

	interrupt	通常来自 I/O	asyn.
	trap	系统调用	syn.
	fault		syn.
	abort		syn

exception 实际上是由核控制子进程的一种方式。

Status of a Process

New
Running
Ready
Stopped
Blocked 等待 I/O
Terminated

Process:

creating processes: fork 0 is child; child pid to parent
-1 on error
exec

reaping : pid_t waitpid(pid_t pid, int *status, int options)

pid: -1, 0, n, -n

status: WIFEXITED(status) & WEXITSTATUS(status)
WIFSIGNALED(status) & WTERMSIG(status)

option: 0: 直到有子进程终止

WNOHANG (no hang): 如果无子进程终止，则立刻返回 0.

pid 的值	含义
> 0	等待进程 ID 等于 pid 的子进程结束。
== 0	等待任何一个与当前进程在同一进程组的子进程结束。
< -1	等待进程组 ID 等于 -pid 的任何一个子进程结束。
== -1	等待任何一个子进程结束 (等价于 wait() 的行为)。

WUNTRACED: 挂起，一直等到子进程 终止或停止

WCONTINUED: 挂起，直到子进程 停止，或自己停止的进程重新开始。

WUNTRACED|WNOHANG: 立刻返回，如果有进程 停止或中止，则返回 PID；如果没有，返回 0

pid_t wait (int * status)
= waitpid (-1, &status, 0)

sleep & pause: sleep: 挂起一段时间
pause: 挂起直到接收到信号

竞争条件：原因：

- Shell一行执行多个命令，当父进程结束，就开始执行下一个命令，不管子进程是否执行完。
- 因此child的数据会与parent的随机交叉
- 先子后父，不会交叉

`int execve (const char *filename, const char *argv[], const char *envp[])`

Signal

发送信号：
硬件
内核
shell命令。
函数：kill, alarm, raise ...

接收信号：从内核态 \rightarrow 用户态中断

`sigHandler_t signal (int signum, sigHandler_t handler)`

阻塞和解除：
sigprocmask
sigemptyset
sigfillset
sigaddset
sigdelset

信号处理：回收子进程 `while (waitpid (-1, NULL, 0) > 0)`

为了避免父子进程竞争，考虑屏蔽 sigchild.

CPU 调度：进程切换的决策

优先目标：

1. Turnaround Time 周转时间

$$T_{turnaround} = T_{completion} - T_{arrival}$$

只考虑周转时间会带来公平性问题

FIFO : 缺：无短的高优先级长。

SJF (最短优先)



STCF (时间片抢占)

2. Response time 响应时间

RR: 时间片队列

$$T_{response} = T_{firstrun} - T_{arrival}$$

3. 高 I/O

MLFQ (Multi-Level Feedback Queue)

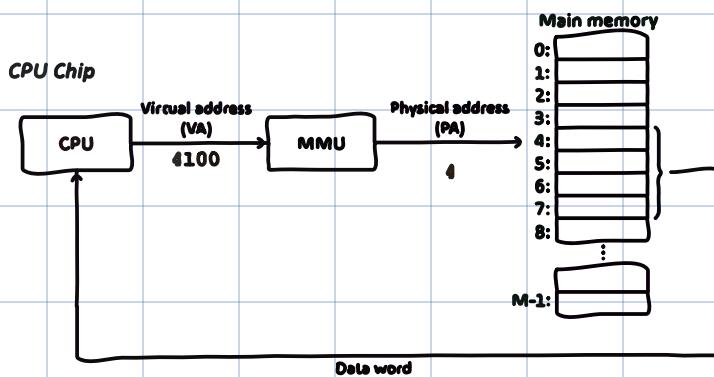
- Rule 1: If Priority(A) > Priority(B), A runs (B doesn't).
- Rule 2: If Priority(A) = Priority(B), A & B run in RR.
- Rule 3: When a job enters the system, it is placed at the highest priority (the topmost queue).
- Rule 4: Once a job uses up its time allotment at a given level (regardless of how many times it has given up the CPU), its priority is reduced (i.e., it moves down one queue).
- Rule 5: After some time period S, move all the jobs in the system to the topmost queue.

随机调度算法: Lottery Scheduling

多核 CPU 调度

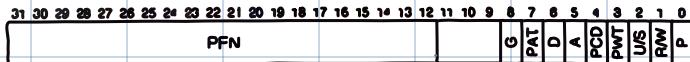
Virtual Memory

A System Using Virtual Addressing



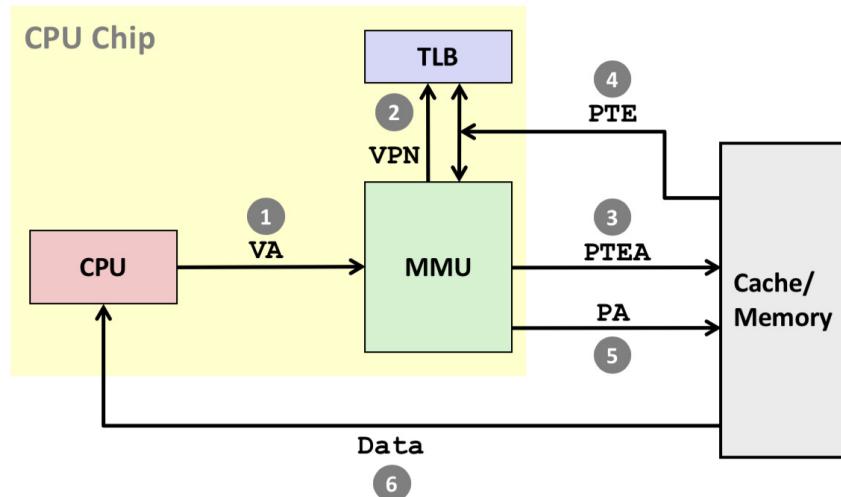
- Used in all modern servers, desktops, and laptops
- One of the great ideas in computer science

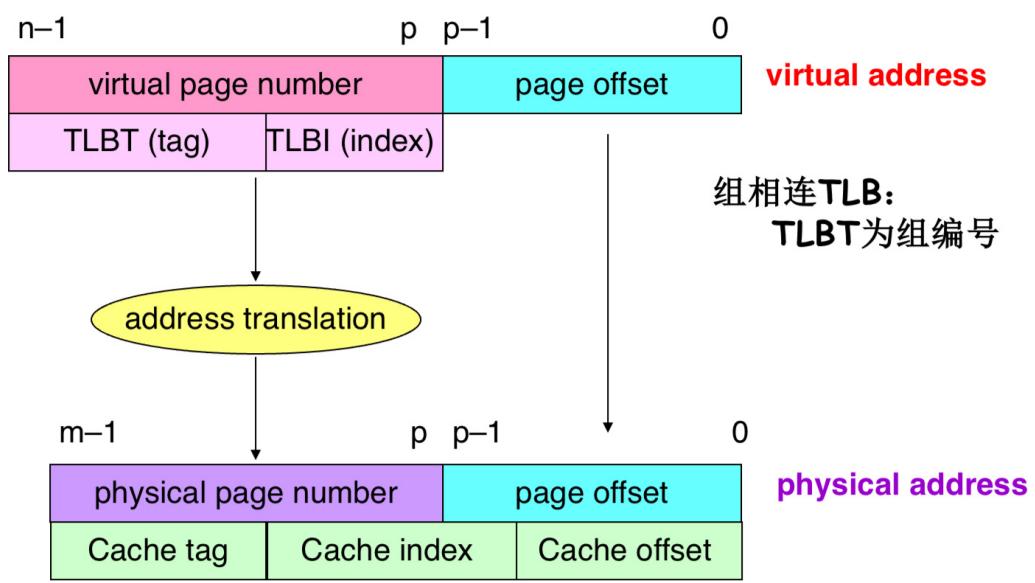
真实的PTE:



- valid bit (P)
- protection bits (R/W) (r/w/e权限)
- present bit (P) (on disk/ in physical memory)
- dirty bit (D) (与disk版本是否一致)
- reference bit/ accessed bit (A) (体现page的流行度, 与replacement有关)
- U/S. user/supervisor, 用户态是否可以访问
- PAT/PWT/PCD/G: 与硬件cache相关

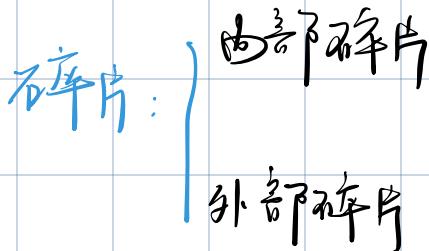
TLB:





多级页表：

动态内存分配



Method 1. *Implicit list* using length-links all blocks

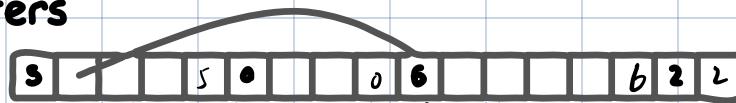


(增加 footer,
方便合并)

e.g.: 由于双字对齐, header 的最后三位可作标志位.

Find a free block: First fit, Next fit, Best fit, Worst fit.

Method 2: *Explicit list* among the free blocks using pointers



Insertion policy: Where in the free list do you put a newly freed block?

LIFO (last-in-first-out) policy

- Insert freed block at the beginning of the free list
- **Pro:** simple and constant time
- **Con:** studies suggest fragmentation is worse than address ordered

Address-ordered policy

- Insert freed blocks so that free list blocks are always in address order:
 $\text{addr}(\text{prev}) < \text{addr}(\text{curr}) < \text{addr}(\text{next})$
- **Con:** requires search
- **Pro:** studies suggest fragmentation is lower than LIFO

Method 3: *Segregated free list*

- Different free lists for different size classes

- #1. 简单分离存储 (Simple Segregated Storage)

每个链表中装大小相同的块, 不分割, 不合并.

- #2. 分离适配 (Segregated Fit)

每个链表中是一个空间范围的 free blocks

- #3. 伙伴系统 (Buddy System)

每个块的大小都是 2 的幂

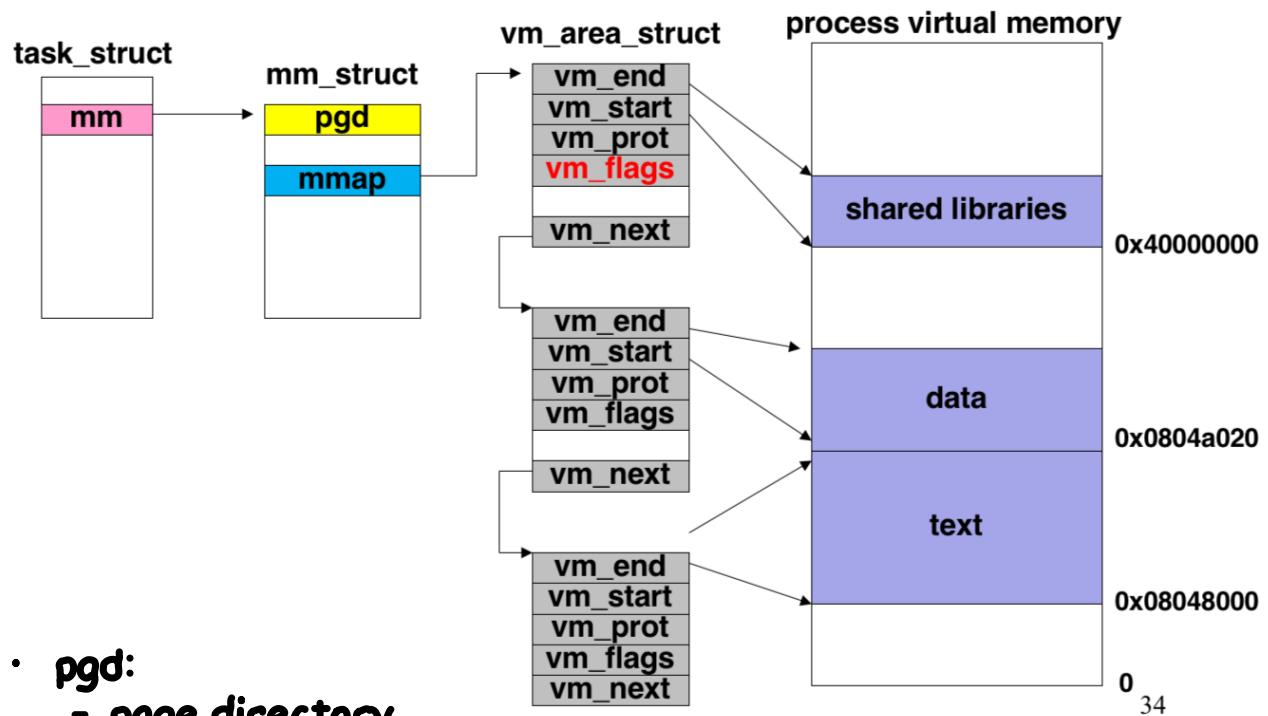
Keeping Track of Free Blocks

Method 4: Blocks sorted by size

- Can use a balanced tree (e.g. Red-Black Tree) with pointers within each free block and the length used as a key

垃圾回收：遍历遍历所有节点，找到所有能到达的块标记之。

内存映射：vm_area_struct:



- **pgd:**
 - page directory address (will be loaded to CR3)
- **vm_prot:**
 - read/write permissions for this area
- **vm_flags**
 - shared with other processes or private to this process

Mmap 防止

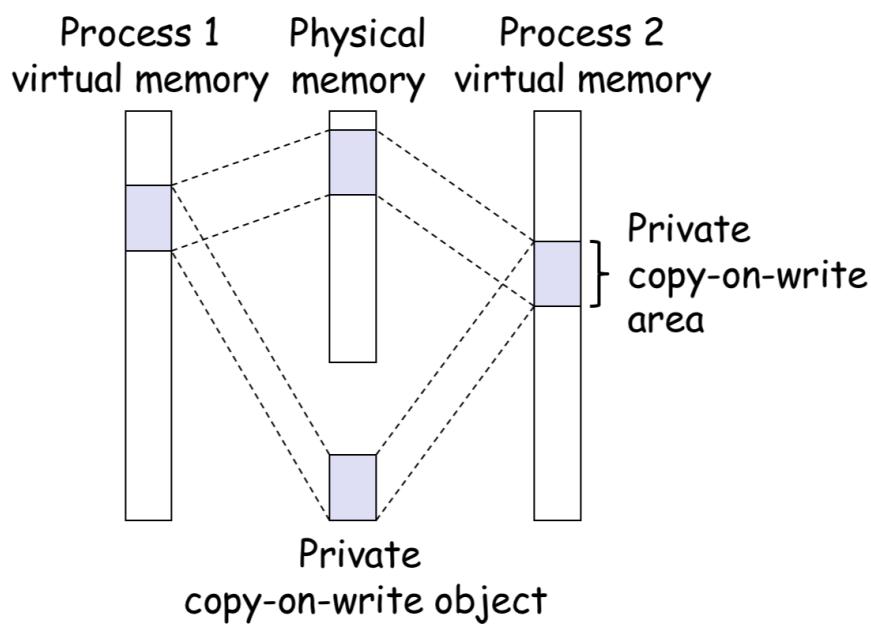
```
void *mmap(void *start, int len, int prot,  
           int flags, int fd, int offset)
```

Map **len** bytes starting at offset **offset** of the file specified by file description **fd**, preferably at address **start** (usually 0 for don't care).

- **prot**: PROT_READ, PROT_WRITE
- **flags**: MAP_PRIVATE, MAP_SHARED

Return a pointer to the mapped area

COW: copy on write



I/O

Devices: 磁盘

寻道时间 等待时间 读取时间

调度算法 { P-scan 内→外 外→内

C-scan 只向内→外读

SSTF (Shortest Seeking Time First)

RAID

(独立磁盘冗余阵列)

RAID 0

无备份

RAID 1

镜像备份

?几个一组是否固定

RAID 5

旋转奇偶校验

JBOD

?

SSD

映射策略

page-level

block-level

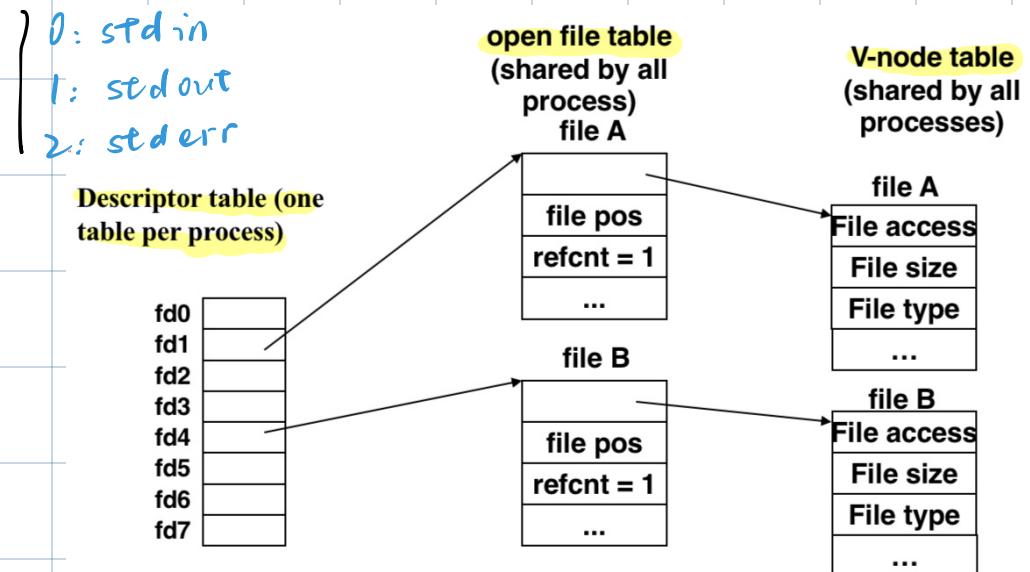
hybrid: log table + data table.

partial merge

full merge

此类型题目不统一，需注意 merge 的要求。

Unix I/O Descriptor table → open file table → V-node table.



接口函数

```
int open(char *filename, int flags, mode_t mode);
    Returns: new file descriptor if OK, -1 on error
```

flags

- O_RDONLY, O_WRONLY, O_RDWR (must have one)
- O_CREAT, O_TRUNC, O_APPEND, O_DIRECT (optional)

mode

- S_IRUSR, S_IWUSR, S_IXUSR
- S_IRGRP, S_IWGRP, S_IXGRP
- S_IROTH, S_IWOTH, S_IXOTH

```
int close(int fd);
    Returns: zero if OK, -1 on error
```

```
#include <unistd.h>

ssize_t read(int fd, void *buf, size_t count);
    returns: number of bytes read if OK,
              0 on EOF, -1 on error

ssize_t write(int fd, const void *buf, size_t count);
    returns: number of bytes written if OK,
              -1 on error
```

读取文件的元数据后.

```
#include <unistd.h>
#include <sys/stat.h>

int stat(const char *filename, struct stat *buf);

int fstat(int fd, struct stat *buf) ;

returns: 0 if OK, -1 on error
```

| 硬链接 不同文件名但相同 inode.

| 软链接 不同文件名, 数据块中存的是指向对应文件的路径字符串.

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

空目录实际上包含一个引用自身和一个引用父目录的条目.

DIR *opendir (const char *name);
returns: 成功则返回处理的指针; 若出错, 则返回NULL

struct dirent *readdir (DIR *dirp);
returns: 成功则返回指向下一个目录项的指针; 没有更多目录或出错则返回NULL

DIR *closedir (DIR *dirp);
returns: 成功则返回0; 若出错, 则返回-1
```

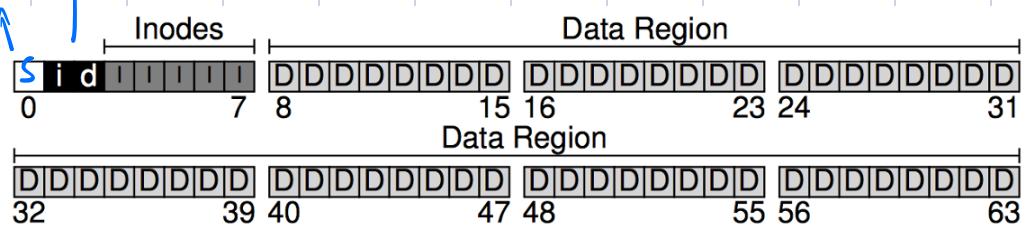
重定向 修改的是 newfd

```
int dup2(int oldfd, int newfd);
returns: nonnegative descriptor if OK, -1 on error
```

File System

bitmap.

super block 1



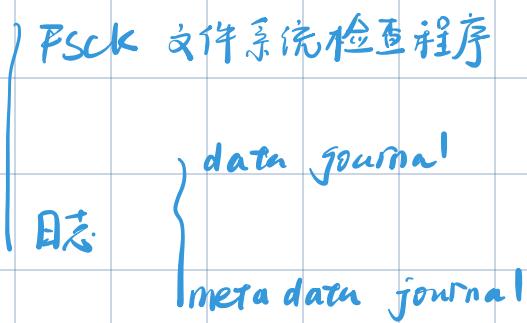
Multi-level index 用多级索引来索引数据块



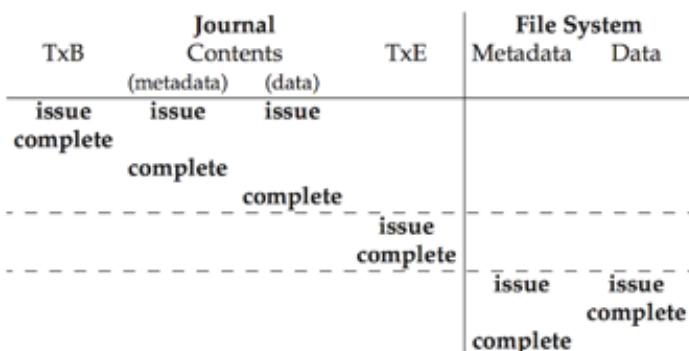
Extent point + length.

linked-based in node 中只包含一个指针

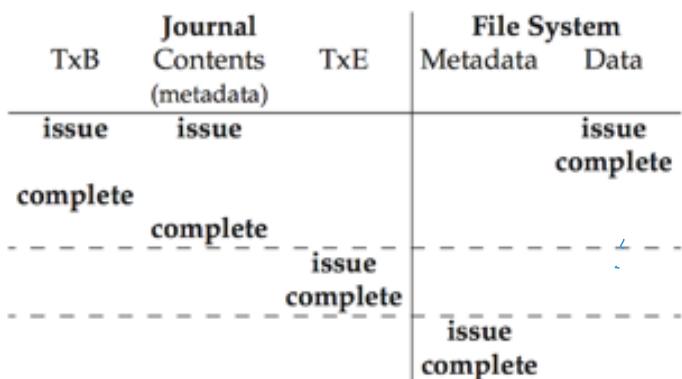
解决崩溃一致性问题的方法：



- 总结时间线: data journaling

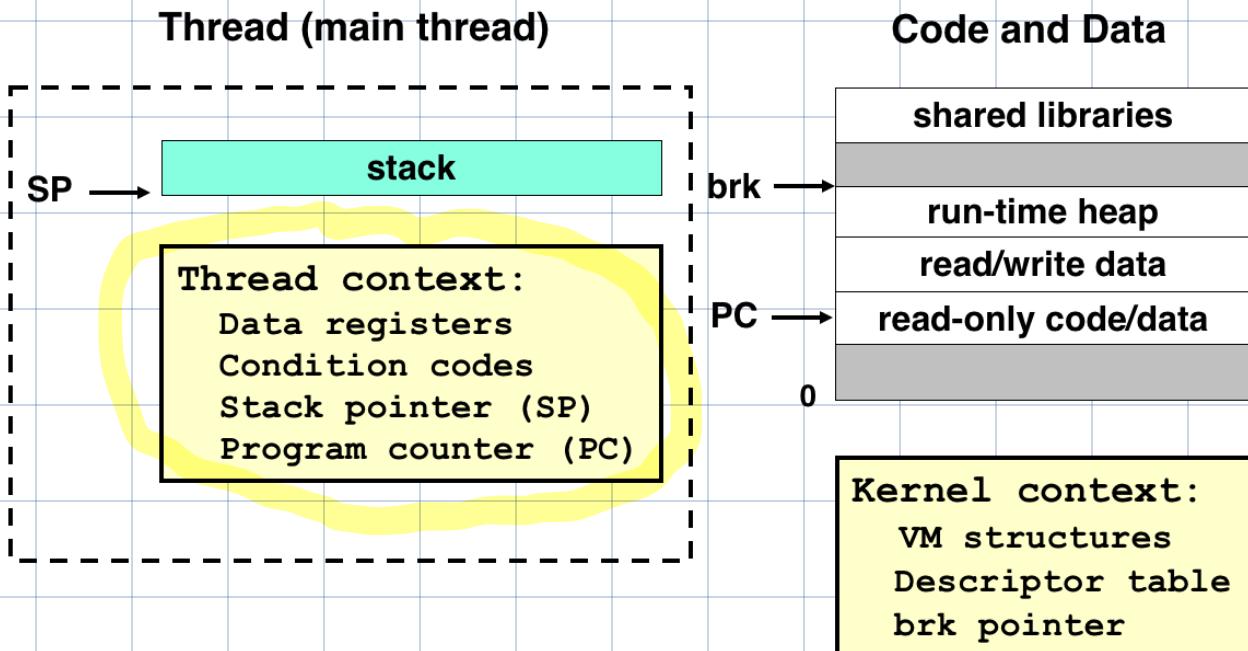


- 总结时间线: metadata journaling



并行编程:

线程的基本概念: 进程 - code . data . heap . kernel context



Posix threads (Pthreads) interface

- Creating and reaping threads

- `pthread_create`
- `pthread_join`

eg: 不 join 必须用线程的 detached 模式。

`pthread_detach (tID)`

- Determining your thread ID

- `pthread_self`

- Terminating threads

- `pthread_cancel` 发送中止请求给某线程
- `pthread_exit` 线程主动退出
- `exit [terminates all threads]` 会导致所有进程退出
- `return [terminates current thread]`

信号量:

```
#include <semaphore.h>

int sem_init(sem_t *sem, 0, unsigned int value);
int sem_wait(sem_t *s); /* P(s) */
int sem_post(sem_t *s); /* V(s) */

#include "csapp.h"

void P(sem_t *s); /* Wrapper function for sem_wait */
void V(sem_t *s); /* Wrapper function for sem_post */
```

西进程图 采分析无锁.

并行编程模型

生产 - 消费者

核心数据结构:

```
struct {
    int *buf;          /* Buffer array */
    int n;             /* Maximum number of slots */
    int front;         /* buf[(front+1)%n] is first item */
    int rear;          /* buf[rear%n] is last item */
    sem_t mutex;       /* protects accesses to buf */
    sem_t slots;       /* Counts available slots */
    sem_t items;        /* Counts available items */
} sbuf_;
```

读写者问题

读优先

读写公平

写优先

线程安全

可重入函数：不涉及任何共享变量

Race，当程序依赖于多个线程的执行顺序时。

没有统一。

并行数据结构.

Sloppy counter: 每个CPU一个lock, 还有一个global lock.
本地累加, 定期汇总.

链表: 一个global lock, 对insert和look-up上锁

队列: head-lock, tail-lock 分别对头尾上锁.

哈希表: 每个桶1个锁 一个桶就是一个链表.

公用

无锁化: CAS.

锁的实现:

Test and set

Compare and Swap.

Load-linked and Store-Conditional

eg: 基本都是依赖一些原子操作