，在页表（也是数组）中找到页的起始地址；将页的起始地址与线性地址中最后12位相加。

目的：

内存节约：如果一级页表中的一个页表条目为空，那么那所指的二级页表就根本不会存在。这表现出一种巨大的潜在节约，因为对于一个典型的程序，4GB虚拟地址空间的大部份都会是未分配的；



32位, PGD = 10bit, PUD = PMD = 0, table = 10bit, offset = 12bit

64位, PUD和PMD ≠ 0

22.请你说一说操作系统中的结构体对齐，字节对齐

参考回答：

1、原因：

- 平台原因（移植原因）：不是所有的硬件平台都能访问任意地址上的任意数据的；某些硬件平台只能在某些地址处取某些特定类型的数据，否则抛出硬件异常。
- 性能原因：数据结构（尤其是栈）应该尽可能地在自然边界上对齐。原因在于，为了访问未对齐的内存，处理器需要作两次内存访问；而对齐的内存访问仅需要一次访问。

2、规则

- 数据成员对齐规则：结构(struct)(或联合(union))的数据成员，第一个数据成员放在offset为0的地方，以后每个数据成员的对齐按照#pragma pack指定的数值和这个数据成员自身长度中，比较小的那个进行。
- 结构(或联合)的整体对齐规则：在数据成员完成各自对齐之后，结构(或联合)本身也要进行对齐，对齐将按照#pragma pack指定的数值和结构(或联合)最大数据成员长度中，比较小的那个进行。
- 结构体作为成员：如果一个结构里有某些结构体成员，则结构体成员要从其内部最大元素大小的整数倍地址开始存储。

3、定义结构体对齐

可以通过预编译命令#pragma pack(n), n=1,2,4,8,16来改变这一系数，其中的n就是指定的“对齐系数”。

4、举例

```
#pragma pack(2)
struct AA {
```



```
int a;          //长度4 > 2 按2对齐; 偏移量为0; 存放位置区间[0,3]
char b;        //长度1 < 2 按1对齐; 偏移量为4; 存放位置区间[4]
short c;       //长度2 = 2 按2对齐; 偏移量要提升到2的倍数6; 存放位置区间[6,7]
char d;        //长度1 < 2 按1对齐; 偏移量为7; 存放位置区间[8]; 共九字节
};
#pragma pack()
```

23.请你说一下虚拟内存置换的方式

参考回答:

比较常见的内存替换算法有: FIFO, LRU, LFU, LRU-K, 2Q。

1、FIFO (先进先出淘汰算法)

思想: 最近刚访问的, 将来访问的可能性比较大。

实现: 使用一个队列, 新加入的页面放入队尾, 每次淘汰队首的页面, 即最先进入的数据, 最先被淘汰。

弊端: 无法体现页面冷热信息

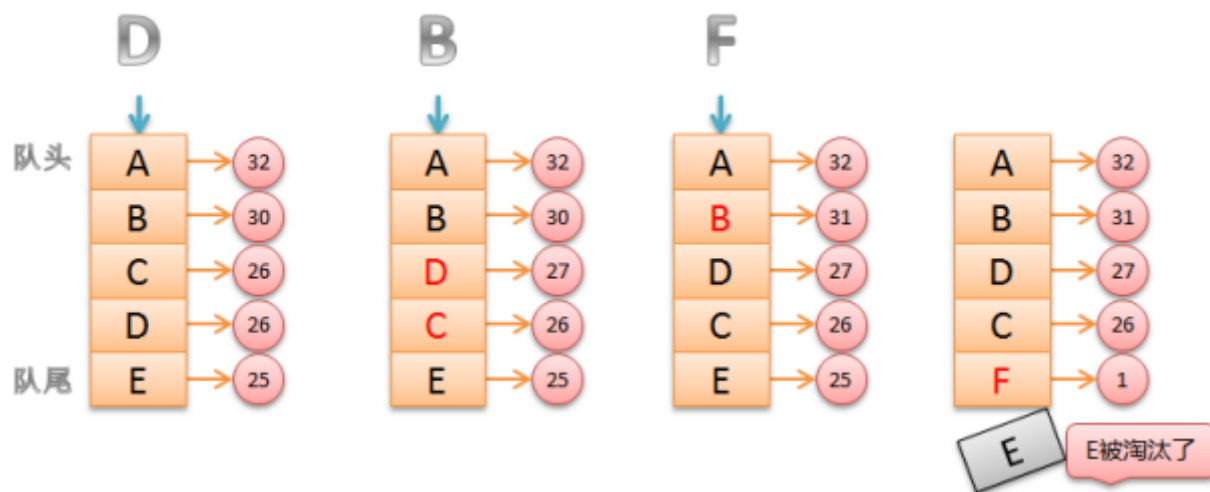
2、LFU (最不经常访问淘汰算法)

思想: 如果数据过去被访问多次, 那么将来被访问的频率也更高。

实现: 每个数据块一个引用计数, 所有数据块按照引用计数排序, 具有相同引用计数的数据块则按照时间排序。每次淘汰队尾数据块。

开销: 排序开销。

弊端: 缓存颠簸。



3、LRU (最近最少使用替换算法)

思想: 如果数据最近被访问过, 那么将来被访问的几率也更高。

实现: 使用一个栈, 新页面或者命中的页面则将该页面移动到栈底, 每次替换栈顶的缓存页面。

优点: LRU算法对热点数据命中率是很高的。

缺陷:

- 1) 缓存颠簸, 当缓存 (1, 2, 3) 满了, 之后数据访问 (0, 3, 2, 1, 0, 3, 2, 1...)。
- 2) 缓存污染, 突然大量偶发性的数据访问, 会让内存中存放大量冷数据。

4、LRU-K (LRU-2、LRU-3)

思想: 最久未使用K次淘汰算法。

LRU-K中的K代表最近使用的次数, 因此LRU可以认为是LRU-1。LRU-K的主要目的是为了解决LRU算法“缓存污染”的问题, 其核心思想是将“最近使用过1次”的判断标准扩展为“最近使用过K次”。

相比LRU，LRU-K需要多维护一个队列，用于记录所有缓存数据被访问的历史。只有当数据的访问次数达到K次的时候，才将数据放入缓存。当需要淘汰数据时，LRU-K会淘汰第K次访问时间距当前时间最大的数据。

实现：

- 1) 数据第一次被访问，加入到访问历史列表；
- 2) 如果数据在访问历史列表里后没有达到K次访问，则按照一定规则（FIFO，LRU）淘汰；
- 3) 当访问历史队列中的数据访问次数达到K次后，将数据索引从历史队列删除，将数据移到缓存队列中，并缓存此数据，缓存队列重新按照时间排序；
- 4) 缓存数据队列中被再次访问后，重新排序；
- 5) 需要淘汰数据时，淘汰缓存队列中排在末尾的数据，即：淘汰“倒数第K次访问离现在最久”的数据。

针对问题：

LRU-K的主要目的是为了解决LRU算法“缓存污染”的问题，其核心思想是将“最近使用过1次”的判断标准扩展为“最近使用过K次”。

5、2Q

类似LRU-2。使用一个FIFO队列和一个LRU队列。

实现：

- 1) 新访问的数据插入到FIFO队列；
- 2) 如果数据在FIFO队列中一直没有被再次访问，则最终按照FIFO规则淘汰；
- 3) 如果数据在FIFO队列中被再次访问，则将数据移到LRU队列头部；
- 4) 如果数据在LRU队列再次被访问，则将数据移到LRU队列头部；
- 5) LRU队列淘汰末尾的数据。

针对问题：LRU的缓存污染

弊端：

当FIFO容量为2时，访问负载是：ABCABCABC会退化为FIFO，用不到LRU。

24.请你说一下多线程，线程同步的几种方式

参考回答：

概念：

进程是对运行时程序的封装，是系统进行资源调度和分配的基本单位，实现了操作系统的并发；线程是进程的子任务，是CPU调度和分派的基本单位，用于保证程序的实时性，实现进程内部的并发；线程是操作系统可识别的最小执行和调度单位。每个线程都独自占用一个虚拟处理器：独自の寄存器组，指令计数器和处理器状态。每个线程完成不同的任务，但是共享同一地址空间（也就是同样的动态内存，映射文件，目标代码等等），打开的文件队列和其他内核资源。

线程间通信的方式：

1、临界区：

通过多线程的串行化来访问公共资源或一段代码，速度快，适合控制数据访问；

2、互斥量 Synchronized/Lock：

采用互斥对象机制，只有拥有互斥对象的线程才有访问公共资源的权限。因为互斥对象只有一个，所以可以保证公共资源不会被多个线程同时访问

3、信号量 Semaphore：

为控制具有有限数量的用户资源而设计的，它允许多个线程在同一时刻去访问同一个资源，但一般需要限制同一时刻访问此资源的最大线程数目。

4、事件(信号), Wait/Notify:

通过通知操作的方式来保持多线程同步，还可以方便的实现多线程优先级的比较操作