# Synchronization Primitives (ch. 30+31+32)

## Operating Systems
### Based on: Three Easy Pieces by Arpaci-Dusseaux

Moshe Sulamy

Tel-Aviv Academic College

# Condition Variables

- Many cases a thread wishes to wait until a certain **condition**
- e.g., waiting for another thread to complete
  - Often called a `join()`

- Shared variable: works, but hugely inefficient

> How should a thread wait for a condition?

# Condition Variables

- **Waiting** on a condition
  - Thread puts itself in a queue until some state of execution
- **Signaling** on a condition
  - Some other thread can wake waiting thread

# Condition Variables

- **Waiting** on a condition
  - Thread puts itself in a queue until some state of execution
- **Signaling** on a condition
  - Some other thread can wake waiting thread

- Name is a bit misleading
  - More of a queue
  - We are responsible for the actual "condition"

# Definitions and Routines

- Declare a condition variable:

```
pthread_cond_t cv;
```

- Operations:

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

- Wait call takes **mutex** as a parameter
  - Caller must be its **owner** (have it locked)
  - Releases the lock, puts caller to sleep
  - On wake up, re-acquires lock and returns

# Parent Waiting For Child

- Two threads:
    - **Parent**:
        - Creates child thread
        - Waits on CV until child completes
    - **Child**:
        - Prints a message ("child")
        - Wakes parent by signaling on CV

# Parent Waiting For Child

```
1  int done = 0;
2  pthread_mutex_t m = PTHREAD_MUTEX_INITIALIZER;
3  pthread_cond_t c = PTHREAD_COND_INITIALIZER;
4
5  void* child(void* arg) {
6      printf("child\n");
7      done = 1;
8      thr_exit();
9      return NULL;
10 }
11 int main(void) {
12     printf("parent: begin\n");
13     pthread_t p;
14     pthread_create(&p, NULL, child, NULL);
15     thr_join();
16     printf("parent: end\n");
17     return 0;
18 }
```

# Parent Waiting For Child

```
1 void thr_exit() {
2     pthread_cond_signal(&c);
3 }
4 void thr_join() {
5     pthread_mutex_lock(&m);
6     pthread_cond_wait(&c, &m);
7     pthread_mutex_unlock(&m);
8 }
```

- Why might this code fail?

# Parent Waiting For Child

```
1 void thr_exit() {
2     pthread_cond_signal(&c);
3 }
4 void thr_join() {
5     pthread_mutex_lock(&m);
6     pthread_cond_wait(&c, &m);
7     pthread_mutex_unlock(&m);
8 }
```

- Why might this code fail?
    - Child runs immediately
    - Will signal, but no thread asleep on CV
    - Parent runs, calls wait and gets stuck
    - Solution?

# Parent Waiting For Child

```
1 void thr_exit() {
2     pthread_cond_signal(&c);
3 }
4 void thr_join() {
5     pthread_mutex_lock(&m);
6     pthread_cond_wait(&c, &m);
7     pthread_mutex_unlock(&m);
8 }
```

- Why might this code fail?
  - Child runs immediately
  - Will signal, but no thread asleep on CV
  - Parent runs, calls wait and gets stuck
  - Solution? use `done` variable

# Parent Waiting For Child

```
1 void thr_exit() {
2     pthread_cond_signal(&c);
3 }
4 void thr_join() {
5     pthread_mutex_lock(&m);
6     if (done == 0)
7         pthread_cond_wait(&c, &m);
8     pthread_mutex_unlock(&m);
9 }
```

- Why might this code fail?

# Parent Waiting For Child

```
1 void thr_exit() {
2     pthread_cond_signal(&c);
3 }
4 void thr_join() {
5     pthread_mutex_lock(&m);
6     if (done == 0)
7         pthread_cond_wait(&c, &m);
8     pthread_mutex_unlock(&m);
9 }
```

- Why might this code fail?
  - Parent calls join, sees `done=0`
  - Interrupted just before wait, context switch to child
  - Child sets done, signal is lost, parent is stuck again
  - Solution?

# Parent Waiting For Child

```
1  void thr_exit() {
2      pthread_cond_signal(&c);
3  }
4  void thr_join() {
5      pthread_mutex_lock(&m);
6      if (done == 0)
7          pthread_cond_wait(&c, &m);
8      pthread_mutex_unlock(&m);
9  }
```

- Why might this code fail?
    - Parent calls join, sees done=0
    - Interrupted just before wait, context switch to child
    - Child sets done, signal is lost, parent is stuck again
    - Solution? hold lock while signaling

```
1  void thr_exit() {
2      pthread_mutex_lock(&m);
3      pthread_cond_signal(&c);
4      pthread_mutex_unlock(&m);
5  }
6  void thr_join() {
7      pthread_mutex_lock(&m);
8      while (done == 0)
9          pthread_cond_wait(&c, &m);
10     pthread_mutex_unlock(&m);
11 }
```

- Additionally, check variable in a loop
  - Condition variable may signal unexpectedly
  - Also crucial for more than 2 threads

# Covering Conditions

- Memory allocator implementation
- Assume zero bytes are free:
    - Thread A calls `allocate(100)`
    - Thread B calls `allocate(10)`
    - Both A and B wait on the condition
    - Thread C calls `free(50)`

# Covering Conditions

- Memory allocator implementation
- Assume zero bytes are free:
  - Thread A calls `allocate(100)`
  - Thread B calls `allocate(10)`
  - Both A and B wait on the condition
  - Thread C calls `free(50)`
    - Which waiting thread wakes up?

# Covering Conditions

- Memory allocator implementation
- Assume zero bytes are free:
  - Thread A calls `allocate(100)`
  - Thread B calls `allocate(10)`
  - Both A and B wait on the condition
  - Thread C calls `free(50)`
    - Which waiting thread wakes up?

- Solution: wake up all waiting threads
  - Use `pthread_cond_broadcast()`
  - Performance cost: too many threads might be woken

# Producer / Consumer

- Also: **bounded buffer** problem
- **Producer**
  - Produces data items
  - Wishes to place items in a buffer
- **Consumer**
  - Grabs items out of the buffer
  - Consumes items in some way
- e.g., web server consumes HTTP requests (work queue)

# Producer / Consumer

- Also used in pipes:
  - `grep foo file.txt | wc -l`
  - `grep` - output lines from `file.txt` containing `foo`
  - `wc -l` - output number of lines from input
  - Shell redirects `grep` standard output to a **pipe**
    - Created by the `pipe` system call
  - Other end connected to standard input of `wc`
  - `grep` - producer, `wc` - consumer
- In-kernel bounded buffer

# Producer / Consumer

```
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  int get() {
10     assert(count == 1);
11     count = 0;
12     return buffer;
13 }
```

# First Attempt

- What are the two problems?

```
1  pthread_cond_t cond;
2  pthread_mutex_t mutex;
3
4  void produce(int i) {
5      pthread_mutex_lock(&mutex);
6      if (count == 1)
7          pthread_cond_wait(&cond, &mutex);
8      put(i);
9      pthread_cond_signal(&cond);
10     pthread_mutex_unlock(&mutex);
11 }
12 int consume() {
13     pthread_mutex_lock(&mutex);
14     if (count == 0)
15         pthread_cond_wait(&cond, &mutex);
16     int tmp = get();
17     pthread_cond_signal(&cond);
18     pthread_mutex_unlock(&mutex);
19     return tmp;
20 }
```

# First Attempt

- More than one consumer / producer

# First Attempt

- More than one consumer / producer
- Two consumers:
  - No `while` on CV in `consume()`
  - Can consume when empty!
  - (and the same for `produce()`)

```
1  int consume() {
2      pthread_mutex_lock(&mutex);
3      while (count == 0)
4          pthread_cond_wait(&cond, &mutex);
5      int tmp = get();
6      pthread_cond_signal(&cond);
7      pthread_mutex_unlock(&mutex);
8      return tmp;
9  }
```

buffer=0

Producer

Consumer 1

Consumer 2

# First Attempt

buffer=0

Producer

Consumer 1

consume

Consumer 2

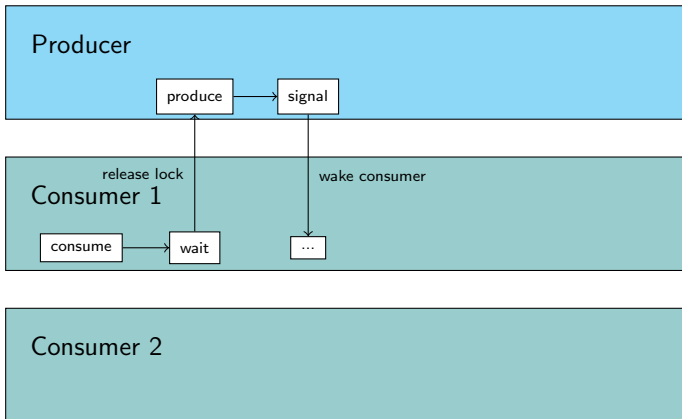# First Attempt

buffer=0

Producer

Consumer 1

consume ⟶ wait

Consumer 2

# First Attempt

buffer=1

Producer

produce → signal

release lock

Consumer 1

consume → wait

Consumer 2

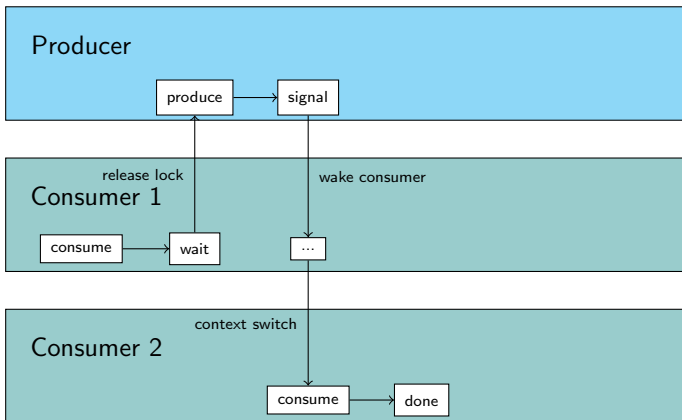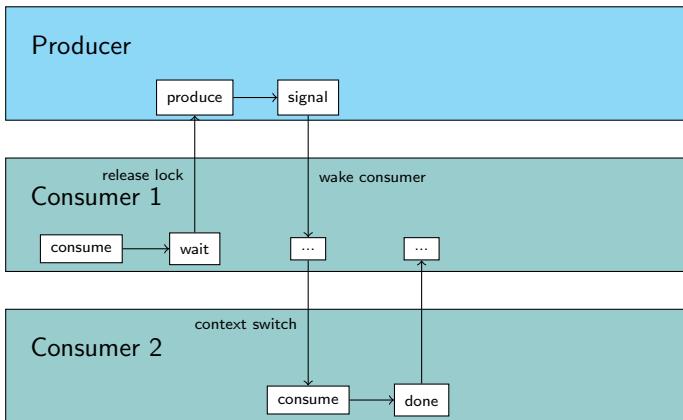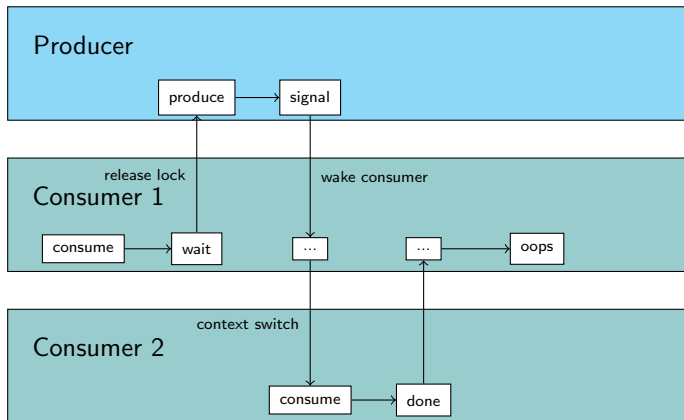# First Attempt

# First Attempt

# First Attempt

# First Attempt

# First Attempt

- Only one condition variable!
  - Consumer might wake another consumer
  - Producer might wake another producer

- Solution?

# First Attempt

- Only one condition variable!
    - Consumer might wake another consumer
    - Producer might wake another producer

- Solution? use **two** condition variables
    - Producer threads wait on `empty`
    - Consumer threads wait on `full`

# Semaphores

- Another **synchronization primitive**
  - Can use **semaphores** as both locks and condition variables

- A **semaphore**
  - Object with an integer value
  - sem_wait(), sem_post()

# Semaphores

- Initialization:

```
#include <semaphore.h>
sem_t s;
sem_init(&s, 0, 1);
```

- Initialized to the value 1 (third argument)
    - Second argument: 0 to share between threads (vs. processes)
    - We will only use 0

# Semaphores

- `sem_wait()`:
    - Decrement semaphore value by one
    - Wait if value is negative

- `sem_post()`:
    - Increment semaphore value by one
    - Wake one waiting thread

- Both operations are performed **atomically**

# Binary Semaphores

- How can we use a semaphore as a lock?

```
1 sem_init(&m, 0, X); // what should X be?
2 ...
3 sem_wait(&m);
4 // critical section
5 sem_post(&m);
```

- What should **X** be?

# Binary Semaphores

- How can we use a semaphore as a lock?

```
1 sem_init(&m, 0, X); // what should X be?
2 ...
3 sem_wait(&m);
4 // critical section
5 sem_post(&m);
```

- What should **X** be?
  - X = 1
  - **Binary semaphore**

# Ordering Primitive

- Similar to **condition variables**
- How can we use a semaphore to wait for an event?

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

- What should **X** be?

# Ordering Primitive

- Similar to **condition variables**
- How can we use a semaphore to wait for an event?

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

- What should **X** be?
  - X=0

# Semaphores

- What value should a semaphore be initialized to?
  - General rules: number of resources
  - Lock: $1 \rightarrow$ can be locked after initialization
  - Ordering: $0 \rightarrow$ nothing to give away at the start

# Producer / Consumer

- How can we implement a **bounded buffer** with semaphores?

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

# First Attempt

- What is the problem?

```
1  sem_t empty; // initialized to MAX
2  sem_t full;  // initialized to 0
3
4  void produce(int value) {
5      sem_wait(&empty);
6      put(value);
7      sem_post(&full);
8  }
9  int consume() {
10     sem_wait(&full);
11     int tmp = get();
12     sem_post(&empty);
13     return tmp;
14 }
```

# First Attempt

Producer 1

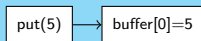Producer 2

Producer 1

put(5)

Producer 2

# First Attempt

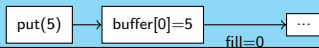# First Attempt

# First Attempt



Producer 1
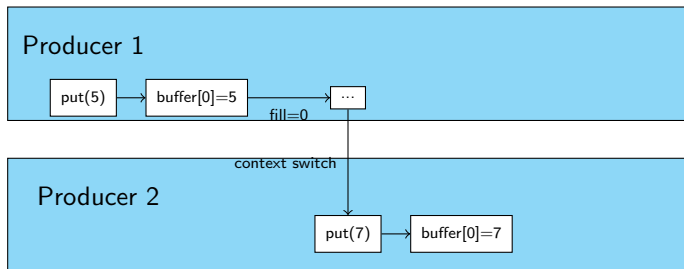
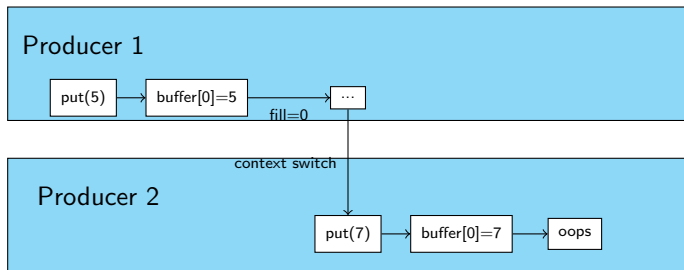put(5) → buffer[0]=5 → ⋯

fill=0

context switch

Producer 2

put(7) → buffer[0]=7

# First Attempt

# With Mutual Exclusion

- What is the problem?

```
1  sem_t mutex; // binary semaphore
2  sem_t empty; // initialized to MAX
3  sem_t full;  // initialized to 0
4
5  void produce(int value) {
6      sem_wait(&mutex);
7      sem_wait(&empty);
8      put(value);
9      sem_post(&full);
10     sem_post(&mutex);
11 }
12 int consume() {
13     sem_wait(&mutex);
14     sem_wait(&full);
15     int tmp = get();
16     sem_post(&empty);
17     sem_post(&mutex);
18     return tmp;
19 }
```

# With Mutual Exclusion

- What is the problem? **deadlock!**
  - Producer waits on empty, holds mutex, consumer can't consume

```c
1  sem_t mutex; // binary semaphore
2  sem_t empty; // initialized to MAX
3  sem_t full;  // initialized to 0
4
5  void produce(int value) {
6      sem_wait(&mutex);
7      sem_wait(&empty);
8      put(value);
9      sem_post(&full);
10     sem_post(&mutex);
11 }
12 int consume() {
13     sem_wait(&mutex);
14     sem_wait(&full);
15     int tmp = get();
16     sem_post(&empty);
17     sem_post(&mutex);
18     return tmp;
19 }
```

# With Mutual Exclusion

- Solution: use mutex around the critical section

```
1  sem_t mutex; // binary semaphore
2  sem_t empty; // initialized to MAX
3  sem_t full;  // initialized to 0
4
5  void produce(int value) {
6      sem_wait(&empty);
7      sem_wait(&mutex);
8      put(value);
9      sem_post(&mutex);
10     sem_post(&full);
11 }
12 int consume() {
13     sem_wait(&full);
14     sem_wait(&mutex);
15     int tmp = get();
16     sem_post(&mutex);
17     sem_post(&empty);
18     return tmp;
19 }
```

# Reader-Writer Locks

- More flexible locking primitive
  - e.g., concurrent operations: inserts and lookups
  - Insert changes state $\rightarrow$ traditional critical section
  - Lookup reads data structure $\rightarrow$ many at once (if no insert)

- **Reader-writer lock**
  - Four operations: acquire/release read/write lock

# Reader-Writer Locks

- A **single writer** can acquire the lock
- Once a reader acquires a **read lock**:
    - **More readers** are allowed to acquire the read lock
    - A writer waits until all readers are finished

# Reader-Writer Locks

- A **single writer** can acquire the lock
- Once a reader acquires a **read lock**:
    - **More readers** are allowed to acquire the read lock
    - A writer waits until all readers are finished

```
1  typedef struct _rwlock_t {
2      sem_t lock;             // binary semaphore
3      sem_t writelock;        // used to allow ONE writer
4      int readers;            // count of readers in CS
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 }
```

# Reader-Writer Locks

```
1  void rwlock_acquire_writelock(rwlock_t* rw) {
2      sem_wait(&rw->writelock);
3  }
4  void rwlock_release_writelock(rwlock_t* rw) {
5      sem_post(&rw->writelock);
6  }
7
8  void rwlock_acquire_readlock(rwlock_t* rw) {
9      sem_wait(&rw->lock);      // CS for readers
10     rw->readers++;
11     if (rw->readers == 1)
12         sem_wait(&rw->writelock); // first reader grabs writelock
13     sem_post(&rw->lock);
14 }
15 void rwlock_release_readlock(rwlock_t* rw) {
16     sem_wait(&rw->lock);      // CS for readers
17     rw->readers--;
18     if (rw->readers == 0)
19         sem_post(&rw->writelock); // last reader releases writelock
20     sem_post(&rw->lock);
21 }
```

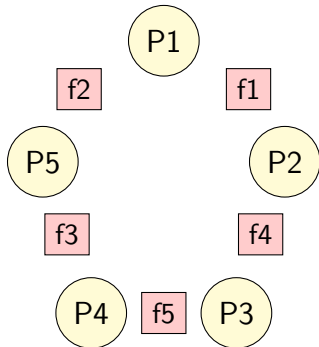# Reader-Writer Locks

- What is the problem?

# Reader-Writer Locks

- What is the problem? **fairness**
    - Easy to **starve** writer
    - How to prevent readers from starving writers?

# Dining Philosophers

- Five philosophers around a table
  - <u>Single</u> fork between each pair
  - Philosophers **think** and **eat**
  - Two forks to eat (left and right)

# Dining Philosophers

- As code:

```
1  while (1) {
2      think();
3      getforks();
4      eat();
5      putforks();
6  }
```

```
1  int left(int p) {
2      return p;
3  }
4  int right(int p) {
5      return (p + 1) % 5;
6  }
```

# Dining Philosophers

- Use a semaphore for each fork:

```
1 void get_forks(int p) {
2     sem_wait(&forks[left(p)]);
3     sem_wait(&forks[right(p)]);
4 }
5 void put_forks(int p) {
6     sem_post(&forks[left(p)]);
7     sem_post(&forks[right(p)]);
8 }
```

- The problem?

# Dining Philosophers

- Use a semaphore for each fork:

```
1 void get_forks(int p) {
2     sem_wait(&forks[left(p)]);
3     sem_wait(&forks[right(p)]);
4 }
5 void put_forks(int p) {
6     sem_post(&forks[left(p)]);
7     sem_post(&forks[right(p)]);
8 }
```

- The problem? **deadlock!**
    - Each philosopher grabs fork on their left
    - All waiting for their right

# Dining Philosophers

- Solution: break the dependency

```
1  void get_forks(int p) {
2      if (p == 4) {
3          sem_wait(&forks[right(p)]);
4          sem_wait(&forks[left(p)]);
5      }
6      else {
7          sem_wait(&forks[left(p)]);
8          sem_wait(&forks[right(p)]);
9      }
10 }
```

# Thread Throttling

- For example: hundreds of threads work in parallel
- Section of code allocates a lot of memory
    - All threads at the same time $\rightarrow$ exceeds physical memory
    - Machine will start thrashing (swapping to and from the disk)

- Solution?

# Thread Throttling

- For example: hundreds of threads work in parallel
- Section of code allocates a lot of memory
  - All threads at the same time $\rightarrow$ exceeds physical memory
  - Machine will start thrashing (swapping to and from the disk)

- Solution? **semaphore**!
  - Initialized to max threads we wish to enter code section
  - Surrounds code section, limits concurrent threads in it

# Implementing Semaphores

- Doesn't maintain invariant: negative value $\rightarrow$ # waiting threads
  - Easier, matches the Linux implementation

```
1   typedef struct __zem_t {
2       int value;
3       pthread_cond_t cond;
4       pthread_mutex_t lock;
5   };
6
7   void zem_init(zem_t* s, int value) {
8       s->value = value;
9       pthread_cond_init(&s->cond);
10      pthread_mutex_init(&s->lock);
11  }
12  void zem_wait(zem_t* s) {
13      pthread_mutex_lock(&s->lock);
14      while (s->value <= 0)
15          pthread_cond_wait(&s->cond, &s->lock);
16      s->value--;
17      pthread_mutex_unlock(&s->lock);
18  }
19  void zem_post(zem_t* s) {
20      pthread_mutex_lock(&s->lock);
21      s->value++;
22      pthread_cond_signal(&s->cond);
23      pthread_mutex_unlock(&s->lock);
24  }
```

# Summary

- **Condition variables**
  - Thread waits until a certain condition
  - `wait()`, `signal()`
  - **Hold lock** while signaling
  - Check value **in a loop**
- **Semaphore**
  - Integer value
  - Decrement on acquire, wait if negative, increment on release
- **Read-write lock**
  - **Single writer** or **multiple readers**
- Dining philosophers
  - Think and eat
- Producer / consumer (bounded buffer)