



进程、线程与纤程

郑晨

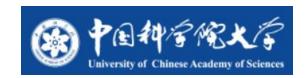
改编声明

- 本课程教学及PPT内容基于上海交通大学并行与分布式系统研究所发布的操作系统课程修改,原课程官网:
 - https://ipads.se.sjtu.edu.cn/courses/os/index.shtml
- 本课程修改人为中国科学院软件研究所,用于国科大操作系统课程教学。





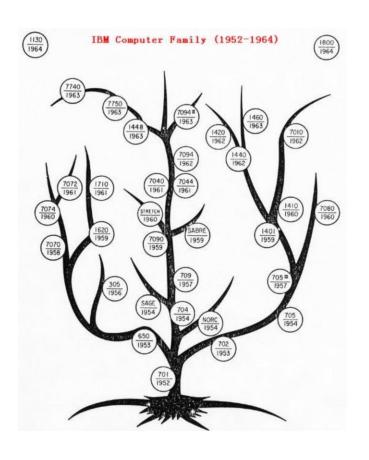




进程机制的引入

进程的诞生和概念 - 进程的状态 - 数据结构 - 基本操作

早期的IBM计算机:批处理作业

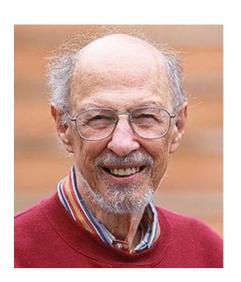


- 常驻监督程序: 用来自动装填纸卡
- 客户为IBM机器开发操作系统
 - 通用汽车、通用电气 /MIT
- 一次处理一个程序
 - 调试时间以天为单位

IBM 7090/7094:分时共享

- · CompatibleTimeSharingSystem(CTSS),在IBM 7090机器上第一次实现了多个程序同时运行,"进程"开始登上历史舞台
- FernandoJ.Corbató
- 1990年图灵奖





单道编程VS多道编程

- · 单道编程:同一时刻只有一个线程
 - MS/DOS, 早期Macintosh, 批处理
 - 简化操作系统设计与实现
 - 去除了并发性 (only one thread accessing resources!)

- ・ 多道编程:同一时刻有多个执行线程
 - Multics, UNIX/Linux, OS/2, Windows NT-8, Mac OSX, Android, iOS
 - 增加了操作系统设计与实现的难度
 - 高并发性

多道编程的挑战

· 需求

- virtual machine abstraction
 - 每个程序都期望独占机器资源
- concurrency
 - 通过程序间并发提高硬件资源的利用
- protection
 - 实现对共享资源的仲裁和独占资源的保护
- Coordination
 - 程序间需要能够彼此通信合作

• 通过进程实现

- 多道并发
- 隔离保护
- 资源共享
- 调度协作

进程定义

· 什么是进程

- 程序 (Program) 执行的一个实例 (instance)

```
Processes: 607 total, 2 running, 605 sleeping, 3679 threads
                                                                        17:21:18
Load Avg: 2.10, 2.89, 3.48 CPU usage: 5.1% user, 6.60% sys, 88.38% idle
SharedLibs: 202M resident, 38M data, 57M linkedit.
MemRegions: 374499 total, 4774M resident, 154M private, 2501M shared.
PhysMem: 16G used (3413M wired), 43M unused.
VM: 3739G vsize, 1991M framework vsize, 22484911(0) swapins, 25909406(0) swapout
Networks: packets: 9383767/12G in, 5335881/3117M out.
Disks: 10326554/244G read, 5007684/171G written.
                    %CPU TIME
                                             #PORT MEM
PID
       COMMAND
                                  #TH
                                                          PURG
                                                                 CMPRS PGRP
       svsmond
                    10.2 13:41.94 3
                                                   3776K+ ØB
                                             29+
                                                                  628K 415
22166 top
                    8.1 00:02.20 1/1
                                                   8712K ØB
                                                                 ØB
                                                                        22166
       kernel_task 7.0 88:39.79 275/4 0
                                                   612M+
                                                          0B
                                                                 0B
                                                                        0
2231
       iTerm2
                    4.5 03:33.03 12
                                                   114M+
                                                          5072K- 45M-
                                                                      2231
                                             563
21098
      Microsoft Po 4.0 03:58.69 38
                                        10
                                             1671+ 382M+
                                                                  85M
                                                                        21098
174
                    3.3 22:37.34 6
       hidd
                                                   4584K
                                                          ØB
                                                                  1508K 174
270
       WindowServer 1.1 02:47:45 9
                                                                  179M 270
                                             6191- 1031M- 121M
       Activity Mon 0.9 05:13.99 5
1294
                                             3192+ 162M+
                                                          0B
                                                                  125M
                                                                       1294
1287
       Google Chrom 0.4 89:09.43 34
                                             1834 811M
                                                          676K
                                                                  355M
                                                                      1287
       mysald
                    0.3 02:50.35 39
                                                   362M
                                                          ØB
                                                                  355M
                                                                      1444
       Google Chrom 0.2 02:05.37 15
1390
                                             187
                                                   78M
                                                                  58M
                                                                       1287
2206
       Google Chrom 0.1 02:20.87 21
                                                                      1287
                                             254
                                                   361M
                                                          0B
                                                                  225M
       thunderbird 0.1 11:04.27 45
                                                                 305M 11453
11453
                                             1730
                                                  627M
                                                          28M
11484
       Google Chrom 0.1 00:38.43 13
                                                                       1287
                                             140
                                                   79M+
                                                          0B
                                                                       1546<mark>9</mark>
15469
       IINA
                    0.1 06:52.55 20
                                             510
                                                   118M
                                                          192K
                                                                  63M
```

用户,程序,进程

- · 用户在系统中有账户
- · 用户加载程序
 - 多个用户可以加载相同程序吗?
 - 一个用户可以加载相同程序的多个实例吗?

- · 进程是程序 (Program) 执行的一个实例 (instance)
 - 两个实例之间的关系是?

程序实例

• 静态变量的地址总是相同的,但其值却是不同的

- 结论: 地址不是绝对的
 - 每个进程有自己独立的内存地址空间
- 优势:
 - Compiler/linker/loader不需要关心具体地址
 - 地址空间可以大于实际内存

进程 VS 程序

· 程序不是进程

- 程序是静态的
 - 代码+数据
- 进程是动态的
 - instruction execution + data + OS resource + more

• 程序vs进程: 不是1对1的映射关系

- 进程不只有代码与数据
- 一个程序可能有多个同时执行的进程
- 一个进程可以运行多个程序

通过进程实现的两种虚拟机制

- · CPU虚拟
 - 通过上下文(Context)环境实现对CPU的虚拟

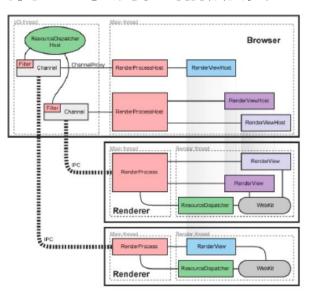
- 内存虚拟
 - 通过地址空间(Address)机制实现对对内存系统的虚拟

- · 两种虚拟机制共同构筑了virtual machine abstraction
 - 多道并发、隔离保护、资源共享、调度协作

Example: Thee Chrome Browser

- · 多个进程: 各个插件、标签页
- · 如果一个页面崩溃,不会使得整个浏览器崩溃





进程的内涵与外延

• 进程是一种操作系统对于执行的抽象

- 是一种执行的实体
- 是系统调度的单位
- 是一个程序的动态执行上下文

· 程序执行的实体

- Register: PC, SP、X1-X31
- Memory: code, data, stack, heap

• 资源分配的载体

- 拥有独立的地址空间:memory (address space), file descriptors, file system context, ...
- 拥有独立I/O 状态:file descriptor table, network sockets

• 程序协调的机制

- 进程间通信IPC、pipe、socket

进程

进程的诞生和概念 - 进程的状态 - 数据结构 - 基本操作

进程:运行中的程序

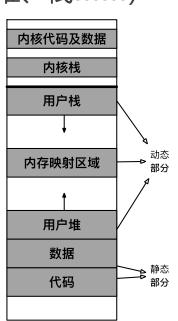
• 进程是计算机程序运行时的抽象

- 静态部分:程序运行需要的代码和数据

- 动态部分:程序运行期间的<mark>状态</mark>(程序计数器、堆、栈.....)

• 进程具有独立的虚拟地址空间

- 每个进程都具有"独占全部内存"的假象
- 内核中同样包含内核栈和内核代码、数据



如何表示进程:进程控制块 (PCB)

- · 每个进程都对应一个元数据, 称为"进程控制块" PCB
 - 进程控制块存储在内核态(为什么?)

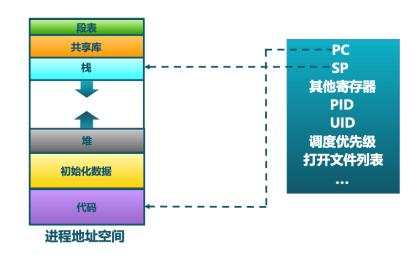
- · 想一想: 进程控制块里至少应该保存哪些信息?
 - 独立的虚拟地址空间
 - 独立的执行上下文

进程的组成

- · 一个进程 (PCB) 包括程序执行过程中的所有状态
 - 内存地址空间
 - 执行程序的代码、数据
 - 执行stack: 包含所有调用的状态
 - Program counter (PC): 指向下一个指令
 - 通用寄存器值
 - 操作系统资源集合: 打开的文件、网络链接等

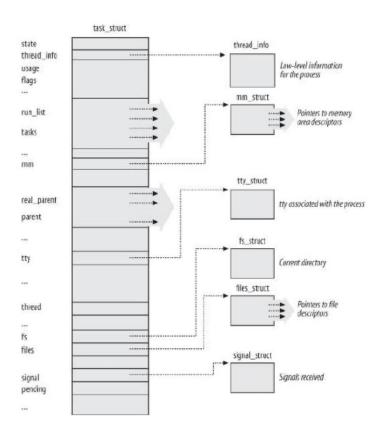
PCB

- 记录进程自身状态,包括
 - 进程状态(new, ready,running, waiting,...)
 - 机器状态(e.g., CPU register, TLB)
 - 调度与统计信息
 - 内存管理信息
 - I/O状态信息



- 将CPU硬件状态从一个进程切换到另一个进程——进程上下文切换(Context Switch)
 - ❖ 一秒钟可以发生100到1000次
- 调度时: 进程放弃CPU给其他进程
 - ❖ 同一时刻只有一个进程在运行
 - ❖ 给重要讲程更多的运行时间

PCB in Linux: task_struct



• Linux 5.5.10中670LOC

task_struct剖析 (1)

```
struct task_struct {
                      /* -1 unrunnable, 0 runnable, >0 stopped */
    volatile long state;
    //...
    long exit_state;
    //...
   进程的状态宏定义:
      #define TASK_RUNNING
      #define TASK_INTERRUPTIBLE
      #define TASK UNINTERRUPTIBLE
      #define __TASK_STOPPED
      #define __TASK_TRACED
      /* in tsk->exit_state */
     #define EXIT_ZOMBIE
     #define EXIT_DEAD
   //...
```

为什么要这么定义?

task_struct剖析 (2)

```
struct task_struct {
   //...
   struct list_head tasks;  //将系统中所有进程通过双向链表链接起来!
   //...
       怎样访问所有的进程呢?
     #define for_each_process(p) \
              for (p = &init task; (p = next task(p)) != &init task; )
     #define next_task(p) \
               list_entry_rcu((p)->tasks.next, struct task_struct, tasks)
```

task_struct剖析 (3)

系统调用 getpid() 返回什么?

Linux系统允许用户使用一个叫做进程标识符的PID来标识进程,PID顺序编号,新创建进程是前一个进程的PID加1,不过PID值有一个上限,达到上限之后再开始循环使用闲置的小PID。

task_struct剖析 (4)

```
struct task_struct {
   //...
   struct task_struct __rcu *real_parent; /* real parent process */
   struct task_struct __rcu *parent; /* recipient of SIGCHLD, wait4() reports */
   struct list head children; /* list of my children */
                                /* linkage in my parent's children list */
   struct list head sibling;
   //...
              进程之间的关系:
                        父子关系
                        兄弟关系
```

current宏定义解析

- · current宏: 一个全局指针,指向当前进程的struct task_struct结构体,即表示当前进程。
- 例如current->pid就能得到当前进程的pid, currentcomm就能得到当前进程的名称。



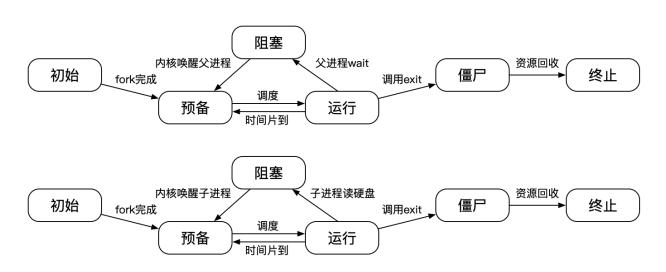
简化模型:单一线程的进程

- 假设每个进程都只有一个线程
 - 历史上确实如此! (如早期的UNIX操作系统)

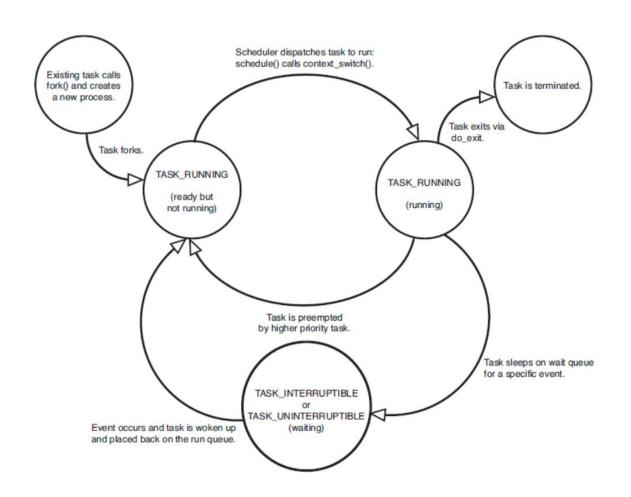
- · 好处: 便于理解操作系统中的各种概念
 - 多线程使操作系统的管理更加复杂
 - 在单一线程的进程中: 线程管理/调度≈进程管理/调度
 - 线程的内容将在后面介绍

进程管理: 即管理进程的生命周期

- 进程自创建到终止可经历多个过程
 - 称为进程状态
- 不同的系统调用和事件会影响进程的状态

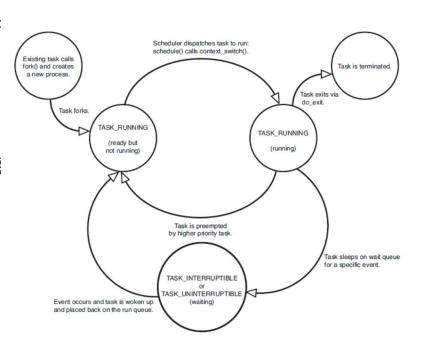


Linux进程5种状态变换流图



Linux进程状态

- TASK_RUNNING(R)
 - 可运行,但不等于正在运行,可能在运行队列中
- TASK_INTERRUPTIBLE(S)
 - 可被中断的睡眠状态,可被信号唤醒
- TASK_UNINTERRUPTIBLE(D)
 - 不可被中断的睡眠状态,俗称"D住了",不能被信号唤醒
- TASK TRACED
 - 被跟踪状态,调试用
- __TASK_STOPPED
 - 停止执行状态,通常在于收到了SIGSTOP等状态



Linux进程状态

Tasks	: 171 tota): 0.1%us 16467276k	(1), 8) (1	1 : 0.1%s :al,	running sy, 0	g, 170 .0%ni, 552k u) sle . 99.	epi 8%i 2	ng, d, (3076)	0 st 0.0%wa 2 4 k fr	opped, , 0.0%hi ee, 171:	0.06, 0.07, 0.05 0 zombie , 0.0%si, 0.0%st 168k buffers 340k cached
PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
14677	voelker	20	0	55548	3232	2364	R	0	0.0	0:00.07	top
24637	voelker	20		86300	6364	1024	S	0	0.0		mosh-server
1	root	20		57812		584	S		0.0	1:26.73	
2	root	20	0	0	0	0	S	0	0.0		kthreadd
3	root	RT	0	0	0	0	S	0	0.0		migration/0
4	root	20	0	0	0	0	S	0	0.0		ksoftirqd/0
5	root	RT	0	0	0	0	S	0	0.0		watchdog/0
	root	RT	0	0	0	0	S	0	0.0		migration/1
7	root	20	0	0	0	0	S	0	0.0		ksoftirqd/1
8	root	RT	0	0	0	0	S	0	0.0	0:00.01	watchdog/1
9	root	RT	0	0	0	0	S	0	0.0		migration/2
10	root	20	0	0	0	0	S	0	0.0	9:44.37	ksoftirqd/2
11	root	RT	0	0	0	0	S	0	0.0	0:00.01	watchdog/2
12	root	RT	0	0	0	0	S	0	0.0	0:18.06	migration/3
13	root	20	0	0	0	0	S	0	0.0	9:01.67	ksoftirqd/3
14	root	RT	0	0	0	0	S	0	0.0	0:00.01	watchdog/3
15	root	20	0	0	0	0	S	0	0.0	2:30.99	events/0

进程创建: fork()

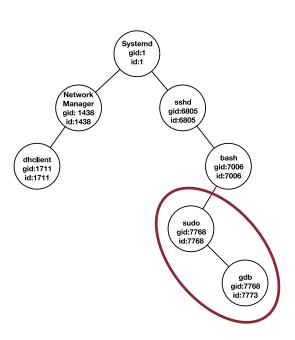
- 语义: 为调用进程创建一个一模一样的新进程
 - 调用进程为**父进程**,新进程为**子进程**
 - 接口简单,无需任何参数

- · fork后的两个进程均为独立进程
 - 拥有不同的进程id
 - 可以并行执行, 互不干扰 (除非使用特定的接口)
 - 父进程和子进程会共享部分数据结构(内存、文件等)

进程树与进程组

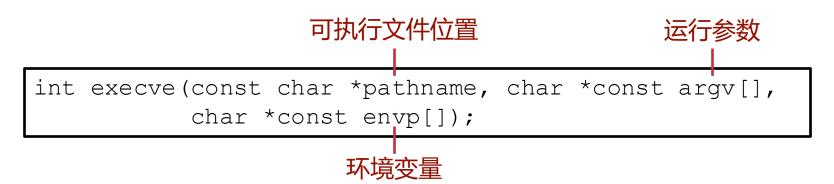
- · fork为进程之间建立了父进程和子进程的关系
 - 进程之间建立了树型结构
 - Linux可使用ps命令查看

- · 多个进程可以属于同一个进程组
 - 子进程默认与父进程属于同一个进程组
 - 可以向同一进程组中的所有进程发送信号
 - 主要用于shell程序中



进程的执行: exec

· 为进程指定可执行文件和参数



- ・ 在fork之后调用
 - exec在载入可执行文件后会重置地址空间

写时拷贝 (Copy-On-Write)

- · 早期的fork实现:将父进程直接拷贝一份
 - 性能差: 时间随占用内存增加而增加
 - 无用功: fork之后如果调用exec, 拷贝的内存就作废了

- · 基本思路: 只拷贝内存映射, 不拷贝实际内存
 - 性能较好: 一条映射至少对应一个4K的页面
 - 调用exec的情况里,减少了无用的拷贝

fork的优缺点分析

· fork的优点

- 接口非常简洁
- 将进程"创建"和"执行"(exec)解耦,提高了灵活度
- 刻画了进程之间的内在关系(进程树、进程组)

· fork的缺点

- 完全拷贝过于粗暴(不如clone)
- 性能差、可扩展性差(不如vfork和spawn)
- 不可组合性 (例如: fork() + pthread())

fork的替代接口

- · vfork: 类似于fork, 但让父子进程共享同一地址空间
 - 优点: 连映射都不需要拷贝, 性能更好
 - 缺点:
 - 只能用在"fork + exec"的场景中
 - 共享地址空间存在安全问题
- · 轶事: vfork的提出最初就是为了解决fork的性能问题
 - 但写时拷贝拯救了fork

Since this function is hard to use correctly from application software, it is recommended to use posix_spawn(3) or fork(2) instead.

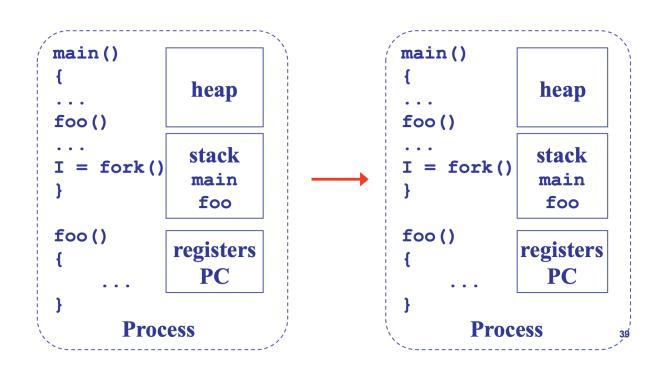
fork的替代接口

- posix_spawn: 相当于fork + exec
 - 优点:可扩展性、性能较好
 - 缺点:不如fork灵活

- · clone: fork的"进阶版",可以选择性地不拷贝内存
 - 优点: 高度可控, 可依照需求调整
 - 缺点:接口比fork复杂,选择性拷贝容易出错

进程API

• Fork() system call 创建当前进程的副本



Fork系统调用

Fork()

- 创建和初始化一个新的PCB
- 创建一个新的地址空间(但,复制父进程内容)
- 初始化地址空间: 复制父进程的整个地址空间内容
- 初始化内核资源(指向父进程使用的资源,如open files)
- 将PCB放于ready队列

Fork return twice

- 将子进程的pid返回父进程,返回 "0" 给子进程

Fork()

```
int main(int argc, char *argv[])
  char *name = argv[0];
  int child pid = fork();
  if (child pid == 0) {
      printf("Child of %s is %d\n", name, getpid());
      return 0;
  } else {
      printf("My child is %d\n", child pid);
      return 0;
```

• 这段代码会输出什么?

Fork()

```
alpenglow (18) ~/tmp> cc t.c
alpenglow (19) ~/tmp> a.out
My child is 486
Child of a.out is 486
```

复制地址空间

Parent

```
child_pid = 486
                                         child_pid = 0
child pid = fork();
                                  child_pid = fork();
                                                              PC
if (child pid == 0) {
                                  if (child pid == 0) {
  printf("child");
                                    printf("child");
                                  } else {
} else {
                                    printf("parent");
  printf("parent");
```

Child

差异

```
child_pid = 486
                                                 child_pid = 0
        child_pid = fork();
                                          child_pid = fork();
        if (child_pid == 0) {
                                          if (child_pid == 0) {
          printf("child");
                                            printf("child");
                                                                      PC
        } else {
                                            else {
                                            printf("parent");
PC
          printf("parent");
```

Parent Child

Fork()

```
alpenglow (18) ~/tmp> cc t.c
alpenglow (19) ~/tmp> a.out
My child is 486
Child of a.out is 486
alpenglow (20) ~/tmp> a.out
Child of a.out is 498
My child is 498
```

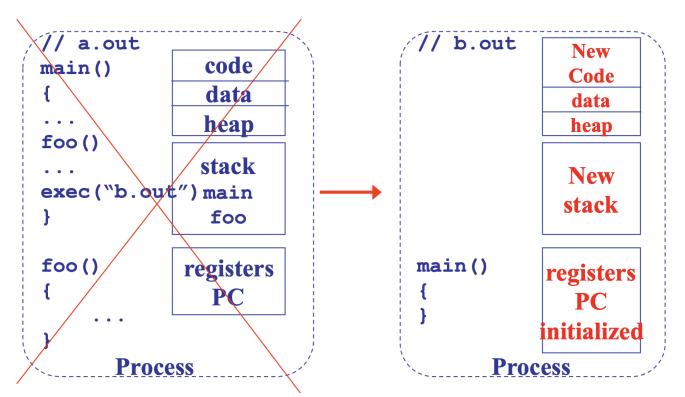
■ 为什么顺序不同?

Why fork()

- · 父子进程协同工作
- 依赖父进程的数据完成任务

```
while (1) {
   int sock = accept();
   if ((child_pid = fork()) == 0) {
        Handle client request and exit
   } else {
        Close socket
   }
}
```

Exec()



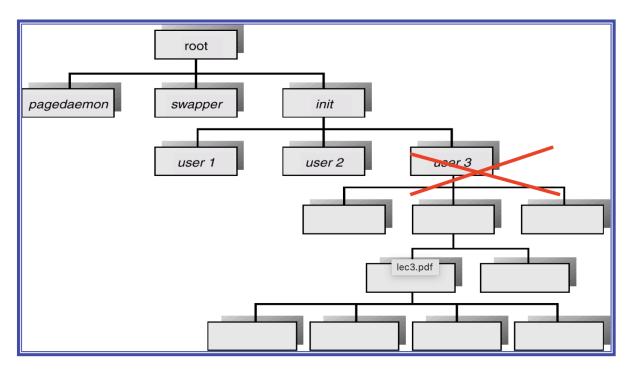
· 执行exec()后,进程的什么没有改变?

Exit()

- · Exit()系统调用
 - 结束进程执行
 - 主要是为了释放资源

- · Exit()执行
 - 终结所有线程(下一节)
 - 关闭所有文件描述符、网络链接
 - 释放占用的内存 (并将内存页换出到磁盘)
 - 删除PCB
- · 为什么需要OS作这些操作?

Linux系统进程树



- · 如果一个父进程先于子进程消失会发生什么?
 - Orphaned children

Orphaned children in Linux

- Orphaned children process will set process #1 as parent
- From Linux 3.4
 - 进程可以触发prctl()系统调用, withPR_SET_CHILD_SUBREAPER
 - Orphaned children proces将会以最近的ancestor process
 为父进程,不一定为process #1

线程

线程的概念 - 线程模型 - 相关数据结构 - 基本操作

为什么需要线程?

- 创建进程的开销较大
 - 包括了数据、代码、堆、栈等

- 进程的隔离性过强
 - 进程间交互:可以通过进程间通信 (IPC),但开销较大

• 进程内部无法支持并行

线程: 更加轻量级的运行时抽象

- 线程只包含运行时的状态
 - 静态部分由进程提供
 - 包括了执行所需的**最小**状态(主要是寄存器和栈)

- · 一个进程可以包含多个线程
 - 每个线程共享同一地址空间(方便数据共享和交互)
 - 允许进程内并行

进阶模型:多线程的进程

· 一个进程可以包含多个线程

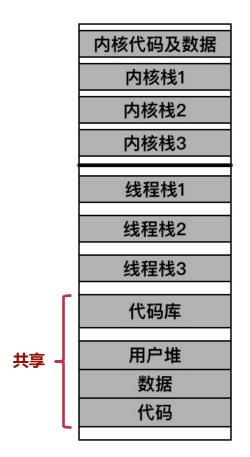
- 一个进程的多线程可以在不同处理器上同时执行
 - 调度的基本单元由进程变为了线程
 - 每个线程都有状态
 - 上下文切换的单位变为了线程

多线程进程的地址空间

· 每个线程拥有自己的栈

内核中也有为线程准备的内核栈

- ・ 其它区域共享
 - 数据、代码、堆……



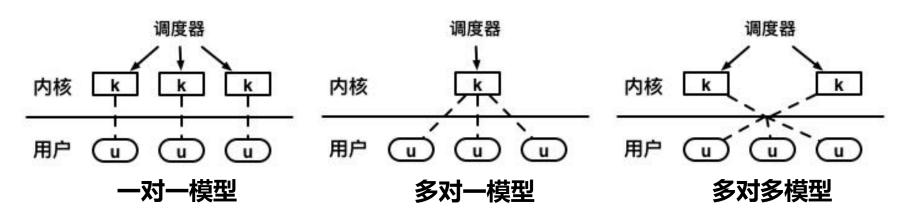
用户态线程与内核态线程

- 根据线程是否受内核管理,可以将线程分为两类
 - 内核态线程: 内核可见, 受内核管理
 - 用户态线程:内核不可见,不受内核直接管理
- 内核态线程
 - 由内核创建,线程相关信息存放在内核中
- · 用户态线程 (纤程)
 - 在应用态创建,线程相关信息主要存放在应用数据中

线程模型

• 线程模型表示了用户态线程与内核态线程之间的联系

- 多对一模型: 多个用户态线程对应一个内核态线程
- 一对一模型: 一个用户态线程对应一个内核态线程
- 多对多模型:多个用户态线程对应多个内核态线程



多对一模型

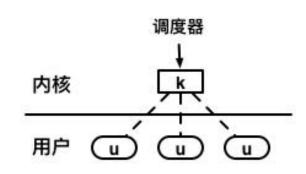
· 将多个用户态线程映射给单一的内核线程

- 优点: 内核管理简单

- 缺点: 可扩展性差, 无法适应多核机器的发展

• 在主流操作系统中被弃用

· 用于各种用户态线程库中



一对一模型

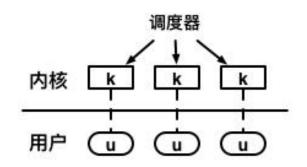
• 每个用户线程映射单独的内核线程

- 优点:解决了多对一模型中的可扩展性问题

- 缺点: 内核线程数量大, 开销大

• 主流操作系统都采用一对一模型

Windows, Linux, OS X......

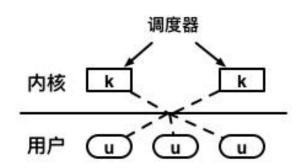


多对多模型(又叫Scheduler Activation)

- N个用户态线程映射到M个内核态线程 (N > M)
 - 优点:解决了可扩展性问题(多对一)和线程过多问题(一对一)
 - 缺点:管理更为复杂

- · Solaris在9之前使用该模型
 - 9之后改为一对一

• 在虚拟化中得到了广泛应用



线程的相关数据结构: TCB

- · 一对一模型的TCB可以分为两部分
- ・ 内核态:与PCB结构类似
 - Linux中进程与线程使用的是同一种数据结构 (task_struct)
 - 上下文切换中会使用
- · 应用态: 可以由线程库定义
 - Linux: pthread结构体
 - Windows: TIB (Thread Information Block)
 - 可以认为是内核TCB的扩展

线程本地存储 (TLS)

- · 不同线程可能会执行相同的代码
 - 线程不具有独立的地址空间, 多线程共享代码段

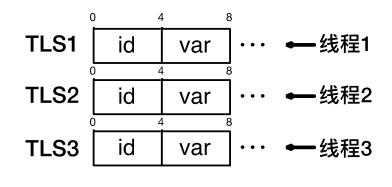
- · 问题: 对于全局变量,不同线程可能需要不同的拷贝
 - 举例:用于标明系统调用错误的errno

· 解决方案: 线程本地存储 (Thread Local Storage)

线程本地存储 (TLS)

- 线程库允许定义每个线程独有的数据
 - thread int id; 会为每个线程定义一个独有的id变量

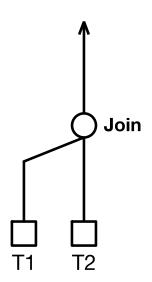
- · 每个线程的TLS结构相似
 - 可通过TCB索引
- ・ TLS寻址模式: 基地址 + 偏移量
 - X86: 段页式 (fs寄存器)
 - RISC-V: 通用寄存器tp(x4)



线程的基本操作:以pthreads为例

- · 创建: pthread_create
 - 内核态: 创建相应的内核态线程及内核栈
 - 应用态: 创建TCB、应用栈和TLS

- ・ 合并: pthread_join
 - 等待另一线程执行完成,并获取其执行结果
 - 可以认为是fork的"逆向操作"



线程的基本操作:以pthreads为例

- ・ 退出: pthread_exit
 - 可设置返回值 (会被pthread_join获取)

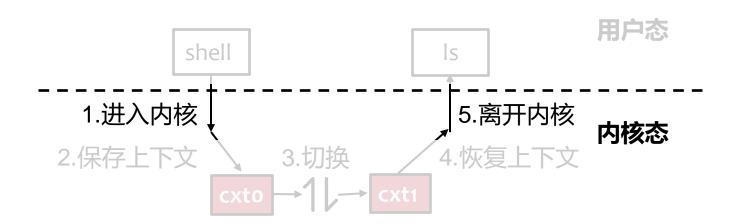
- ・ 暂停: pthread_yield
 - 立即暂停执行,出让CPU资源给其它线程
 - 好处:可以帮助调度器做出更优的决策

上下文切换

上下文组成 - 上下文切换

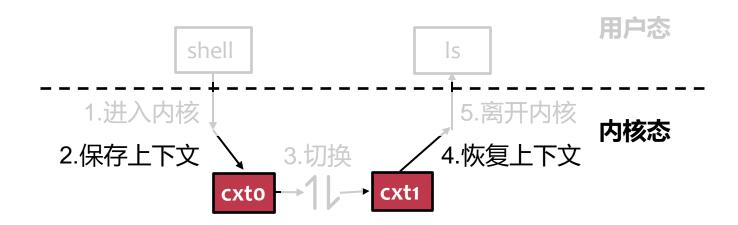
如何实现上下文切换?

• 进程怎样切换到内核中执行? 如何切换回用户态?



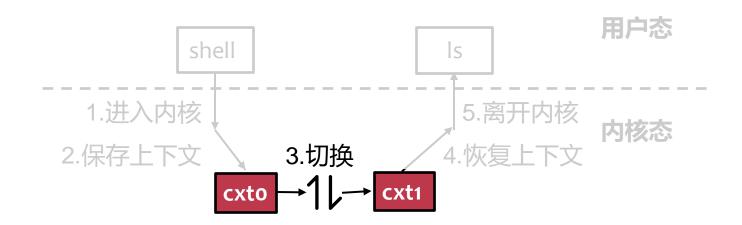
如何实现上下文切换?

- 进程怎样切换到内核中执行? 如何切换回用户态?
- · 如何对上下文进行保存和恢复?



如何实现上下文切换?

- 进程怎样切换到内核中执行? 如何切换回用户态?
- 如何对上下文进行保存和恢复?
- 如何实现关键的切换步骤?



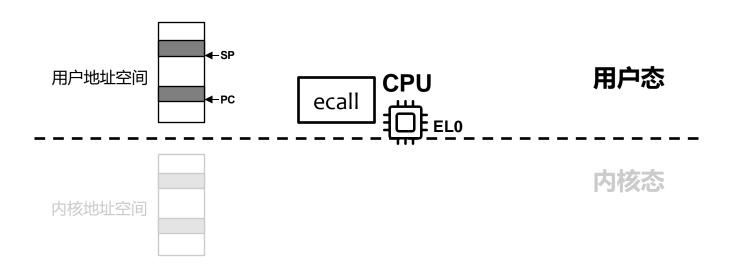
进程上下文的组成 (RISC-V)

- 进程上下文需要包含哪些内容?
 - 常规寄存器: X0-X31
 - 程序计数器 (PC)
 - 系统寄存器,在不同的特权模式下有相对应的寄存器

· 思考:为什么进程上下文只需要保存寄存器信息,而不用保存内存?

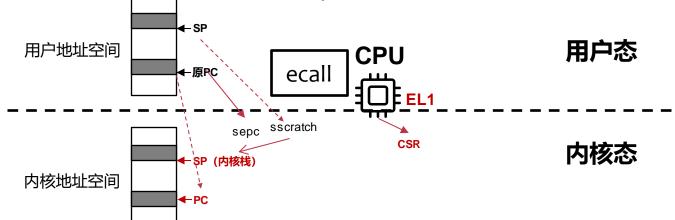
进程的内核态执行: 切换到内核态

· RISC-V提供了硬件支持,使进程切换到内核态执行

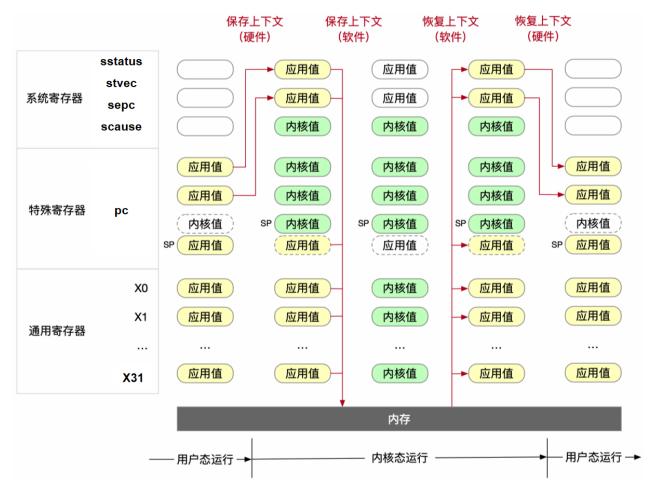


进程的内核态执行: 切换到内核态

- · RISC-V提供了硬件支持,使进程切换到内核态执行
 - 原特权级写入sstatus的 SPP等字段
 - sepc 会被修改为 Trap 处理完成后会执行的下一条指令的地址。
 - scause/stval 分别会被修改成 Trap 的原因以及相关的附加信息。
 - CPU 会跳转到 stvec 所设置的 Trap 处理入口地址,并将当前特权级设置为S,然后从Trap 处理入口地址处开始执行。



回顾: 用户态/内核态切换时的处理器状态变化



回顾: 处理器在切换过程中的任务

- 1. 将发生异常事件的指令地址保存在sepc中
- 2. 将异常事件的原因保存在scause/stval
 - 例如,是执行svc指令导致的,还是访存缺页导致的
- 3. 将处理器的当前状态保存在sstatus寄存器的SPP等字段
- 4. 将引发缺页异常的内存地址保存在stvec中
- 5. 栈寄存器不再使用SP (用户态栈寄存器) , 开始使用SP (内核态栈寄存器)
 - 内核态栈寄存器,可以使用sscratch来保存当前进程task struct的地址,异常发生时取出
- 6. 修改sstatus寄存器中的特权级标志位,设置为内核态
- 7. 找到异常处理函数的入口地址,并将该地址写入PC,开始运行操作系统
 - 根据stvec寄存器中保存的异常向量表基地址,以及发生异常事件的类型确定

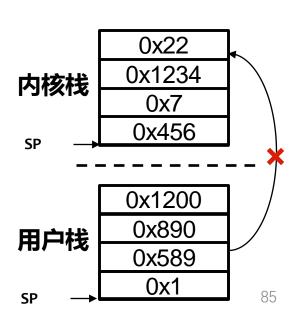
进程的内核态执行: 内核栈

· 为什么需要"又一个栈" (内核栈)?

- 进程在内核中依然执行代码,有读写临时数据的需求
- 进程在用户态和内核态的数据应该相互隔离, 增强安全性

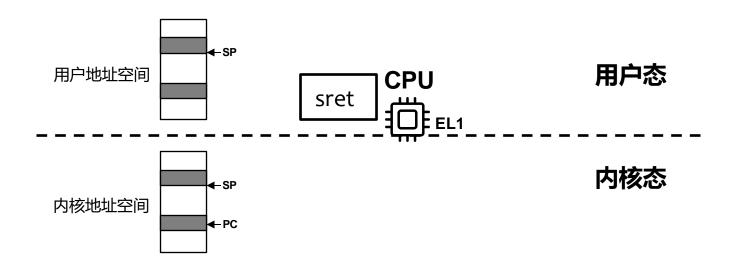
· RISC-V实现:一个栈指针寄存器

- RISC-V只有通用寄存器SP,需要保存恢复
- 从用户态进入时:使用sscratch存储内核栈地 址作为中转



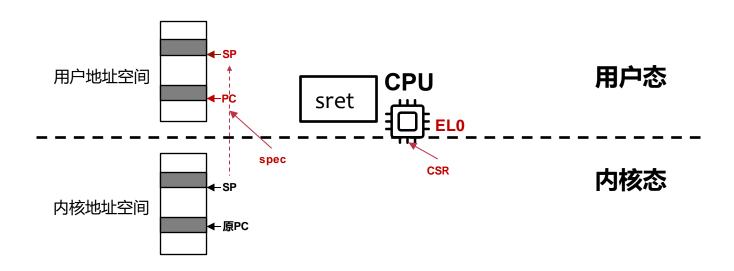
进程的内核态执行:返回用户态

· 进入内核态的"逆过程", RISC-V同样提供了硬件支持

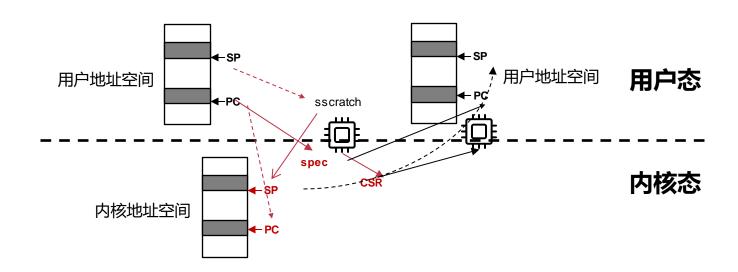


进程的内核态执行:返回用户态

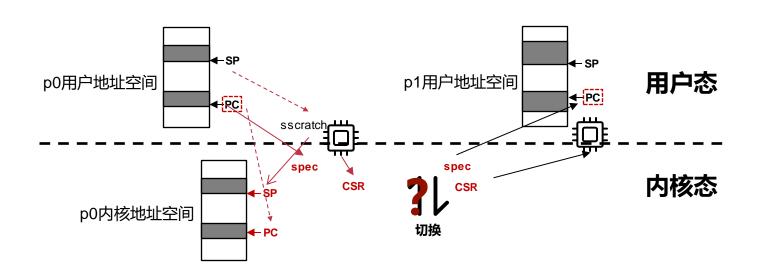
- · 进入内核态的"逆过程", RISC-V同样提供了硬件支持
 - CPU 会将当前的特权级按照 sstatus 的 SPP 字段设置为 U;
 - CPU 会跳转到 sepc 寄存器指向的那条指令,然后继续执行



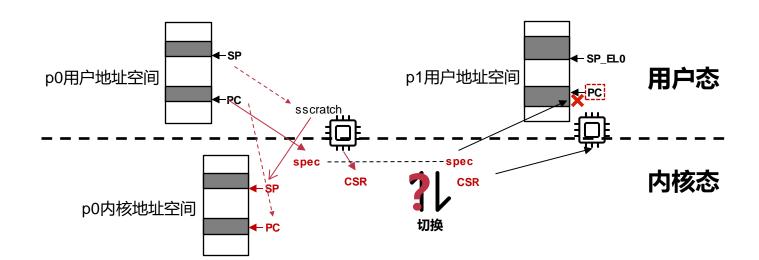
· 通过支持内核/用户态切换,可以实现进程进入内核态并返回(比如system call)



- · 但是该机制并不足以实现上下文切换
 - 不同进程地址空间不同,使用的寄存器值也不同(如PC)

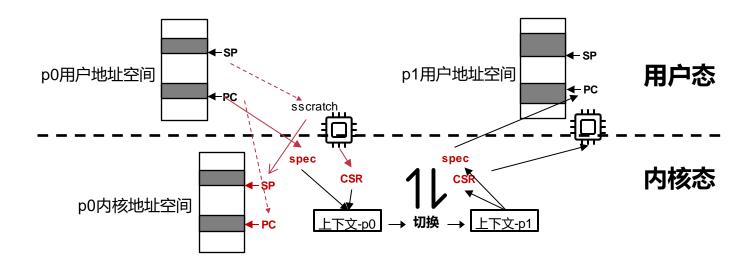


- 但是该机制并不足以实现上下文切换
 - 不同进程地址空间不同,使用的寄存器值也不同(如PC)
 - 但是:**寄存器只有一个!**直接恢复会导致错误



• 但是该机制并不足以实现上下文切换

- 不同进程地址空间不同,使用的寄存器值也不同(如PC)
- 但是:**寄存器只有一个!**直接恢复会导致错误
- 解决方法:保存**上下文**(寄存器)到内存,用于之后恢复



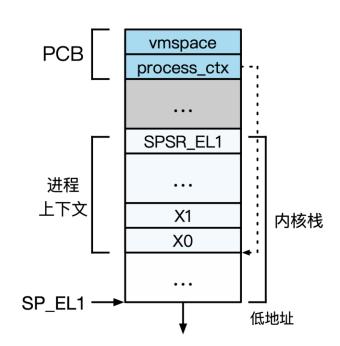
上下文与其他内核数据结构

· 与进程相关的三种内核数据结构: PCB、上下文、内核栈

高地址

· PCB保存指向上下文的引用

• 上下文的位置固定在内核栈底部



上下文保存与恢复

REG_S x5, PT_T0(sp)

save_from_x6_to_x31

发生异常时进入 handle_exception, 第一步先保存上下文 /* save all GPs except x1 ~ x5 */ .macro save_from_x6_to_x31 SYM_CODE_START(handle_exception) REG_S x6, PT_T1(sp) REG_S x7, PT_T2(sp) REG_S x8, PT_S0(sp) * If coming from userspace, preserve the user thread pointer and load REG_S x9, PT_S1(sp) * the kernel thread pointer. If we came from the kernel, the scrotch REG_S $\times 10$, PT_A0(sp) REG_S x11, PT_A1(sp) * register will contain 0, and we should continue on the current TP. REG_S x12, PT_A2(sp) REG_S $\times 13$, PT_A3(sp) REG_S x14, PT_A4(sp) csrrw tp, CSR_SCRATCH, tp $REG_S \times 15$, $PT_A5(sp)$ REG_S $\times 16$, PT_A6(sp) bnez tp. _save_context REG S x17, PT A7(sp) REG_S \times 18, PT_S2(sp) REG S x19, PT_S3(sp) REG_S x20, PT_S4(sp) REG_S x21, PT_S5(sp) REG_S x22, PT_S6(sp) **REG_S x23**, PT_S7(sp) _save_context: REG_S x24, PT_S8(sp) REG S x25. PT S9(sp) REG_S sp, TASK_TI_USER_SP(tp) REG_S x26, PT_S10(sp) **REG_S x27**, PT_S11(sp) REG_L sp, TASK_TI_KERNEL_SP(tp) REG_S $\times 28$, PT_T3(sp) addi sp, sp, -(PT_SIZE_ON_STACK) REG_S x29, PT_T4(sp) REG_S $\times 30$, PT_T5(sp) REG_S x1, PT_RA(sp) REG_S $\times 31$, PT_T6(sp) .endm $REG_S x3$, $PT_GP(sp)$

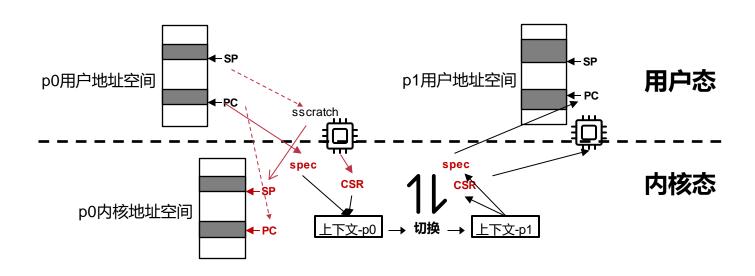
arch/riscv/include/asm/asm.h

上下文保存与恢复

```
SYM_CODE_START_NOALIGN(ret_from_exception)
       REG_L s0, PT_STATUS(sp)
        csrw CSR_STATUS, a0
        csrw CSR_EPC, a2
                                  恢复上下文类似,只是反过
                                  来操作
        REG_L x1, PT_RA(sp)
        REG_L x3, PT_GP(sp)
        REG_L x4, PT_TP(sp)
        REG_L x5, PT_T0(sp)
        restore_from_x6_to_x31
        REG_L x2, PT_SP(sp)
```

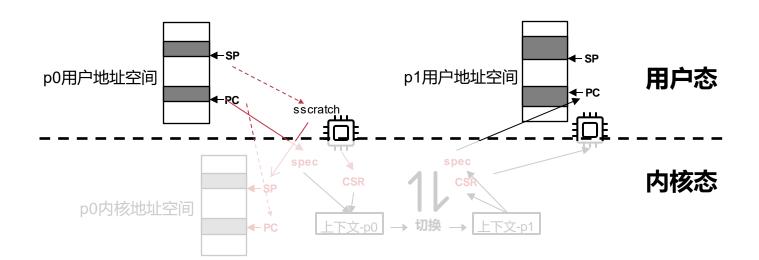
万事俱备,只欠切换!

- · 上下文保存/恢复机制为进程间切换奠定了基础
- 思考: 最关键的切换包含哪些步骤呢?



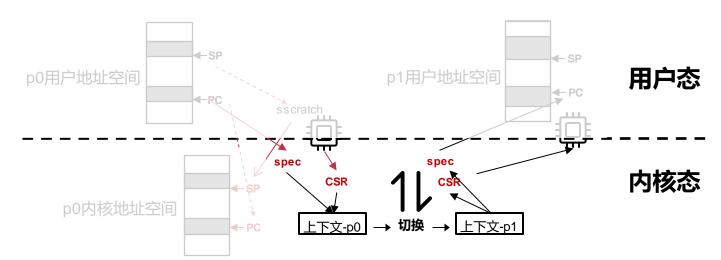
万事俱备, 只欠切换!

- · 上下文保存/恢复机制为进程间切换奠定了基础
- 思考: 最关键的切换包含哪些步骤呢?
 - 关键1:如何切换到p1的地址空间?



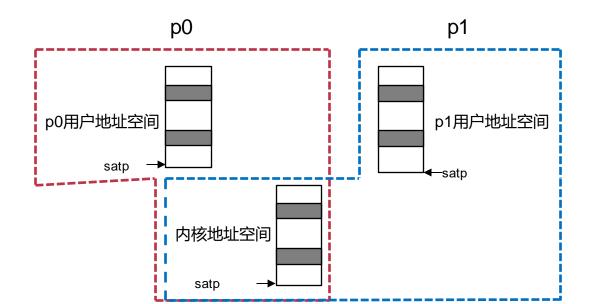
万事俱备, 只欠切换!

- · 上下文保存/恢复机制为进程间切换奠定了基础
- 思考: 最关键的切换包含哪些步骤呢?
 - 关键1:如何切换到p1的地址空间?
 - 关键2:如何切换到已经存储的p1上下文并进行恢复?



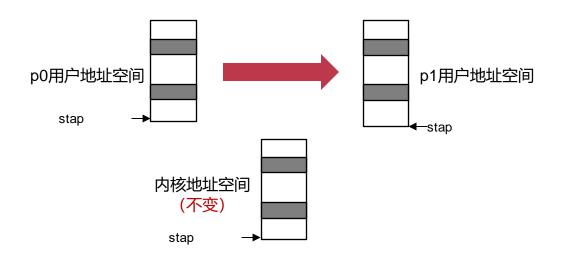
步骤1: 地址空间的切换

- · 回顾: RISC-V的地址空间管理
 - 内核与用户态地址空间分开管理
 - 用户地址空间独有,内核地址空间共享



步骤1: 地址空间的切换

- · 回顾: RISC-V的地址空间管理
 - 内核与用户态地址空间分开管理
 - 用户地址空间独有,内核地址空间共享
 - 因此, 只需要实现用户地址空间切换即可



步骤1: 地址空间的切换

- RISC-V中的实现: switch_mm()
 - 获取p1的PCB, 并获取其vmspace, 最后设置satp

```
switch mm (struct mm struct *prev, struct mm struct *next,
             struct task struc<u>t *</u>task)
static void set mm noasid(struct mm struct *mm)
        csr write(CSR_SATP, virt to pfn(mm->pgd) | satp_mode);
                                   转为物理地址
        local flush tlb all();
    static inline void local flush tlb all (void)
                      volatile ("sfence.vma" : : : "memory");
                            刷新TLB(为什么?
```

步骤2: 如何切换到p1的上下文?

· 回顾: 当p1保存上下文时, 其内核栈应该是怎样的?



SP指向上下文起始地址

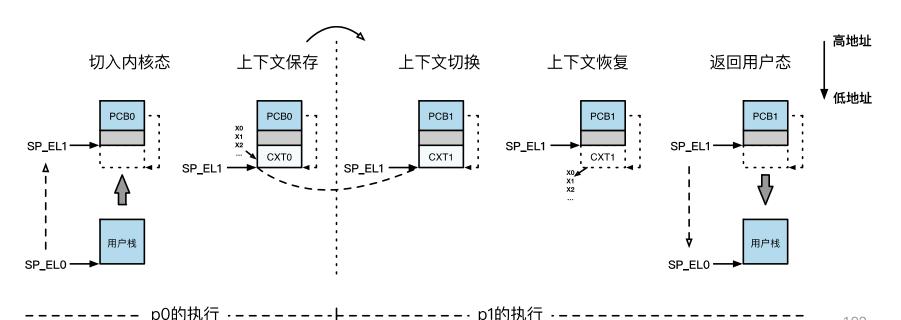
101

· p0/p1共享内核地址空间,因此直接切换内核栈指针即可

```
context switch(struct rq *rq, struct task struct *prev,
                   struct task struct *next, struct rq flags *rf){
   /* Here we just switch the register
      state and the stack. */
                                                                             CSR
    switch to (prev, next, prev);
                                                                                      p1上下文
    barrier();
                                                                               X1
                                                                               XΩ
                                                                             CSR
       REG_L s11, TASK_THREAD_S11_RA(a4)
                                                                                      女才土0g
                                                                               X1
       /* The of set of thread_info in tag
                                                                               X0
     move tp, a1
   SYM_FUNC_END(__switch_to)
```

总结: 上下文切换栈变化全过程

- 共涉及两次权限等级切换、三次栈切换
- 内核栈的切换是线程切换执行的"分界点"



102

纤程

纤程的概念 - 编程模型 - Windows和编程语言支持

一对一线程模型的局限

- · 复杂应用: 对调度存在更多需求
 - 生产者消费者模型: 生产者完成后, 消费者最好马上被调度
 - 内核调度器的信息不足,无法完成及时调度

- · "短命"线程:执行时间亚毫秒级 (如处理web请求)
 - 内核线程初始化时间较长,造成执行开销
 - 线程上下文切换频繁,开销较大

纤程 (用户态线程)

· 比线程更加轻量级的运行时抽象

- 不单独对应内核线程
- 一个内核线程可以对应多个纤程(多对一)

• 纤程的优点

- 不需要创建内核线程,开销小
- 上下文切换快(不需要进入内核)
- 允许用户态自主调度,有助于做出更优的调度决策

Linux对于纤程的支持: ucontext

- · 每个ucontext可以看作一个用户态线程
 - makecontext: 创建新的ucontext
 - setcontext:纤程上下文切换
 - getcontext:保存当前的ucontext

纤程的例子: 生产者 - 消费者

消费者 生产者 void consume() { void produce() { buf[++cnt] = rand(); process(buf[cnt]); setcontext(&cxt2); setcontext(&cxt1); 主纤程 makecontext(&cxt1, produce, ...); makecontext(&cxt2, consume, ...); setcontext(&cxt1);

setcontext的代码片段

```
1 ENTRY (__setcontext)
     // 恢复被调用者保存的通用寄存器
     ldp x18, x19, [x0, register_offset + 18 * SZREG]
     ldp x20, x21, [x0, register\_offset + 20 * SZREG]
     ldp x22, x23, [x0, register_offset + 22 * SZREG]
     1dp \times 24, \times 25, [\times 0, register offset + 24 \times SZREG]
 8
     //恢复用户栈
10
     ldr x2, [x0, sp_offset]
11
    mov sp, x2
     // 恢复浮点寄存器及参数
12
13
    // 恢复 PC 并返回
14
15 ldr x16, [x0, pc_offset]
     br x16
16
```

108

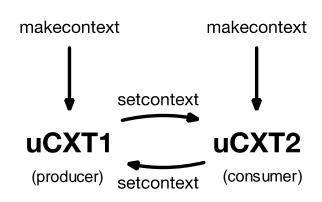
从例子看纤程的优势

• 纤程切换及时

- 当生产者完成任务后,可直接用户态切换到消费者
- 对该线程来说是最优调度(内核调度器和难做到)

· 高效上下文切换

- 切换不进入内核态,开销小
- 即时频繁切换也不会造成过大开销



Windows对于纤程的支持: Fiber库

· 与ucontext类似的编程模型

- createFiber: 创建新的纤程
- SwitchToFiber: 纤程切换

· 支持纤程本地存储 (FLS)

- Fiber Local Storage
- 当一个内核线程对应单个纤程时,FLS与TLS结构相同
- 当一个内核线程对应多个纤程时,TLS可分裂为多个FLS

程序语言中对纤程的支持: 协程

- · 许多高级程序语言都对协程提供了支持
 - go、python、lua......
 - C++自20开始也支持了协程
- · 协程也拥有状态 (新生/暂停/终止/执行)
 - 核心操作: yield (使协程暂停执行)、resume (继续执行)

