Design philosophy of operating systems (IV)

Hung-Wei Tseng



Recap: Each process has a separate virtual memory space





Virtually, every process seems to have a processor, but only a few of them are physically executing.

Recap: The basic process API of UNIX

- fork
- wait
- exec
- exit



Recap: How to implement redirection in shell

- Say, we want to do ./a > b.txt
- fork
- The forked code opens b.txt
- The forked code dup the file descriptor to stdin/stdout
- The forked code closes b.txt
- exec("./a", NULL)

code int pid, fd; char cmd[2048], prompt = "myshell\$" while(gets(cmd) != NULL) { ((pid = fork()) == 0) { fd = open("b.txt", O_RDWR | O_CREAT, S_IRUSR | S_IWUSR); dup2(fd, stdout); close(fd); execv("./a",NULL); else printf("%s ",prompt); The shell can respond to next input static data heap stack

Homework for you: Think about the case when your fork is equivalent to fork+exec()

char cmd[2048], prompt = "myshell\$" while(gets(cmd) != NULL) { if ((pid = fork()) == 0) {		
<pre>fd = open("b.txt", O_RDWR O_CREA S_IWUSR); dup2(fd, stdout); close(fd); execv("./a",NULL); } else printf("%s ",prompt); }</pre>		
static data		





- The hardware is changing
 - Multiprocessors
 - Networked computing
- The software

be built and future development of UNIX-like systems for new architectures can continue. The computing environment for which Mach is targeted spans a wide class of systems, providing basic support for large, general purpose multiprocessors, smaller multiprocessor networks and individual workstations (see

- The demand of extending an OS easily
- Repetitive but confusing mechanisms for similar stuffs

As the complexity of distributed environments and multiprocessor architectures increases, it becomes increasingly important to return to the original UNIX model of consistent interfaces to system facilities. Moreover, there is a clear need to allow the underlying system to be transparently extended to allow user-state processes to provide services which in the past could only be fully integrated into UNIX by adding code to the operating system kernel.

Make UNIX great again!



- Mach: A New Kernel Foundation For UNIX Development (cont.)
- Taxonomy of Kernels
- Synchronization

Mach: A New Kernel Foundation For UNIX Development

Mike Accetta, Robert Baron, William Bolosky, David Golub, Richard Rashid, Avadis Tevanian, **Michael Young Computer Science Department, Carnegie Mellon University**

Tasks/processes



Intel Sandy Bridge







Main memory is eventually shared among processor



Tasks/processes



Threads



Tasks/Processes and threads

- How many of the following regarding the comparison of parallelizing computation tasks using processes and threads is/are correct?

 - The context switch and creation overhead of processes is higher

 you have to change page tables, warm up TLBs, warm up caches, create a new memory space ...

 The overhead of exchanging data among different computing tasks for the same applications is higher in process model
 - ③ The demand of memory usage is higher when using processes
 - The security and isolation guarantees are better achieved using processes 4

- A. 0 B. 1
- C. 2
- D. 3



- separate address, it's not easy to access data from another process

Case study: Chrome v.s. Firefox





What's in the kernel?

- How many of the following Mach features/functions are implemented in the kernel?
 - ① I/O device drivers
 - ² File system
 - ③ Shell
 - ④ Virtual memory management
 - A. 0
 - B. 1
 - C. 2
 - D. 3

E. 4



What's in the kernel?

- How many of the following Mach features/functions are implemented in the kernel?
 - ① I/O device drivers
 - ² File system
 - ③ Shell
 - ④ Virtual memory management
 - A. 0
 - B. 1
 - C. 2
 - D. 3

E. 4



What's in the kernel?

- How many of the following Mach features/fun implemented in the kernel?
 - ① I/O device drivers
 - ² File system
 - ③ Shell
 - ④ Virtual memory management
 - A. 0
 - B. 1
 - C. 2
 - D. 3 E. 4

Policies

Mechanisms

Functionality:





Policy? Mechanisms?

- How many of the following terms belongs to "policies"?
 - Least-recently used (LRU)
 - ² First-in, first-out
 - ③ Paging
 - ④ Preemptive scheduling
 - ⑤ Capability
 - A. 0
 - B. 1
 - C. 2
 - D. 3

E. 4



Policy? Mechanisms?

- How many of the following terms belongs to "policies"?
 - Least-recently used (LRU)
 - ② First-in, first-out
 - ③ Paging
 - ④ Preemptive scheduling
 - ⑤ Capability
 - A. 0
 - B. 1
 - C. 2
 - D. 3

E. 4





Policy? Mechanisms?

- How many of the following terms belongs to "policies"?
 - ① Least-recently used (LRU)
 - ² First-in, first-out
 - ③ Paging
 - ④ Preemptive scheduling
 - ⑤ Capability
 - A. 0
 - B. 1
 - C. 2
 - D. 3
 - E. 4

— Policy

- Policy

- Mechanism
- Mechanism
- Mechanism



Whys v.s. whats

How many pairs of the "why" and the "what" in Mach are correct?

	Why	
(1)	Support for multiprocessors	Threads
(2)	Networked computing	Messages/Ports
(3)	OS Extensibility	Microkernel/Obj
(4)	Repetitive but confusing mechanisms	Messages/Ports
Α.	0	
В.	1	
C.	2	
D.	3	
Ε.	4	

What

ject-oriented design

Types of kernels

- What type of kernels does the UNIX described in Dennis M. Ritchie's paper belong to?
 - A. Microkernel the kernel only provides a minimal set of services/ mechanisms including memory management, multitasking and interprocess communication
 - B. Monolithic the kernel implements every function that cannot be in a user-space library: device drivers, scheduler, memory handling, file systems, network stacks
 - C. Modular the kernel provides a basic set of functions like microkernels, but allows load/unload kernel modules if necessary D. Layered kernel — the kernel follows strict layered design that lowerorder module cannot interact with higher-order modules

Types of kernels

- What type of kernels does the UNIX described in Dennis M. Ritchie's paper belong to?
 - A. Microkernel the kernel only provides a minimal set of services/ mechanisms including memory management, multitasking and interprocess communication
 - B. Monolithic the kernel implements every function that cannot be in a user-space library: device drivers, scheduler, memory handling, file systems, network stacks
 - C. Modular the kernel provides a basic set of functions like microkernels, but allows load/unload kernel modules if necessary D. Layered kernel — the kernel follows strict layered design that lowerorder module cannot interact with higher-order modules



Types of kernels

- What type of kernels does the UNIX described in Dennis M. Ritchie's paper belong to?
 - A. Microkernel the kernel only provides a minimal set of services/ mechanisms including memory management, multitasking and interprocess communication Mach, Nucleus
 - B. Monolithic the kernel implements every function that cannot be in a user-space library: device drivers, scheduler, memory handling, file systems, network stacks **Old UNIX**
 - C. Modular the kernel provides a basic set of functions like microkernels, but allows load/unload kernel modules if necessary Linux, Windows, MacOS, FreeBSD D. Layered kernel — the kernel follows strict layered design that lower-
 - order module cannot interact with higher-order modules THE

Types of Kernels



Only mechanisms are in the kernel

Original UNIX

Nucleus, Mach

Linux, Windows, **MacOS**

Why not microkernels?

- Although Mach's design strongly influenced modern operating systems, why most modern operating systems do not adopt the design of microkernels?
 - A. Microkernels are more difficult to extend than monolithic kernels
 - B. Microkernels are more difficult to maintain than monolithic kernels
 - C. Microkernels are less stable than monolithic kernels
 - D. Microkernels are not as competitive as monolithic kernels in terms of application performance
 - E. Microkernels are less flexible than monolithic kernels



Why not microkernels?

- Although Mach's design strongly influenced modern operating systems, why most modern operating systems do not adopt the design of microkernels?
 - A. Microkernels are more difficult to extend than monolithic kernels
 - B. Microkernels are more difficult to maintain than monolithic kernels
 - C. Microkernels are less stable than monolithic kernels
 - D. Microkernels are not as competitive as monolithic kernels in terms of application performance
 - E. Microkernels are less flexible than monolithic kernels





Why not microkernels?

- Although Mach's design strongly influenced modern operating systems, why most modern operating systems do not adopt the design of microkernels?
 - A. Microkernels are more difficult to extend than monolithic kernels
 - B. Microkernels are more difficult to maintain than monolithic kernels
 - C. Microkernels are less stable than monolithic kernels
 - D. Microkernels are not as competitive as monolithic kernels in terms of application performance
 - E. Microkernels are less flexible than monolithic kernels



Context switches!

The impact of Mach

- Threads
- Extensible operating system kernel design
- Strongly influenced modern operating systems
 - Windows NT/2000/XP/7/8/10
 - MacOS

С

developer.apple.com/library/archive/documentation/Darwin/Conceptual/KernelProgramming/Mach/Mach.html

Documentation Archive

Table of Contents

- About This Document
- Keep Out
- Kernel Architecture Overview
- The Early Boot Process
- Security Considerations
- Performance Considerations
- Kernel Programming Style
- Mach Overview
- Memory and Virtual Memory
- Mach Scheduling and Thread Interfaces
- Bootstrap Contexts
- I/O Kit Overview
- BSD Overview
- File Systems Overview Network Architecture
- Boundary Crossings
- Synchronization Primitives
- Miscellaneous Kernel Services
- Kernel Extension Overview
- Building and Debugging Kernels
- Bibliography Revision History Glossary

Mach Overview

The fundamental services and primitives of the OS X kernel are based on Mach 3.0. Apple has modified and extended Mach to better meet OS X functional and p Mach 3.0 was originally conceived as a simple, extensible, communications microkernel. It is capable of running as a stand-alone kernel, with other traditional o

networking stacks running as user-mode servers.

However, in OS X, Mach is linked with other kernel components into a single kernel address space. This is primarily for performance; it is much faster to make a messages or do remote procedure calls (*RPC*) between separate tasks. This modular structure results in a more robust and extensible system than a monolithic l microkernel.

Thus in OS X, Mach is not primarily a communication hub between clients and servers. Instead, its value consists of its abstractions, its extensibility, and its flex

- object-based APIs with communication channels (for example, ports) as object references
- highly parallel execution, including preemptively scheduled threads and support for SMP.
- a flexible scheduling framework, with support for real-time usage
- a complete set of IPC primitives, including messaging, RPC, synchronization, and notification
- support for large virtual address spaces, shared memory regions, and memory objects backed by persistent store.
- proven extensibility and portability, for example across instruction set architectures and in distributed environments.
- security and resource management as a fundamental principle of design; all resources are virtualized

Mach Kernel Abstractions

Mach provides a small set of abstractions that have been designed to be both simple and powerful. These are the main kernel abstractions:

- Tasks. The units of resource ownership; each task consists of a virtual address space, a port right namespace, and one or more threads. (Similar to a process.)
- Threads. The units of CPU execution within a task.
- Address space. In conjunction with memory managers, Mach implements the notion of a sparse virtual address space and shared memory.
- Memory objects. The internal units of memory management. Memory objects include named entries and regions; they are representations of potentially persi-
- Ports. Secure, simplex communication channels, accessible only via send and receive capabilities (known as port rights).
- IPC. Message queues, remote procedure calls, notifications, semaphores, and lock sets.
- Time. Clocks, timers, and waiting.

30

Kernel Programming Guide



Thread programming & synchronization



The virtual memory of multithreaded applications



Bounded-Buffer Problem

- Also referred to as "producer-consumer" problem
- Producer places items in shared buffer
- Consumer removes items from shared buffer





38	2	15
----	---	----

We need to control accesses to the buffer!



int buffer[BUFF_SIZE]; // shared global

int item = buffer[out]; out = (out + 1) % BUFF_SIZE;

// do something w/ item

Solving the "Critical Section Problem"

- 1. Mutual exclusion at most one process/thread in its critical section
- Progress a thread outside of its critical section cannot block another thread from entering its critical section
 Fairness — a thread cannot be postponed indefinitely from
- Fairness a thread cannot be postponed in entering its critical section
- 4. Accommodate nondeterminism the solution should work regardless the speed of executing threads and the number of processors

Use locks

• t		int buffer[BUFF_SIZE volatile unsigned in		
<pre>int</pre>	<pre>main(int argc, char *argv[]) { pthread_t p; printf("parent: begin\n"); // init here Pthread_create(&p, NULL, child, NULL int in = 0; while(TRUE) { int item =; Pthread_mutex_lock(&lock); buffer[in] = item; in = (in + 1) % BUFF_SIZE; Pthread_mutex_unlock(&lock); } printf("parent: end\n"); return 0; </pre>);	<pre>void *child(void int out = 0; printf("chil while(TRUE) Pthread_ int item out = (o Pthread_ // do so } return NULL;</pre>	

]; // shared global t lock = 0;

1 *arg) {
.d\n");
{
.mutex_lock(&lock);
n = buffer[out];
out + 1) % BUFF_SIZE;
.mutex_unlock(&lock);
omething w/ item

How to implement lock/unlock

int buffer[BUFF_SIZE]; // shared global volatile unsigned int lock = 0;

```
int main(int argc, char *argv[]) {
   pthread_t p;
   printf("parent: begin\n");
   // init here
   Pthread_create(&p, NULL, child, NULL);
                                                }
   int in = 0;
                                                return NULL;
   while(TRUE) {
                                            }
      int item = ...;
      Pthread_mutex_lock(&lock);
      buffer[in] = item;
                                       while (*lock == 1) // TEST (lock)
      in = (in + 1) \% BUFF_SIZE;
                                     ; // spin
      Pthread_mutex_unlock(&lock);
                                       }
                                   }
   printf("parent: end\n");
   return 0;
}
                                    {
                                       *lock = 0;
                                   }
```

- void *child(void *arg) {
 - printf("child\n");

int out = 0;

while(TRUE) {

```
Pthread_mutex_lock(&lock);
int item = buffer[out];
out = (out + 1) \% BUFF_SIZE;
Pthread_mutex_unlock(&lock);
// do something w/ item
```

void Pthread_mutex_lock(volatile unsigned int *lock) {

void Pthread_mutex_unlock(volatile unsigned int *lock)

- How many of the following can the naive implementation guarantee for the producer-consumer problem?
 - At most one process/thread in its critical section
 - ② A thread outside of its critical section cannot block another thread from entering its critical section
 - ③ A thread cannot be postponed indefinitely from entering its critical section
 - The solution should work regardless the speed of executing threads and the (4) number of processors
 - A. 0
 - B. 1
 - C. 2
 - D. 3
 - E. 4



volatile unsigned int *lock) { / TEST (lock)

/ SET (lock)

k(volatile unsigned int *lock) {

- How many of the following can the naive implementation guarantee for the producer-consumer problem?
 - At most one process/thread in its critical section
 - ② A thread outside of its critical section cannot block another thread from entering its critical section
 - ③ A thread cannot be postponed indefinitely from entering its critical section
 - The solution should work regardless the speed of executing threads and the (4) number of processors
 - A. 0
 - B. 1
 - C. 2
 - D. 3
 - E. 4





volatile unsigned int *lock) { / TEST (lock)

/ SET (lock)

k(volatile unsigned int *lock) {

int buffer[BUFF_SIZE]; // shared global volatile unsigned int lock = 0;

```
while(TRUE) {
                                                      Pthread_mutex_lock(&lock);
int main(int argc, char *argv[]) {
                                                      int item = buffer[out];
    pthread_t p;
                                                      out = (out + 1) \% BUFF_SIZE;
    printf("parent: begin\n");
                                                      Pthread_mutex_unlock(&lock);
    // init here
                                                      // do something w/ item
    Pthread_create(&p, NULL, child, NULL);
    int in = 0;
                                                  return NULL;
   while(TRUE) {
                                              }
       int item = ...;
       Pthread_mutex_lock(&lock);
                                     void Pthread_mutex_lock(volatile unsigned int *lock) {
       buffer[in] = item;
                                         while (*lock == 1) // TEST (lock)
       in = what if context switch
                                     // spin
      Pthread_mutex_unlock.x.pc
happens here?
                                        }
                                     }
    printf("parent: end\n");
   return 0;
                                     void Pthread_mutex_unlock(volatile unsigned int *lock)
}
                                     {
                                         *lock = 0;
                                     }
```

- void *child(void *arg) {
 - printf("child\n");

int out = 0;

 How many of the following can the naive implementation guarantee for the producer-consumer problem?



- X At most one process/thread in its critical section
- ② A thread outside of its critical section cannot block another thread from entering its critical section
- ③ A thread cannot be postponed indefinitely from entering its critical section
- The solution should work regardless the speed of executing threads and the (4)number of processors
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4



volatile unsigned int *lock) { / TEST (lock)

/ SET (lock)

k(volatile unsigned int *lock) {

int buffer[BUFF_SIZE]; // shared global volatile unsigned int lock = 0;

```
while(TRUE) {
                                                    Pthread_mutex_lock(&lock);
int main(int argc, char *argv[ what if the thread
                                                  ____int item = buffer[out];
    pthread_t p;
                                                  out = (out + 1) % BUFF_SIZE;
    printf("parent: begin\n")crashes/halts here?
                                                    Pthread_mutex_unlock(&lock);
    // init here
                                                    // do something w/ item
    Pthread_create(&p, NULL, child, NULL);
    int in = 0;
                                                return NULL;
    while(TRUE) {
                                             }
       int item = ...;
       Pthread_mutex_lock(&lock);
                                    void Pthread_mutex_lock(volatile unsigned int *lock) {
      buffer[in] = item;
                                        while (*lock == 1) // TEST (lock)
       in = what if context switch
                                     // spin
      Pthread_mutex_unlock.ter?
                                        }
                                    }
    printf("parent: end\n");
    return 0;
                                    void Pthread_mutex_unlock(volatile unsigned int *lock)
}
                                    Ł
                                        *lock = 0;
                                    }
```

- void *child(void *arg) {
 - printf("child\n");

int out = 0;

 How many of the following can the naive implementation guarantee for the producer-consumer problem?



- X At most one process/thread in its critical section
- ② A thread outside of its critical section cannot block another thread from entering its critical section



- A thread cannot be postponed indefinitely from entering its critical section
- The solution should work regardless the speed of executing threads and the (4) number of processors
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4



volatile unsigned int *lock) { / TEST (lock)

/ SET (lock)

k(volatile unsigned int *lock) {

int buffer[BUFF_SIZE]; // shared global volatile unsigned int lock = 0;

```
while(TRUE) {
                                                    Pthread_mutex_lock(&lock);
int main(int argc, char *argv[ what if the thread
                                                  __int item = buffer[out];
    pthread_t p;
                                                   out = (out + 1) % BUFF_SIZE;
    printf("parent: begin\n")crashes/halts here?
                                                    Pthread_mutex_unlock(&lock);
    // init here
                                                    // do something w/ item
    Pthread_create(&p, NULL, child, NULL);
    int in = 0;
                                                return NULL;
    while(TRUE) {
                                             }
       int item = ...;
       Pthread_mutex_lock(&lock);
                                    void Pthread_mutex_lock(volatile unsigned int *lock) {
       buffer[in] = item;
                                        in = what if context switch
                                      // spin
      Pthread_mutex_unlock.a.e?
                                       *lock = 1; Iock as 0 at this point
    }
                                    }
    printf("parent: end\n");
    return 0;
                                    void Pthread_mutex_unlock(volatile unsigned int *lock)
}
                                    \mathbf{I}
                                        *lock = 0;
                                    }
```

void *child(void *arg) {

```
printf("child\n");
```

int out = 0;

 How many of the following can the naive implementation guarantee for the producer-consumer problem?



- X At most one process/thread in its critical section
- ② A thread outside of its critical section cannot block another thread from entering its critical section



A thread cannot be postponed indefinitely from entering its critical section



- The solution should work regardless the speed of executing threads and the number of processors
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

```
while (*lock == 1) // TEST (lock)
  ; // spin
    *lock = 1;
}
    *lock = 0;
}
```



void Pthread_mutex_lock(volatile unsigned int *lock) {

// SET (lock)

void Pthread_mutex_unlock(volatile unsigned int *lock) {

int buffer[BUFF_SIZE]; // shared global volatile unsigned int lock = 0;

```
while(TRUE) {
int main(int argc, char *argv[ what if the thread
                                                 ____int item = buffer[out];
   pthread_t p;
    printf("parent: begin\n")crashes/halts here?
    // init here
    Pthread_create(&p, NULL, child, NULL);
    int in = 0;
                                               return NULL;
   while(TRUE) {
                                            }
      int item = ...;
      Pthread_mutex_lock(&lock);
                                   void Pthread_mutex_lock(volatile unsigned int *lock) {
      buffer[in] = item;
                                       in = what if context switch
                                    // spin
      Pthread_mutex_unlock.te?
                                       *lock = 1; Iock as 0 at this point
                                   }
    printf("parent: end\n");
   return 0;
                                   void Pthread_mutex_unlock(volatile unsigned int *lock)
}
                                   Ł
                                       *lock = 0;
                                   }
```

- void *child(void *arg) {
 - printf("child\n");

int out = 0;

Pthread_mutex_lock(&lock); out = (out + 1) % BUFF_SIZE; Pthread_mutex_unlock(&lock); // do something w/ item

coherence cache misses? page fault?

 How many of the following can the naive implementation guarantee for the producer-consumer problem?



- X At most one process/thread in its critical section
- X A thread outside of its critical section cannot block another thread from entering its critical section



X A thread cannot be postponed indefinitely from entering its critical section



The solution should work regardless the speed of executing threads and the number of processors





volatile unsigned int *lock) { / TEST (lock)

/ SET (lock)

k(volatile unsigned int *lock) {

Announcement

- Reading quizzes due next Tuesday
 - Welcome new friends! will drop a total of 6 reading quizzes for the quarter
 - Attendance count as 4 reading guizzes
 - We plan to have a total of 11 reading guizzes
- Office Hour links are inside Google Calendar events
 - https://calendar.google.com/calendar/u/0/r? cid=ucr.edu_b8u6dvkretn6kq6igunlc6bldg@group.calendar.google.com
 - Different links from lecture ones
 - We cannot share through any public channels so that we can better avoid Zoom bombing
- We will make both midterm and final exams online this quarter
 - Avoid the uncertainty of COVID-19
 - Avoid high-density in the classroom (only sits 60 and we have 59 for now) during examines

Computer Science & Engineering





