



哈爾濱工業大學

HARBIN INSTITUTE OF TECHNOLOGY



操作系统

Operating Systems

刘川意 教授

liuchuanyi@hit.edu.cn

哈尔滨工业大学（深圳）

2021年10月

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



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-threaded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**, a piece of code that accesses a *shared* resource, usually a variable or data structure
- **Race Condition** 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.
- **Indeterminate**, 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.
- **Mutual Exclusion** primitives, guarantee that only a single thread ever enters a critical section, thus avoiding races, and resulting in deterministic program outputs

Review: Thread

- 轻量化执行环境, new abstraction for a single running process
- 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**.



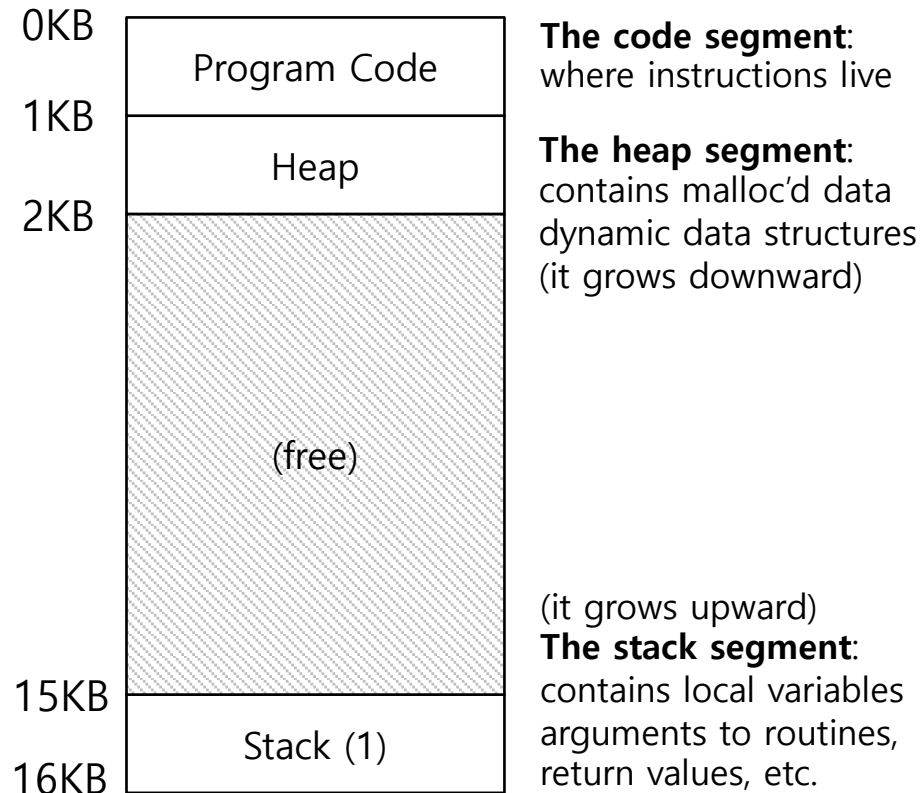
Context switch between threads

- Each thread has its own program counter and set of registers.
 - One or more **thread control blocks(TCBs)** are needed to store the state of each thread.

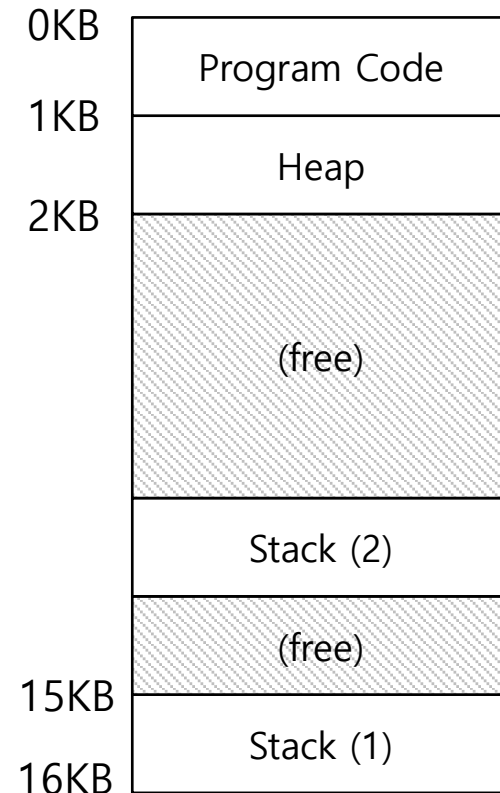
- When switching from running one (T1) to running the other (T2),
 - The register state of T1 be saved.
 - The register state of T2 restored.
 - The **address space remains** the same.

The stack of the relevant thread

- There will be **one stack per thread**.



**A Single-Threaded
Address Space**



**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",
31.               cnt);
32.     else
33.         printf("OK cnt=%ld\n", cnt);
34.     exit(0);
35. }
```

badcnt.c

```
38. /* Thread routine */
39. void *thread(void *vargp)
40. {
41.     long i, niters =
42.         *((long *)vargp);
43.
44.     for (i = 0; i < niters; i++)
45.         cnt++;
46.
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

Assembly Code for Counter Loop

■ 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 | j1 | 10a <thread+0x1d> | | |
| 12b: | b8 00 00 00 00 | mov | \$0x0,%eax | | |
| 130: | 5d | pop | %rbp | | |
| 131: | c3 | retq | | | |

H_i

L_i

U_i

S_i

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++;
```

00000000000000957 <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 | | |

怎么计算?

H_i

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,

| OS | Thread1 | Thread2 | (after instruction) | | |
|-----------|-------------------------|---------------------|---------------------|------|---------|
| | | | PC | %eax | counter |
| | before critical section | | 100 | 0 | 50 |
| | mov 0x8049a1c, %eax | | 105 | 50 | 50 |
| | add \$0x1, %eax | | 108 | 51 | 50 |
| interrupt | save T1's state | | | | |
| | 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 | save T2's state | | | | |
| | 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.
 - Multiple threads executing critical section can result in a race condition.
 - Need to support **atomicity** for critical sections (**mutual exclusion**)



Locks

- Ensure that any such critical section executes as if it were a single atomic instruction (**execute a series of instructions atomically**).

```
1  lock_t mutex;  
2  . . .  
3  lock(&mutex);  
4  balance = balance + 1;  
5  unlock(&mutex);
```

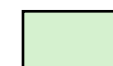
→ Critical section

Concurrent Execution(并发执行)

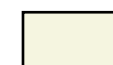
■ *Key idea:* In general, any sequentially consistent interleaving is possible, but some give an unexpected result!

- I_i denotes that thread i executes instruction I
- $\%rdx_i$ is the content of $\%rdx$ in thread i 's context

| 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

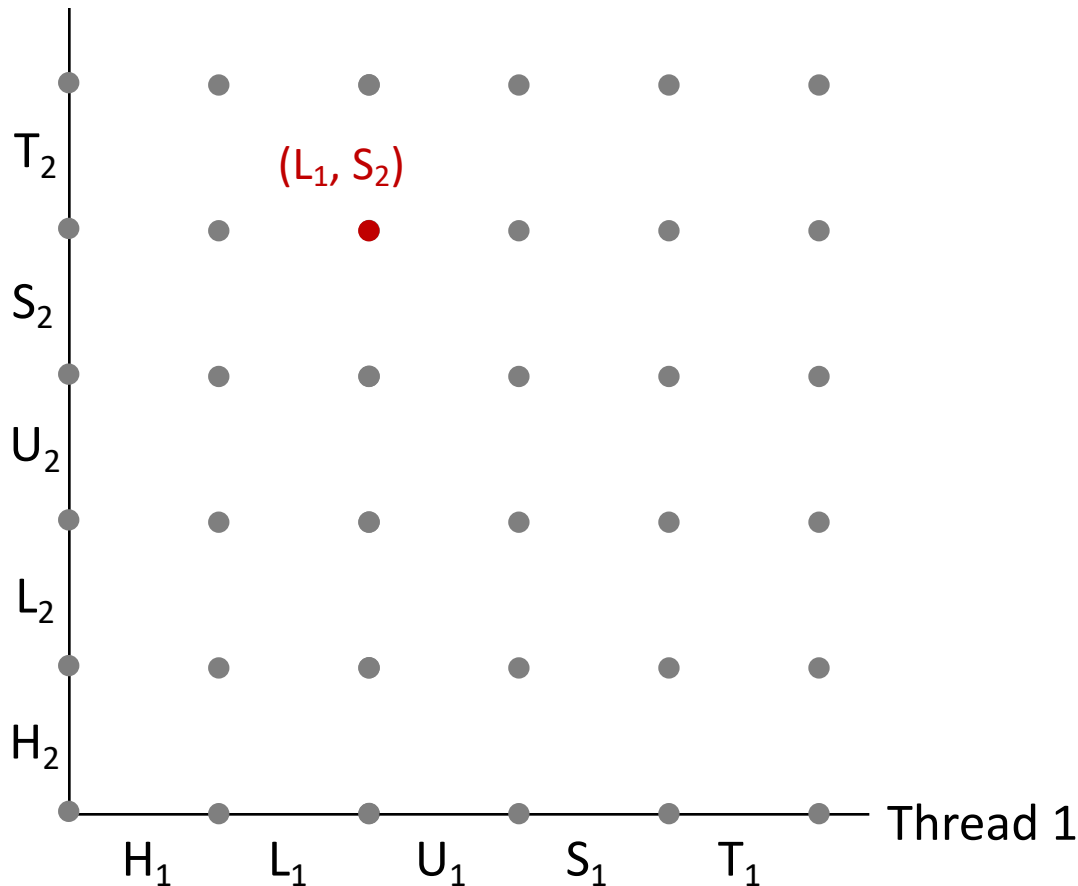
| i (thread) | instr _i | %rdx ₁ | %rdx ₂ | cnt |
|------------|--------------------|-------------------|-------------------|-----|
| 1 | H ₁ | - | - | 0 |
| 1 | L ₁ | 0 | - | 0 |
| 1 | U ₁ | 1 | - | 0 |
| 2 | H ₂ | - | - | 0 |
| 2 | L ₂ | - | 0 | 0 |
| 1 | S ₁ | 1 | - | 1 |
| 1 | T ₁ | 1 | - | 1 |
|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

Thread 2



A *progress graph* depicts the discrete *execution state space* of concurrent threads.

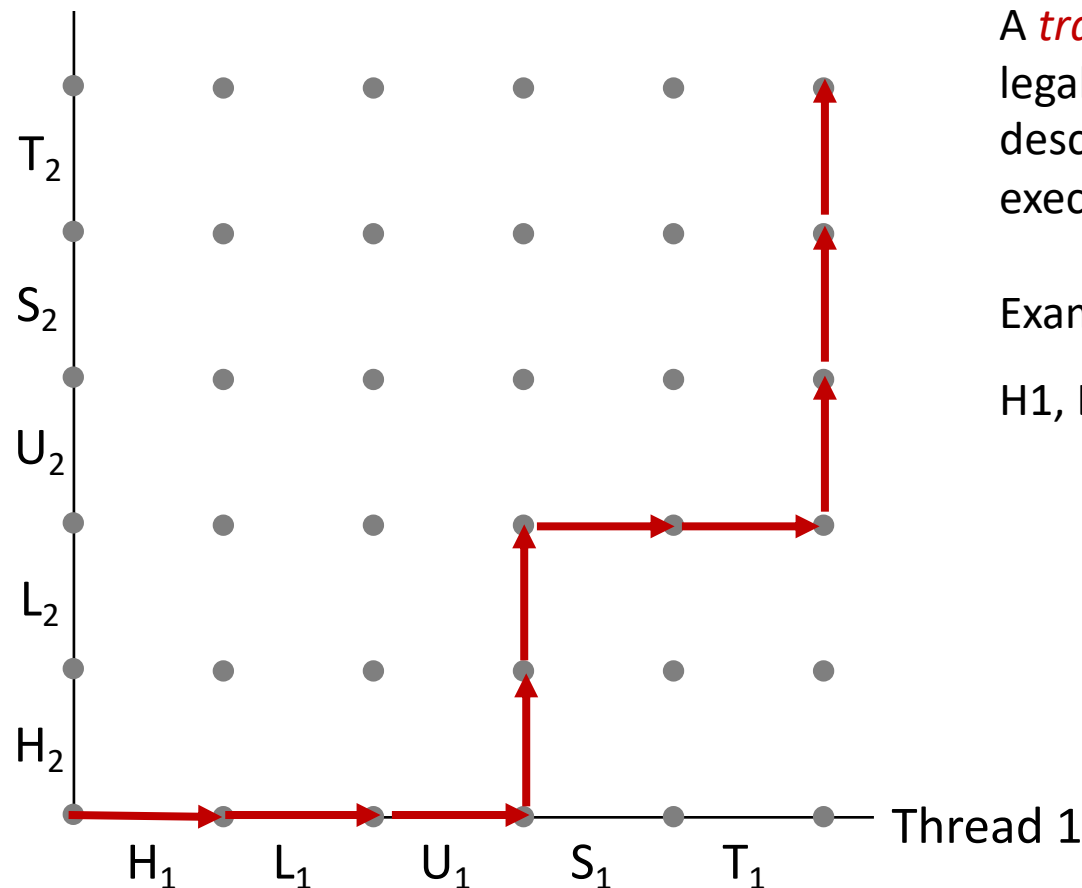
Each axis corresponds to the sequential order of instructions in a thread.

Each point corresponds to a possible *execution state* ($Inst_1, Inst_2$).

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

Trajectories in Progress Graphs

Thread 2

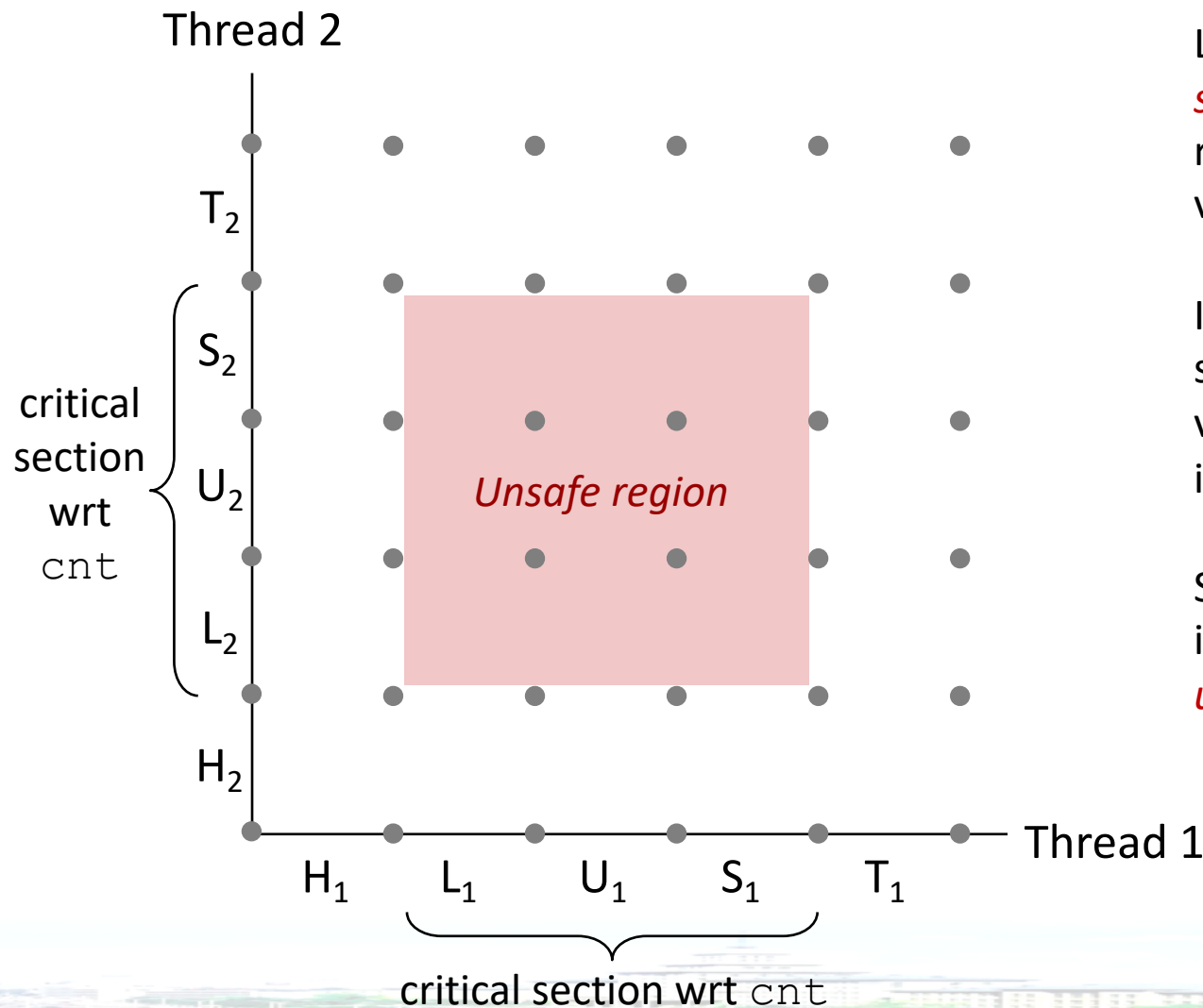


A *trajectory* (軌道) is a sequence of legal state transitions (轉換) that describes one possible concurrent execution of the threads.

Example:

$H_1, L_1, U_1, H_2, L_2, S_1, T_1, U_2, S_2, T_2$

Critical Sections and Unsafe Regions

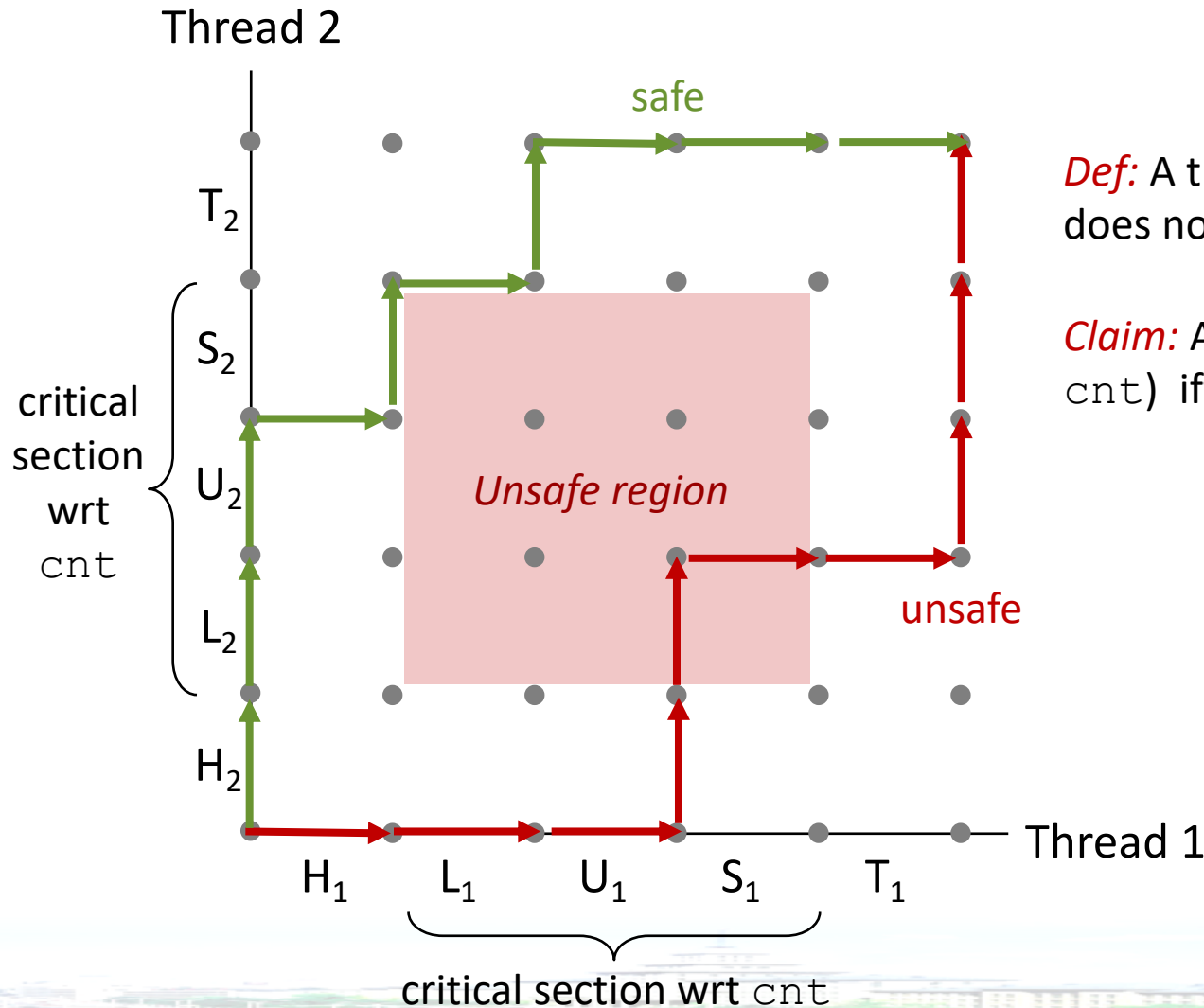


L , U , and S form a **critical section** (临界区) with respect to the shared variable `cnt`

Instructions in critical sections (wrt some shared variable) should not be interleaved

Sets of states where such interleaving occurs form **unsafe regions**

Critical Sections and Unsafe Regions



Def: A trajectory (轨道) is *safe* iff it does not enter any unsafe region

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

Enforcing Mutual Exclusion

- *Question:* How can we guarantee a safe trajectory?
- Answer: We must **synchronize** (同步) the execution of the threads so that they can never have an unsafe trajectory.
 - i.e., need to guarantee **mutually exclusive access** (互斥地访问) for each critical section.
- Classic solution:
 - Semaphores (信号量) (Edsger Dijkstra)
- Other approaches (out of our scope)
 - Mutex and condition variables (Pthreads)
 - Monitors (Java)

【例题1】子进程运行结束会向父进程发送_____信号。

【例题2】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：SIGCHLD/17号 C

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

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 the state of 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.

lock()原语的语义(semantics)

- lock()
 - **Try to** acquire the lock.
 - If no other thread holds the lock, the thread will **acquire** the lock.
 - **Enter** the *critical section*.
 - ▶ This thread is said to be the owner of the lock.
 - Other threads are *prevented from* entering the critical section while the first thread that holds the lock is in there.

Pthread Locks - mutex

- The name that the POSIX library uses for a lock.
 - Used to provide **mutual exclusion** 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).



Lock如何实现？

- Efficient locks provided mutual exclusion at **low cost**.
- Building a lock need some help from the **hardware** and the **OS**.



如何评价lock原语？

■ Mutual exclusion 正确性

- Does the lock work, preventing multiple threads from entering *a critical section*?

■ Fairness 公平性

- Does each thread contending for the lock get a fair shot at acquiring it once it is free? (**Starvation**)

■ Performance 性能

- The time overheads added by using the lock

Controlling Interrupts 基于中断控制的锁实现

■ Disable Interrupts for critical sections

- 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.
- ▶ Do not work on **multiprocessors** 多处理器体系结构这种方式不work
- ▶ Code that masks or unmask interrupts be executed *slowly* by modern CPUs

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:** Spin-waiting 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`.
- *Simultaneously* **update**(setting) said value to `new`.
- This sequence of operations is **performed atomically**.

A Simple Spin Lock 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 a preemptive scheduler.
- 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:** yes

- The spin lock only allows a single thread to entry the critical section.

■ **Fairness:** no

- Spin locks don't provide any fairness guarantees.
- Indeed, a thread spinning may spin *forever*.

■ **Performance:**

- In the single CPU, performance overheads can be quire *painful*.
- If the number of threads roughly equals the number of CPUs, spin locks work *reasonably well*.



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.
 - ▶ **success**: return 1 and update the value at `ptr` to `value`.
 - ▶ **fail**: the value at `ptr` is not updates and 0 is returned.

Load-Linked and Store-Conditional (Cont.)

```
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 fetch-and add.
 - Ensure progress for all threads. → fairness

```
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.

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

如何解决“自旋空转”，办法1: Just Yield

- When you are going to spin, **give up the CPU** to another thread.
 - OS system call moves the caller from the *running state* to the *ready state*.
 - The cost of a **context switch** can be substantial and the **starvation** problem still exists.

```
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 waiting to enter the lock.
- `park()`
 - Put a calling thread to sleep
- `unpark(threadID)`
 - Wake a particular thread as designated by `threadID`.



Using Queues: Sleeping Instead of Spinning

```
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



Using Queues: Sleeping Instead of Spinning

```
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.)

Futex (Cont.)

```
16         if (v >= 0)
17             continue;
18         futex_wait(mutex, v);
19     }
20 }
21
22 void mutex_unlock(int *mutex) {
23     /* Adding 0x80000000 to the counter results in 0 if and only if
24        there are not other interested threads */
25     if (atomic_add_zero(mutex, 0x80000000))
26         return;
27     /* There are other threads waiting for this mutex,
28        wake one of them up */
29     futex_wake(mutex);
30 }
```

Linux-based Futex Locks (Cont.)

Two-Phase Locks

- A two-phase lock realizes that **spinning can be useful** if the lock *is about to* be released.
 - **First phase**
 - ▶ The lock spins for a while, *hoping that* it can acquire the lock.
 - ▶ If the lock is not acquired during the first spin phase, a second phase is entered,
 - **Second phase**
 - ▶ The caller is put to sleep.
 - ▶ The caller is only woken up when the lock becomes free later.

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

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.





Example: Concurrent Counter without Lock

■ 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

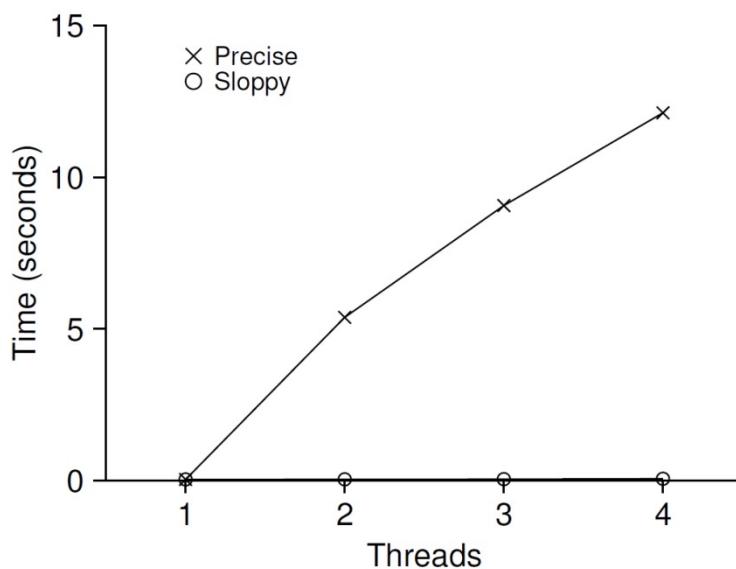


- acquired when calling a routine manipulating the data structure.

```
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    }
16
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.

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, on per CPU core
 - A single **global counter**
 - There are **locks**:
 - ▶ 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*.
 - ▶ Thus counter updates are **scalable**.
 - The local values are periodically transferred to the global counter.
 - ▶ Acquire the global lock
 - ▶ Increment it by the local counter's value
 - ▶ The local counter is then reset to zero.

The basic idea of sloppy counting (Cont.)

- How often the local-to-global transfer occurs is determined by a threshold, S (sloppiness).
 - The smaller S :
 - ▶ The more the counter behaves like the *non-scalable counter*.
 - The bigger S :
 - ▶ The more scalable the counter.
 - ▶ The further off the global value might be from the *actual count*.

Sloppy counter example

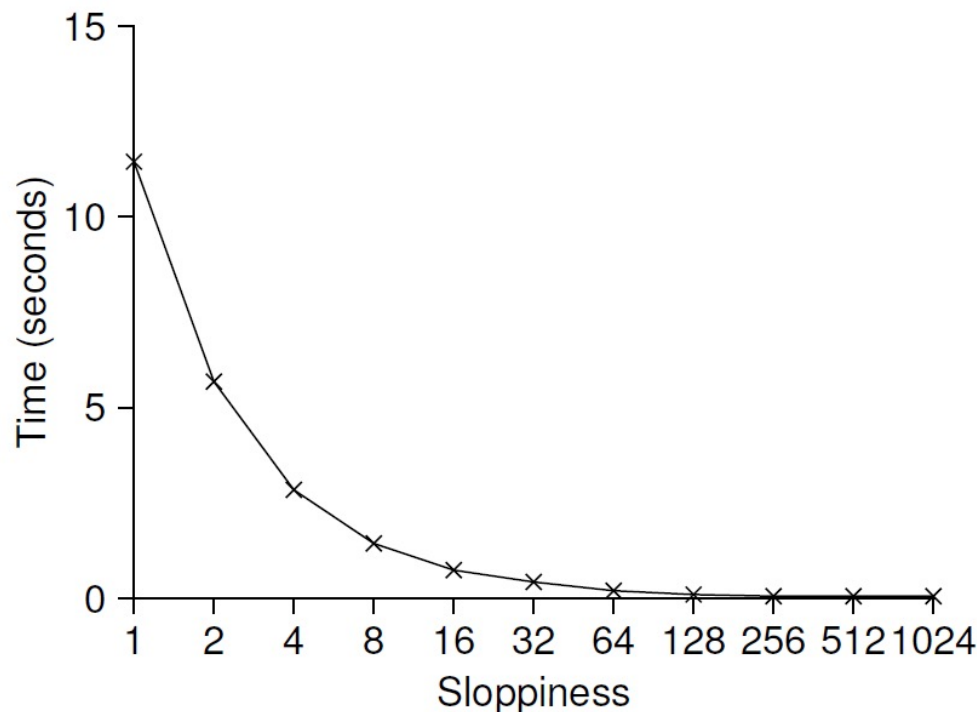
■ Tracing the Sloppy Counters

- 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 | $5 \rightarrow 0$ | 1 | 3 | 4 | 5 (from L_1) |
| 7 | 0 | 2 | 4 | $5 \rightarrow 0$ | 10 (from L_4) |

Importance of the threshold value S

- Each four threads increments a counter 1 million times on four CPUs.
 - Low $S \rightarrow$ Performance is **poor**, The global count is always quite **accurate**.
 - High $S \rightarrow$ Performance is **excellent**, The global count **lags**.



Scaling Sloppy Counters

Sloppy Counter Implementation

```
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->threshold = 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
```

Sloppy Counter Implementation (Cont.)

```
(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

```
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.)

- The code **acquires** a lock in the insert routine upon entry.
- The code **releases** the lock upon exit.
 - If `malloc()` happens to *fail*, the code must also release the lock before failing the insert.
 - This kind of exceptional control flow has been shown to be **quite error prone**.
 - **Solution:** The lock and release *only surround* the actual critical section in the insert code

Concurrent Linked List: Rewritten

```
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
```



Concurrent Linked List: Rewritten (Cont.)

(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)
 - Add **a lock per node** of the list instead of having a single lock for the entire list.
 - 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.
 - ▶ However, in practice, the overheads of acquiring and releasing locks for each node of a list traversal is *prohibitive*.

Michael and Scott Concurrent Queues

- 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.

- Add a dummy node
 - Allocated in the queue initialization code
 - Enable the separation of head and tail operations

Concurrent Queues (Cont.)

```
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.)
```

Concurrent Queues (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     }
```

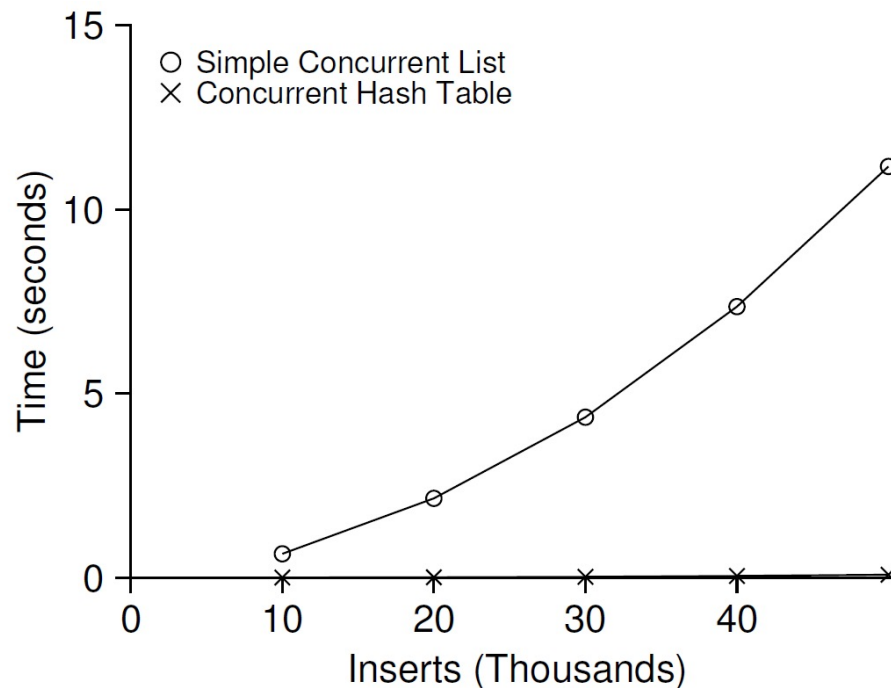
Concurrent Hash Table

- 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*.



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.**

Concurrent Hash Table



```
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     }
```

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

Condition Variables 条件变量的引入

- There are many cases where a thread wishes to check 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
```

Parent waiting fore child: Spin-based Approach

```
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 inefficient 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.
 - **Signaling** on the condition
 - ▶ Some other thread, *when it changes said state*, can wake one of those waiting threads and allow them to continue.

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 mutex 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.

Parent waiting for Child: Use a condition variable

```
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
```


Parent waiting for Child: Use a condition variable

```
(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 waiting for Child: Use a condition variable

■ 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 no thread asleep 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 interrupted 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.

The Producer/Consumer (Bound Buffer) Problem

■ Producer

- Produce data items
- Wish to place data items in a buffer

■ Consumer

- Grab data items out of the buffer consume them in some way

■ Example: Multi-threaded web server

- A *producer* puts HTTP requests in to a work queue
- *Consumer threads* take requests out of this queue and process them



Bounded buffer

- A bounded buffer is used when you pipe the output of one program into another.
 - Example: `grep foo file.txt | wc -l`
 - ▶ The `grep` process is the producer.
 - ▶ The `wc` process is the consumer.
 - ▶ Between them is an in-kernel bounded buffer.
 - Bounded buffer is Shared resource → **Synchronized access** is required.

The Put and Get Routines (Version 1)

```
1      int buffer;  
2      int count = 0;           // initially, empty  
3  
4      void put(int value) {  
5          assert(count == 0);  
6          count = 1;  
7          buffer = value;  
8      }  
9  
10     int get() {  
11         assert(count == 1);  
12         count = 0;  
13         return buffer;  
14     }
```

- Only put data into the buffer when `count` is zero.
 - ▶ i.e., when the buffer is *empty*.
- Only get data from the buffer when `count` is one.
 - ▶ i.e., when the buffer is *full*.



Producer/Consumer Threads (Version 1)

```
1      void *producer(void *arg) {
2          int i;
3          int loops = (int) arg;
4          for (i = 0; i < loops; i++) {
5              put(i);
6          }
7      }
8
9      void *consumer(void *arg) {
10         int i;
11         while (1) {
12             int tmp = get();
13             printf("%d\n", tmp);
14         }
15     }
```

- **Producer** puts an integer into the shared buffer loops number of times.
- **Consumer** gets the data out of that shared buffer.

Producer/Consumer: Single CV and If Statement



```
1      cond_t cond;
2      mutex_t mutex;
3
4      void *producer(void *arg) {
5          int i;
6          for (i = 0; i < loops; i++) {
7              Pthread_mutex_lock(&mutex);           // p1
8              if (count == 1)                       // p2
9                  Pthread_cond_wait(&cond, &mutex); // p3
10             put(i);                                // p4
11             Pthread_cond_signal(&cond);            // p5
12             Pthread_mutex_unlock(&mutex);          // p6
13         }
14     }
15
16     void *consumer(void *arg) {
17         int i;
18         for (i = 0; i < loops; i++) {
19             Pthread_mutex_lock(&mutex);             // c1
20             if (count == 0)                         // c2
21                 Pthread_cond_wait(&cond, &mutex);  // c3
22             int tmp = get();                         // c4
23             Pthread_cond_signal(&cond);             // c5
24             Pthread_mutex_unlock(&mutex);           // c6
25             printf("%d\n", tmp);
26         }
27     }
```

If we have more than one of producer and consumer?

Thread Trace: Broken Solution (Version 1)

| T_{c1} | State | T_{c2} | State | T_p | State | Count | Comment |
|----------|---------|----------|---------|-------|---------|-------|------------------------|
| c1 | Running | | Ready | | Ready | 0 | |
| c2 | Running | | Ready | | Ready | 0 | |
| c3 | Sleep | | Ready | | Ready | 0 | Nothing to get |
| | Sleep | | Ready | p1 | Running | 0 | |
| | Sleep | | Ready | p2 | Running | 0 | |
| | Sleep | | Ready | p4 | Running | 1 | Buffer now full |
| | Ready | | Ready | p5 | Running | 1 | T_{c1} awoken |
| | Ready | | Ready | p6 | Running | 1 | |
| | Ready | | Ready | p1 | Running | 1 | |
| | Ready | | Ready | p2 | Running | 1 | |
| | Ready | | Ready | p3 | Sleep | 1 | Buffer full; sleep |
| | Ready | c1 | Running | | Sleep | 1 | T_{c2} sneaks in ... |
| | Ready | c2 | Running | | Sleep | 1 | |
| | Ready | c4 | Running | | Sleep | 0 | ... and grabs data |
| | Ready | c5 | Running | | Ready | 0 | T_p awoken |
| | Ready | c6 | Running | | Ready | 0 | |
| c4 | Running | | Ready | | Ready | 0 | Oh oh! No data |

Thread Trace: Broken Solution (Version 1)

- The problem arises for a simple reason:
 - After the producer woke T_{c1} , but before T_{c1} ever ran, the state of the bounded buffer *changed by* T_{c2} .
 - There is no guarantee that when the woken thread runs, the state will still be as desired → Mesa semantics.
 - ▶ Virtually every system ever built employs *Mesa semantics*.
 - Hoare semantics provides a stronger guarantee that the woken thread will run immediately upon being woken.



Producer/Consumer: Single CV and While

```
1      cond_t cond;
2      mutex_t mutex;
3
4      void *producer(void *arg) {
5          int i;
6          for (i = 0; i < loops; i++) {
7              Pthread_mutex_lock(&mutex);           // p1
8              while (count == 1)                   // p2
9                  Pthread_cond_wait(&cond, &mutex); // p3
10             put(i);                               // p4
11             Pthread_cond_signal(&cond);           // p5
12             Pthread_mutex_unlock(&mutex);         // p6
13         }
14     }
15
16     void *consumer(void *arg) {
17         int i;
18         for (i = 0; i < loops; i++) {
19             Pthread_mutex_lock(&mutex);           // c1
20             while (count == 0)                   // c2
21                 Pthread_cond_wait(&cond, &mutex); // c3
22             int tmp = get();                      // c4
23             Pthread_cond_signal(&cond);           // c5
24             Pthread_mutex_unlock(&mutex);         // c6
25             printf("%d\n", tmp);
26         }
27     }
```



Thread Trace: Broken Solution (Version 2)

| T_{c1} | State | T_{c2} | State | T_p | State | Count | Comment |
|----------|---------|----------|---------|-------|---------|-------|---------------------------------------|
| c1 | Running | | Ready | | Ready | 0 | |
| c2 | Running | | Ready | | Ready | 0 | |
| c3 | Sleep | | Ready | | Ready | 0 | Nothing to get |
| | Sleep | c1 | Running | | Ready | 0 | |
| | Sleep | c2 | Running | | Ready | 0 | |
| | Sleep | c3 | Sleep | | Ready | 0 | Nothing to get |
| | Sleep | | Sleep | p1 | Running | 0 | |
| | Sleep | | Sleep | p2 | Running | 0 | |
| | Sleep | | Sleep | p4 | Running | 1 | Buffer now full |
| | Ready | | Sleep | p5 | Running | 1 | T_{c1} awoken |
| | Ready | | Sleep | p6 | Running | 1 | |
| | Ready | | Sleep | p1 | Running | 1 | |
| | Ready | | Sleep | p2 | Running | 1 | |
| | Ready | | Sleep | p3 | Sleep | 1 | Must sleep (full) |
| c2 | Running | | Sleep | | Sleep | 1 | Recheck condition |
| c4 | Running | | Sleep | | Sleep | 0 | T_{c1} grabs data |
| c5 | Running | | Ready | | Sleep | 0 | Oops! Woke T_{c2} |

Thread Trace: Broken Solution (Version 2) (Cont.)

| T_{c1} | State | T_{c2} | State | T_p | State | Count | Comment |
|----------|---------|----------|---------|-------|-------|-------|---------------------|
| ... | ... | ... | ... | ... | ... | ... | (cont.) |
| c6 | Running | | Ready | | Sleep | 0 | |
| c1 | Running | | Ready | | Sleep | 0 | |
| c2 | Running | | Ready | | Sleep | 0 | |
| c3 | Sleep | | Ready | | Sleep | 0 | Nothing to get |
| | Sleep | c2 | Running | | Sleep | 0 | |
| | Sleep | c3 | Sleep | | Sleep | 0 | Everyone asleep ... |

- ◆ A consumer should not wake other consumers, only producers, and vice-versa.

The single Buffer Producer/Consumer Solution



■ Use **two** condition variables and while

- **Producer** threads wait on the condition `empty`, and signals `fill`.
- **Consumer** threads wait on `fill` and signal `empty`.

```
1      cond_t empty, fill;
2      mutex_t mutex;
3
4      void *producer(void *arg) {
5          int i;
6          for (i = 0; i < loops; i++) {
7              Pthread_mutex_lock(&mutex);
8              while (count == 1)
9                  Pthread_cond_wait(&empty, &mutex);
10             put(i);
11             Pthread_cond_signal(&fill);
12             Pthread_mutex_unlock(&mutex);
13         }
14     }
15
16     void *consumer(void *arg) {
17         int i;
18         for (i = 0; i < loops; i++) {
19             Pthread_mutex_lock(&mutex);
20             while (count == 0)
21                 Pthread_cond_wait(&fill, &mutex);
22             int tmp = get();
23             Pthread_cond_signal(&empty);
24             Pthread_mutex_unlock(&mutex);
25             printf("%d\n", tmp);
26         }
27     }
```

The Final Producer/Consumer Solution

- More **concurrency** and **efficiency** → Add more buffer slots.
 - Allow concurrent production or consuming to take place.
 - Reduce context switches.

```
1      int buffer[MAX];
2      int fill = 0;
3      int use = 0;
4      int count = 0;
5
6      void put(int value) {
7          buffer[fill] = value;
8          fill = (fill + 1) % MAX;
9          count++;
10     }
11
12     int get() {
13         int tmp = buffer[use];
14         use = (use + 1) % MAX;
15         count--;
16         return tmp;
17     }
```

The Final Put and Get Routines

The Final Producer/Consumer Solution (Cont.)

```
1      cond_t empty, fill;
2      mutex_t mutex;
3
4      void *producer(void *arg) {
5          int i;
6          for (i = 0; i < loops; i++) {
7              Pthread_mutex_lock(&mutex);           // p1
8              while (count == MAX)                 // p2
9                  Pthread_cond_wait(&empty, &mutex); // p3
10             put(i);                               // p4
11             Pthread_cond_signal(&fill);           // p5
12             Pthread_mutex_unlock(&mutex);         // p6
13         }
14     }
15
16     void *consumer(void *arg) {
17         int i;
18         for (i = 0; i < loops; i++) {
19             Pthread_mutex_lock(&mutex);           // c1
20             while (count == 0)                     // c2
21                 Pthread_cond_wait(&fill, &mutex); // c3
22             int tmp = get();                       // c4
23             Pthread_cond_signal(&empty);          // c5
24             Pthread_mutex_unlock(&mutex);         // c6
25             printf("%d\n", tmp);
26         }
27     }
```

Covering Conditions

- Assume there are zero bytes free
 - Thread T_a calls `allocate(100)`.
 - Thread T_b calls `allocate(10)`.
 - Both T_a and T_b wait on the condition and go to sleep.
 - Thread T_c calls `free(50)`.

Which waiting thread should be woken up?

Covering Conditions (Cont.)

```
1      // how many bytes of the heap are free?
2      int bytesLeft = MAX_HEAP_SIZE;
3
4      // need lock and condition too
5      cond_t c;
6      mutex_t m;
7
8      void *
9      allocate(int size) {
10         Pthread_mutex_lock(&m);
11         while (bytesLeft < size)
12             Pthread_cond_wait(&c, &m);
13         void *ptr = ...;           // get mem from heap
14         bytesLeft -= size;
15         Pthread_mutex_unlock(&m);
16         return ptr;
17     }
18
19     void free(void *ptr, int size) {
20         Pthread_mutex_lock(&m);
21         bytesLeft += size;
22         Pthread_cond_signal(&c);   // whom to signal??
23         Pthread_mutex_unlock(&m);
24     }
```

Covering Conditions (Cont.)

- Solution (Suggested by Lampson and Redell)
 - Replace `pthread_cond_signal()` with `pthread_cond_broadcast()`
 - `pthread_cond_broadcast()`
 - ▶ Wake up **all waiting threads**.
 - ▶ Cost: too many threads might be woken.
 - ▶ Threads that shouldn't be awake will simply wake up, re-check the condition, and then go back to sleep.

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

Semaphore: A definition

■ An object with an integer value

- We can manipulate 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 shared between *threads in the same process*.

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 suspend execution waiting for a subsequent post.
- When negative, the value of the semaphore is equal to the number of waiting threads.

Semaphore: Interact with semaphore (Cont.)

■ 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  
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 | <i>Interrupt; Switch → T1</i> | Ready | | Running |
| 0 | | Ready | call sem_wait() | Running |
| -1 | | Ready | decrement sem | Running |
| -1 | | Ready | (sem < 0) → sleep | sleeping |
| -1 | | Running | <i>Switch → T0</i> | 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 | <i>Interrupt; Switch → T1</i> | 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 |

Semaphores As Condition Variables

```
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**.

Thread Trace: Parent Waiting For Child (Case 1)

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

| Value | Parent | State | Child | State |
|-------|--|----------|--|---------|
| 0 | Create(Child) | Running | <i>(Child exists; is runnable)</i> | Ready |
| 0 | call <code>sem_wait()</code> | Running | | Ready |
| -1 | decrement sem | Running | | Ready |
| -1 | $(sem < 0) \rightarrow \text{sleep}$ | sleeping | | Ready |
| -1 | <i>Switch</i> \rightarrow <i>Child</i> | sleeping | child runs | Running |
| -1 | | sleeping | call <code>sem_post()</code> | Running |
| 0 | | sleeping | increment sem | Running |
| 0 | | Ready | wake(Parent) | Running |
| 0 | | Ready | <code>sem_post()</code> returns | Running |
| 0 | | Ready | <i>Interrupt; Switch</i> \rightarrow <i>Parent</i> | Ready |
| 0 | <code>sem_wait()</code> 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 | <i>Interrupt; switch→Child</i> | Ready | child runs | Running |
| 0 | | Ready | call <code>sem_post()</code> | Running |
| 1 | | Ready | increment sem | Running |
| 1 | | Ready | wake (nobody) | Running |
| 1 | | Ready | <code>sem_post()</code> returns | Running |
| 1 | parent runs | Running | <i>Interrupt; Switch→Parent</i> | Ready |
| 1 | call <code>sem_wait()</code> | Running | | Ready |
| 0 | decrement sem | Running | | Ready |
| 0 | (sem<0)→awake | Running | | Ready |
| 0 | <code>sem_wait()</code> 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.

■ **Consumer:** `get()` interface

- Wait for a buffer to become *filled* before using it.

```
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 }
```



The Producer/Consumer (Bounded-Buffer) Problem

```
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

The Producer/Consumer (Bounded-Buffer) Problem

```
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`.
 - ▶ It means that the old data there is overwritten.
- We've forgotten here is **mutual exclusion**.
 - ▶ The filling of a buffer and incrementing of the index into the buffer is a **critical section**.

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 }
15
(Cont.)
```

Adding Mutual Exclusion (Incorrectly)

A Solution: Adding Mutual Exclusion

```
(Cont.)
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)

A Solution: Adding Mutual Exclusion (Cont.)



- 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 still holds the mutex!
 - 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
(Cont.)
```

Adding Mutual Exclusion (Correctly)

Finally, A Working Solution

```
(Cont.)
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  }
27
28  int main(int argc, char *argv[]) {
29      // ...
30      sem_init(&empty, 0, MAX); // MAX buffers are empty to begin with ...
31      sem_init(&full, 0, 0);    // ... and 0 are full
32      sem_init(&mutex, 0, 1);   // mutex=1 because it is a lock
33      // ...
34  }
```

Adding Mutual Exclusion (Correctly)

Reader-Writer Locks

- Imagine a number of concurrent list operations, including **inserts** and simple **lookups**.
 - **insert:**
 - ▶ Change the state of the list
 - ▶ A traditional critical section makes sense.
 - **lookup:**
 - ▶ Simply *read* the data structure.
 - ▶ As long as we can guarantee that no insert is on-going, we can allow many lookups to proceed **concurrently**.

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

A Reader-Writer Locks

- Only **a single writer** can acquire the lock.
- Once a reader has acquired **a read lock**,
 - **More readers** will be allowed to acquire the read lock too.
 - A writer will have to wait until all readers are finished.

```
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     ...
```


A Reader-Writer Locks (Cont.)

```
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 }
```

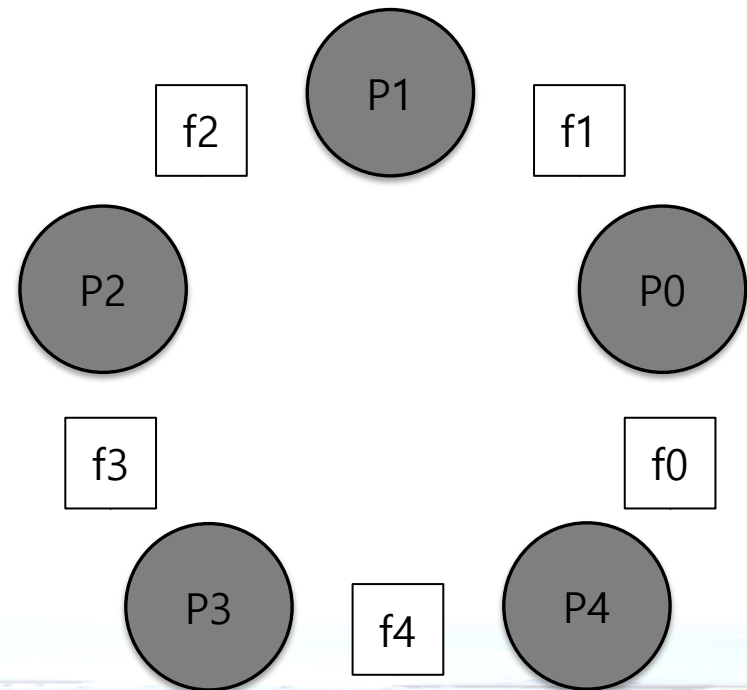
A Reader-Writer Locks (Cont.)

- The reader-writer locks have **fairness problem**.
 - It would be relatively easy for reader to **starve writer**.
 - How to prevent more readers from entering the lock once a writer is waiting?



The Dining Philosophers

- Assume there are five “**philosophers**” sitting around a table.
 - Between each pair of philosophers is a single fork (five total).
 - The philosophers each have times where they **think**, and don’t need any forks, and times where they **eat**.
 - In order to *eat*, a philosopher needs **two forks**, both the one on their *left* and the one on their *right*.
 - **The contention for these forks.**



The Dining Philosophers (Cont.)

■ Key challenge

- There is **no deadlock**.
- **No** philosopher **starves** and never gets to eat.
- **Concurrency** 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 \rightarrow call `left(p)`.
- ▶ Philosopher p wishes to refer to the fork on their right \rightarrow call `right(p)`.

The Dining Philosophers (Cont.)

- 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.
 - ▶ Each will be stuck *holding one fork* and waiting for another, *forever*.

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.**

How To Implement Semaphores

■ Build our own version of semaphores called Zemaphores

```
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 }
21 ...
```

How To Implement Semaphores (Cont.)

```
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  }
```

- Semaphore don't maintain the invariant that *the value of the semaphore*.
 - ▶ The value never be lower than zero.
 - ▶ This behavior is **easier** to implement and **matches** the current Linux implementation.

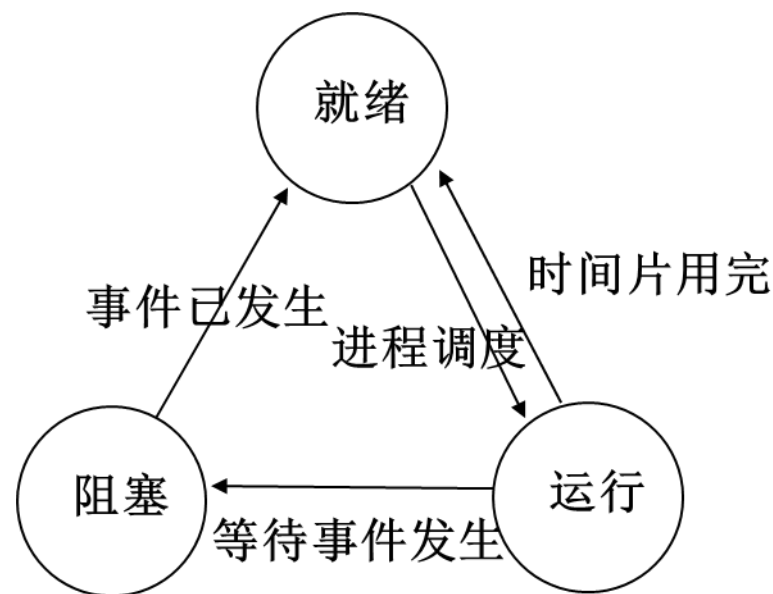
例题

【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。