# 操作系统

# Operating Systems

**刘 川 意　教授**

liuchuanyi@hit.edu.cn

哈尔滨工业大学 (深圳)

2025年9月

# Module 5: Concurrency & Synchronization
## 并发与同步

1. **Concurrency(并发) Introduction**

2. **Locks**

3. **基于Lock的并发数据结构**

4. **Condition Variables 条件变量**

5. **Semaphore 信号量**

6. **常见并发问题**

7. **基于事件的并发**

# Concurrency为什么放到OS中讲?

- **History** **!**

    - OS Kernel是第一个并发程序，如：write()的设计，中断对shared structures的影响 （page tables, process lists, file system structures, and virtually every kernel data structure has to be carefully accessed)

    - 很多并发处理技术是在OS中发明和实现的

    - multi-threads进程中，应用程序也需要考虑并发

# 并发相关的重要术语

- **Critical Section(临界区)**, a piece of code that accesses a *shared* resource, usually a variable or data structure. 临界区是一段访问共享资源（通常是变量或数据结构）的代码。

- **Race Condition(条件竞争)** arises if multiple threads of execution enter the critical section at roughly the same time; both attempt to update the shared data structure, leading to a surprising (and perhaps undesirable) outcome. 如果多个正在执行的线程同时进入临界区，则会出现条件竞争；这些线程都试图更新共享的数据结构，会造成意料之外的结果。

- **An Indeterminate(不确定的)** program consists of one or more race conditions; the output of the program varies from run to run, depending on which threads ran when. The outcome is thus not **deterministic**, something we usually expect from computer systems. 不确定的进程由一个或多个条件竞争组成；程序每一次运行的输出都有可能不同，具体取决于每个线程运行的时间。

- **Mutual Exclusion(互斥)** primitives(原语), guarantee that only a single thread ever enters a critical section, thus avoiding races, and resulting in deterministic program outputs. 互斥原语保证只有单个线程进入临界区，从而避免竞争，并导致确定性的进程输出。

# Review: Thread

- 轻量化执行环境，new abstraction for <u>a single running process</u>

- Multi-threaded 程序的特点：
    - A multi-threaded program has more than one point of execution.
    - Multiple PCs (Program Counter)
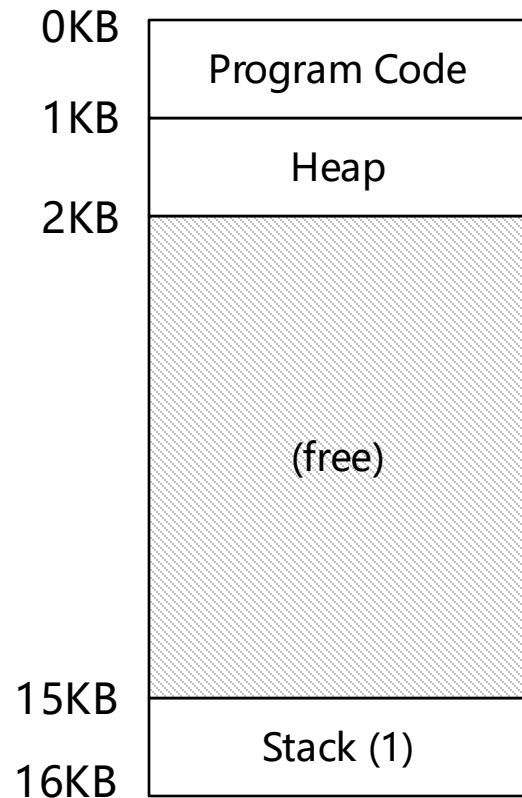    - They share the same address space.

# Context switch between threads

- Each thread has its own <u>program counter</u> and <u>set of registers</u>. 每个线程拥有独立的PC和寄存器。

  - One or more **thread control blocks(TCBs)** are needed to store the state of each thread. 需要一个或多个线程控制块（TCB）来存储每个线程的状态。

- When switching from running one (T1) to running the other (T2),

  - The register state of T1 be saved.

  - The register state of T2 restored.

  - The address space remains the same. 地址空间保持不变。

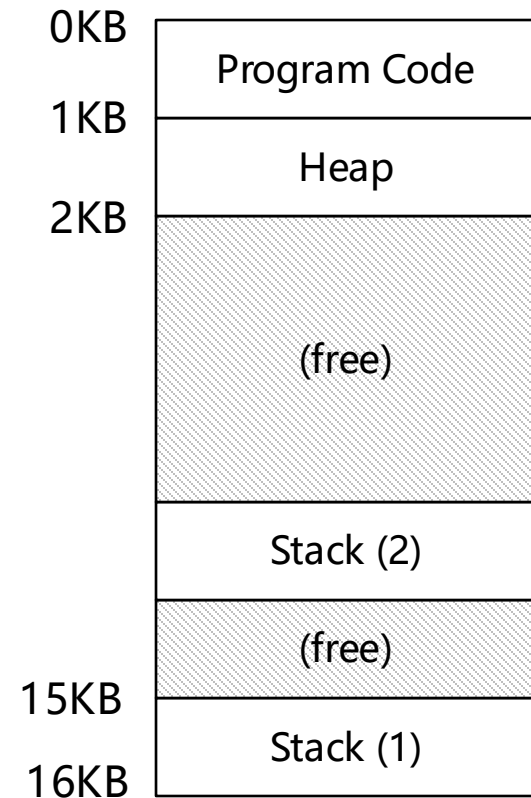# The stack of the relevant thread

- There will be one stack per thread.

|  | A Single-Threaded Address Space | |
| --- | --- | --- |
| 0KB | Program Code | |
| 1KB | Heap | |
| 2KB | | |
| | (free) | |
| 15KB | | |
| | Stack (1) | |
| 16KB | | |

**The code segment**:
where instructions live

**The heap segment**:
contains malloc'd data(malloc 动态申请的内存数据)
dynamic data structures(动态的数据结构)
(it grows downward)

(it grows upward)
**The stack segment**:
contains local variables(局部变量), arguments to routines(函数参数), return values(返回值), etc.

**A Single-Threaded Address Space**

| | Two threaded Address Space | |
| --- | --- | --- |
| 0KB | Program Code | |
| 1KB | Heap | |
| 2KB | | |
| | (free) | |
| | Stack (2) | |
| | (free) | |
| 15KB | | |
| | Stack (1) | |
| 16KB | | |

**Two threaded Address Space**

# `badcnt.c`: Improper Synchronization

```
7.    /* Global shared variable */
8.    volatile long cnt = 0; /* Counter */

17.   int main(int argc, char **argv)
18.   {
19.     long niters;
20.     pthread_t tid1, tid2;

21.     niters = atoi(argv[1]);
22.     Pthread_create(&tid1, NULL,
23.       thread, &niters);
24.     Pthread_create(&tid2, NULL,
25.       thread, &niters);
26.     Pthread_join(tid1, NULL);
27.     Pthread_join(tid2, NULL);

28.     /* Check result */
29.     if (cnt != (2 * niters))
30.       printf("BOOM! cnt=%ld\n", cnt);
31.     else
32.       printf("OK cnt=%ld\n", cnt);
33.     exit(0);
34.   }
```

badcnt.c

```
38.   /* Thread routine */
39.   void *thread(void *vargp)
40.   {
41.     long i, niters =
42.           *((long *)vargp);

44.     for (i = 0; i < niters; i++)
45.       cnt++;

47.     return NULL;
48.   }
```

```
[zs_cao@localhost conc]$ ./badcnt 10000
OK cnt=20000
[zs_cao@localhost conc]$ ./badcnt 10000
BOOM! cnt=17302
[zs_cao@localhost conc]$ ./badcnt 10000
OK cnt=20000
```

线程并发执行的问题

# Assembly Code for Counter Loop

- 编译:
  - gcc –s badcnt.c –o badcnt.s
  - vim badcnt.s

```
for (i = 0; i < niters; i++)
     cnt++;
```

```
 94        movq    %rdi, -24(%rbp)
 95        movq    -24(%rbp), %rax
 96        movq    (%rax), %rax
 97        movq    %rax, -8(%rbp)
 98        movq    $0, -16(%rbp)
 99        jmp     .L6
100 .L7:
101        movq    cnt(%rip), %rax
102        addq    $1, %rax
103        movq    %rax, cnt(%rip)
104        addq    $1, -16(%rbp)
105 .L6:
106        movq    -16(%rbp), %rax
107        cmpq    -8(%rbp), %rax
108        jl      .L7
109        movl    $0, %eax
110        popq    %rbp
```

$H_i$ : Head

$L_i$ : Load cnt
$U_i$ : Update cnt
$S_i$ : Store cnt

$T_i$ : Tail

☐ 汇编:

   ☐ gcc –c badcnt.s –o badcnt.o

   ☐ objdump -dx  badcnt.o

```
for (i = 0; i < niters; i++)
     cnt++;
```

```
00000000000000ed <thread>:
  ed:    55                         push    %rbp
  ee:    48 89 e5                   mov     %rsp,%rbp
  f1:    48 89 7d e8                mov     %rdi,-0x18(%rbp)
  f5:    48 8b 45 e8                mov     -0x18(%rbp),%rax
  f9:    48 8b 00                   mov     (%rax),%rax
  fc:    48 89 45 f8                mov     %rax,-0x8(%rbp)
 100:    48 c7 45 f0 00 00 00       movq    $0x0,-0x10(%rbp)
 107:    00
 108:    eb 17                      jmp     121 <thread+0x34>
 10a:    48 8b 05 00 00 00 00       mov     0x0(%rip),%rax        # 111
<thread+0x24>
                        10d: R_X86_64_PC32        cnt-0x4
 111:    48 83 c0 01                add     $0x1,%rax
 115:    48 89 05 00 00 00 00       mov     %rax,0x0(%rip)        # 11c
<thread+0x2f>
                        118: R_X86_64_PC32        cnt-0x4
 11c:    48 83 45 f0 01             addq    $0x1,-0x10(%rbp)
 121:    48 8b 45 f0                mov     -0x10(%rbp),%rax
 125:    48 3b 45 f8                cmp     -0x8(%rbp),%rax
 129:    7c df                      jl      10a <thread+0x1d>
 12b:    b8 00 00 00 00             mov     $0x0,%eax
 130:    5d                         pop     %rbp
 131:    c3                         retq
```

$\left.\right\} H_i$

$L_i$

$\left.\right\} U_i$

$S_i$

$\left.\right\} T_i$

# Assembly Code for Counter Loop

□ 链接:

- □ gcc –o badcnt.c –o badcnt -lpthread
- □ objdump -d badcnt

```
for (i = 0; i < niters; i++)
    cnt++;
```

```
0000000000000957 <thread>:
 957:    55                       push   %rbp
 958:    48 89 e5                 mov    %rsp,%rbp
 95b:    48 89 7d e8              mov    %rdi,-0x18(%rbp)
 95f:    48 8b 45 e8              mov    -0x18(%rbp),%rax
 963:    48 8b 00                 mov    (%rax),%rax
 966:    48 89 45 f8              mov    %rax,-0x8(%rbp)
 96a:    48 c7 45 f0 00 00 00     movq   $0x0,-0x10(%rbp)
 971:    00
 972:    eb 17                    jmp    98b <thread+0x34>
 974:    48 8b 05 b5 06 20 00     mov    0x2006b5(%rip),%rax
   # 201030 <cnt>
 97b:    48 83 c0 01              add    $0x1,%rax
 97f:    48 89 05 aa 06 20 00     mov    %rax,0x2006aa(%rip)
   # 201030 <cnt>
 986:    48 83 45 f0 01           addq   $0x1,-0x10(%rbp)
 98b:    48 8b 45 f0              mov    -0x10(%rbp),%rax
 98f:    48 3b 45 f8              cmp    -0x8(%rbp),%rax
 993:    7c df                    jl     974 <thread+0x1d>
 995:    b8 00 00 00 00           mov    $0x0,%eax
 99a:    5d                       pop    %rbp
 99b:    c3                       retq
```

$H_i$

怎么计算?

$L_i$
$U_i$
$S_i$

$T_i$

# Race condition

- 把上述示例简化一下：
  - counter = counter + 1 (default is 50)
  - We expect the result is 52. However,

|  |  |  | (after instruction) | | |
| --- | --- | --- | --- | --- | --- |
| OS | Thread1 | Thread2 | PC | %eax | counter |
|  | before critical section |  | 100 | 0 | 50 |
|  | mov 0x8049a1c, %eax |  | 105 | 50 | 50 |
|  | add $0x1, %eax |  | 108 | 51 | 50 |
| interrupt<br>save T1's state<br>restore T2's state |  |  | 100 | 0 | 50 |
|  |  | mov 0x8049a1c, %eax | 105 | 50 | 50 |
|  |  | add $0x1, %eax | 108 | 51 | 50 |
|  |  | mov %eax, 0x8049a1c | 113 | 51 | 51 |
| interrupt<br>save T2's state<br>restore T1's state |  |  | 108 | 51 | 50 |
|  | mov %eax, 0x8049a1c |  | 113 | 51 | **51** |

# Critical section

- A piece of code that accesses a shared variable and must not be concurrently executed by more than one thread. 一段访问共享变量的代码不能并发地被超过一个线程执行。

  - Multiple threads executing critical section can result in a race condition. 多线程执行临界区代码会引起条件竞争。

  - Need to support **atomicity** for critical sections (**mutual exclusion**). 需要支持访问临界区的原子性（互斥）。

# Critical section

【例题】下列对临界区的论述中，正确的是（　）。

A. 临界区是指进程中用于实现进程互斥的那段代码。

B. 临界区是指进程中用于实现进程同步的那段代码。

C. 临界区是指进程中用于实现进程通信的那段代码。

D. 临界区是指进程中用于访问临界资源的那段代码。

答案：D

解析：【PPT第4页】**Critical Section(临界区)**, a piece of code that accesses a *shared* resource, usually a variable or data structure.

临界区是一段访问共享资源（通常是变量或数据结构）的代码。

# Critical section

【例题】下列准则中，实现临界区互斥机制必须遵循的是（ ）。

I. 两个进程不能同时进入临界区

II. 允许进程访问空闲的临界资源

III. 进程等待进入临界区的时间是有限的

IV. 不能进入临界区的执行态进程立即放弃CPU

A. 仅I、IV B. 仅II、III

C. 仅I、II、III D. 仅I、III、IV

答案：C

IV选项，不一定必须满足，因为某些机制（如自旋锁）允许忙等待。

# Critical section

【例题】两个旅行社甲和乙为旅客到某航空公司订飞机票，形成互斥资源的是（ ）。

A. 旅行社

B. 航空公司

C. 飞机票

D. 旅行社与航空公司

答案：C

解析：一张飞机票不能售给不同的旅客，因此飞机票是互斥资源，其他因素只是为完成飞机票订票的中间过程，与互斥资源无关。

# Locks

- Ensure that any such critical section executes as if it were a single atomic instruction (**execute a series of instructions atomically**). 保证任意临界区代码像原子操作一样执行（即原子地执行一系列指令）。

```
1    lock_t mutex;
2    . . .
3    lock(&mutex);
4    balance = balance + 1;       →  Critical section
5    unlock(&mutex);
```
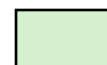
# Concurrent Execution(并发执行)

□ *Key idea:* In general, any sequentially consistent interleaving is possible, but some give an unexpected result! 通常，任何顺序的交错执行都是有可能的，但有些会给出意想不到的结果！

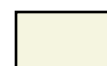  □ $I_i$ denotes that thread i executes instruction I. $I_i$ 表示线程i执行指令I。

  □ %rdx$_i$ is the content of %rdx in thread i's context. %rdx$_i$是线程i中%rdx的值。

| i (thread) | instr$_i$ | %rdx$_1$ | %rdx$_2$ | cnt |
|------------|-----------|----------|----------|-----|
| 1 | H$_1$ | - | - | 0 |
| 1 | L$_1$ | 0 | - | 0 |
| 1 | U$_1$ | 1 | - | 0 |
| 1 | S$_1$ | 1 | - | 1 |
| 2 | H$_2$ | - | - | 1 |
| 2 | L$_2$ | - | 1 | 1 |
| 2 | U$_2$ | - | 2 | 1 |
| 2 | S$_2$ | - | 2 | 2 |
| 2 | T$_2$ | - | 2 | 2 |
| 1 | T$_1$ | 1 | - | 2 |

Thread 1 critical section

Thread 2 critical section

*OK*

# Concurrent Execution (cont)

☐ Incorrect ordering: two threads increment the counter, but the result is 1 instead of 2. 两个线程同时增加counter的值，但结果是1而不是2。
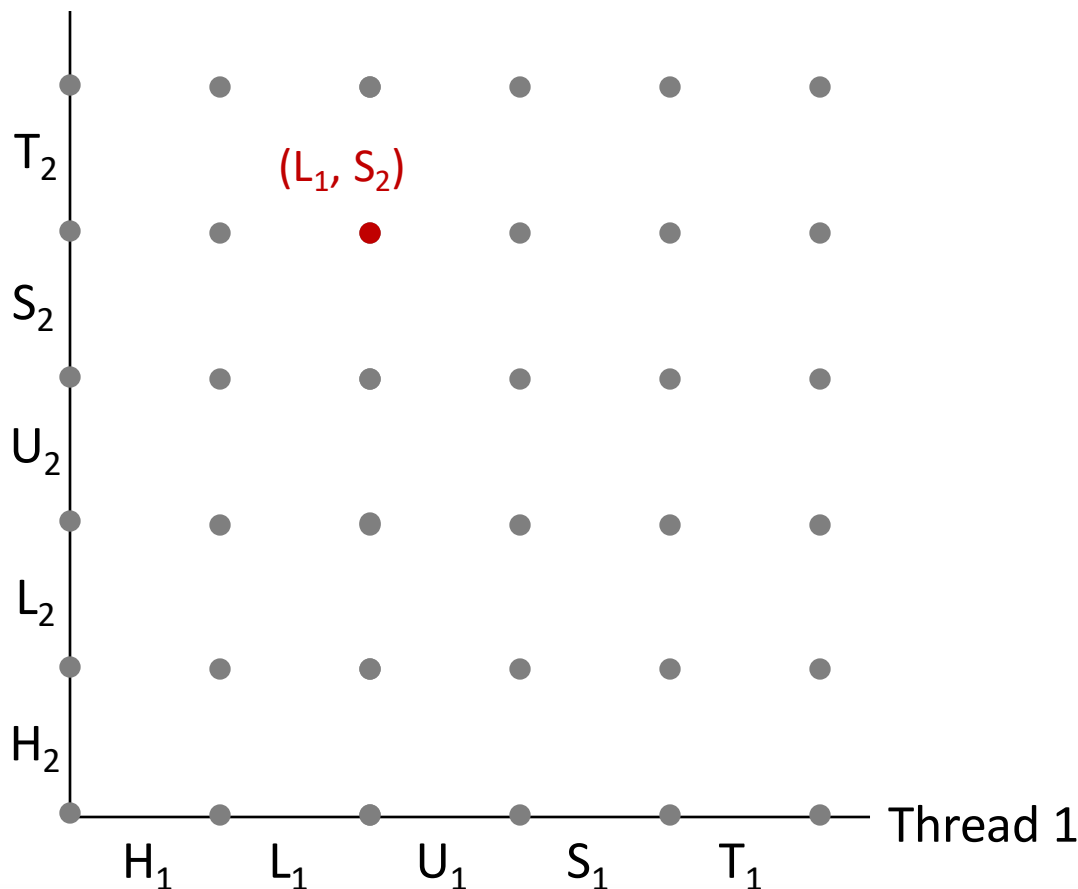
| i (thread) | $instr_i$ | $\%rdx_1$ | $\%rdx_2$ | cnt |
|:---:|:---:|:---:|:---:|:---:|
| 1 | $H_1$ | - | - | 0 |
| 1 | $L_1$ | 0 | - | 0 |
| 1 | $U_1$ | 1 | - | 0 |
| 2 | $H_2$ | - | - | 0 |
| 2 | $L_2$ | - | 0 | 0 |
| 1 | $S_1$ | 1 | - | 1 |
| 1 | $T_1$ | 1 | - | 1 |
| 2 | $U_2$ | - | 1 | 1 |
| 2 | $S_2$ | - | 1 | 1 |
| 2 | $T_2$ | - | 1 | 1 |

S1应该在L2之前执行

*Oops!*

# Progress Graphs （进度图）

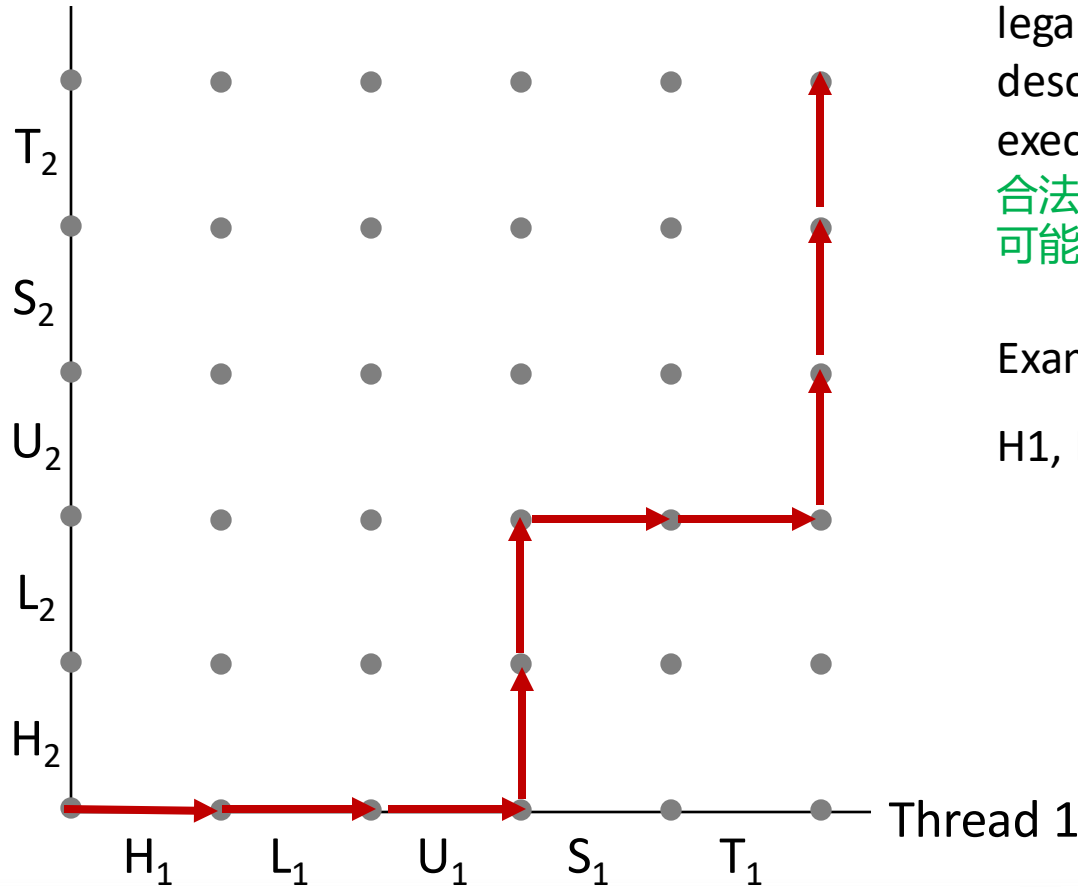A *progress graph* depicts the discrete *execution state space* of concurrent threads. 进度图描述了并发线程的离散执行状态空间。

Each axis corresponds to the sequential order of instructions in a thread. 每个轴表示一个进程的指令执行顺序。

Each point corresponds to a possible *execution state* ($Inst_1$, $Inst_2$). 每个点表示一个可能的执行状态。

E.g., ($L_1$, $S_2$) denotes state where thread 1 has completed $L_1$ and thread 2 has completed $S_2$.

Thread 2

$T_2$

($L_1$, $S_2$)

$S_2$

$U_2$

$L_2$

$H_2$

$H_1$   $L_1$   $U_1$   $S_1$   $T_1$   Thread 1

Thread 2

T$_2$

S$_2$

U$_2$

L$_2$
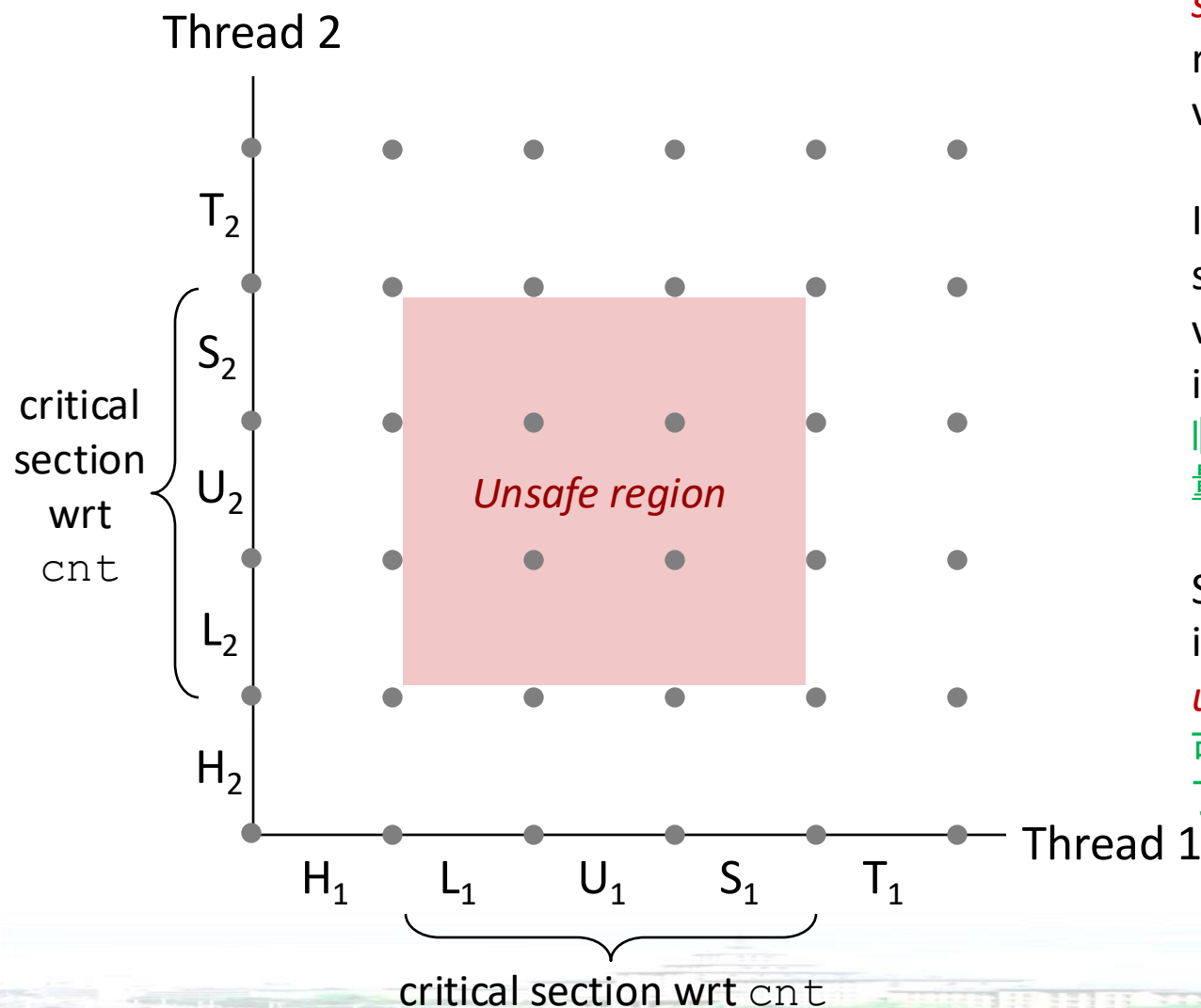
H$_2$

H$_1$   L$_1$   U$_1$   S$_1$   T$_1$   Thread 1

A *trajectory* （轨道） is a sequence of legal state transitions（转换） that describes one possible concurrent execution of the threads. 轨道是一系列合法的状态转换，用于描述线程的一个可能的并发执行序列。

Example:

H1, L1, U1, H2, L2,  S1, T1, U2, S2, T2

L, U, and S form a *critical section（临界区）* with respect to the shared variable `cnt`

Instructions in critical sections (wrt some shared variable) should not be interleaved
临界区中的指令（写入共享变量）不能被交错执行

Sets of states where such interleaving occurs form *unsafe regions*
可能出现交错的状态集合形成了*不安全区域*

**Thread 2**

$T_2$

$S_2$

critical section wrt `cnt`

$U_2$

$L_2$

$H_2$

*Unsafe region*

$H_1$    $L_1$    $U_1$    $S_1$    $T_1$    **Thread 1**

critical section wrt `cnt`

Def: A trajectory(轨道) is *safe* iff it does not enter any unsafe region
一个轨道是安全的当且仅当它不进入任何不安全区域。

Claim: A trajectory is correct (wrt `cnt`) iff it is safe

# Enforcing Mutual Exclusion

- *Question:* How can we guarantee a safe trajectory? 如何确保得到一个安全的轨道?

- Answer: We must *synchronize* （同步） the execution of the threads so that they can never have an unsafe trajectory. 我们需要对进程的执行进行同步，保证进程不存在不安全的轨道。
    - i.e., need to guarantee *mutually exclusive access （互斥地访问）* for each critical section. 即需要保证对每个临界区的**互斥访问**。

- Classic solution:
    - Semaphores（信号量） (Edsger Dijkstra)

- Other approaches (out of our scope)
    - Mutex and condition variables (Pthreads)
    - Monitors (Java)

1. **Concurrency Introduction**

2. <span style="color:red">**Locks**</span>

3. **基于Lock的并发数据结构**

4. **Condition Variables 条件变量**

5. **Semaphore 信号量**

6. **常见并发问题**

7. **基于事件的并发**

# Locks: The Basic Idea

- Ensure that any **critical section** executes as if it were a single atomic instruction. "全部或都不"

  - Eg. update of a shared variable

  ```
  balance = balance + 1;
  ```

  - Add some code around the critical section

  ```
  1    lock_t mutex; // some globally-allocated lock 'mutex'
  2    …
  3    lock(&mutex);
  4    balance = balance + 1;
  5    unlock(&mutex);
  ```

# Lock变量

- Lock variable holds <u>the state of</u> the lock.
  - **available** (or **unlocked** or **free**)
    - ▸ No thread holds the lock.

  - **acquired** (or **locked** or **held**)
    - ▸ Exactly one thread holds the lock and presumably is in a critical section. 有且只有一个进程拥有锁，且这个进程很可能处在临界区中。

# lock()原语的语义（semantics）

- `lock()`

  - **Try to** acquire the lock.

  - If <u>no other thread holds</u> the lock, the thread will **acquire** the lock.

  - **Enter** the *critical section*.

    - This thread is said to be <u>the owner of</u> the lock.

  - Other threads are *prevented from* entering the critical section while the first thread that holds the lock is in there. 当一个拥有锁的进程进入一个临界区时，其他线程不能进入这个临界区。

# Pthread Locks - mutex

- The name that the POSIX library uses for a <u>lock</u>.

  - Used to provide mutual exclusion(互斥) between threads.

```
1   pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
2
3   Pthread_mutex_lock(&lock); // wrapper for pthread_mutex_lock()
4   balance = balance + 1;
5   Pthread_mutex_unlock(&lock);
```

  - We may be using *different locks* to protect *different variables* → Increase concurrency (a more **fine-grained** approach). 我们可以使用不同的锁来保护不同的变量 → 提升并发性（更细粒度的方法）。

# Lock如何实现?

- Efficient locks provided mutual exclusion at low cost. 高效的锁提供最低开销的互斥。

- Building a lock need some help from the **hardware** and the **OS**. 构造一个锁需要硬件和操作系统的协助。

# 如何评价lock原语?

- **Mutual exclusion 正确性**
  - Does the lock work, preventing multiple threads from entering *a critical section*? 这个锁是否能保证多线程不能同时进入临界区?

- **Fairness 公平性**
  - Does each thread contending for the lock get a fair shot at acquiring it once it is free? (Starvation) 每个争夺锁的线程在锁被释放后，是否都能获得公平的获取锁机会？（饥饿）

- **Performance 性能**
  - The time overheads added by using the lock 使用锁带来的时间开销

# Controlling Interrupts 基于中断控制的锁实现

- **Disable Interrupts** for critical sections 在临界区禁用中断

  - One of the earliest solutions used to provide mutual exclusion

  - Invented for single-processor systems.

  ```
  1   void lock() {
  2       DisableInterrupts();
  3   }
  4   void unlock() {
  5       EnableInterrupts();
  6   }
  ```

  - Problem:

    ▸ Require too much *trust* in applications

      – Greedy (or malicious) program could monopolize the processor. 贪婪的（或恶意的）程序会独占处理器。

    ▸ Do not work on multiprocessors 多处理器体系结构这种方式不work

    ▸ Code that masks or unmasks interrupts be executed *slowly* by modern CPUs 对中断mask或unmask的指令在现代CPU上执行速度很慢

# Why hardware support needed?

☐ **First attempt**: Using a *flag* denoting whether the lock is held or not.

   ☐ The code below has problems.

```
1    typedef struct __lock_t { int flag; } lock_t;
2
3    void init(lock_t *mutex) {
4        // 0 → lock is available, 1 → held
5        mutex->flag = 0;
6    }
7
8    void lock(lock_t *mutex) {
9        while (mutex->flag == 1)  // TEST the flag
10                ;  // spin-wait (do nothing)
11        mutex->flag = 1;  // now SET it !
12    }
13
14   void unlock(lock_t *mutex) {
15        mutex->flag = 0;
16   }
```

# Why hardware support needed? (Cont.)

- **Problem 1**: No Mutual Exclusion (assume `flag=0` to begin)

| Thread1 | Thread2 |
|---|---|

```
call lock()
while (flag == 1)
interrupt: switch to Thread 2
```

```
                                    call lock()
                                    while (flag == 1)
                                    flag = 1;
                                    interrupt: switch to Thread 1
```

- **Problem 2**: <u>Spin-waiting</u> wastes time waiting for another thread.

```
flag = 1;  // set flag to 1 (too!)
```

- So, we need an atomic instruction supported by Hardware!
  - ***test-and-set* instruction**, also known as ***atomic exchange***

# 基于Test-and-set硬件指令实现

- An instruction to support the creation of simple locks

```
1    int TestAndSet(int *ptr, int new) {
2        int old = *ptr;          // fetch old value at ptr
3        *ptr = new;      // store 'new' into ptr
4        return old;      // return the old value
5    }
```

- **return**(testing) old value pointed to by the `ptr`. 返回ptr指向的旧的值。
- *Simultaneously* **update**(setting) said value to `new`. 同步地将ptr指向的值设置为 `new`。
- This sequence of operations is performed atomically. 这一系列操作是原子地执行的。

# A Simple Spin Lock(自旋锁) using test-and-set

```
1    typedef struct __lock_t {
2        int flag;
3    } lock_t;
4
5    void init(lock_t *lock) {
6        // 0 indicates that lock is available,
7        // 1 that it is held
8        lock->flag = 0;
9    }
10
11   void lock(lock_t *lock) {
12       while (TestAndSet(&lock->flag, 1) == 1)
13               ;           // spin-wait
14   }
15
16   void unlock(lock_t *lock) {
17       lock->flag = 0;
18   }
```

- **Note**: To work correctly on *a single processor*, it requires <u>a preemptive scheduler</u>.
- 在单处理器体系结构中，需要OS kernel实现抢占式调度策略来支持

# Evaluating Spin Locks

□ **Correctness(正确性)**: yes

　□ The spin lock only allows a single thread to entry the critical section. 自旋锁只允许一个进程进入临界区。

□ **Fairness(公平性)**: no

　□ Spin locks <u>don't provide any fairness</u> guarantees. Spin locks不能为公平性提供任何保证。

　□ Indeed, a thread spinning may spin *forever*. 实际上，一个自旋锁可能永远自旋（即原地"打转"）。

□ **Performance(性能)**:

　□ In the single CPU, performance overheads can be quire *painful*. 在单核CPU中，性能开销非常大。

　□ If the number of threads roughly equals the number of CPUs, spin locks work *reasonably well*. 如果线程数约等于 CPU 数，则自旋锁效果很好。

# 基于Compare-And-Swap硬件指令实现

- Test whether the value at the address(`ptr`) is equal to `expected`.

  - *If so*, update the memory location pointed to by `ptr` with the `new` value.

  - *In either case*, return the actual value at that memory location.

```
1    int CompareAndSwap(int *ptr, int expected, int new) {
2        int actual = *ptr;
3        if (actual == expected)
4                *ptr = new;
5        return actual;
6    }
```

**Compare-and-Swap hardware atomic instruction (C-style)**

```
1    void lock(lock_t *lock) {
2        while (CompareAndSwap(&lock->flag, 0, 1) == 1)
3                ; // spin
4    }
```

**Spin lock with compare-and-swap**

# Compare-And-Swap (Cont.)

- C-callable x86-version of compare-and-swap

```
1    char CompareAndSwap(int *ptr, int old, int new) {
2        unsigned char ret;
3
4        // Note that sete sets a 'byte' not the word
5        __asm__ __volatile__ (
6                " lock\n"
7                " cmpxchgl %2,%1\n"
8                " sete %0\n"
9                : "=q" (ret), "=m" (*ptr)
10               : "r" (new), "m" (*ptr), "a" (old)
11               : "memory");
12       return ret;
13   }
```

# Load-Linked and Store-Conditional

```
1    int LoadLinked(int *ptr) {
2        return *ptr;
3    }
4
5    int StoreConditional(int *ptr, int value) {
6      if (no one has updated *ptr since the LoadLinked to this address) {
7        *ptr = value;
8        return 1; // success!
9      } else {
10       return 0; // failed to update
11     }
12   }
```

**Load-linked And Store-conditional**

☐ The store-conditional *only succeeds* if no intermittent store to the address has taken place. 仅当ptr不被修改时，store-conditional会成功（即返回1）。

▸ **success**: return 1 and <u>update</u> the value at `ptr` to `value`.

▸ **fail**: the value at `ptr` is <u>not updates</u> and 0 is returned.

```
1    void lock(lock_t *lock) {
2        while (1) {
3                while (LoadLinked(&lock->flag) == 1)
4                        ; // spin until it's zero
5                if (StoreConditional(&lock->flag, 1) == 1)
6                        return; // if set-it-to-1 was a success: all done
7                                otherwise: try it all over again
8        }
9    }
10
11   void unlock(lock_t *lock) {
12       lock->flag = 0;
13   }
```

**Using LL/SC To Build A Lock**

```
1    void lock(lock_t *lock) {
2        while (LoadLinked(&lock->flag)||!StoreConditional(&lock->flag, 1))
3                ; // spin
4    }
```

**A more concise form of the `lock()` using LL/SC**

# Fetch-And-Add

- Atomically increment a value while returning the old value at a particular address.

```
1    int FetchAndAdd(int *ptr) {
2        int old = *ptr;
3        *ptr = old + 1;
4        return old;
5    }
```

**Fetch-And-Add Hardware atomic instruction (C-style)**

# Ticket Lock

- **Ticket lock** can be built with <u>fetch-and add</u>.
  - Ensure progress for all threads. → fairness

```c
1    typedef struct __lock_t {
2        int ticket;
3        int turn;
4    } lock_t;
5
6    void lock_init(lock_t *lock) {
7        lock->ticket = 0;
8        lock->turn = 0;
9    }
10
11   void lock(lock_t *lock) {
12       int myturn = FetchAndAdd(&lock->ticket);
13       while (lock->turn != myturn)
14               ; // spin
15   }
16   void unlock(lock_t *lock) {
17       FetchAndAdd(&lock->turn);
18   }
```

# So Much Spinning

☐ Hardware-based spin locks are simple and they work.

☐ In some cases, these solutions can be quite inefficient.

    ☐ Any time a thread gets caught *spinning*, it **wastes an entire time slice** doing nothing but checking a value. 每当线程自旋时，它会浪费整个时间片，除了检查锁的值以外什么也不做。

> **How To Avoid *Spinning*?**
> **We'll need OS Support too!**

- When you are going to spin, give up the CPU to another thread. 当线程将要自旋时，直接让出CPU。

  - OS system call moves the caller from the *running state* to the *ready state*. OS通过系统调用将线程从运行态变为就绪态。

  - The cost of a **context switch** can be substantial and the **starvation** problem still exists. 上下文切换的开销很大，且饥饿问题依然存在。

```
1   void init() {
2       flag = 0;
3   }
4
5   void lock() {
6       while (TestAndSet(&flag, 1) == 1)
7           yield(); // give up the CPU
8   }
9
10  void unlock() {
11      flag = 0;
12  }
```

**Lock with Test-and-set and Yield**

# 办法2: Using Queues: Sleeping, not Spinning

- **Queue** to keep track of which threads are <u>waiting</u> to enter the lock. 用队列来追踪正在等待锁的线程。

- `park()`
  - Put a calling thread to sleep 将调用线程置于睡眠状态

- `unpark(threadID)`
  - Wake a particular thread as designated by `threadID`. 唤醒特定threadID的进程

```
1    typedef struct __lock_t { int flag; int guard; queue_t *q; } lock_t;
2
3    void lock_init(lock_t *m) {
4        m->flag = 0;
5        m->guard = 0;
6        queue_init(m->q);
7    }
8
9    void lock(lock_t *m) {
10       while (TestAndSet(&m->guard, 1) == 1)
11           ; // acquire guard lock by spinning
12       if (m->flag == 0) {
13           m->flag = 1; // lock is acquired
14           m->guard = 0;
15       } else {
16           queue_add(m->q, gettid());
17           m->guard = 0;
18           park();
19       }
20   }
21   …
```

**Lock With Queues, Test-and-set, Yield, And Wakeup**

```
22  void unlock(lock_t *m) {
23      while (TestAndSet(&m->guard, 1) == 1)
24          ; // acquire guard lock by spinning
25      if (queue_empty(m->q))
26          m->flag = 0; // let go of lock; no one wants it
27      else
28          unpark(queue_remove(m->q)); // hold lock (for next thread!)
29      m->guard = 0;
30  }
```

**Lock With Queues, Test-and-set, Yield, And Wakeup (Cont.)**

# Two-Phase Locks

- A two-phase lock realizes that spinning can be useful if the lock *is about to* be released. 二阶段锁意识到如果锁即将释放，则自旋可能是有用的。

  - **First phase**

    - The lock spins for a while, *hoping that* it can acquire the lock.自旋一段时间，希望可以获取到锁定资源

    - If the lock is not acquired during the first spin phase, <u>a second phase</u> is entered.如果在第一阶段自旋期间没有获取锁定资源，则进入第二阶段

  - **Second phase**

    - The caller is put to sleep.

    - The caller is only woken up when the lock becomes free later.

1. **Concurrency Introduction**

2. **Locks**

3. **基于Lock的并发数据结构**

4. **Condition Variables 条件变量**

5. **Semaphore 信号量**

6. **常见并发问题**

7. **基于事件的并发**

# Lock-based Concurrent Data structure

- Adding locks to a data structure makes the structure **thread safe**.

  - How locks are added determine both the correctness and performance of the data structure.

- Simple but not scalable

```
1       typedef struct __counter_t {
2               int value;
3       } counter_t;
4
5       void init(counter_t *c) {
6               c->value = 0;
7       }
8
9       void increment(counter_t *c) {
10              c->value++;
11      }
12
13      void decrement(counter_t *c) {
14              c->value--;
15      }
16
17      int get(counter_t *c) {
18              return c->value;
19      }
```
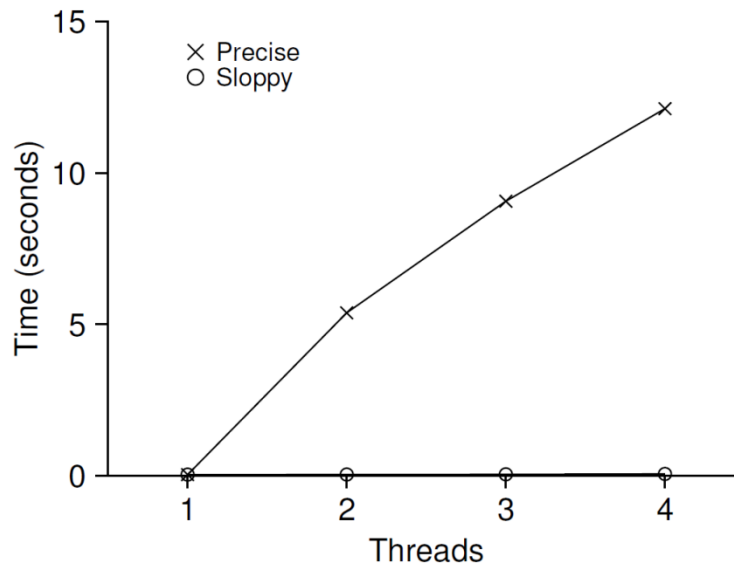
# Add a single lock

☐ acquired when calling a routine manipulating the data structure.

```c
1       typedef struct __counter_t {
2               int value;
3               pthread_lock_t lock;
4       } counter_t;
5
6       void init(counter_t *c) {
7               c->value = 0;
8               Pthread_mutex_init(&c->lock, NULL);
9       }
10
11      void increment(counter_t *c) {
12              Pthread_mutex_lock(&c->lock);
13              c->value++;
14              Pthread_mutex_unlock(&c->lock);
15 }
17      void decrement(counter_t *c) {
18              Pthread_mutex_lock(&c->lock);
19              c->value--;
20              Pthread_mutex_unlock(&c->lock);
21      }
22
23      int get(counter_t *c) {
24              Pthread_mutex_lock(&c->lock);
25              int rc = c->value;
26              Pthread_mutex_unlock(&c->lock);
27              return rc;
28      }
```

# The performance cost of the simple approach

- Each thread updates a single shared counter.
    - Each thread updates the counter one million times.
    - iMac with four Intel 2.7GHz i5 CPUs.

**Performance of
Traditional vs. Sloppy Counters**
(Threshold of Sloppy, *S*, is set to 1024)

**Synchronized counter scales poorly.**

同步计数器的扩展性差

# Perfect Scaling

☐ Even though more work is done, it is **done in parallel**.

☐ The time taken to complete the task is *not increased*.

# Sloppy counter

- The sloppy counter works by representing …

  - A single **logical counter** via numerous local physical counters, <u>on per CPU core</u> 在每个CPU核心上通过多个本地物理计数器来表示单个逻辑计数器

  - A single **global counter**

  - There are **lock**s:

    - ▸ One for each local counter and one for the global counter

- Example: on a machine with four CPUs

  - Four local counters

  - One global counter

# The basic idea of sloppy counting

- When a thread running on a core wishes to increment the counter.

  - It increment its local counter.

  - Each CPU has its own local counter:

    - Threads across CPUs can update local counters *without contention*. 跨CPU的线程可以更新本地计数器而不会发生争用

    - Thus counter updates are scalable.

  - The local values are periodically transferred to the global counter.局部值会定期传输到全局计数器

    - Acquire the global lock

    - Increment it by the local counter's value

    - The local counter is then reset to zero.

- How often the local-to-global transfer occurs is determined by a threshold, $S$ (sloppiness). 本地到全局传输发生的频率由阈值S决定。

  - The smaller $S$:

    - The more the counter behaves like the *non-scalable counter*.S越小计数器的行为越像不可扩展的计数器

  - The bigger $S$:

    - The more scalable the counter.

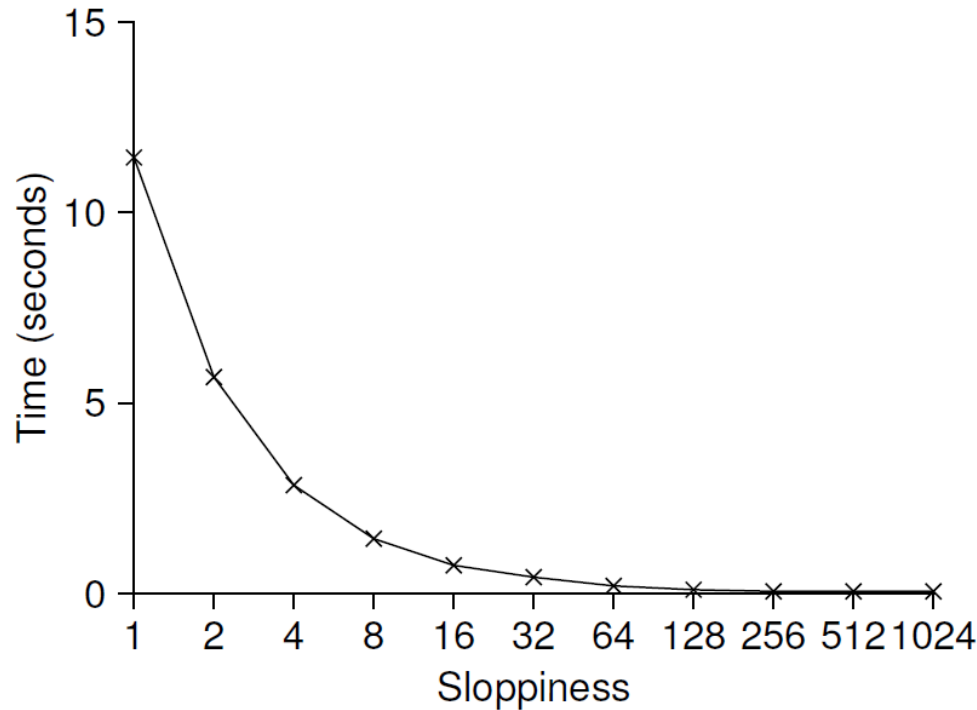    - The further off the global value might be from the *actual count*.全局值可能离计数器的实际值偏差越大

# Sloppy counter example

- Tracing the Sloppy Counters
  - The threshold S is set to 5.
  - There are threads on each of 4 CPUs
  - Each thread updates their local counters $L_1 \dots L_4$.

| Time | $L_1$ | $L_2$ | $L_3$ | $L_4$ | G |
|------|-------|-------|-------|-------|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 1 | 1 | 0 |
| 2 | 1 | 0 | 2 | 1 | 0 |
| 3 | 2 | 0 | 3 | 1 | 0 |
| 4 | 3 | 0 | 3 | 2 | 0 |
| 5 | 4 | 1 | 3 | 3 | 0 |
| 6 | 5 → 0 | 1 | 3 | 4 | 5 (from $L_1$) |
| 7 | 0 | 2 | 4 | 5 → 0 | 10 (from $L_4$) |

# Importance of the threshold value *s*

- Each four threads increments a counter 1 million times on four CPUs.

    - Low S → Performance is **poor**, The global count is always quire **accurate**.

    - High S → Performance is **excellent**, The global count **lags**.全局计数滞后



**Scaling Sloppy Counters**

# Sloppy Counter Implementation

```c
1        typedef struct __counter_t {
2            int global;                    // global count
3            pthread_mutex_t glock;         // global lock
4            int local[NUMCPUS];            // local count (per cpu)
5            pthread_mutex_t llock[NUMCPUS]; // ... and locks
6            int threshold;         // update frequency
7        } counter_t;
8
9        // init: record threshold, init locks, init values
10       //       of all local counts and global count
11       void init(counter_t *c, int threshold) {
12           c->thres hold = threshold;
13
14           c->global = 0;
15           pthread_mutex_init(&c->glock, NULL);
16
17           int i;
18           for (i = 0; i < NUMCPUS; i++) {
19               c->local[i] = 0;
20               pthread_mutex_init(&c->llock[i], NULL);
21           }
22       }
23
```

```
(Cont.)
24      // update: usually, just grab local lock and update local
amount
25      //          once local count has risen by 'threshold', grab
global
26      //          lock and transfer local values to it
27      void update(counter_t *c, int threadID, int amt) {
28          pthread_mutex_lock(&c->llock[threadID]);
29          c->local[threadID] += amt;          // assumes amt > 0
30          if (c->local[threadID] >= c->threshold) { // transfer
to global
31              pthread_mutex_lock(&c->glock);
32              c->global += c->local[threadID];
33              pthread_mutex_unlock(&c->glock);
34              c->local[threadID] = 0;
35          }
36          pthread_mutex_unlock(&c->llock[threadID]);
37      }
38
39      // get: just return global amount (which may not be perfect)
40      int get(counter_t *c) {
41          pthread_mutex_lock(&c->glock);
42          int val = c->global;
43          pthread_mutex_unlock(&c->glock);
44          return val;          // only approximate!
45      }
```

# Concurrent Linked Lists

```c
1          // basic node structure
2        typedef struct __node_t {
3                int key;
4                struct __node_t *next;
5        } node_t;
6
7        // basic list structure (one used per list)
8        typedef struct __list_t {
9                node_t *head;
10               pthread_mutex_t lock;
11       } list_t;
12
13       void List_Init(list_t *L) {
14               L->head = NULL;
15               pthread_mutex_init(&L->lock, NULL);
16       }
17
(Cont.)
```

# Concurrent Linked Lists

```
(Cont.)
18      int List_Insert(list_t *L, int key) {
19              pthread_mutex_lock(&L->lock);
20              node_t *new = malloc(sizeof(node_t));
21              if (new == NULL) {
22                      perror("malloc");
23                      pthread_mutex_unlock(&L->lock);
24          return -1; // fail
26          new->key = key;
27          new->next = L->head;
28          L->head = new;
29          pthread_mutex_unlock(&L->lock);
30          return 0; // success
31  }
32      int List_Lookup(list_t *L, int key) {
33              pthread_mutex_lock(&L->lock);
34              node_t *curr = L->head;
35              while (curr) {
36                      if (curr->key == key) {
37                              pthread_mutex_unlock(&L->lock);
38                              return 0; // success
39                      }
40                      curr = curr->next;
41              }
42              pthread_mutex_unlock(&L->lock);
43              return -1; // failure
44      }
```

# Concurrent Linked Lists (Cont.) 并发链表

- The code **acquires** a lock in the insert routine upon entry.在进入时获取插入操作例程中的锁

- The code **releases** the lock upon exit.在退出时释放锁

  - If `malloc()` happens to *fail*, the code must also <u>release the lock</u> before failing the insert.如果malloc失败，必须在插入操作失败前释放锁

  - This kind of exceptional control flow has been shown to be quite error prone.这种异常的控制流已被证明非常容易出错

  - **Solution**: The lock and release *only surround* the actual critical section in the insert code锁定和释放仅围绕插入操作代码的关键部分

```
1          void List_Init(list_t *L) {
2                  L->head = NULL;
3                  pthread_mutex_init(&L->lock, NULL);
4          }
5
6          void List_Insert(list_t *L, int key) {
7                  // synchronization not needed
8                  node_t *new = malloc(sizeof(node_t));
9                  if (new == NULL) {
10                         perror("malloc");
11                         return;
12                 }
13                 new->key = key;
14
15                 // just lock critical section
16                 pthread_mutex_lock(&L->lock);
17                 new->next = L->head;
18                 L->head = new;
19                 pthread_mutex_unlock(&L->lock);
20         }
21
```

```
(Cont.)
22      int List_Lookup(list_t *L, int key) {
23              int rv = -1;
24              pthread_mutex_lock(&L->lock);
25              node_t *curr = L->head;
26              while (curr) {
27                      if (curr->key == key) {
28                              rv = 0;
29                              break;
30                      }
31                      curr = curr->next;
32              }
33              pthread_mutex_unlock(&L->lock);
34              return rv; // now both success and failure
35      }
```

# Scaling Linked List

- Hand-over-hand locking (lock coupling) 锁耦合

  - Add **a lock per node** of the list instead of having a single lock for the entire list.为列表的每个节点添加一个锁，而不是为整个列表添加一个锁

  - When traversing the list,

    - First grabs the next node's lock.

    - And then releases the current node's lock.

  - Enable a high degree of concurrency in list operations.在列表操作中启用高度并发

    - However, in practice, <u>the overheads of </u>acquiring and releasing locks for each node of a list traversal is *prohibitive*.但在实践中，为列表遍历每个节点获取和释放锁的开销难以接受

- There are two locks.

    - One for the **head** of the queue.

    - One for the **tail**.

    - The goal of these two locks is to enable concurrency of *enqueue* and *dequeue* operations.这两个锁的目标是启用入队和出队操作的并发性

- Add a dummy node

    - Allocated in the queue initialization code 在队列初始化代码中分配

    - Enable the separation of head and tail operations隔离头部和尾部的操作

```
1          typedef struct __node_t {
2                  int value;
3                  struct __node_t *next;
4          } node_t;
5
6          typedef struct __queue_t {
7                  node_t *head;
8                  node_t *tail;
9                  pthread_mutex_t headLock;
10                 pthread_mutex_t tailLock;
11         } queue_t;
12
13         void Queue_Init(queue_t *q) {
14                 node_t *tmp = malloc(sizeof(node_t));
15                 tmp->next = NULL;
16                 q->head = q->tail = tmp;
17                 pthread_mutex_init(&q->headLock, NULL);
18                 pthread_mutex_init(&q->tailLock, NULL);
19         }
20
(Cont.)
```

```
(Cont.)
21      void Queue_Enqueue(queue_t *q, int value) {
22              node_t *tmp = malloc(sizeof(node_t));
23              assert(tmp != NULL);
24
25              tmp->value = value;
26              tmp->next = NULL;
27
28              pthread_mutex_lock(&q->tailLock);
29              q->tail->next = tmp;
30              q->tail = tmp;
31               pthread_mutex_unlock(&q->tailLock);
32      }
33      int Queue_Dequeue(queue_t *q, int *value) {
34              pthread_mutex_lock(&q->headLock);
35              node_t *tmp = q->head;
36              node_t *newHead = tmp->next;
37              if (newHead == NULL) {
38                      pthread_mutex_unlock(&q->headLock);
39                      return -1; // queue was empty
40              }
41              *value = newHead->value;
42              q->head = newHead;
43              pthread_mutex_unlock(&q->headLock);
44              free(tmp);
45              return 0;
46      }
```
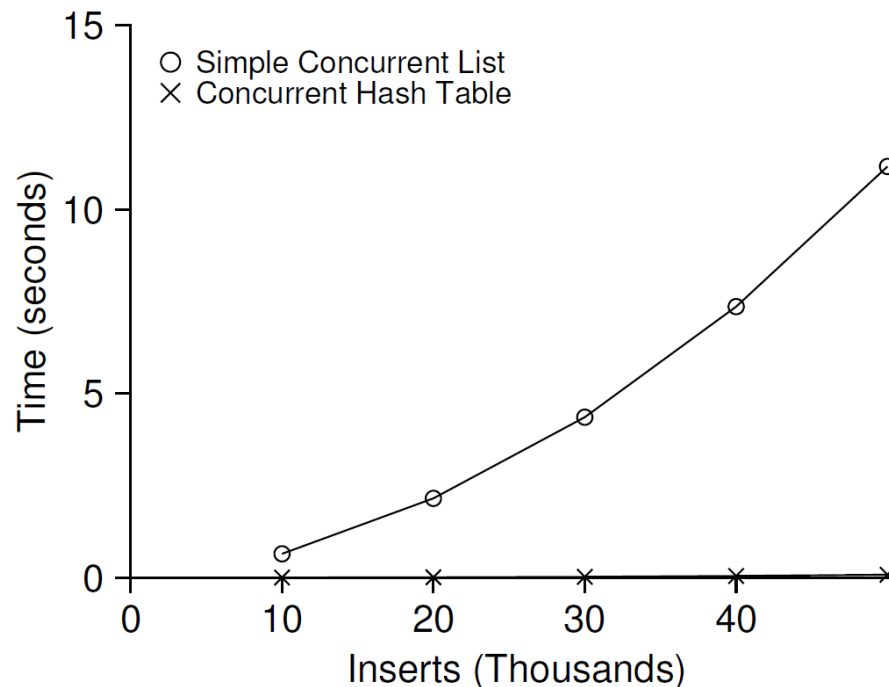
# Concurrent Hash Table

- Focus on a simple hash table

  - The hash table does not resize.

  - Built using the concurrent lists

  - It uses a lock per hash bucket each of which is represented by *a list*.
    它对每个哈希桶使用一个锁，每个哈希桶都由一个列表表示。

# Performance of Concurrent Hash Table

- From 10,000 to 50,000 concurrent updates from each of four threads.
  - iMac with four Intel 2.7GHz i5 CPUs.



**The simple concurrent hash table scales magnificently.**

```
1          #define BUCKETS (101)
2
3          typedef struct __hash_t {
4                  list_t lists[BUCKETS];
5          } hash_t;
6
7          void Hash_Init(hash_t *H) {
8                  int i;
9                  for (i = 0; i < BUCKETS; i++) {
10                         List_Init(&H->lists[i]);
11                 }
12         }
13
14         int Hash_Insert(hash_t *H, int key) {
15                 int bucket = key % BUCKETS;
16                 return List_Insert(&H->lists[bucket], key);
17         }
18
19         int Hash_Lookup(hash_t *H, int key) {
20                 int bucket = key % BUCKETS;
21                 return List_Lookup(&H->lists[bucket], key);
22         }
```

1. **Concurrency Introduction**

2. **Locks**

3. **基于Lock的并发数据结构**

4. **Condition Variables 条件变量**

5. **Semaphore 信号量**

6. **常见并发问题**

7. **基于事件的并发**

# Condition Variables 条件变量的引入

- There are many cases where a thread wishes to <u>check</u> whether a **condition** is true before continuing its execution. 一个线程需要检查另一个的状态，并据此决定自己是否继续执行

- Example:
  - A parent thread might wish to check whether a child thread has *completed*.
  - This is often called a `join()`.

# Condition Variables (Cont.)

**A Parent Waiting For Its Child**

```
1          void *child(void *arg) {
2              printf("child\n");
3              // XXX how to indicate we are done?
4              return NULL;
5          }
6
7          int main(int argc, char *argv[]) {
8              printf("parent: begin\n");
9              pthread_t c;
10             Pthread_create(&c, NULL, child, NULL); // create child
11             // XXX how to wait for child?
12             printf("parent: end\n");
13             return 0;
14         }
```

**What we would like to see here is:**

```
parent: begin
child
parent: end
```

```
1          volatile int done = 0;
2
3      void *child(void *arg) {
4          printf("child\n");
5          done = 1;
6          return NULL;
7      }
8
9      int main(int argc, char *argv[]) {
10         printf("parent: begin\n");
11         pthread_t c;
12         Pthread_create(&c, NULL, child, NULL); // create child
13         while (done == 0)
14             ; // spin
15         printf("parent: end\n");
16         return 0;
17     }
```

☐ This is hugely <u>inefficient</u> as the parent spins and **wastes CPU time**.

# How to wait for a condition

- Condition variable 本质上是一个队列及对该队列的操作原语

  - **Waiting** on the condition

    - An explicit queue that threads can put themselves on when some state of execution is not as desired. 当某些执行状态不符合要求时，Waiting让该线程将自己放入对应的显式队列中

  - **Signaling** on the condition

    - Some other thread, *when it changes said state*, can wake one of those waiting threads and allow them to continue. 其他一些线程在改变所述状态时，可以唤醒其中一个等待线程并允许它们继续

# Definition and Routines

- Declare condition variable

```
pthread cond t c;
```

  - Proper initialization is required.

- Operation (the POSIX calls)

```
pthread_cond_wait(pthread_cond_t *c, pthread_mutex_t *m);   // wait()
pthread_cond_signal(pthread_cond_t *c);                      // signal()
```

  - The wait() call takes a <u>mutex</u> as a parameter.
    - The wait() call release the lock and put the calling thread to sleep.
    - When the thread wakes up, it must re-acquire the lock.

```
1           int done = 0;
2           pthread_mutex_t m = PTHREAD_MUTEX_INITIALIZER;
3           pthread_cond_t c = PTHREAD_COND_INITIALIZER;
4
5       void thr_exit() {
6               Pthread_mutex_lock(&m);
7               done = 1;
8               Pthread_cond_signal(&c);
9               Pthread_mutex_unlock(&m);
10      }
11
12      void *child(void *arg) {
13              printf("child\n");
14              thr_exit();
15              return NULL;
16      }
17
18      void thr_join() {
19              Pthread_mutex_lock(&m);
20              while (done == 0)
21                      Pthread_cond_wait(&c, &m);
22              Pthread_mutex_unlock(&m);
23      }
24
```

```
(cont.)
25      int main(int argc, char *argv[]) {
26              printf("parent: begin\n");
27              pthread_t p;
28              Pthread_create(&p, NULL, child, NULL);
29              thr_join();
30              printf("parent: end\n");
31              return 0;
32      }
```

Create the child thread and continues running itself.

- **Parent:**
  - Create the child thread and continues running itself.
  - Call into `thr_join()` to wait for the child thread to complete.
    - Acquire the lock
    - Check if the child is done
    - Put itself to sleep by calling `wait()`
    - Release the lock
- **Child:**
  - Print the message "child"
  - Call `thr_exit()` to wake the parent thread
    - Grab the lock
    - Set the state variable `done`
    - Signal the parent thus waking it.

# The importance of the state variable done

```
1        void thr_exit() {
2                Pthread_mutex_lock(&m);
3                Pthread_cond_signal(&c);
4                Pthread_mutex_unlock(&m);
5        }
6
7        void thr_join() {
8                Pthread_mutex_lock(&m);
9                Pthread_cond_wait(&c, &m);
10               Pthread_mutex_unlock(&m);
11       }
```

**`thr_exit()` and `thr_join()` without variable `done`**

- Imagine the case where the *child runs immediately*.
  - The child will signal, but there is <u>no thread asleep</u> on the condition.
  - When the parent runs, it will call wait and be stuck.
  - No thread will ever wake it.

# Another poor implementation

```
1        void thr_exit() {
2                done = 1;
3                Pthread_cond_signal(&c);
4        }
5
6        void thr_join() {
7                if (done == 0)
8                        Pthread_cond_wait(&c);
9        }
```

- The issue here is a subtle **race condition**.
  - The parent calls `thr_join()`.
    - The parent checks the value of `done`.
    - It will see that it is 0 and try to go to sleep.
    - *Just before* it calls wait to go to sleep, the parent is <u>interrupted</u> and the child runs.
  - The child changes the state variable `done` to 1 and signals.
    - But no thread is waiting and thus no thread is woken.
    - When the parent runs again, it sleeps forever.

# Module 5: Concurrency & Synchronization
## 并发与同步

1. **Concurrency Introduction**

2. **Locks**

3. **基于Lock的并发数据结构**

4. **Condition Variables 条件变量**

5. **<span style="color:red">Semaphore 信号量</span>**

6. **常见并发问题**

7. **基于事件的并发**

# Semaphore: A definition

- An object with an integer value

  - We can manipulate(操作)with two routines; `sem_wait()` and `sem_post()`.

  - Initialization

```
1    #include <semaphore.h>
2    sem_t s;
3    sem_init(&s, 0, 1); // initialize s to the value 1
```

  - ▸ Declare a semaphore `s` and initialize it to the value 1
  - ▸ The second argument, 0, indicates that the semaphore is <u>shared</u> between *threads in the same process*.第二个参数0表示信号量在同一个进程中的线程之间共享

# Semaphore: wait原语

- `sem_wait()`

```
1  int sem_wait(sem_t *s) {
2      decrement the value of semaphore s by one
3      wait if value of semaphore s is negative
4  }
```

- 减1，若大于0，则返回，否则挂起等待被post唤醒

- If the value of the semaphore was *one* or *higher* when called `sem_wait()`, **return right away**.

- It will cause the caller to <u>suspend execution</u> waiting for a subsequent post.它将导致调用者暂停执行，等待后续的post操作

- When negative, the value of the semaphore is equal to the number of waiting threads.当为负数时，信号量的值等于等待线程的数量

# Semaphore: post原语

- sem_post()

```
1   int sem_post(sem_t *s) {
2       increment the value of semaphore s by one
3       if there are one or more threads waiting, wake one
4   }
```

- 加1，若有等待线程，则唤醒一个

- Simply **increments** the value of the semaphore.

- If there is a thread waiting to be woken, **wakes** one of them up.

# Binary Semaphores (Locks)

- What should **x** be?
  - The initial value should be **1**.

```
1   sem_t m;
2   sem_init(&m, 0, X); // initialize semaphore to X; what should X be?
3
4   sem_wait(&m);
5   //critical section here 临界区
6   sem_post(&m);
```

# Thread Trace: Single Thread Using A Semaphore

| Value of Semaphore | Thread 0 | Thread 1 |
|:---:|:---|:---:|
| 1 | | |
| 1 | call sema_wait() | |
| 0 | sem_wait() returns | |
| 0 | (crit sect) | |
| 0 | call sem_post() | |
| 1 | sem_post() returns | |

# Thread Trace: Two Threads Using A Semaphore

| Value | Thread 0 | State | Thread 1 | State |
|---|---|---|---|---|
| 1 | | Running | | Ready |
| 1 | call sem_wait() | Running | | Ready |
| 0 | sem_wait() retruns | Running | | Ready |
| 0 | (crit set: begin) | Running | | Ready |
| 0 | *Interrupt; Switch → T1* | Ready | | Running |
| 0 | | Ready | call sem_wait() | Running |
| -1 | | Ready | decrement sem | Running |
| -1 | | Ready | (sem < 0)→sleep | sleeping |
| -1 | | Running | *Switch → T0* | sleeping |
| -1 | (crit sect: end) | Running | | sleeping |
| -1 | call sem_post() | Running | | sleeping |
| 0 | increment sem | Running | | sleeping |
| 0 | wake(T1) | Running | | Ready |
| 0 | sem_post() returns | Running | | Ready |
| 0 | *Interrupt; Switch → T1* | Ready | | Running |
| 0 | | Ready | sem_wait() retruns | Running |
| 0 | | Ready | (crit sect) | Running |
| 0 | | Ready | call sem_post() | Running |
| 1 | | Ready | sem_post() returns | Running |

```
1    sem_t s;
2
3    void *
4    child(void *arg) {
5        printf("child\n");
6        sem_post(&s); // signal here: child is done
7        return NULL;
8    }
9
10   int
11   main(int argc, char *argv[]) {
12       sem_init(&s, 0, X); // what should X be?
13       printf("parent: begin\n");
14       pthread_t c;
15       pthread_create(c, NULL, child, NULL);
16       sem_wait(&s); // wait here for child
17       printf("parent: end\n");
18       return 0;
19   }
```

**A Parent Waiting For Its Child**

```
parent: begin
child
parent: end
```

**The execution result**

- What should **x** be?

  ▸ The value of semaphore should be set to is **0**.

☐ The parent call `sem_wait()` before the child has called `sem_post()`.

| Value | Parent | State | Child | State |
|---|---|---|---|---|
| 0 | `Create(Child)` | Running | *(Child exists; is runnable)* | Ready |
| 0 | `call sem_wait()` | Running | | Ready |
| -1 | `decrement sem` | Running | | Ready |
| -1 | `(sem < 0)→sleep` | sleeping | | Ready |
| -1 | *Switch→Child* | sleeping | `child runs` | Running |
| -1 | | sleeping | `call sem_post()` | Running |
| 0 | | sleeping | `increment sem` | Running |
| 0 | | Ready | `wake(Parent)` | Running |
| 0 | | Ready | `sem_post() returns` | Running |
| 0 | | Ready | *Interrupt; Switch→Parent* | Ready |
| 0 | `sem_wait() retruns` | Running | | Ready |

# Thread Trace: Parent Waiting For Child (Case 2)

- The child runs to completion before the parent call `sem_wait()`.

| Value | Parent | State | Child | State |
|-------|--------|-------|-------|-------|
| 0 | `Create(Child)` | Running | *(Child exists; is runnable)* | Ready |
| 0 | *Interrupt; switch→Child* | Ready | `child runs` | Running |
| 0 | | Ready | `call sem_post()` | Running |
| 1 | | Ready | ` increment sem` | Running |
| 1 | | Ready | ` wake(nobody)` | Running |
| 1 | | Ready | `sem_post() returns` | Running |
| 1 | `parent runs` | Running | *Interrupt; Switch→Parent* | Ready |
| 1 | `call sem_wait()` | Running | | Ready |
| 0 | ` decrement sem` | Running | | Ready |
| 0 | ` (sem<0)→awake` | Running | | Ready |
| 0 | `sem_wait() retruns` | Running | | Ready |

# The Producer/Consumer (Bounded-Buffer) Problem

- **Producer**: `put()` interface

  - Wait for a buffer to become *empty* in order to put data into it.等待缓冲区变为空，以便将数据放入其中

- **Consumer**: `get()` interface

  - Wait for a buffer to become *filled* before using it.等待缓冲区被填满后再使用它。

```c
1   int buffer[MAX];
2   int fill = 0;
3   int use = 0;
4
5   void put(int value) {
6       buffer[fill] = value;    // line f1
7       fill = (fill + 1) % MAX;  // line f2
8   }
9
10  int get() {
11      int tmp = buffer[use];   // line g1
12      use = (use + 1) % MAX;   // line g2
13      return tmp;
14  }
```

```
1     sem_t empty;
2     sem_t full;
3
4     void *producer(void *arg) {
5         int i;
6         for (i = 0; i < loops; i++) {
7                 sem_wait(&empty);        // line P1
8                 put(i);                  // line P2
9                 sem_post(&full);         // line P3
10        }
11    }
12
13    void *consumer(void *arg) {
14        int i, tmp = 0;
15        while (tmp != -1) {
16                sem_wait(&full);         // line C1
17                tmp = get();             // line C2
18                sem_post(&empty);        // line C3
19                printf("%d\n", tmp);
20        }
21    }
22    …
```

**First Attempt: Adding the Full and Empty Conditions**

```
21   int main(int argc, char *argv[]) {
22       // …
23       sem_init(&empty, 0, MAX);         // MAX buffers are empty to begin with…
24       sem_init(&full, 0, 0);            // … and 0 are full
25       // …
26   }
```

**First Attempt: Adding the Full and Empty Conditions (Cont.)**

- Imagine that `MAX` is greater than 1 .

  - If there are multiple producers, race condition can happen at line *f1*. 如果有多个生产者，竞态条件可能发生在line f1

  - It means that the old data there is overwritten.这意味着旧数据被覆盖了

- We've forgotten here is **mutual exclusion**.

  - The filling of a buffer and incrementing of the index into the buffer is a critical section.缓冲区的填充和向缓冲区增加索引是临界段

# A Solution: Adding Mutual Exclusion

```
1    sem_t empty;
2    sem_t full;
3    sem_t mutex;
4
5    void *producer(void *arg) {
6        int i;
7        for (i = 0; i < loops; i++) {
8                sem_wait(&mutex);        // line p0 (NEW LINE)
9                sem_wait(&empty);        // line p1
10               put(i);                  // line p2
11               sem_post(&full);         // line p3
12               sem_post(&mutex);        // line p4 (NEW LINE)
13       }
14   }

16   void *consumer(void *arg) {
17       int i;
18       for (i = 0; i < loops; i++) {
19               sem_wait(&mutex);        // line c0 (NEW LINE)
20               sem_wait(&full);         // line c1
21               int tmp = get();         // line c2
22               sem_post(&empty);        // line c3
23               sem_post(&mutex);        // line c4 (NEW LINE)
24               printf("%d\n", tmp);
25       }
26   }
```

**Adding Mutual Exclusion (Incorrectly)**

- Imagine two thread: one producer and one consumer.
    - The consumer **acquire** the `mutex` (line c0).
    - The consumer **calls** `sem_wait()` on the full semaphore (line c1).
    - The consumer is **blocked** and **yield** the CPU.
        - ▸ The consumer <u>still holds the mutex</u>!
    - The producer **calls** `sem_wait()` on the binary `mutex` semaphore (line p0).
    - The producer is now **stuck** waiting too. a classic deadlock.

# Finally, A Working Solution

```
1    sem_t empty;
2    sem_t full;
3    sem_t mutex;
4
5    void *producer(void *arg) {
6        int i;
7        for (i = 0; i < loops; i++) {
8                sem_wait(&empty);        // line p1
9                sem_wait(&mutex);        // line p1.5 (MOVED MUTEX HERE…)
10               put(i);                  // line p2
11               sem_post(&mutex);        // line p2.5 (… AND HERE)
12               sem_post(&full);         // line p3
13       }
14   }
15
16    void *consumer(void *arg) {
17       int i;
18       for (i = 0; i < loops; i++) {
19               sem_wait(&full);         // line c1
20               sem_wait(&mutex);        // line c1.5 (MOVED MUTEX HERE…)
21               int tmp = get();         // line c2
22               sem_post(&mutex);        // line c2.5 (… AND HERE)
23               sem_post(&empty);        // line c3
24               printf("%d\n", tmp);
25       }
26   }
```

**Adding Mutual Exclusion** **(Correctly)**

# Finally, A Working Solution

```
(Cont.)

27  int main(int argc, char *argv[]) {
28      // …
29      sem_init(&empty, 0, MAX); // MAX buffers are empty to begin with …
30      sem_init(&full, 0, 0);    // ... and 0 are full
31      sem_init(&mutex, 0, 1);   // mutex=1 because it is a lock
32      // …
33  }
```

# Reader-Writer Locks

- Imagine a number of concurrent list operations, including **inserts** and simple **lookups**.想象一下许多并发的列表操作，包括插入和简单的查找

  - **insert:**
    - Change the state of the list更改列表的状态
    - A traditional <u>critical section</u> makes sense.传统的临界区是有意义的。

  - **lookup:**
    - Simply *read* the data structure.只需读取数据结构
    - As long as we can guarantee that no insert is on-going, we can allow many lookups to proceed concurrently.只要能够保证没有插入正在进行，就可以允许多个查找同时进行

> **This special type of lock is known as a reader-write lock.**

# A Reader-Writer Locks

- Only **a single writer** can acquire the lock.只有一个写者可以获得锁

- Once a reader has acquired a read lock,一旦读者获得了读锁

  - **More readers** will be allowed to acquire the read lock too.更多的读者也将被允许获得读锁。

  - A writer will <u>have to wait</u> until all readers are finished.写者必须等到所有读者都看完。

```
1   typedef struct _rwlock_t {
2       sem_t lock;        // binary semaphore (basic lock)
3       sem_t writelock;   // used to allow ONE writer or MANY readers
4       int readers;       // count of readers reading in critical section
5   } rwlock_t;
6
7   void rwlock_init(rwlock_t *rw) {
8       rw->readers = 0;
9       sem_init(&rw->lock, 0, 1);
10      sem_init(&rw->writelock, 0, 1);
11  }
12
13  void rwlock_acquire_readlock(rwlock_t *rw) {
14      sem_wait(&rw->lock);
15      …
```

```
15          rw->readers++;
16          if (rw->readers == 1)
17                  sem_wait(&rw->writelock); // first reader acquires writelock
18          sem_post(&rw->lock);
19  }
20
21  void rwlock_release_readlock(rwlock_t *rw) {
22          sem_wait(&rw->lock);
23          rw->readers--;
24          if (rw->readers == 0)
25                  sem_post(&rw->writelock); // last reader releases writelock
26          sem_post(&rw->lock);
27  }
28
29  void rwlock_acquire_writelock(rwlock_t *rw) {
30          sem_wait(&rw->writelock);
31  }
32
33  void rwlock_release_writelock(rwlock_t *rw) {
34          sem_post(&rw->writelock);
35  }
```

# A Reader-Writer Locks (Cont.)

- The reader-writer locks have fairness(公平) problem.

  - It would be relatively easy for reader to **starve writer**.对于读者来说，饿死写者是相对容易的

  - How to <u>prevent</u> more readers from entering the lock once a writer is waiting?在写者等待时，如何防止更多的读者进入锁？

# The Dining Philosophers

- Assume there are five "**philosophers**" sitting around a table.
  - Between each pair of philosophers is <u>a single fork</u> (five total).每一对哲学家之间都有一个叉(总共五个叉)
  - The philosophers each have times where they **think**, and don't need any forks, and times where they **eat**.每个哲学家都有他们思考的时候，他们不需要叉子，也有他们吃饭的时候。
  - In order to *eat*, a philosopher needs two forks, both the one on their *left* and the one on their *right*.为了吃饭，哲学家需要两把叉子，一把在左边，一把在右边
  - **The contention(争用) for these forks.**

# The Dining Philosophers (Cont.)

- ☐ Key challenge
  - ☐ There is **no deadlock()**.
  - ☐ **No** philosopher **starves** and never gets to eat. 没有一个哲学家挨饿不吃东西
  - ☐ **Concurrency(并发)** is high.

```
while (1) {
        think();
        getforks();
        eat();
        putforks();
}
```

**Basic loop of each philosopher**

```
// helper functions
int left(int p) { return p; }

int right(int p) {
        return (p + 1) % 5;
}
```

**Helper functions (Downey's solutions)**

- ‣ Philosopher `p` wishes to refer to the fork on their left → call `left(p)`.
- ‣ Philosopher `p` wishes to refer to the fork on their right → call `right(p)`.

- We need some **semaphore**, one for each fork: `sem_t forks[5]`.

```
1    void getforks() {
2        sem_wait(forks[left(p)]);
3        sem_wait(forks[right(p)]);
4    }
5
6    void putforks() {
7        sem_post(forks[left(p)]);
8        sem_post(forks[right(p)]);
9    }
```

The `getforks()` and `putforks()` Routines (Broken Solution)

- Deadlock occur!

  ‣ If each philosopher happens to **grab the fork on their left** before any philosopher can grab the fork on their right.如果每一位哲学家碰巧在任何一位哲学家抓住他们右边的叉子之前，抓住了他们左边的叉子

  ‣ Each will be stuck *holding one fork* and waiting for another, *forever*. 每个人都会被卡住，拿着一个叉子，永远等着另一个

# A Solution: Breaking The Dependency

- Change how forks are acquired.

  - Let's assume that philosopher 4 acquire the forks in a *different order*.

```
1    void getforks() {
2        if (p == 4) {
3                sem_wait(forks[right(p)]);
4                sem_wait(forks[left(p)]);
5        } else {
6                sem_wait(forks[left(p)]);
7                sem_wait(forks[right(p)]);
8        }
9    }
```

  ▸ There is no situation where each philosopher grabs one fork and is stuck waiting for another. **The cycle of waiting is broken**.不存在每一位哲学家抓住一个叉子而等待另一个的情况。等待的循环被打破。

# How To Implement Semaphores

☐ Build our own version of semaphores called Zemaphores

```c
1    typedef struct __Zem_t {
2        int value;
3        pthread_cond_t cond;
4        pthread_mutex_t lock;
5    } Zem_t;
6
7    // only one thread can call this
8    void Zem_init(Zem_t *s, int value) {
9        s->value = value;
10       Cond_init(&s->cond);
11       Mutex_init(&s->lock);
12   }
13
14   void Zem_wait(Zem_t *s) {
15       Mutex_lock(&s->lock);
16       while (s->value <= 0)
17           Cond_wait(&s->cond, &s->lock);
18       s->value--;
19       Mutex_unlock(&s->lock);
20   }

22   void Zem_post(Zem_t *s) {
23       Mutex_lock(&s->lock);
24       s->value++;
25       Cond_signal(&s->cond);
26       Mutex_unlock(&s->lock);
27   }
```

# How To Implement Semaphores (Cont.)

- Zemaphore don't maintain the invariant that *the value of* the semaphore. Zemaphore不保持信号量值不变

  - The value <u>never be lower than zero</u>.该值不得低于0

  - This behavior is **easier** to implement and **matches** the current Linux implementation.这种行为更易于实现，并与当前的Linux实现相匹配

# 例题

【例题】对于两个并发进程，设互斥信号量为mutex（初值为1），若 mutex=-1，则（　）。

A. 表示没有进程进入临界区

B. 表示有一个进程进入临界区

C. 表示有两个进程进入临界区

D. 表示有一个进程进入临界区，另一个进程在等待进入

答案：D

P操作：进入临界区，mutex减1。

V操作：释放临界区，mutex加1。

mutex初值为1，此时mutex=-1说明已经有一个进程在临界区，另一个在等待。

# 例题

【例题】当 V 操作唤醒一个等待进程时，被唤醒进程变为（ ）态。

A．运行

B．阻塞

C．就绪

D．完成

答案：C

等待唤醒的进程处于阻塞态，被唤醒后进入就绪态。只有就绪进程能获得处理器资源，被唤醒的进程并不能直接转化为运行态。



进程：状态转换

# 例题

【例题】有三个进程共享同一程序段，而每次只允许两个进程进入该程序段，若用 PV 操作同步机制，则信号量 S 的取值范围是（ ）

A．2,1,0,-1

B．3,2,1,0

C．2,1,0,-1,-2

D．1,0,-1,-2

答案：A

因为每次允许两个进程进入该程序段，信号量最大值取2（否则三个进程可以同时进入程序段）。至多有三个进程申请，则信号量最小为-1。

【例题】若一个信号量的初值为3，经过多次PV操作后当前值为-1，这表示等待进入临界区的进程数是（ ）。

A. 1

B. 2

C. 3

D. 4

答案：A

信号量的初值为3，表示可以有3个进程进入临界区。如果经过多次PV操作后，信号量的当前值为-1，则表示当前有4个进程尝试进入临界区，其中3个进程已进入，1个进程在等待。

【2013统考真题】某博物馆最多可容纳500人同时参观，有一个出入口，该出入口一次仅允许一人通过，参观者的活动描述如下：

```
cobegin
        参观者进程i;
        {
                ...
                进门;
                ...
                参观;
                ...
                出门;
                ...
        }
coend
```

请添加必要的信号量和 P, V[或 wait(), signal()]操作，以实现上述过程中的互斥与同步。要求写出完整的过程，说明信号量的含义并赋值。

解答：
出入口一次仅允许一个人通过，设置互斥信号量mutex，初值为1。博物馆最多可同时容纳500人，因此设置信号量empty，初值为500。

```
Semaphore empty = 500; // 博物馆可以容纳的最多人数
Semaphore mutex = 1; // 用于出入口资源的控制
cobegin
参观者进程i:
{
        ...
        P(empty); // 可容纳人数减1
        P(mutex); // 互斥使用门
        进门;
        V(mutex);
        参观;
        P(mutex); // 互斥使用门
        出门;
        V(mutex);
        V(empty); // 可容纳人数增加1
        ...
}
coend
```

# 例题

【2011统考真题】某银行提供1个服务窗口和10个供顾客等待的座位。顾客到达银行时，若有空座位，则到取号机上领取一个号，等待叫号。取号机每次仅允许一位顾客使用。当营业员空闲时，通过叫号选取一位顾客，并为其服务。顾客和营业员的活动过程描述如下：

```
cobegin
{
        process 顾客i
        {
                从取号机获取一个号码;
                等待叫号;
                获取服务;
        }
}
```

```
{
        process 营业员
        {
                while(TRUE)
                {
                        叫号;
                        为顾客服务;
                }
        }
}
coend
```

请添加必要的信号量和 P, V[或 wait(), signal()]操作，实现上述过程中的互斥与同步。要求写出完整的过程，说明信号量的含义并赋值。

# 例题

解答：
互斥资源：取号机（一次只有一位顾客领号），因此设置互斥信号量mutex。
同步问题：顾客需要获得空座位等待叫号。营业员空闲时，将选取一位顾客并为其服务。是否有空座位影响等待顾客的数量，是否有顾客决定了营业员是否能开始服务，因此分别设置信号量empty和full来实现这一同步关系。
另外，顾客获得空座位后，需要等待叫号和被服务。顾客与营业员就服务何时开始又构成了一个同步关系，定义信号量service来完成这一同步过程。

```
semaphore empty = 10; // 空座位数量
semaphore full = 0; // 已占座位的数量
semaphore mutex = 1; // 互斥使用取号机
semaphore service = 0; // 等待叫号（当前是否正在服务）

cobegin
{
    Process 顾客 i{
        P(empty); // 等空位
        P(mutex); // 申请使用取号机
        取号;
        V(mutex); // 取号结束
        V(full); // 通知营业员有新顾客
        P(service); // 等待营业员叫号
        接收服务
    }
                                {
                                    Process 营业员{
                                        while(True){
                                            P(full); // 没有顾客则休息
                                            V(empty); // 离开座位
                                            V(service); // 叫号
                                            为顾客服务;
                                        }
                                    }
                                }
                                coend
}
```

# 信号量问题解题步骤

## 1.找出问题中所有同步与互斥的关系

- 互斥
  - 找到进程竞争的临界资源
  - 抓住"仅允许"或类似词汇
  - 博物馆出入口每次仅允许一人通过、取号机每次仅允许一位顾客使用
- 同步
  - 不同进程对资源合作处理
  - 博物馆内容纳500人、空座位10个

## 2.确定信号量个数及每个信号量的初值

- 互斥
  - 用1个信号量，代表资源是否被互斥使用
  - 一般初值为0或1
- 同步
  - 用1或2个信号量，1个信号量用于判断是否资源为empty或full，2个信号量分别判断资源是否为empty和full
  - 一般初值为题目中给出的特定数值（容纳500人、空座位10个）或0

## 3.用类似程序的语言描述算法

# Module 5: Concurrency & Synchronization
## 并发与同步

1. **Concurrency Introduction**

2. **Locks**

3. **基于Lock的并发数据结构**

4. **Condition Variables 条件变量**

5. **Semaphore 信号量**

6. **常见并发问题**

7. **基于事件的并发**

# Common Concurrency Problems

- More recent work focuses on studying other types of common concurrency bugs. 而不是deadlock 越来越多近期的工作专注于研究其它类型的常见并发问题，而不是死锁

  - Take a brief look at some example concurrency problems found in real code bases. 简单看一下在实际代码库中发现的一些并发问题事例

# What Types Of Bugs Exist?

- Focus on four major open-source applications
  - MySQL, Apache, Mozilla, OpenOffice.

| Application | What it does | Non-Deadlock | Deadlock |
|---|---|---|---|
| MySQL | Database Server | 14 | 9 |
| Apache | Web Server | 13 | 4 |
| Mozilla | Web Browser | 41 | 16 |
| Open Office | Office Suite | 6 | 2 |
| **Total** | | **74** | **31** |

**Bugs In Modern Applications**

# Non-Deadlock Bugs

- Make up a majority of concurrency bugs. **非死锁类bug占并发bug的大头！**

- Two major types of non deadlock bugs:

  - Atomicity violation 违反原子性

  - Order violation 违反顺序性

# Atomicity-Violation Bugs

- The desired **serializability** among multiple memory accesses is *violated*.
(违反原子性错误）违反了多个内存访问之间所要求的序列化性

  - Simple Example found in MySQL:

    - Two different threads access the field `proc_info` in the `struct thd`.

```
1    Thread1::
2    if(thd->proc_info){
3        …
4        fputs(thd->proc_info , …);
5        …
6    }
7
8    Thread2::
9    thd->proc_info = NULL;
```

# Atomicity-Violation Bugs (Cont.)

☐ **Solution**: Simply add locks around the shared-variable references.

```
1    pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
2
3    Thread1::
4    pthread_mutex_lock(&lock);
5    if(thd->proc_info){
6        …
7        fputs(thd->proc_info , …);
8        …
9    }
10   pthread_mutex_unlock(&lock);
11
12   Thread2::
13   pthread_mutex_lock(&lock);
14   thd->proc_info = NULL;
15   pthread_mutex_unlock(&lock);
```

# Order-Violation Bugs

- The desired order between two memory accesses is <u>flipped</u>. 两个内存访问之间所要求的顺序颠倒了

  - i.e., **A** should always be executed before **B**, but the order is not enforced during execution. 进程A应该始终在进程B之前执行，但是在执行过程中这个顺序并不是强制执行的

  - **Example**:

```
1    Thread1::
2    void init(){
3        mThread = PR_CreateThread(mMain, …);
4    }
5
6    Thread2::
7    void mMain(…){
8        mState = mThread->State
9    }
```

  - ▸ The code in Thread2 seems to assume that the variable `mThread` has already been *initialized* (and is not `NULL`).

# Order-Violation Bugs (Cont.)

☐ **Solution**: Enforce ordering using condition variables

```
1    pthread_mutex_t mtLock = PTHREAD_MUTEX_INITIALIZER;
2    pthread_cond_t mtCond = PTHREAD_COND_INITIALIZER;
3    int mtInit = 0;
4
5    Thread 1::
6    void init(){
7        …
8        mThread = PR_CreateThread(mMain,…);
9
10       // signal that the thread has been created.
11       pthread_mutex_lock(&mtLock);
12       mtInit = 1;
13       pthread_cond_signal(&mtCond);
14       pthread_mutex_unlock(&mtLock);
15       …
16   }
17
18   Thread2::
19   void mMain(…){
20       …
```

```
21          // wait for the thread to be initialized …
22          pthread_mutex_lock(&mtLock);
23          while(mtInit == 0)
24                  pthread_cond_wait(&mtCond, &mtLock);
25          pthread_mutex_unlock(&mtLock);
26
27          mState = mThread->State;
28          …
29  }
```

# Deadlock Bugs

```
Thread 1:          Thread 2:

lock(L1);          lock(L2);

lock(L2);          lock(L1);
```

- The presence of a cycle
  - `Thread1` is holding a lock `L1` and waiting for another one, `L2`.
  - `Thread2` that holds lock `L2` is waiting for `L1` to be release.

哈尔滨工业大学（深圳）
HARBIN INSTITUTE OF TECHNOLOGY, SHENZHEN

■ **看一个实际的例子**

■ **现在分析这个例子**

  ■ **竞争使用资源: 道路**

  ■ **A占有道路1，又要请求道路2，B占有…**

  ■ **形成了无限等待**

# 死锁概念(Deadlock)

■ **死锁: <span style="color:red">多个进程（线程）因循环等待资源而造成无法执行的现象。</span>**

占有 → 进程A → 等待

资源1 资源2

等待 → 进程B → 占有

■ **死锁会造成进程（线程）无法执行**

■ **<span style="color:red">死锁会造成系统资源的极大浪费(资源无法释放)</span>**

- **多个进程因等待资源才造成死锁**
- **资源: 进程在完成其任务过程所需要的所有对象**
    - CPU、内存、磁盘块、外设、文件、信号量 …
- **显然有些资源不会造成死锁，而有些会**
    - 只读文件是不会造成进程等待的，也就不会死锁
    - 打印机一次只能让一个进程使用，就会造成死锁

    **称为互斥访问资源**

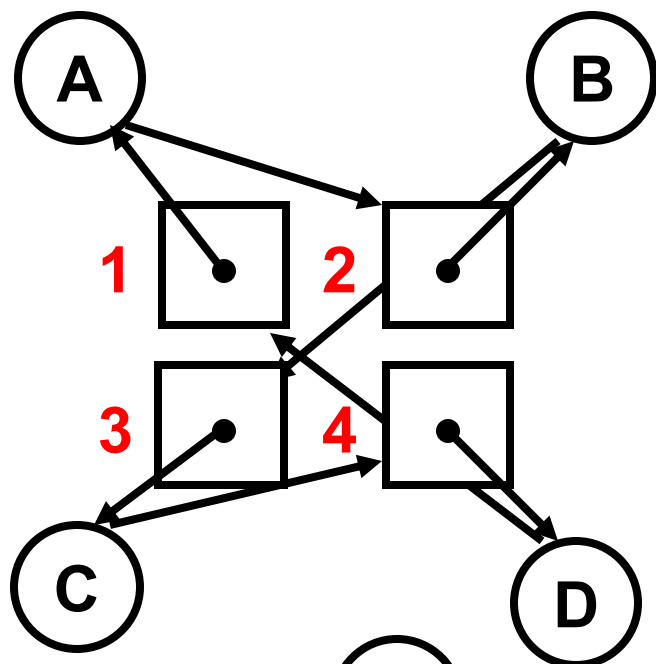    - **显然，资源互斥访问是死锁的必要条件**

- **资源请求需要形成环路等待才死锁！如何描述这种等待关系？**

# 资源分配图

- **资源分配图模型**

    - **一个进程集合$\{P_1, P_2, \ldots, P_n\}$**

    - **一资源类型集合$\{R_1, R_2, \ldots, R_m\}$**

    - **资源类型$R_i$有$W_i$个实例**

    - **资源请求边：有向边$P_i \rightarrow R_j$**

    - **资源分配边：有向边$R_i \rightarrow P_k$**



记号

**存在环路：**

1→A→2→B→3→C→4→D→1

**产生死锁**

**存在环路：**

P$_1$→R$_1$→P$_3$→R$_2$→P$_1$

**但并不死锁，仍可继续执行**

占有　等待

# 死锁的4个必要条件

- ## 互斥使用(Mutual exclusion)
  - ### 至少有一个资源互斥使用

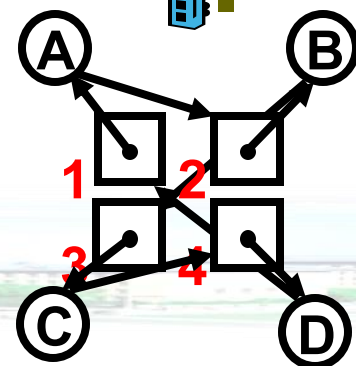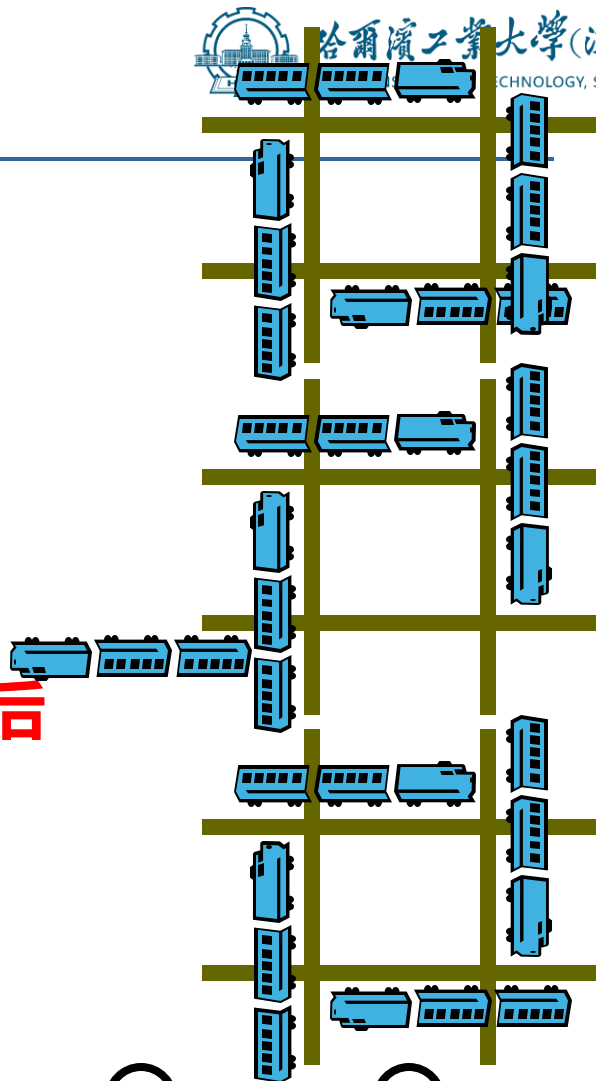- ## 不可抢占(No preemption)
  - ### 资源只能自愿放弃，如车开走以后

- ## 请求和保持(Hold and wait)
  - ### 进程必须占有资源，再去申请

- ## 循环等待(Circular wait)
  - ### 在资源分配图中存在一个环路

# Conditional for Deadlock

☐ <u>Four conditions</u> need to hold for a deadlock to occur.

| Condition | Description |
|---|---|
| **Mutual Exclusion** | **Threads claim exclusive control of resources that they require.** 进程要求对所需要的资源进行排它性控制 |
| **Hold-and-wait** | **Threads hold resources allocated to them while waiting for additional resources.** 进程在等待其它资源时保留已分配给它们的资源 |
| **No preemption** | **Resources cannot be forcibly removed from threads that are holding them.** 不能强制地移除进程所持有的资源 |
| **Circular wait** | **There exists a circular chain of threads such that each thread holds one more resources that are being requested by the next thread in the chain.** 存在一种进程资源的循环等待链，链中每一个进程已获得的资源同时被链中下一个进程所请求 |

☐ If any of these four conditions are not met, **deadlock cannot occur**.

【例题】某计算机系统中有 8 台打印机，由 K 个进程竞争使用，每个进程最多需要 3 台打印机。该系统可能发生死锁的 K 的最小值是（  )。

A．2

B．3

C．4

D．5

答案：C

考虑最极端情况，因为每个进程最多需要3台打印机，若每个进程已经占有了2台打印机，则只要还有多的打印机，总能满足一个进程达到3台的条件，然后顺利执行。所以将8台打印机分给 K 个进程，每个进程有2台打印机，这就是极端情况，K = 4。

# 死锁处理方法概述

- **死锁预防**　　"no smoking"，预防火灾
  - 破坏死锁的必要条件

- **死锁避免**　　检测到煤气超标时，自动切断电源
  - 检测每个资源请求，如果造成死锁就拒绝

- **死锁检测+恢复**　　发现火灾时，立刻拿起灭火器
  - 检测到死锁出现时，剥夺一些进程的资源

- **死锁忽略**　　在太阳上可以对火灾全然不顾
  - 就好像没有出现死锁一样

- **破坏互斥使用**
  - 资源的固有特性，通常**无法破除，如打印机**

- **破除不可抢占**
  - 如果一个进程占有资源并申请另一个不能立即分配的资源，那么**已分配资源就可被抢占（即持有不用即可抢占）**
  - 如果申请的资源得到满足，则抢占其他资源一次性分配给该进程
  - 只对状态能保存和恢复的资源(如CPU，内存空间)有效，对打印机等外设不适用
    - **实例：两个进程使用串口，都要读串口，数据不同不可恢复。**

■ **破除请求和保持**

   ■ 在进程执行前，**一次性申请所有需要的资源**

   ■ 缺点1: 需要预知未来，编程困难

   ■ 缺点2: 许多资源分配后很长时间后才使用，资源利用率低

- **破除循环等待**

  - 对资源类型进行排序，**资源申请必须按序进行**
  - **例如：** 所有的进程必须先申请磁盘驱动，再申请打印机，再….，如同日常交通中的单行道

  - 缺点: 如果编程时就需考虑，用户会觉得很别扭；可能需要释放某些资源(申请序号小的资源)，进程可能会无法执行

- **总之，破除死锁的必要条件会引入不合理因素，实际中很少使用。**

【例题】一次性分配所有资源的方法可以预防死锁的发生，它破坏了死锁4个必要条件中的(  )。

A．互斥使用（Mutual exclusion）

B．不可抢占（No preemption）

C．请求和保持（Hold-and-wait）

D．循环等待（Circular wait）

答案：C

解析：一次性分配所有资源的方法是当进程需要资源时，一次性提出所有的请求，若请求的所有资源满足则分配，只要有一项不满足，就不分配任何资源。这种分配方式不会部分地占有资源，因此破坏了"请求和保持"（进程在等待其它资源时保留已分配给它们的资源）。

不死锁就成了问题的核心！

- **思想**: 判断此次请求**是否造成死锁**

    若会造成死锁，则拒绝该请求

- **安全状态定义**：如果系统中的所有进程存在一个可完成的执行序列$P_1$，…$P_n$，则称系统处于安全状态

都能执行完成当然就不死锁

- **安全序列**：上面的执行序列$P_1$，…$P_n$

如何找到这样的序列？

# 死锁避免之银行家算法

一个银行家：目前手里只有1亿
第A个开发商：已贷款15亿，资金紧张还需3亿。
第B个开发商：已贷款5亿，还需贷款1亿，运转良好能收回。
第C个开发商：已贷款2亿，欲贷款18亿
… …

**开发商B还钱，再借给A，则可以继续借给C**

银行家当前可用的资金（**Available**）？可以利用的资金，即可用的加上能收回的共有多少（**work**）？各个开发商已贷款——已分配的资金（**Allocation**）？各个开发商还需要贷款（**need**）

# 死锁避免之银行家算法

## ■ 安全序列P$_1$，…P$_n$应该满足的性质：

Pi(1≤i≤n)需要资源 ≤ 剩余资源 + 分配给Pj(1≤j<i)资源

```
1．Banker()
2.  int  n,m;   //系统中进程总数n和资源种类总数m
3.  int  Available[m]; //资源当前可用总量
4.  int  Allocation[n][m];
5.         //当前给分配给每个进程的各种资源数量
6.  int  Need[n][m];
7.         //当前每个进程还需分配的各种资源数量
8.  int  Work[m]; //当前可分配的资源，包括可收回的
9.  bool  Finish[n]; //进程是否结束
```

## ■ 安全状态判定（思路）：

①初始化设定：

Work = Available （*动态记录当前可（收回）分配资源*）

Finish[i]=false （*设定所有进程均未完成*）

②查找这样的进程$P_i$（*未完成但目前剩余资源可满足其需要，这样的进程是能够完成的*）：

a）Finish[i] = =false        b）Need[i] ≤ Work

如果没有这样的进程$P_i$，则跳转到第④步

③（*若有则*）$P_i$一定能完成，并归还其占用的资源，即：

a）Finish[i] = true        b）Work = Work +Allocation[i]

GOTO 第②步，继续查找

④如果所有进程$P_i$都是能完成的，即Finish[i]=ture

则系统处于*安全状态*，否则系统处于*不安全状态*

哈尔滨工业大学(深圳)
HARBIN INSTITUTE OF TECHNOLOGY, SHENZHEN

■ **当前状态:**

|  | **Allocation** | **Need** | **Available** |
|---|---|---|---|
|  | $A\ B\ C$ | $A\ B\ C$ | $A\ B\ C$ |
| P0 | ~~0 1 0~~ | ~~7 4 3~~ | 3 3 2 |
| P1 | ~~2 0 0~~ | ~~1 2 2~~ | |
| P2 | ~~3 0 2~~ | ~~6 0 0~~ | |
| P3 | ~~2 1 1~~ | ~~0 1 1~~ | |
| P4 | ~~0 0 2~~ | ~~4 3 1~~ | |

Work=[3  3  2]

$P_1$ Work=[5  3  2]

$P_3$ Work=[7  4  3]

$P_4$ Work=[7  4  5]

$P_0$ Work=[7  5  5]

$P_2$ Work=[10  5  7]

■ **安全序列是<$P_1$, $P_3$, $P_4$, $P_0$ ,$P_2$>**

■ **安全序列是唯一的吗?**

```
1.bool Found;
2.Work = Available; Finish = false;
3.while(true){
4.    Found = false; //是否为安全序列找到一个新进程
5.    for(i=1; i<=n; i++){
6.      if(Finish[i]==false && Need[i]<=Work){
7.        Work = Work + Allocation[i];
8.        Finish[i] = true;
9.        printf("%d->",i);//输出安全序列
10.       Found = true;
11.     }
12.   } 没有安全序列或已经找到
13.   if(Found==false)break;
14.}
15.for(i=1;i<=n;i++)
16.  if(Finish[i]==false)
17.    return "deadlock";
```

$$T(n)=O(mn^2)$$

最好情形：安全状态就是$P_1$-$P_n$
最坏情形：$P_n$-$P_1$

# 死锁避免之资源请求算法

**思想：可用的资源可以满足某个进程的资源请求，则分配，然后寻找安全序列，找到，分配成功，找不到，已分配资源收回。**

```
1. extern Banker();
2. int Request[m]; /*进程Pi的资源申请*/
3. if(Request>Need[i]) return "error";
4. if(Request>Available) sleep();
5. Available=Available-Request;
6. Allocation[i]=Allocation[i]+Request;
7. Need[i]=Need[i]-Request;
8.                    /*先将资源分配给Pi*/
9. if(Banker()=="deadlock")
10.          /*调用银行家算法判定是否会死锁*/
11.    拒绝Request;/*若算法判定deadlock则拒绝请求, 资源回滚*/
```

**哈尔滨工业大学(深圳)**
HARBIN INSTITUTE OF TECHNOLOGY, SHENZHEN

■ $P_1$申请资源(1,0,2)

| | Allocation | | | Need | | | Available | | |
|---|---|---|---|---|---|---|---|---|---|
| | A | B | C | A | B | C | A | B | C |
| P0 | 0 | 1 | 0 | 7 | 4 | 3 | 2 | 3 | 0 |
| P1 | 3 | 0 | 2 | 0 | 2 | 0 | | | |
| P2 | 3 | 0 | 2 | 6 | 0 | 0 | | | |
| P3 | 2 | 1 | 1 | 0 | 1 | 1 | | | |
| P4 | 0 | 0 | 2 | 4 | 3 | 1 | | | |

| | Allocation | | | Need | | | Available | | |
|---|---|---|---|---|---|---|---|---|---|
| | A | B | C | A | B | C | A | B | C |
| P0 | 0 | 1 | 0 | 7 | 4 | 3 | 3 | 3 | 2 |
| P1 | 2 | 0 | 0 | 1 | 2 | 2 | | | |
| P2 | 3 | 0 | 2 | 6 | 0 | 0 | | | |
| P3 | 2 | 1 | 1 | 0 | 1 | 1 | | | |
| P4 | 0 | 0 | 2 | 4 | 3 | 1 | | | |

■ 序列<$P_1$, $P_3$, $P_2$, $P_4$, $P_0$>是安全的

■ 此次申请允许

# 死锁避免之资源请求实例(2)

- ## $P_0$再申请$(0,2,0)$

| | Allocation | Need | Available |
|---|---|---|---|
| | A B C | A B C | A B C |
| P0 | 0 3 0 | 7 2 3 | 2 1 0 |
| P1 | 3 0 2 | 0 2 0 | |
| P2 | 3 0 2 | 6 0 0 | |
| P3 | 2 1 1 | 0 1 1 | |
| P4 | 0 0 2 | 4 3 1 | |

| | Allocation | Need | Available |
|---|---|---|---|
| | A B C | A B C | A B C |
| P0 | 0 1 0 | 7 4 3 | 2 3 0 |
| P1 | 3 0 2 | 0 2 0 | |
| P2 | 3 0 2 | 6 0 0 | |
| P3 | 2 1 1 | 0 1 1 | |
| P4 | 0 0 2 | 4 3 1 | |

- 进程$P_0$, $P_1$, $P_2$, $P_3$, $P_4$一个也没法执行，死锁进程组

- 此次申请被拒绝

# 银行家算法讨论：

- 每个进程进入系统时必须告知所需资源的最大数量
  **对应用程序员要求高**

- 安全序列寻找算法（安全状态判定算法）计算**时间复杂度为$O(mn^2)$，过于复杂**

- 若每次资源请求都要调用银行家算法，耗时过大，
  **系统效率降低**

- 采用此算法，存在情况：当前有资源可用，尽管可能很快就会释放，由于会使整体进程处于不安全状态，而不被分配，致使**资源利用率大大降低**

- **基本原因: 每次申请都执行O(mn²),效率低**
- **对策: 只要可用资源足够,则分配, <span style="color:red">发现问题再处理</span>**
  - **定时检测或者当发现资源利用率低时检测**

```
1.bool Found;
2.int Request[n][m];
3.Work = Available; Finish = false;
4.if Allocation[i] != 0: Finish[i] = false;
5.else: Finish[i] = true;
6.while(true){
7.    Found = false; //是否为安全序列找到一个新进程
8.    for(i=1; i<=n; i++){
9.      if(Finish[i]==false && Request[i]<=Work){
10.          Work = Work + Allocation[i];
11.          Finish[i] = true;
12.          Found = true;
13.      }
14.    }
15.    if(Found==false)break;
16.}
17.for(i=1;i<=n;i++) {
18.  if(Finish[i]==false) {
19.    deadlock = deadlock + {i}; return "deadlock";
20.  }
21.}
```

**//对银行家算法进行改进**

**对于无分配资源的进程,不论其是否获得请求资源,则认为其是完成的**

# 死锁检测+恢复: 死锁恢复

- ■ **终止进程　选谁终止？**
  - ■ **优先级？占用资源多的？…**

- ■ **剥夺资源　进程需要回滚 (rollback)**
  - ■ **回滚点的选取？如何回滚？…**

# 鸵鸟算法（死锁忽略）

- **死锁预防？**
  - **引入太多不合理因素…**

- **死锁避免？**

  - **每次申请都执行银行家算法$O(mn^2)$，效率太低**

- **死锁检测+恢复？**

  - **还要执行银行家算法$O(mn^2)$，且恢复并不容易**

- **鸵鸟算法: 对死锁不做任何处理……**

  - **死锁出现时，手动干预——重新启动**
  - **死锁出现不是确定的，避免死锁付出的代价毫无意义**
  - **有趣的是大多数操作系统都用它，如UNIX和Windows**

- **进程竞争资源 ⇒ 有可能形成循环竞争 ⇒ 死锁**

- **死锁需要处理 ⇒ 死锁分析 ⇒ 死锁的必要条件**

- **死锁处理 ⇒ 预防、避免、检测+恢复、忽略**

- **死锁预防: 破除必要条件 ⇒ 引入了不合理因素**

- **死锁避免: 用银行家算法找安全序列 ⇒ 效率太低**

- **死锁检测恢复: 银行家算法找死锁进程组并恢复 ⇒ 实现较难**

- **死锁忽略: 就好像没有死锁 ⇒ 现在用的最多**

**任何思想、概念、技术的主流都会
随着时间而改变，操作系统尤为明显!**

【例题】在下列死锁的解决方法中，属于死锁预防策略的是(  )。

A．银行家算法

B．资源有序分配算法

C．死锁检测算法

D．资源分配图化简算法

答案：B

A．属于**死锁避免**

B．破坏了循环等待，**破坏死锁的必要条件**属于**死锁预防**

C．属于**死锁检测**

D．属于**死锁检测**

# 信号处理

【例题】异步信号安全的函数要么是可重入的，要么不能被信号处理程序中断，包括I/O函数（）

A.printf          B.sprintf          C.write          D.malloc

**答案:**

**C**