

OPERATING

SYSTEM

Study Notes

Zhijie Xia
Study Notes
github.com/zhijie-os

Contents

1 Virtualization	5
1.1 Segmentation And Paging	5
1.1.1 Segmentation	5
1.1.2 Free Space Management	6
1.1.3 Advanced Page Tables	7
1.2 Swapping	7
1.2.1 Swapping Mechanisms	7
1.2.2 Swapping Policies	7
2 Cucurrency	11
2.1 Cucurrency and Thread	11
2.1.1 Similarities between Processes and Threads	11
2.1.2 Difference between Processes and Threads	11
2.1.3 Benefit of using threads	11
2.1.4 Problem with threads: Race Condition	12
2.1.5 Thread API	13
2.1.6 Locks and Building one	14
2.1.7 Locked Data Structures	18
2.1.8 Conditional Variables	21
2.1.9 Semaphore	26
2.1.10 Deadlock bugs	31
2.1.11 Event-Based Concurrency	33
3 Persistence	35
3.1 I/O devices	35
3.2 Hard Disk Drives	38
3.2.1 RAID	39
3.3 Files and Directories	41
3.3.1 File System Implementation	45
3.3.2 The Inode	46

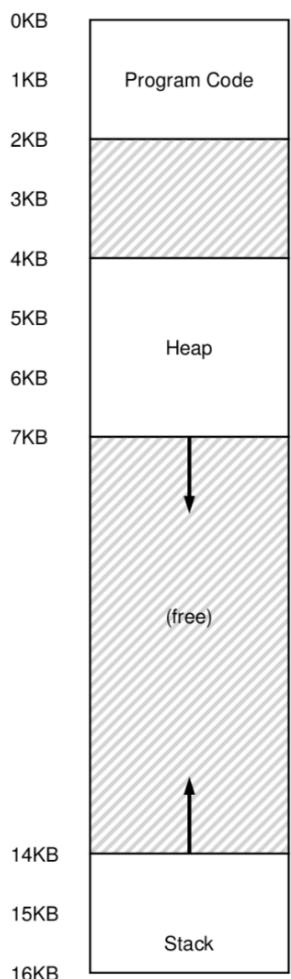
Chapter 1

Virtualization

1.1 Segmentation And Paging

1.1.1 Segmentation

Segmentation is compiler's view about memory.



In canonical address space, we have three logically-different segments:

1. Code/Text
2. Heap
3. Stack

Base/Bound registers

Each segments would have a pair of base and bound registers.
We would use base and bound/limit registers to translate address.

$$\text{Physical Address} = \text{Base Address} + \text{Offset}$$

Bound register is used to check boundary. If offset is greater than bound register value \Rightarrow Segmentation fault.

Note: the stack grows in opposite direction

Referring Segments

We would chop up the address space into two parts:

1. 2 bits: indicate the segments
2. the rest: offset within the segments



Assume 8-bit address space in use:

1. 00xxxxxx: this is invalid, we would address 1/4 less because we don't use 00.
2. 01xxxxxx: would be in the Code segment.
3. 10xxxxxx: would be in the Heap segment.
4. 11xxxxxx: would be in the Stack segment.

Support for sharing

It is good idea to share the code segment, but how the OS supports it?

There would be protection bits to indicate whether the segment can be used for.

Segment	Base	Size (max 4K)	Grows Positive?	Protection
Code ₀₀	32K	2K	1	Read-Execute
Heap ₀₁	34K	3K	1	Read-Write
Stack ₁₁	28K	2K	0	Read-Write

OS support

So support segmentation, the OS

1. Context Switch: The segment registers must be saved and restored.
2. Support growth or shrinkage of a segment: Management of free space, Compacting physical memory and Free-list management algorithms.

1.1.2 Free Space Management

Low-level Mechanisms

Splitting and coalescing:

Basic Strategies

1.1.3 Advanced Page Tables

1.2 Swapping

1.2.1 Swapping Mechanisms

Use hard disk drive to stash portions of address spaces that currently aren't in great demand, so that we can support programs that take more memories than RAM has.

Swap Space

Swap Space: Reserved space on the disk for moving pages back and forth.

And OS can read from and write to swap space in page-sized units. Also OS needs to know the exact address of the swap space in order to quickly swap pages.

The Present Bit

Inside of page table entry, there is a present bit.

When the present bit is set, that indicates the page is in the physical memory.

If the present bit is off, the page is not in the memory. When a TLB miss resulting in a page table entry and found the present bit is off, it is a page fault (indicates the demanding page is not in the physical memory).

The Page Fault

The OS would invoke Page-Fault Handler to deal with page fault.

1. The OS would look into PTE(page table entry) to find the address to fetch.
2. After completing disk I/O, the OS updates PTE to mark the page as present.
3. Update the PFN(physical frame number) of the PTE to the in-memory location of the fetched-page.
4. Retry the instruction.

When the memory is full

When page-fault and the memory is full, the OS needs to kick out a page to place a new page in. Therefore, page-replacement policy is needed.

1.2.2 Swapping Policies

Cache Management

Main memory holds some subset of all the pages of ongoing processes \Rightarrow main memory is a cache for virtual memory pages.

The goal is to minimize the number of cache misses.

Average Memory Access Time

$$AMAT = T_M + (P_{miss} \times T_D)$$

Where,

1. T_M : the cost of accessing memory
2. T_D : the cost of accessing disk
3. P_{miss} : the probability of cache miss

Assume that $T_M = 100$ nanoseconds, $T_D = 10$ milliseconds .

If the hit rate is 0.9, the AMTA would be $100\text{ns} + 0.1 \cdot 10\text{ms} = 1.0001 \text{ ms}$.

However, when the hit rate is 0.99, the AMTA would be 10.01 microsecs which is 100x faster. That is, the performance is heavily based on the hit rate \rightleftharpoons swapping policy matters.

The Optimal Replacement Policy

The optimal replacement policy leads to the fewest number of misses overall.

Belady (a person) showed that a simple policy that leads to optimal: The page that would be accessed furthest in the future is the optimal policy (This is very like shortest job first CPU scheduling,i.e, impossible!)

Access	Hit/Miss?	Evict	Resulting Cache State
0	Miss		0
1	Miss		0, 1
2	Miss		0, 1, 2
0	Hit		0, 1, 2
1	Hit		0, 1, 2
3	Miss	2	0, 1, 3
0	Hit		0, 1, 3
3	Hit		0, 1, 3
1	Hit		0, 1, 3
2	Miss	3	0, 1, 2
1	Hit		0, 1, 2

First three accesses are misses, and the misses are called **cold-start misses**.

When Access 3 at the first time, the optimal policy decides to evict 2 because 0 and 1 would be accessed before 2.

However, future is unpredictable, an another approach is needed.

FIFO

Old friend, First In First Out.

FIFO is very simple to implement.

Access	Hit/Miss?	Evict	Resulting Cache State
0	Miss	First-in→	0
1	Miss	First-in→	0, 1
2	Miss	First-in→	0, 1, 2
0	Hit	First-in→	0, 1, 2
1	Hit	First-in→	0, 1, 2
3	Miss	0	First-in→ 1, 2, 3
0	Miss	1	First-in→ 2, 3, 0
3	Hit		2, 3, 0
1	Miss	2	First-in→ 3, 0, 1
2	Miss	3	First-in→ 0, 1, 2
1	Hit		First-in→ 0, 1, 2

FIFO could often do poorly because it cannot determine the importance of a page.

Random

Simple to implement, and do well if the distribution of accessing page is uniform distributed. However, it is unlikely for accessing page to follow a particular distribution.

Random is better than FIFO, and a bit worst than optimal.

LRU

If page is accessed in the near past, it is likely to be accessed in near future.

Some historically-based algorithms are used. LFU: Least-Frequently-Used, and LRU: Least-Recently-Used.

Access	Hit/Miss?	Evict	Resulting Cache State
0	Miss		LRU→ 0
1	Miss		LRU→ 0, 1
2	Miss		LRU→ 0, 1, 2
0	Hit		LRU→ 1, 2, 0
1	Hit		LRU→ 2, 0, 1
3	Miss	2	LRU→ 0, 1, 3
0	Hit		LRU→ 1, 3, 0
3	Hit		LRU→ 1, 0, 3
1	Hit		LRU→ 0, 3, 1
2	Miss	0	LRU→ 3, 1, 2
1	Hit		LRU→ 3, 2, 1

The LRU policy works well to matching the optimal.

Implement Historical Algorithms

Using LRU, the system needs to count the least- and most-recently used which is a lot of work. Bad implementation would led heavy performance penalty.

Even adding timestamp for every process accessing, it is unlikely to scan the all pages to find the absolute least recently used page.

Approximating LRU would be a solution:

1. Need a use bit: use bit is set to 1 when the page is accessed.
2. Use a clock algorithm: a clock pointer points to each particular page, if the use bit is 1, the OS clear the use bit and the pointer points to the next page. If found a page with use bit 0, replace it. The worst case is looping through the entire set of pages for one circle.

Dirty Pages

If a page is modified while in the memory, it is dirty. It would cost a lot to evict dirty page since the page must be written to the disk first. Therefore, most algorithms would favor to evict clean pages over dirty pages.

To support the behavior, the hardware includes a modified bit. The bit is set when the page is modified and cleared when written to disk.

Other Policies

Page selection policy determines when to bring a page into the memory. For most pages, OS would use **demanding paging**, which means the OS brings the page into memory when it is accessed. Or OS would predict which page would be accessed in the future and bring it to the memory, this is called **prefetching**.

Another policy determines how the OS writes pages out to disk. **Clustering** is a behavior that OS buffers the changes and write out to disk in one write.

Thrashing

Thrashing: When the demanding of pages exceed the available physical memory, the system would constantly being paging.

Linux would ran **out-of-memory killer** to choose some memory-intensive process and kill them.

Chapter 2

Cucurrency

2.1 Cucurrency and Thread

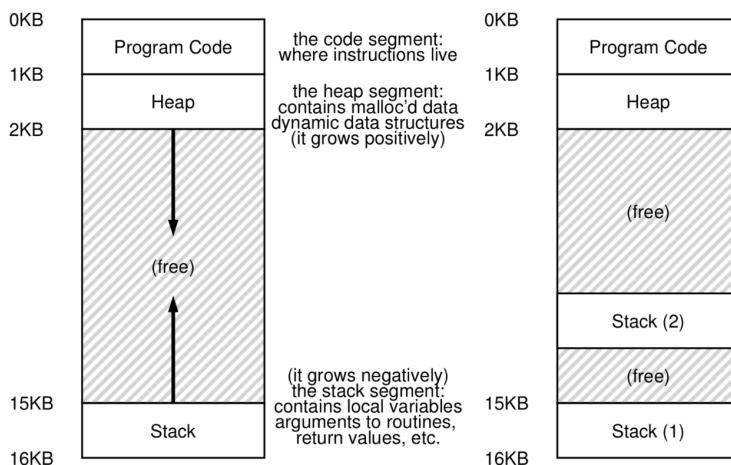
Thread: very much like process, except threads share the same address space and thus can access the same data.

2.1.1 Similarities between Processes and Threads

A thread has a program counter(PC), own set of registers. When switching from thread T1 to thread T2, a context switch must take place. There would be thread control blocks (TCB) like PCB to store the state of each thread.

2.1.2 Difference between Processes and Threads

1. The address space of threads within a single process is the same.
2. Multi-threaded Address Spaces has different structure: one stack per thread called *thread-local* storage.



2.1.3 Benefit of using threads

Parallelism

Run works parallelly.

Avoid Blocking

Avoid blocking program progress due to slow I/O; while one thread in the program waits, the CPU scheduler can switch to other ready threads.

Threading enables overlap of I/O with other activities within a single program.

Why not use processes instead?

Threads make it easy to share data, and often used to cooperate with other threads to finish tasks.

Processes are more sound choice for logically separate tasks when little sharing of data structures in memory is needed.

2.1.4 Problem with threads: Race Condition

The execution sequence of threads is indeterministic.

Create two threads to update on the same global variable with the same function.

```

11 void *mythread(void *arg) {
12     char *letter = arg;
13     int i; // stack (private per thread)
14     printf("%s: begin [addr of i: %p]\n", letter, &i); //threads would share the same data
15     for (i = 0; i < max; i++) {
16         // at here, it would like
17
18         // ldr x1,counter
19         // add x1,x1,1 ---- if the context switch happens here, it could cause problem.
20         // str x1,counter
21         counter = counter + 1; // shared: only one
22     }
23     printf("%s: done\n", letter);
24     return NULL;
25 }
26

```

The problem can be:

The image shows three terminal sessions (01:46:32, 01:46:43, and 01:46:46) running on a Linux system. Each session runs the command `./thread_counter 200000`. The output shows two threads, A and B, updating a shared counter. In each session, the final value of the counter is different (196611, 219206, and 225131 respectively), demonstrating the indeterministic nature of thread execution due to context switching.

```

[01:46:32|zhijie@zhijieLinux:[Processes_Threads] ➔ ./thread_counter 200000
main: begin [counter = 0] [ed918070]
A: begin [addr of i: 0x7f11f6c74e3c]
B: begin [addr of i: 0x7f11f6473e3c]
B: done
A: done
main: done
[counter: 196611]
[should: 400000]
[01:46:43|zhijie@zhijieLinux:[Processes_Threads] ➔ ./thread_counter 200000
main: begin [counter = 0] [57dbc070]
A: begin [addr of i: 0x7f28e9ff3e3c]
B: begin [addr of i: 0x7f28e97f2e3c]
A: done
B: done
main: done
[counter: 219206]
[should: 400000]
[01:46:46|zhijie@zhijieLinux:[Processes_Threads] ➔ ./thread_counter 200000
main: begin [counter = 0] [76b52070]
A: begin [addr of i: 0x7f21177dee3c]
B: begin [addr of i: 0x7f2116fdde3c]
A: done
B: done
main: done
[counter: 225131]
[should: 400000]

```

Assembly Code

In ARMx8 Assembly:

counter = counter + 1 is equivalent to

1. ldr x1,[counter]
2. add x1, x1, 1
3. str x1, [counter]

The work flow of causing problem

1. Thread A loads counter into x_1 , say $x_1=50$.
2. Context switch happens, and switch to Thread B.
3. Now thread B loads counter into x_1 , $x_1=50$.
4. Thread B increase x_1 by 1, $x_1=51$.
5. Thread B stores x_1 back to counter, $counter=51$.
6. Context switch happens, and switch back to Thread A.
7. Context Switch restores x_1 for A,i.e, $x_1=50$. And A won't load counter to x_1 again
8. Thread A increase x_1 by 1, $x_1=51$.
9. Thread A stores x_1 back to counter, $counter=51$.
10. Thus, counter is set to 51 twice, although it should be 52 after the flow.

Critical Section

Critical Section: A piece of code that accesses a shared variable, and must not be concurrently executed by more than one thread.

Mutual Exclusion

Mutual Exclusion: if one thread is executing within the critical section, the others will be prevented from doing so.

Race Condition

Multiple threads of execution enter the critical section at roughly the same time; both attempt to update the shared data structure, leading unexpected outcome.

Atomicity

Atomic operation: grouped actions to be executed in one scheduling, i.e, the operation won't be interrupted.

For example, $x_1 = x_1 + x_2$ could be done in one single step with hardware support, instead of load, add, and store.

It is desired to support atomicity for critical sections.

2.1.5 Thread API

Lock

In POSIX library, a lock needs to be initialized

```
1 int rc = pthread_mutex_init(&lock, NULL);
2 assert(rc == 0); // always check success!
```

Also, a thread can acquire a lock, and release a lock.

```
1 int pthread_mutex_lock(pthread_mutex_t *mutex);
2 int pthread_mutex_unlock(pthread_mutex_t *mutex);
```

Condition Variable

Condition variables are useful when some kind of signaling must take place between threads if one thread is waiting for another to do something before it continues.

```
1 int pthread_cond_wait(pthread_cond_t *cond, pthread_mutex_t *mutex);
2 int pthread_cond_signal(pthread_cond_t *cond);
```

A thread must hold a lock to call either `wait()` or `signal()`. `pthread_cond_wait()`, puts the calling thread to sleep. `pthread_cond_signal()` awakes the waiting thread.

The reason that `pthread_cond_wait()` takes two parameter is because it needs to specify which thread to give the lock to. When `pthread_cond_wait()` the calling thread release the lock and pass it to another thread.

2.1.6 Locks and Building one

Basic

Lock is used around the critical section. It is a global variable that either available or acquired and exactly one thread can hold it at a time.

`mutex` in POSIX means ***mutual exclusion*** between threads.

Evaluating Locks

Basic criteria:

1. Mutual exclusion: A lock must provide mutual exclusion, i.e, the lock should preventing multiple threads from entering a critical section.
2. Fairness: Prevent starving a lock.
3. Performance: How many overhead would be added to use the lock.

Lock by Controlling Interrupts

One of the earliest implementation of lock is disable interrupts.

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

This approach works since it assumes mutual exclusion, and it is very simple.
However, it has flaws:

1. Privileged action: malicious process would disable the interrupts and never enable interrupts again.
2. Interrupts would get lost: for example, I/O interrupts would get lost and some processes which are waiting on those interrupts cannot move forward.
3. Inefficient approach: It is very costly to enable/disable interrupts.
4. No support for multiprocessors.

A Fail Attempt: Using a Flag

```

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 }
```

Correctness Problem: Interleaving would give more locks than just one.

Thread 1	Thread 2
call lock ()	
while (flag == 1)	
interrupt: switch to Thread 2	
	call lock ()
	while (flag == 1)
	flag = 1;
	interrupt: switch to Thread 1
flag = 1; // set flag to 1 (too!)	

Performance Problem: While loops is valid instruction that would use CPU, it is very likely a thread which acquiring the lock spends its timeslot to loop. This behavior is called **busy-waiting** or **spin-waiting**.

Test-and-Set

Test-and-Set is atomic instruction supported by the hardware. It both gets and sets the value in a register/address.

```

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

With **Test-and-Set**, we can build a correct lock.

```

1  typedef struct __lock_t {
2  |   int flag;
3  } lock_t;
4
5  void init(lock_t *lock) {
6  |   // 0: lock is available, 1: lock is held
7  |   lock->flag = 0;
8  }
9
10 void lock(lock_t *lock) {
11 |   while (TestAndSet(&lock->flag, 1) == 1)
12 |   |   ; // spin-wait (do nothing)
13 }
14
15 void unlock(lock_t *lock) {
16 |   lock->flag = 0;
17 }

```

However, the Test-and-Set approach doesn't guarantee

1. Fairness: There is no intelligence invoked to provide fairness.
2. Performance: It is painful on a single CPU. Acceptable on multiple CPUs, because the CPU scheduler would switch the waiting thread out after its timeslot.

Compare-And-Swap

Compare-And-Swap: Test whether the value equals; if so, update the memory value. Finally, it would return the original value.

```

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

```

With **Compare-And-Swap**, it is possible to build a spin clock

```

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

```

Fetch-And-Add

```

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

```

And a spin clock can be built, a *ticket lock*

```

1  typedef struct __lock_t {
2      int ticket;
3      int turn;
4  } lock_t;
5
6
7  void lock_init(lock_t *lock) {
8      lock->ticket = 0;
9      lock->turn = 0;
10 }
11
12 void lock(lock_t *lock) {
13     int myturn = FetchAndAdd(&lock->ticket);
14     while (lock->turn != myturn)
15     | ; // spin
16 }
17
18 void unlock(lock_t *lock) {
19     | | lock->turn = lock->turn + 1;
20 }

```

The advantage of *ticket lock* is that the approach "remember" the requests of lock, i.e, it is like to pick a ticket that has number on it. And once the thread picks its ticket, all it needs to do is to wait and be called.

⇒ this approach guarantees fairness.

Other approaches are like fighting with each other for a single ticket.

Spin locks and hardware limitation

With the extra hardware supports, we can now build spin locks. However the inefficiency is like a disaster.

Suppose each threads executes the same amount of time, with N threads and a single lock. The actual work done is $\frac{1}{N}$, and $\frac{N-1}{N}$ is useless busy-waiting.

The hardware cannot solve everything, a smart software needs to be introduced in OS.

Just yield, Baby

```

void lock() {
    while (TestAndSet(&flag, 1) == 1)
        yield(); // give up the CPU
}

```

`yield()` is an operating system primitive which a thread can call when it wants to give up the CPU and let another thread to run, i.e, descheduling the calling thread and move it from running to ready.

Efficiency Problem: Suppose there are N threads and a single lock and the scheduler is taking round robin , N-1 yield would be called and only 1 critical section instruction would be executed. That is still bad.

Fairness Problem: If the scheduler is not using round robin, a thread would be picked consecutive to call `yield()` which introduce the possiblity of starving a process.

Using Queues: Sleeping Instead of Spinning

There are some controls needed over which thread next gets to acquire the lock after the current holder release it.

Some OS support is need ⇒ A queue to keep track of which threads are waiting to acquire the lock.

`park()` and `unpark()`.

In Solaris: *park()* to put a calling thread to sleep and *unpark(threadID)* to wake a particular thread as designated by *threadID*.

When a thread tries to acquire the lock, the OS would park() to put the thread into sleeping and awakes it by calling unpark(threadID) when the lock is free.

```

1  typedef struct __lock_t {
2  |    int flag;
3  |    int guard;
4  |    queue_t *q;
5 } lock_t;
6
7 void lock_init(lock_t *m) {
8 |    m->flag = 0;
9 |    m->guard = 0;
10 |   queue_init(m->q);
11 }
```

Flag indicates whether the lock is available, *guard* is used to ensure atomicity within the lock()/unlock().

```

13 void lock(lock_t *m) {
14 |    while (TestAndSet(&m->guard, 1) == 1)
15 |        ; //acquire guard lock by spinning
16 |    if (m->flag == 0) {
17 |        m->flag = 1; //lock is acquired
18 |        m->guard = 0;
19 |    } else {
20 |        queue_add(m->q, gettid());
21 |        m->guard = 0;
22 |        park();
23 |    }
24 }
25
26 void unlock(lock_t *m) {
27 |    while (TestAndSet(&m->guard, 1) == 1)
28 |        ; //acquire guard lock by spinning
29 |    if (queue_empty(m->q))
30 |        m->flag = 0; // let go of lock; no one wants it
31 |    else
32 |        unpark(queue_remove(m->q)); // hold lock
33 |                                         // (for next thread!)
34 |    m->guard = 0;
35 }
```

The guard is like another lock for the lock()/unlock(), one thread needs to hold the guard to acquire the lock or put itself into sleep and release the guard.

Advantage:

1. Small spinning/waste.
2. Fairness by using the queue.

Futex

A linux approach.... Kinda complicated, would take a look when studying linux.

2.1.7 Locked Data Structures

We can use locks in some common data structures to make the structures thread safe.

Cucurrent Counters

```

class counter{
    pthread_mutex_t lock;
public:
    int value;

    counter(){
        init();
    }

    void init(){
        value = 0;
        pthread_mutex_init(&lock,NULL);
    }

    void increment(){
        pthread_mutex_lock(&lock);
        value++;
        pthread_mutex_unlock(&lock);
    }

    void decrement(){
        pthread_mutex_lock(&lock);
        value--;
        pthread_mutex_unlock(&lock);
    }

    int getValue(){
        pthread_mutex_lock(&lock);
        int rc = value;
        pthread_mutex_unlock(&lock);
        return rc;
    }
};

```

The lock is acquired and released automatically as you call and return from object methods, like a monitor.

However, this approach scales poorly.

Approximate counters works scably. The idea is to give each thread a local counter, each thread updates on its local counter, and when the counter reaches the threshold. It acquire the global lock and accumulate its local counter to the global counter.

Concurrent Linked Lists

```

class Node {
public:
    int key;
    Node *next;
};

class List {
    Node *head;
    pthread_mutex_t lock;

public:
    void List_Init() {
        head = NULL;
        pthread_mutex_init(&lock, NULL);
    }

    void List_Insert(int key) {

        Node *toInsert = new Node;
        toInsert->key = key;

        pthread_mutex_lock(&lock);

        // actuall critical section
        toInsert->next = head;
        head = toInsert;

        pthread_mutex_unlock(&lock);
    }

    Node* List_Lookup(int key) {
        pthread_mutex_lock(&lock);
        Node *iter = head;
        while (iter != NULL) {
            if (iter->key == key) {
                pthread_mutex_unlock(&lock);
                return iter;
            }
            iter = iter->next;
        }
        pthread_mutex_unlock(&lock);
        return NULL;
    }
};

```

However, this concurrent linked list is not scalable.

Hand-over-hand locking : Instead of having a single lock for the entire list, add a lock per node of the list. When traversing the list, the code first grabs the next node's lock and then release the current node's lock.(Potentially Deadlock). Also, practically it is hard to make such a structure faster than the simple single lock approach, as the overheads of acquiring and releasing locks for each node is prohibitive.

Concurrent Queues

It is always standard to add a big lock to entire structure.

Here is a more interesting concurrent queue designed by Michael and Scott: A lock for the head, and a lock for the tail.

```

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 head_lock, tail_lock;
10 }
11
12 void Queue_Init(queue_t *q) {
13   node_t *tmp = malloc(sizeof(node_t));
14   tmp->next = NULL;
15   q->head = q->tail = tmp;
16   pthread_mutex_init(&q->head_lock, NULL);
17   pthread_mutex_init(&q->tail_lock, NULL);
18 }
19
20 void Queue_Enqueue(queue_t *q, int value) {
21   node_t *tmp = malloc(sizeof(node_t));
22   assert(tmp != NULL);
23   tmp->value = value;
24   tmp->next = NULL;
25
26   pthread_mutex_lock(&q->tail_lock);
27   q->tail->next = tmp;
28   q->tail = tmp;
29   pthread_mutex_unlock(&q->tail_lock);
30 }
31
32 int Queue_Dequeue(queue_t *q, int *value) {
33   pthread_mutex_lock(&q->head_lock);
34   node_t *tmp = q->head;
35   node_t *new_head = tmp->next;
36   if (new_head == NULL) {
37     pthread_mutex_unlock(&q->head_lock);
38     return -1; // queue was empty
39   }
40   *value = new_head->value;
41   q->head = new_head;
42   pthread_mutex_unlock(&q->head_lock);
43   free(tmp);
44   return 0;
45 }

```

Concurrent Hash

A chaining hashtable is just an array of linked list. Since we have concurrent linked list, we could just use concurrent linked list to construct a hashtable. No more work needed.

```

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 int Hash_Insert(hash_t *H, int key) {
14   return List_Insert(&H->lists[key % BUCKETS], key);
15 }
16
17 int Hash_Lookup(hash_t *H, int key) {
18   return List_Lookup(&H->lists[key % BUCKETS], key);
19 }

```

2.1.8 Conditional Variables

There are many cases that we want to check whether a condition is true before continuing execute. For example, a parent thread would like to wait its children to finish execution using `pthread_join()`.

A global conditional variable can be used, and parent spins to wait until the global variable becomes true. However, we would like another approach that puts parent into sleep instead of spinning.

Definition and Routines

A **condition variable** is an explicit **queue** that threads can put themselves on when some condition, which is not as desired by waiting on the condition.

From my understanding, if a thread A is waiting for a state, it would put it into the queue that every thread in the queues wants the same state. Thread B can make the state happen, and then thread B would wake up one or more of those threads in the queue, including A.

Syntax and Usage

In pthread library, `pthreadcond_t name;` can be used to create condition variable.

And a condition variable has two operations: `wait()` and `signal.`

```
1 pthread_cond_wait(pthread_cond_t *c, pthread_mutex_t *m);
2 pthread_cond_signal(pthread_cond_t *c);
```

`wait()` is like "Hey, I need something and put me into sleep until that thing is available."

`signal()` is like "Ok, I make the thing you wanted, and I am going to wake you up."

Note that `wait()` takes a lock as second parameter, and that is because it assumes the lock is held by the caller when calling `wait()` and it would release the lock and be put into sleep (automatically). The `wait()` won't return unless the lock is reacquired by the caller.

```
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
```

But, why! Why! Why we need the state variable?

```
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 }
```

Imagine, child thread executes upon the creation. It would signal when the queue is empty and exit. When the control is returned to the parent, parent would still call `wait()` and no one is going to wake it up. A **If** instead a **While** would probably makes much sense.

TIP: Always Hold the Lock While Signaling

Mandatory: Hold the Lock While Calling Wait.

Producer/Consumer (Bounded Buffer)

Problem Statement: Producers generate data items and place them in the buffer; consumers grab said items from the buffer and consume them in some way.

A real world example: Web server; A producer puts HTTP requests into a work queue. While, the consumer threads

```

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 }
```

We use int for simplicity. This buffer only has one slot for item. `put()` should set `count` to 1 when the slot is empty and `get()` should set `count` to 0 when the slot is filled.

A broken approach:

```

1 int loops; // must initialize somewhere...
2 cond_t cond;
3 mutex_t mutex;
4
5 void *producer(void *arg) {
6     int i;
7     for (i = 0; i < loops; i++){
8         Pthread_mutex_lock(&mutex);           // p1
9         if (count == 1)                      // p2
10            Pthread_cond_wait(&cond, &mutex); // p3
11         put(i);                          // p4
12         Pthread_cond_signal(&cond);       // p5
13         Pthread_mutex_unlock(&mutex);    // p6
14     }
15 }
16
17 void *consumer(void *arg) {
18     int i;
19     for(i = 0; i < loops; i++) {
20         Pthread_mutex_lock(&mutex);           // c1
21         if (count == 0)                      // c2
22            Pthread_cond_wait(&cond, &mutex); // c3
23         int tmp = get();                  // c4
24         Pthread_cond_signal(&cond);       // c5
25         Pthread_mutex_unlock(&mutex);    // c6
26         printf("%d\n", tmp);
27     }
28 }
```

This works only when there is exactly one producer and one consumer, since the condition variable is set in a way that producer and consumer take turns. However, this is not useful anyway. Suppose there are two consumer c1,c2 and one producer p1.

Upon the creation of the c1, it runs and check the buffer is empty it waits. Then p1 runs and put item into the buffer and it call signal. Now c2 passed the if statement, and both c1,c2 call `get()`.

You see here, the signal only tells at the moment the world is desired and you can wake up. But when

waiting one wakes up, it could find the work has already changed once again, i.e, not desired.

A "better" broken approach: use while instead of if

```

1 int loops; // must initialize somewhere...
2 cond_t cond;
3 mutex_t mutex;
4
5 void *producer(void *arg) {
6     int i;
7     for (i = 0; i < loops; i++) {
8         Pthread_mutex_lock(&mutex);           // p1
9         if (count == 1)                      // p2
10            Pthread_cond_wait(&cond, &mutex); // p3
11            put(i);                         // p4
12            Pthread_cond_signal(&cond);      // p5
13            Pthread_mutex_unlock(&mutex);    // p6
14    }
15 }
16
17 void *consumer(void *arg) {
18     int i;
19     for(i = 0; i < loops; i++) {
20         Pthread_mutex_lock(&mutex);           // c1
21         while (count == 0)                  // c2
22             Pthread_cond_wait(&cond, &mutex); // c3
23         int tmp = get();                  // c4
24         Pthread_cond_signal(&cond);        // c5
25         Pthread_mutex_unlock(&mutex);      // c6
26         printf("%d\n", tmp);
27     }
28 }
```

This would address the issue in last approach, c1 and c2 can safely consume the product. When c1 wakes up, it checks whether the world is desired again (because of while).

This "always use while loops" is called *Mesa semantics*.

However, the above approach is still broken due to the nondeterministic signaling.

1. both c1 and c2 sleep when they find there is item in the buffer.
2. p1 puts the item in the buffer and signal one consumer to wake up, and go to sleep
3. assume c1 wakes up, and cosumed the one item.
4. c1 signals on the condition variable, and went to sleep.
5. c2 wakes up, and check the buffer is empty, back to sleep.
6. all p1,c1,c2 are sleeping

This is due to there is one condition variable, and the signaling is nondeterministic. A consumer should never wakes up a consumer, and a producer should never wake up a producer.

Correct solution:

Producer wakes up consumer, and consumer wakes producer. Thus, two condition variables needed.

```

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 }
28

```

Note that two condition variables are in using: *empty* and *fill*, if is used to replace while in producer part.

When producer filled the buffer, it signals full "Hey, we have item in the buffer". And it waits for empty "Well, call me when the buffer is empty".

When consumer emptied the buffer, it signals empty "Hey, producer. The buffer is empty, please fill". And it waits for full "Also, call me when it is filled, I need more".

General solution:

We let producer produces multiple items when wakes up. And let consumer consumes multiple items when wakes up.

With single pair of producer and consumer, this approach is more efficient as it reduces context switches; with multiple producers or consumers, it even allows concurrent producing or consuming to take place.

```

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

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);
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 }
28

```

Covering Conditions: Broadcast

```

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 * allocate(int size) {
9     Pthread_mutex_lock(&m);
10    while (bytesLeft < size)
11        Pthread_cond_wait(&c, &m);
12    void *ptr = ...; // get mem from heap
13    bytesLeft -= size;
14    Pthread_mutex_unlock(&m);
15    return ptr;
16 }
17
18 void free(void *ptr, int size) {
19     Pthread_mutex_lock(&m);
20     bytesLeft += size;
21     Pthread_cond_signal(&c); // whom to signal??
22     Pthread_mutex_unlock(&m);
23 }

```

Problem Scenario:

There are two threads Ta,Tb. Ta asks for 10 bytes and Tb asks for 100 bytes. At the creations of Ta and Tb, there is 0 bytes free memory. Therefore Ta and Tb sleep. When Tf frees up 50 bytes, Tf signaled Tb which needs 100 bytes. Tb would sleep again, and no memory would be allocated even Ta's request can be fulfilled.

Solution:

Instead of *pthreadcondsignal()*, we can use *pthreadcondbroadcast*.

However, the cost is very noticeable since all threads would be waken up.

2.1.9 Semaphore

Definition:

A semaphore is an object with an integer value that we can manipulate with two routines; in POSIX, these routines are *sem_wait()* and *sem_post()*.

The initial value of the semaphore determines its behavior, it is necessary to initialize the semaphore to some value.

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

After initialize the semaphore, one can interact with it using *sem_wait()* and *sem_post*.

```
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 }
5
6 int sem_post(sem_t *s) {
7     increment the value of semaphore s by one
8     if there are one or more threads waiting, wake one
9 }
10
```

The value of the semaphore is equal to the number of, when negative,waiting threads.

And the initial value is the number of resources that can be give away at start.

Binary Semaphore (Locks)

```
1 sem_t m;
2 sem_init(&m, 0, X); // i
3
4 sem_wait(&m);
5 // critical section here
6 sem_post(&m);
```

This is essential just a lock, when semaphore is initialized to 1.

However, this is useless since the lock already serves the purpose.

Producer/Consumer (Bound Buffer):Semaphore approach

Note that semaphore works like a combination of lock and condition variable.

```
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 }
15
```

```

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
23 int main(int argc, char *argv[]) {
24     // ...
25     sem_init(&empty, 0, MAX); // MAX are empty
26     sem_init(&full, 0, 0);   // 0 are full
27     // ...
28 }
```

However, there would be a race condition at the put() and get():

1. Pa gets to run get() first, it first puts item in slot i.
2. just before Pa increment the slot counter, it gets interrupted.
3. Pb still have slot counter at slot i, and Pb overwrites slot i.

That is , there is no mutual exclusion for put() and get(), but only signaling whether one can put or get.

```

1 void *producer(void *arg) {
2     int i;
3     for (i = 0; i < loops; i++) {
4         sem_wait(&mutex); // Line P0 (NEW LINE)
5         sem_wait(&empty); // Line P1
6         put(i);           // Line P2
7         sem_post(&full); // Line P3
8         sem_post(&mutex); // Line P4 (NEW LINE)
9     }
10 }
11 void *consumer(void *arg) {
12     int i;
13     for (i = 0; i < loops; i++) {
14         sem_wait(&mutex); // Line C0 (NEW LINE)
15         sem_wait(&full); // Line C1
16         int tmp = get(); // Line C2
17         sem_post(&empty); // Line C3
18         sem_post(&mutex); // Line C4 (NEW LINE)
19         printf("%d\n", tmp);
20     }
21 }
22
```

Now, the revised solution use semaphore (like a lock) to add mutual exclusion.

However, it still doesn't work because deadlock since semaphore don't increment itself when sleep

The consumers and producers share a single lock.

1. Consumer C1 runs first, the use sem-wait(&mutex) to required the lock, mutex = 0.
2. C1 continue, and calls sem-wait(&full), however there is no item in the slot. And C1 sleeps.
3. When any other producers or consumers try to do something, they don't have the lock, any sem-wait(&mutex) would decrease the mutex to negative value and sleeps forever.

```

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

```

Now a truly working vision,

To resolve the deadlock, simply reduce the scope of lock at this case.

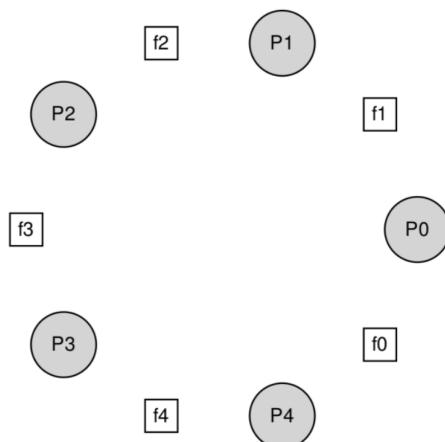
Reader-Writer Locks

For concurrent list operations, we need to guarantee that no two insertion races, but it is ok for multiple readers to lookup concurrently.

reader-write lock: The implementation is in reader-writer.cpp

Dining Philosophers

Problem statement: There are five “philosophers” sitting around a table. Between each pair of philosophers is a single fork (and thus, 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



Here is the behavior for a philosopher:

```
while (1) {
    think();
    get_forks(p);
    eat();
    put_forks(p);
}
```

Some help functions needed

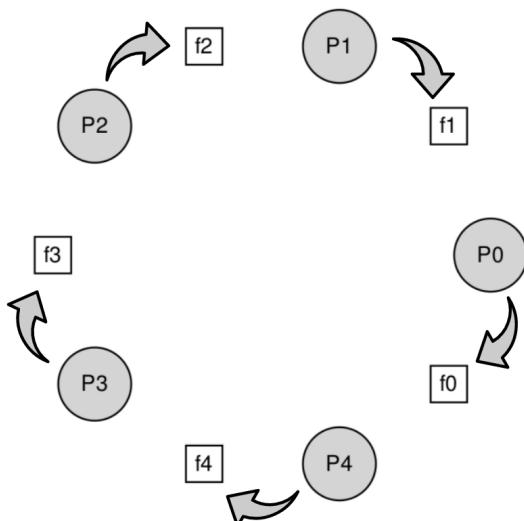
```
int left(int p){
    return p;
}

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

Broken solution:

```
// potentially deadlock: if every philosopher grabs the left one
// that is , all threads makes to sem_wait(&forks[ left(p)]) and no further
void get_forks(int p){
    sem_wait(&forks[left(p)]);
    sem_wait(&forks[right(p)]);
}

void put_forks(int p){
    sem_post(&forks[left(p)]);
    sem_post(&forks[right(p)]);
}
```



Dijkstra solution: break the dependency

Make one philosopher behaves different than other ones.

```
void get_forks(int p){
    // letting philosopher 4 behaves differently , there would be deadlock
    // since there won't be a circle now.
    if(p==4){
        sem_wait(&forks[right(p)]);
```

```

        sem_wait(&forks [ left (p) ]);
    } else {
        sem_wait(&forks [ left (p) ]);
        sem_wait(&forks [ right (p) ]);
    }
}

```

Thread Throttling

How can a programmer prevent "too many" threads from doing something at once and bogging the system down?

Answer is **Throttling**: decide a threshold, use semaphore to limit the number of threads concurrently executing.

How implement semaphore

```

typedef struct Semaphore{
    int value;
    pthread_cond_t cond;
    pthread_mutex_t lock;
}Semaphore;

void sem_init(sem *s, int value){
    s-> value = value;
    cond_init(&s->cond);
    mutex_init(&s->lock);
}

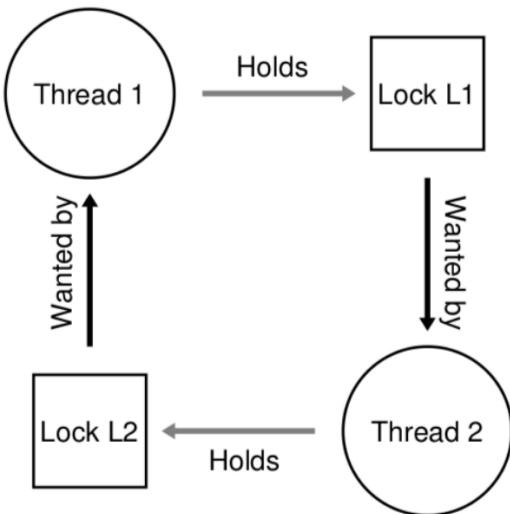
void sem_wait(sem *s){
    pthread_mutex_lock(&s->lock);

    // however, this value would never go negative
    // easy to implement, linux's approach
    while(s->value <=0){
        pthread_cond_wait(&s->cond,&s->lock );
    }
    s->value = s->value - 1;
    pthread_mutex_unlock(&s->lock );
}

void sem_post(sem *s){
    pthread_mutex_lock(&s->lock );
    s->value = s->value - 1;
    pthread_cond_signal(&s->cond);
    pthread_mutex_unlock(&s->lock );
}

```

2.1.10 Deadlock bugs



Four conditions need to be hold for a deadlock to occur:

1. **Mutual exclusion:** Threads claim exclusive control of resources.
2. **Hold-and-wait:** Threads hold the resource and wait for other necessary resource, wont let go.
3. **No preemption:** Resources cannot be forcibly removed from the threads that are holding them.
4. **Circular wait:** There exists a circular chain of threads such that each threads holds one or more resources that are being requested by the next thread in the chain.

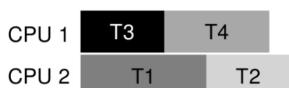
Prevention

Make one of the four condition false all the time to prevent deadlock. However, this is not practical.

Deadlock avoidance via Scheduling

Instead of deadlock prevention, in some scenarios deadlock avoidance is preferable. Avoidance requires some global knowledge of which locks various threads might grab during their execution, and subsequently, schedules said threads in a way as to guarantee no deadlock can occur.

	T1	T2	T3	T4
L1	yes	yes	no	no
L2	yes	yes	yes	no



From the images, it is possible to schedule threads to avoid deadlock
if we need the resources that each thread wants before hand.

	T1	T2	T3	T4
L1	yes	yes	yes	no
L2	yes	yes	yes	no



Also, the scheduler could procure length-considerable scheduling.

Thus, this approach can be only used in limited environment, for example, in an embedded system where one has full knowledge of the entire set of tasks must be run and the locks that they need.

Also, such approach can limit the degree of concurrency.

Detection And Recovery

One strategy can allow deadlock to occur occasionally, and then take some action once it has been detected.

A deadlock detector runs periodically, building a resource graph and checking it for cycles.

2.1.11 Event-Based Concurrency

Why event-based concurrency:

1. Managing concurrency correctly in multi-threaded applications are challenging.
2. Multi-threaded application, the developer has little or no control over what is scheduled at a given moment in time.

Given the difficulty of building a general-purpose scheduler that works well in all cases for all workloads, one cannot expect the scheduler schedule optimally everytime. Thus, we have **event-based concurrency** to address the problem.

The basic idea: An Event Loop

One can simply wait for something to occur; when it does, you check what type of event it is and do the small amount of work it requires, which may include issuing I/O requests, or scheduling other events for future handling.

An example can be event-based server. Such applications are based around on an **event loop**.

```
while(1){
    events = getEvents();
    for( each event in events){
        processEvent(e); // event handler
    }
}
```

When event handler processes an event, it is the only work being done. It is just like a scheduler schedule the event. Thus, we can decide which event to handle next to take over control on which work should be scheduled next.

But how? How can the event-based server determines which events are taking place?

Important API: select() or poll()

In most systems, an API either called select() or poll() system calls to receive events.

The API enables a program to check whether there is any incoming I/O that should be attended to.
select()

```
int select( int nfds ,
            fd_set *restrict readfds ,
            fd_set *restrict writefds ,
            fd_set *restrict errorfds ,
            struct timeval *restrict timeout );
```

select() lets you check whether descriptors can be read from or written to. The former lets a server determine that a new packet has arrived. The latter lets the service know when to reply(i.e, the outbound queue is not full).

Also, there is "timeout" interval. Common usage is to set it to NULL, so that select() would be blocked indefinitely, until a descriptor is ready; another common usage is to set it to 0, so that select() return immediately.

Blocking vs Non-Blocking interfaces

Blocking(synchronous) interfaces do all the work before returning to the caller; non-blocking(asynchronous) interfaces begin some work but return immediately, and let the remaining work done in background.

Non-blocking interfaces are essential in the event-based approach.

Problem: Blocking System Calls

There is one issue with event-based programming: issue blocking system call.

Without the nature of overlapping I/O and other computations that multi-threaded programming has, event-based systems can only block until the call completes, i.e, no solution for the problem.

Thus, we make a rule for event-based systems : no blocking calls are allowed.

Asynchronous I/O

To overcome the issue of blocking calls in event-based system, modern operating system have introduced new ways to issue I/O requests to the disk system, referred to generically as asynchronous I/O.

The interfaces allow issuing I/O request and return immediately. Additional interfaces are also available for the system to decide whether an I/O request is fulfilled.

```
struct aiocb { // AIO control block
    int             aio_fildes;      /* File descriptor */
    off_t           aio_offset;     /* File offset */
    volatile void * aio_buf;       /* Location of buffer */
    size_t          aio_nbytes;    /* Length of transfer */
    int             aio_reqprio;   /* Request priority offset */
    struct sigevent aio_sigevent; /* Signal number and value */
    int             aio_lio_opcode; /* Operation to be performed */
};

int aio_read(struct aiocb *aiocbp);

// check whether the request referred to by aiocbp has completed.
// If yes, 0 is returned; otherwise, EINPROGRESS is returned.
int aio_error(const struct aiocb *aiocbp);
```

Problem: State Management

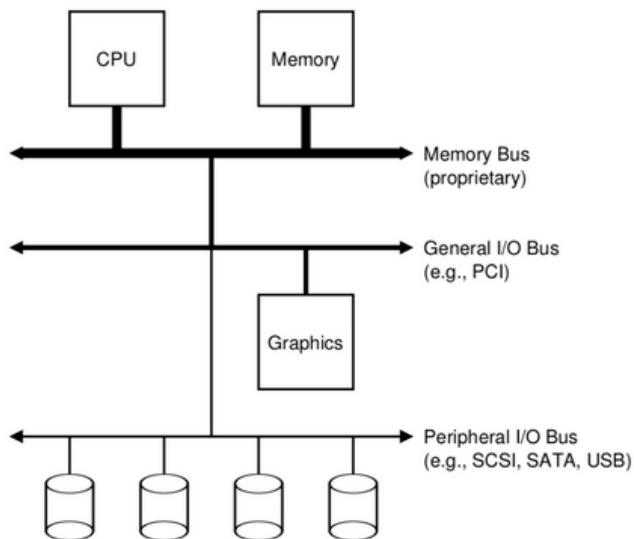
Another problem with event-based approach is that such code is commonly more complicated to write. For example, when an event handler issues an asynchronous I/O, it must package up some program state for the next event handler to use when the I/O finally completes.

Chapter 3

Persistence

3.1 I/O devices

Prototypical System Architecture

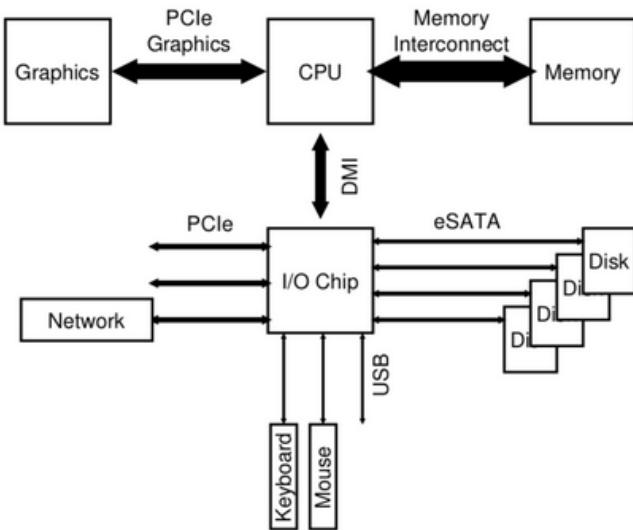


CPU and Memory are connected to the **Memory Bus**. Other devices like graphic card are connected to the system via a general **I/O bus**, which should be PCI(or its derivatives). There is even lower bus called **peripheral bus**, such that **SCSI,SATA** or **USB** can be attached here. Those connect slow devices to the system, including disks, mice and keyboards.

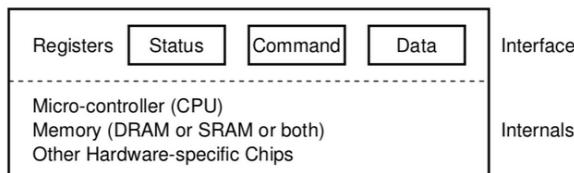
The idea is that: components that demand high performance are nearer the CPU. Lower performance components are further away.

Modern System Architecture

Modern systems use specialized chipsets and faster point-to-point interconnects to improve performance.



Canonical Device



First, each device has its interface to communicate with other device. The interface includes Registers, Buffers.

Second, each device has internal processors/chips to do the processing. Some complex devices may even have CPU and RAM. This internal is also referred as internal.

Cononical Protocol

1. Status Register: indicates the current status of the device.
2. Command Register: to tell device to perform a certain task.
3. Data Register: pass data to the device or get data from the device.

An example of interaction between OS and the device

```

While (STATUS==BUSY)
; // wait until the device is not busy

Write data to DATA register
Write command to COMMAND register
    (starts the device and executes the command)

While (STATUS==BUSY)
; // wait until device is done with your request

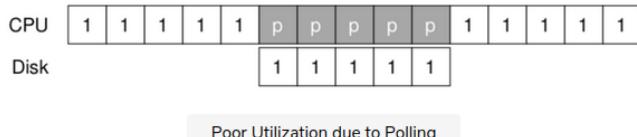
```

Although this approach works correctly, it is inefficient. It wastes a great deal of CPU time waiting for the STATUS becomes not busy.

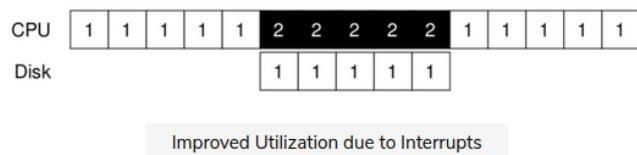
Lower CPU overhead with Interrupts

Instead of polling the device repeatedly, the OS can issue an request, put the calling process to sleep, and context switch to another task. When the device fulfills the request, it raise hardware interrupts to tell the job is done.

Here is a comparison between polling approach and interrupts approach.



Interrupt approach allows for overlap.



Problem with Interrupts

For example, imagine a device that performs its tasks very quickly: the first poll usually finds the device to be done with a task. Using an interrupt, in this case, will actually slow down the system because switching to another process, handling the interrupt, and switching back to the issuing process is expensive.

Thus, if a device is fast, it may be best to poll. If it is slow, interrupts, which allow overlap, are best.

If the speed of the device is not known, it will be best to use a **hybrid(two-phased)** that polls for a little while and then, if the device is not yet finished, uses interrupts.

Also, in networking or other I/O intensive systems, it is better to avoid interrupts. When a huge stream of incoming packets each generate an interrupt, it is possible for the OS to **livelock**-only process interrupts.

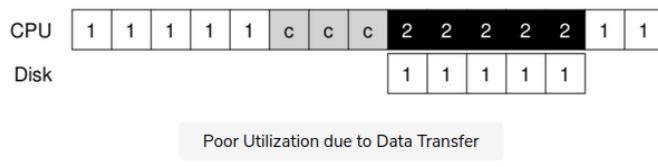
Coalescing

A device waits for a bit before delivering the interrupt to the CPU \Rightarrow collect more interrupts to handle in one delivery.

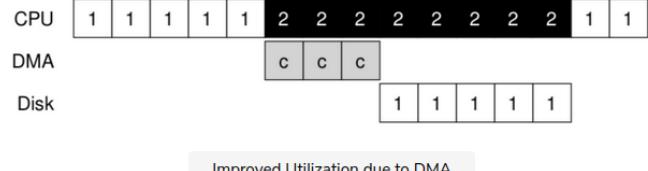
DMA: more efficient data movement

Suppose, a program wants to write data into disk, it has to copy data from memory to disk DATA explicitly, one word at a time.

The writing is marked as 'c'.



However, with Direct Memory Access(DMA). The CPU can tell DMA engine where the data exists in the memory, how much data to copy, and which device to send it to. DMA controller will do all the write, and issues interrupt when finished.



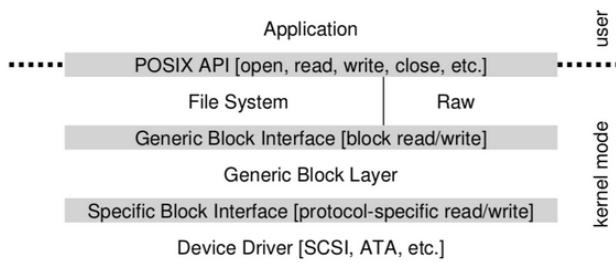
Explicit I/O instructions

On x86, **in** and **out** instructions can be used to communicate with specified devices. Such instructions are usually **privileged**.

Memory-mapped I/O

With this approach, the hardware makes device registers available as if they were memory locations. To access a particular register, the OS issues a load or store to the address.

Device Driver



Device driver provide abstraction, and OS can communicate with different devices without knowing detailed implementation.

Problem with device driver

1. rich features lost in the encapsulation, since uniform is required.
2. 70 percent of linux kernel reveal is device drivers.
3. drivers are written by "amateurs"

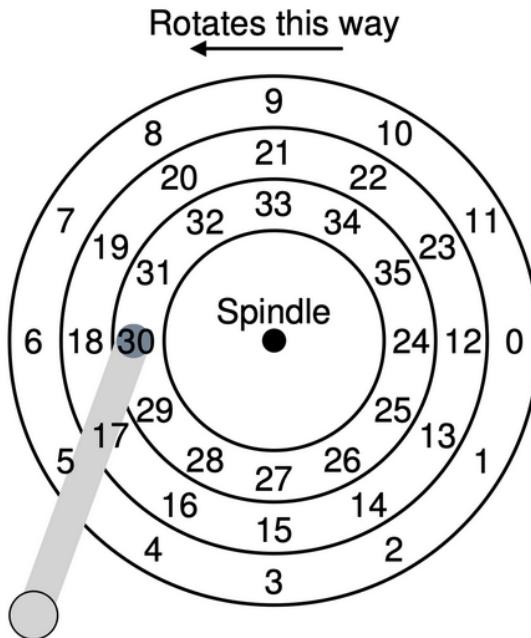
3.2 Hard Disk Drives

Interface

The drive consists of a large number of sections (512-byte blocks), each can be written or read. Sectors are numbered from 0 to n-1 on a disk with n sectors. Thus, the disk can be viewed as an array of sectors. Disk [n][512].

One can assume that accessing blocks in a contiguous chunk is fastest access mode.

Seek Time



Seek: move the disk arm to the correct track, very expensive.
A seek has many phases

1. acceleration phase as the disk arm gets moving.
2. coasting as the arm is moving at full speed.
3. deceleration as the arm slows down.
4. settling as the head is positioned over the correct track.

The **settling time** is quite significant, e.g., 0.5 to 2ms.

After seeking, we have to wait the sector to pass the head, that is **rotation delay**.

HDD I/O time

$$T_{I/O} = T_{seek} + T_{rotation} + T_{transfer}$$

Disk scheduling

Unlike CPU scheduling, the job size is known in the disk scheduling. Thus, the scheduler would try to follow the **principle of SJF** in its operation.

1. SSTF: shortest seek time first. Potential starvation and seek changes when fulfill a request (NBF:nearest-block-first is a makeup).
2. Elevator(SCAN/C-SCAN): Move head back and forth across the disk, servicing requests in order across the tracks.
3. SPTF: shortest positioning time first. This depends on the seek time and rotation time, if the seek time is faster, it favors less rotation. Otherwise, it favors less seeks.

3.2.1 RAID

RAID: Redundant Array of Inexpensive Disks.

Externally, RAID looks just like a disk. Internally, the RAID is complex, consisting of multiple disks, memory, and some processors to manage them.

Advantages of RAID

1. Performance: use multiple disk in parallel can speed up I/O
2. Capacity
3. Reliability: with some redundancy, it would prevent a single point failure.
4. **Transparency:** the system sees the RAID as a gigant disk.

Fault model

RAID is designed to detect and recover from certain kinds of disk faults.

The fail-stop fault model: in this model, disk can either be working or failed. When a disk has failed, we assume it is permanently lost.

Assume failed disk can be easily detected.

Evaluation of RAID

Three critera

1. Capacity: given N disks with B blocks, how much space is available to clients?
2. Reliability: how many disk faults that a design can tolerate?
3. Performance: before evaluating performance, we present a set of typical workloads that one should consider.

Three important RAID designs:

1. RAID LEVEL 0 (striping)
2. RAID LEVEL 1 (mirroring)
3. RAID LEVEL 4/5 (parity-based redundancy)

RAID LEVEL 0

There is no redundancy at level 0(striping), simply round robin fashion to spread the blocks across the disks.

Disk 0	Disk 1	Disk 2	Disk 3
0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

Evaluation:

1. Perfomance: it is designed to extract most parallelism from array. Sequential workload is way higher than random workload.
2. Capacity: striping is perfect at this case, 100 percent capacity can be used.
- 3.

Throughput : N^*S and N^*R

RAID LEVEL 1

This is mirroring and get a copy of all the data (in separate disks) to tolerate disk fault.

Disk 0	Disk 1	Disk 2	Disk 3
0	0	1	1
2	2	3	3
4	4	5	5
6	6	7	7

For read, mirroring can read either copy, but for write, it must update contents in both copies.

1. Capacity: The usable capacity for client is $(N*B)/2$
2. Reliability: This can tolerate 1 disk failure for certain, and $N/2$ disk failure by chance.
- 3.

Throughput : $N*S/2$, $N*R/2$

RAID LEVEL 4

Use **parity**, we can ensure reliability without huge capacity penalty as mirroring has.

Disk 0	Disk 1	Disk 2	Disk 3	Disk 4
0	1	2	3	P0
4	5	6	7	P1
8	9	10	11	P2
12	13	14	15	P3

Use an extra disk 5 to store the parity which is calculated from the striping of disk 1-4.

It turned out that XOR is excellent parity. Invariant: $C1 \text{ xor } C2 \text{ xor } C3 \text{ xor } C4 = P$, reconstruction: $P \text{ xor } C1 \text{ xor } C2 \text{ xor } C3 = C4$.

1. Capacity: $(N-1)*B$
2. Reliability: It only tolerate a single disk failure.
3. Performance: $(N-1)*R$, $(N-1)*S$

However, to maintain the parity is very costly and it is the bottleneck. Yes, we can write parallelly in the disk 1-4, but we cannot update the parity parallelly.

RAID LEVEL 5

We can use **Rotating Parity** to address the writing to parity issue.

Disk 0	Disk 1	Disk 2	Disk 3	Disk 4
0	1	2	3	P0
5	6	7	P1	4
10	11	P2	8	9
15	P3	12	13	14
P4	16	17	18	19

3.3 Files and Directories

File: linear array of bytes that one can read or write. Each file has a low-level name called inode for historical reason.

Directory: directory also has a low-level name, and it contains a list of pairs $\{\text{file name}, \text{inode}\}$, $\{\text{directory name}, \text{inode}\}$ (we can have a directory tree/hierarchy at this case).

The File System Interface

```
int fd = open("foo", O_CREAT|O_WRONLY|O_TRUNC,
              S_IRUSR|S_IWUSR);

// returns a file descriptor.
```

A file descriptor is just an integer, private per process, and is used in UNIX systems to access file. Think it as a capacity, a handle that gives you the power to perform certain operations. Or a pointer to an object of type file.

File Descriptors

File descriptors 0,1,2 are taken

1. standard input
2. standard output
3. standard error

If a process opens first file, it would certainly be file descriptor 3.

Each process maintains an array of file descriptors, each of which refers to an entry in the system-wide open file table. Each entry in this table tracks which underlying file the descriptor refers to, the current offset, and other relevant details such as whether the file is readable or writable.

Read and Write, non-sequentially

Each read and write implicitly updates the offset.

System Calls	Return Code	Current Offset
fd = open("file", O_RDONLY);	3	0
read(fd, buffer, 100);	100	100
read(fd, buffer, 100);	100	200
read(fd, buffer, 100);	100	300
read(fd, buffer, 100);	0	300
close(fd);	0	-

lseek() explicitly change the offset.

System Calls	Return Code	Current Offset
fd = open("file", O_RDONLY);	3	0
lseek(fd, 200, SEEK_SET);	200	200
read(fd, buffer, 50);	50	250
close(fd);	0	-

```
// only change the offset for the next I/O operation
off_t lseek(int fildes, off_t offset, int whence);

/*
If whence is SEEK_SET, the offset is set to offset bytes.
If whence is SEEK_CUR, the offset is set to its current
location plus offset bytes.
If whence is SEEK_END, the offset is set to the size of
the file plus offset bytes.
*/
```

```
// stores in the process's structure
// file table entry
struct file {
    int ref;
    char readable;
    char writable;
```

```

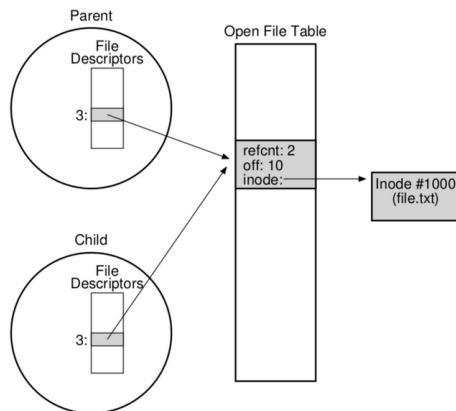
    struct inode *ip;
    uint off;
};

// open file table
struct {
    struct spinlock lock;
    struct file file[NFILE];
} ftable;

```

Shared File Table Entries: fork() and dup()

With fork () the child and parent share the same open file table entry,s since they shared the same file descriptor.



Note that there is an **reference count** in the entry - refcnt. When a file table entry is shared, its reference count is incremented, only all processes close the file, the entry would be removed.

Also, there is **dup()** which allows a process to create a new file descriptor that refers to the same underlying open file as an existing descriptor. Basically copy the file descriptor.

Writing Immediately with fsync()

write() only tells the file system to write the data into disk sometimes in the future. The system would buffer those writes in memory for performance reason.

For some systems, it is ideal to write data immediately. fsync() is such API that supports the functionality. The fsync() routine returns once all of these writes are complete.

```

int fd = open("foo", O_CREAT|O_WRONLY|O_TRUNC,
              S_IRUSR|S_IWUSR);
assert(fd > -1);
int rc = write(fd, buffer, size);
assert(rc == size);
rc = fsync(fd);
assert(rc == 0);

```

The above sequence doesn't guarantee everything that one might expect; in some cases, it is required to fsync() the directory that contains foo. Because foo might be the newly created file, without this additional fsync(), this might lead to application-level bugs.

Renaming Files

```

int fd = open("foo.txt.tmp", O_WRONLY|O_CREAT|O_TRUNC,
              S_IRUSR|S_IWUSR);
write(fd, buffer, size); // write out new version of file

```

```
fsync(fd);
close(fd);

// rename the file
rename("foo.txt.tmp", "foo.txt");
```

First, rename() is atomic call with respect to system crashes. If the system crashed, it would either be renamed the old name or new name.

Metadata of Files

```
struct stat {
    dev_t      st_dev;        //ID of device containing file
    ino_t      st_ino;        //inode number
    mode_t     st_mode;       //protection
    nlink_t    st_nlink;      //number of hard links
    uid_t      st_uid;        //user ID of owner
    gid_t      st_gid;        //group ID of owner
    dev_t      st_rdev;       //device ID (if special file)
    off_t      st_size;       //total size, in bytes
    blksize_t  st_blksize;    // blocksize for filesystem I/O
    blkcnt_t   st_blocks;     // number of blocks allocated
    time_t     st_atime;      //time of last access
    time_t     st_mtime;      //time of last modification
    time_t     st_ctime;      //time of last status change
};
```

You should just think of an inode as a persistent data structure kept by the file system that has information like we see above inside of it.

All inode reside on disk; a copy of active ones are usually cached in the memory to speed up access.

Removing Files

unlink() is the API to remove a file. But why doesn't it called remove() or delete().

Making Directories

mkdir

Reading and Deleting Directories

opendir(), readdir(), closedir().

For each directory entry in the readdir(dir).

```
struct dirent {
    char d_name[256];          // filename
    ino_t d_ino;                // inode number
    off_t d_off;                // offset to the next dirent
    unsigned short d_reclen;    // length of this record
    unsigned char d_type;       // type of file
};
```

rmdir() has the requirement that the directory be empty.

Hard Link

link() creates another name in the directory you are creating the link to, and refers it to the same inode number of the original file. It only gives you another file name that refers to the same file.

unlink() unlinks a file with the inode, and decrease the reference count(link count) within the inode. When the reference count reaches zero does the file system also free the inode and related data blocks, and thus truly "deleted" the file.

Symbolic Link

Hard links are limited because you cannot create one to a directory.
 Symbolic points to the pathname of file.

Permission Bits and Access Control Lists

Unix has the permission bits like "-rw-r-r-". (owner, group, other)
 Access Control Lists specify a list of who can and cannot read a set of files.

Making and Mounting a File System

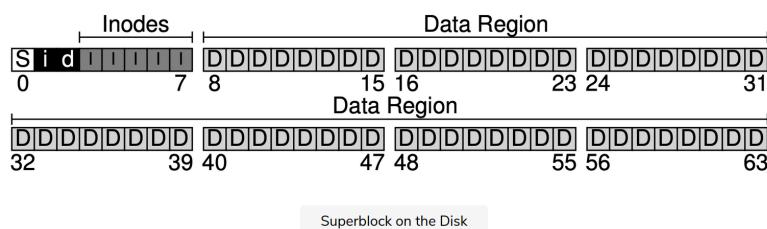
To make a file system, most file systems provide a tool, usually referred to as mkfs (pronounced "make fs"), that performs exactly this task.

However, once such a file system is created, it needs to be made accessible within the uniform file-system tree. This task is achieved via the mount program (which makes the underlying system call mount() to do the real work). What mount does, quite simply is take an existing directory as a target mount point and essentially paste a new file system onto the directory tree at that point.

3.3.1 File System Implementation

A unreal file system that services for introduction purpose: VSFS (very simple file system).
 To understand a file system, think about two different aspects of them:

1. Data Structure: use arrays of blocks and other objects.
2. Access Method: open(), read() and write()



Blocks

First, divide the disk into blocks. The commonly-used size is 4 KB. Assume 64-bit block here.

Data region

The portion that actually stores the user data.

Inodes(I)

There is inode table such that lists all inodes in the filesystem.

Allocation Structures

There is some way to track whether inodes or data blocks are free or allocated.

1. free list
2. bitmap: used in the graph

Superblock

The superblock contains information about the particular file system, including, for example, how many inodes and data blocks are in the file system.

When the system is mounted, the OS would read the superblock to initialize various parameters, and then attach the volume to the file-system tree.

3.3.2 The Inode

INode is short for **index node**, the historical name given to it in UNIX, and possibly earlier systems, used because these nodes were originally arranged in an array, and the array indexed into when accessing a particular inode.

Each inode is implicitly referred to by a number (called **i-number**),

The Inode Table (Closeup)								
	iblock 0	iblock 1	iblock 2	iblock 3	iblock 4			
Super	i-bmap	d-bmap	0 1 2 3 16 17 18 19 32 33 34 35 48 49 50 51 64 65 66 67	4 5 6 7 20 21 22 23 36 37 38 39 52 53 54 55 68 69 70 71	8 9 10 11 24 25 26 27 40 41 42 43 56 57 58 59 72 73 74 75	12 13 14 15 28 29 30 31 44 45 46 47 60 61 62 63 76 77 78 79		
0KB	4KB	8KB	12KB	16KB	20KB	24KB	28KB	32KB

Inside each inode is virtually all of the **metadata** you need about a file:

1. file type
2. file size
3. the number of blocks allocated to it
4. protection information
5. timestamps (when the file was created, modified or last accessed)
6. the location on the disk

Size	Name	What is this inode field for?
2	mode	can this file be read/written/executed?
2	uid	who owns this file?
4	size	how many bytes are in this file?
4	time	what time was this file last accessed?
4	ctime	what time was this file created?
4	mtime	what time was this file last modified?
4	dtime	what time was this inode deleted?
2	gid	which group does this file belong to?
2	links_count	how many hard links are there to this file?
4	blocks	how many blocks have been allocated to this file?
4	flags	how should ext2 use this inode?
4	osd1	an OS-dependent field
60	block	a set of disk pointers (15 total)
4	generation	file version (used by NFS)
4	file_acl	a new permissions model beyond mode bits
4	dir_acl	called access control lists

Simplified Ext2 Inode

Multi-level index

The inode could contain **direct pointers** to the data blocks. However, this approach is very limited since the file could be very large, and it requires a gigantic table of direct pointers to refer it.

To resolve this, **indirect pointers**(a pointer to a table of pointers) are used to implement multi-level indexing.

