1. What are some of the differences between a processor running in *privileged mode* (also called *kernel mode*) and user mode? Why are the two modes needed? In user-mode: • CPU control registers are inaccessible. • CPU management instructions are inaccessible. • Part's of the address space (containing kernel code and data) are inaccessible. Some device memory and registers (or ports) are inaccessible. The two modes of operation are required to ensure that applications (running in user-mode) cannot bypass, circumvent, or take control of the operating system.

**Tutorial Week 2** 

**Questions and Answers** 

**Operating Systems Intro** 

2. What are the two main roles of an Operating System? 1) It provides a high-level abstract machine for programmers (hides the details of the hardware) 2) It is a resource manager that divides resources amongst competing programs or users according to some system policy. 3. Given a high-level understanding of file systems, explain how a file system fulfills the two roles of an operating

system? At the level of the hardware, storage involves low-level controller hardware and storage devices that store blocks of data at many locations in the store. The OS filesystem abstracts above all these details

and provides an interface to store, name and organise arbitrary unstructured data. The filesystem also arbitrates between competing processor by managing allocated and free space on

the storage device, in addition to enforcing limits on storage consumption (e.g. quotas).

4. Which of the following instructions (or instruction sequences) should only be allowed in kernel mode? 1. Disable all interrupts. 2. Read the time of day clock.

3. Set the time of day clock.

4. Change the memory map. 5. Write to the hard disk controller register. 6. Trigger the write of all buffered blocks associated with a file back to disk (fsync).

1,3,4,5 need to be restricted to kernel mode.

OS system call interface 5. The following code contains the use of typical UNIX process management system calls: fork(), execl(), exit() and getpid(). If you are unfamiliar with their function, browse the man pages on a UNIX/Linux machine get an overview, e.g: man fork

Answer the following questions about the code below.

for (i = 1; i <= FORK DEPTH; i++) {

/\* we're in the parent process after successfully forking a child \*/

c. What is the process id of /bin/echo?

#include <sys/types.h> #include <unistd.h> #include <stdlib.h> #include <stdio.h>

#define FORK DEPTH 3

my\_pid = getpid();

r = fork();

 $if (r > 0) {$ 

 $}$  else if (r == 0) {

my\_pid = getpid();

exit(1);

exit(1);

} else { /\* r < 0 \*/</pre>

}

} } }

6.

main()

}

int fd; int len; ssize\_t r;

if (fd < 0) {

exit(1);

if (r < 0) {

exit(1);

close(fd);

#include <sys/types.h> #include <sys/stat.h> #include <fcntl.h> #include <stdio.h> #include <stdlib.h> #include <string.h>

main() {

> int fd; int len; ssize\_t r; off\_t off;

if (fd < 0) {

exit(1);

 $if (r < 0) {$ 

exit(1);

if (off < 0) {

exit(1);

if (r < 0) {

exit(1);

close(fd);

observe the output.

open file?

b. 3

9. Compile and run the following code.

#include <unistd.h> #include <stdlib.h> #include <stdio.h> #include <errno.h>

> r = chdir(".."); $if (r < 0) {$

> > exit(1);

1. read() 2. printf() 3. memcpy() 4. open() 5. strncpy()

functions.

Ready.

multithreaded?

x = x + 1;

condition exists.

void foo()

int j;

i = i + 1;j = j + 1;

dependent on i).

void inc mem(int \*iptr)

\*iptr = \*iptr + 1;

critical section.

/\* random stuff\*/

/\* more random stuff \*/

int i;

{

}

{

}

**Critical sections** 

**Processes and Threads** 

Events that cause transitions:

• Running threads = 0 or 1.

stores, then the second thread stores).

• Blocked = N - Running - Ready • Ready = N - Running - Blocked

perror("Eek!");

perror("Double eek!");

{

int r;

a. What do the following code do?

c. In what directory does /bin/ls run in? Why?

r = execl("/bin/ls","/bin/ls",NULL);

10. On UNIX, which of the following are considered system calls? Why?

a. What is strace doing?

you expect to see.

}

}

}

if (i == FORK DEPTH) {

independent after forking.

not change.

a. What does the following code do?

#include <sys/types.h> #include <sys/stat.h> #include <fcntl.h> #include <stdio.h> #include <stdlib.h> #include <string.h>

main()

int i, r; pid\_t my\_pid;

a. What is the value of i in the parent and child after fork.

b. What is the value of my\_pid in a parent after a child updates it?

e. How many times is *Hello World* printed when FORK\_DEPTH is 3?

d. Why is the code after execl not expected to be reached in the normal case?

f. How many processes are created when running the code (including the first process)?

printf("Parent process %d forked child process %d\n",my\_pid, r);

/\* run /bin/echo if we are at maximum depth, otherwise continue loop \*/

a. The child is a new independent process that is a copy of the parent. i in the child will have

b. my\_pid in a parent is not updated by any action of the child and the child and parent are

d. A successful execl results in the current code being replaced. execl does not return if it

c. execl replaces the *content* of a running process with specified executable. The process id does

/\* We're in the child process, so update my\_pid \*/

/\* we never expect to get here, just bail out \*/

whatever the value was in the parent at the point of forking.

succeeds as there is no previous code to return to. e. Hello World is printed 4 times if the FORK\_DEPTH is 3. f. There are 8 processes involved in the execution of the code.

fd = open("testfile", O\_WRONLY | O\_CREAT, 0600);

printf("Attempting to write %d bytes\n",len);

/\* just ungracefully bail out \*/

perror("File open failed");

r = write(fd, teststr, len);

perror("File write failed");

printf("Wrote %d bytes\n", (int) r);

len = strlen(teststr);

b. In addition to O WRONLY, what are the other 2 ways one can open a file?

c. What open return in fd, what is it used for? Consider success and failure in your answer.

char teststr[] = "The quick brown fox jumps over the lazy dog.\n";

a. The code writes a string to a file. It will create a new file if needed (O CREAT).

file related systems cases to identify the file to operate on.

char teststr[] = "The quick brown fox jumps over the lazy dog.\n";

7. The following code is a variation of the previous code that writes twice.

b. What is lseek() doing that is affecting the final file size? c. What over options are there in addition to SEEK\_SET?.

fd = open("testfile2", O\_WRONLY | O\_CREAT, 0600);

printf("Attempting to write %d bytes\n",len);

/\* just ungracefully bail out \*/

perror("File open failed");

r = write(fd, teststr, len);

perror("File write failed");

off = lseek(fd, 5, SEEK\_SET);

r = write(fd, teststr, len);

perror("File lseek failed");

perror("File write failed");

printf("Wrote %d bytes\n", (int) r);

a. 50 bytes. For each open file, the operating system keeps track of the current offset within the file. The current offset is where the next read or write will start from. The current offset is

second write begins from offset 5, writes 45 bytes, giving 50 bytes in total in the file.

8. Compile either of the previous two code fragments on a UNIX/Linux machine and run strace ./a.out and

b. Without modifying the above code to print fd, what is the value of the file descriptor used to write to the

c. printf does not appear in the system call trace. What is appearing in it's place? What's happening here?

to the system call. There are a lot of system calls at the beginning of a trace related to

c. printf is a library function that creates a buffer based on the string specification that it is passed. The buffer is then written to the console using write() to file descriptor 1.

b. After the program runs, the current working directory of the shell is the same. Why?

in the directory hierarchy), and then runs 1s to list the directory.

current working directory of the parent process remains the same.

of every process, so ls runs in the current directory of child process.

a. strace is printing a trace of all system calls invoked by a process, together with the arguments

dynamically loading code libraries. Towards the end of the trace you will see the system calls

a. The code sets the current working directory of the process to be the parent directory (one higher

b. The shell forks a child process that runs the code. Each process has its own current working directory, so the code above changes the current working directory of the child process, the

c. exec replaces the content of the child process with 1s, not the environment the child process

1 and 4 are system calls, 2 is a C library functions which can call write(), 3 and 5 a simply library

11. In the three-state process model, what do each of the three states signify? What transitions are possible between

the process is not runnable as it is waiting for some event prior to continuing execution.

The three states are: Running, the process is currently being executed on the CPU; Ready, the process is ready to execute, but has not yet been selected for execution by the dispatcher; and *Blocked* where

Possible transitions are Running to Ready, Ready to Running, Running to Blocked, and Blocked to

• Running to Ready: timeslice expired, yield, or higher priority process becomes ready.

• Running to Blocked: A requested resource (file, disk block, printer, mutex) is unavailable, so the

• Blocked to Ready: a resource has become available, so all processes blocked waiting for the

12. Given N threads in a uniprocessor system. How many threads can be running at the same point in time? How many

13. Compare reading a file using a single-threaded file server and a multithreaded file server. Within the file server, it takes 15 msec to get a request for work and do all the necessary processing, assuming the required block is in the main memory disk block cache. A disk operation is required for one third of the requests, which takes an additional 75 msec during which the thread sleeps. How many requests/sec can a server handled if it is single threaded? If it is

In the single-threaded case, the cache hits take 15 msec and cache misses take 90 msec. The weighted average is  $2/3 \times 15 + 1/3 \times 90$ . Thus the mean request takes 40 msec and the server can do 25 per second. For a multithreaded server, all the waiting for the disk is overlapped, so every request takes 15

14. The following fragment of code is a single line of code. How might a race condition occur if it is executed

concurrently by multiple threads? Can you give an example of how an incorrect result can be computed for x.

The single code statement is compiled into multiple machine instructions. The memory location

15. The following function is called by multiple threads (potentially concurrently) in a multi-threaded program. Identify the critical section(s) that require(s) mutual exclusion. Describe the race condition or why no race

corresponding to x is loaded into a register, incremented, and then stored back to memory. During the interval between the load and store in the first thread, another thread may perform a load, increment, and store, and when control passes back to the first thread, the results of the second are overwritten are them overwritten. Another outcome would be for the results of the first to be overwritten by the second (as the first thread loads, increments, then the second thread loads, increments, then the first thread

There is no race condition on j, since it is a local variable per thread. However, i is a variable shared

Whether \*iptr = \*iptr + 1 forms a critical section depends on the scope of the pointer passed to inc mem. If the pointer points to a local variable, then there is no race. If the pointer points to a shared global variable there is potential for a race, and thus the increment would become a critical section.

between threads. Thus i = i + 1 would form a critical section (assuming no random stuff is

16. The following function is called by threads in a multi-thread program. Under what conditions would it form a

each of the states, and what causes a process (or thread) to undertake such a transition?

• Ready to Running: Dispatcher chose the next thread to run.

resource now become ready to continue execution.

msec, and the server can handle 66 2/3 requests per second.

process is blocked waiting for the resource to become available.

threads can be *ready* at the same time? How many threads can be *blocked* at the same time?

runs in. The current working directory is part of the environment that the OS manages on behalf

c. See the man page for details on SEEK CUR and SEEK END

usually at the location of offset of the end of the previous read or write. So one would expect the file size to be 90 bytes after two 45 byte writes, except for lseeks interference (see below). b. lseek sets the current offset to a specific location in the file. The lseek in the code moves the current offset from 45 bytes (after the initial write) to 5 bytes from the start of the file. The

printf("Wrote %d bytes\n", (int) r);

len = strlen(teststr);

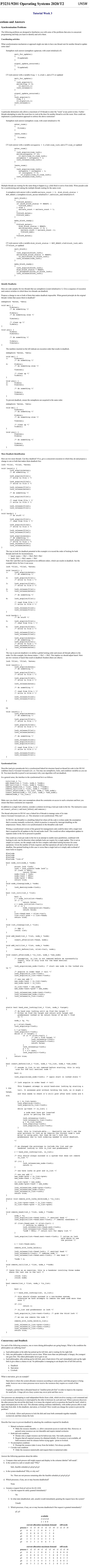
a. How big is the file (in bytes) after the two writes?

b. The other ways of opening a file are read-only (O\_RDONLY) and read-write (O\_RDWR). c. In case of failure fd is set to -1 to signify an error. In the case of success, fd is set to a file

descriptor (an integer) that becomes a handle to the file. The file descriptor is used in the other

r = execl("/bin/echo","/bin/echo","Hello World",NULL);

/\* Eek, not expecting to fail, just bail ungracefully \*/



decoded. Rather than discard this potentially useful work, the architecture rules state that the instruction after a branch is always executed before the instruction at the target of the branch. 2. The goal of this question is to have you reverse engineer some of the C compiler function calling convention (instead of reading it from a manual). The following code contains 6 functions that take 1 to 6 integer arguments. Each function sums its arguments and returns the sum as a the result. #include <stdio.h> /\* function protoypes, would normally be in header files \*/ int arg1(int a); int arg2(int a, int b); int arg3(int a, int b, int c); int arg4(int a, int b, int c, int d); int arg5(int a, int b, int c, int d, int e ); int arg6(int a, int b, int c, int d, int e, int f); /\* implementations \*/ int arg1(int a)

**Tutorial Week 4** 

The pipeline structure of the MIPS CPU means that when a jump instruction reaches the "execute" phase and a new program counter is generated, the instruction after the jump will already have been

**Questions and Answers** 

1. What is a branch delay?

R3000 and assembly

return a + b + c + d + e;

return a + b + c + d + e + f;

int arg6(int a, int b, int c, int d, int e, int f)

/\* do nothing main, so we can compile it \*/

03e00008

00801021

03e00008

00851021

00851021

03e00008

00461021

00852021

00861021

03e00008

00471021

00852021

00863021

00c73821

8fa20010

03e00008

00e21021

00852021

00863021

00c73821

8fa20010

00000000

00e22021

8fa20014

03e00008

00821021

03e00008

00001021

b. Why is there no stack references in arg2?

f. Why does arg5 and arg6 reference the stack?

}

}

}

int main()

of for the sake of clarity).

004000f0 <arg1>:

004000f8 <arg2>: 4000f8:

00400100 <arg3>:

0040010c <arq4>:

0040011c <arg5>:

00400134 <arg6>:

4000f0:

4000f4:

4000fc:

400100:

400104:

400108:

40010c:

400110:

400114:

400118:

40011c:

400120:

400124:

400128:

40012c:

400130:

400134:

400138:

40013c:

400140:

400144:

400148:

40014c:

400150:

400154:

40015c:

00400158 <main>: 400158:

c. What does jr ra do?

to store them.

arg1 returns.

void reverse\_print(char \*s)

reverse print(s+1);

reverse\_print(teststr);

write(STDOUT\_FILENO,s,1);

a. Describe what each line in the code is doing.

27bdffe8

afbf0014

afb00010

80820000

0000000

10400007

00808021

0c10003c

24840001

24040001

02002821

0c1000af

24060001

8fbf0014

8fb00010

03e00008

27bd0018

004000f0 <reverse\_print>:

4000f0:

and s0

4000f4:

400100:

400104:

400110:

40011c:

400120:

400124:

function

400128:

40012c:

400130:

**Threads** 

byte).

<reverse print+0x34>

Call reverse print

pointer to write.

Call write function

Return to the caller.

5. Compare cooperative versus preemptive multithreading?

and select a ready thread to run next.

nop

overwritten.

if (\*s != '\0') {

off for the sake of clarity).

004000f0 <reverse print>:

} }

}

int main()

4000f0:

4000f4:

4000f8:

4000fc:

400100:

400104:

400108:

40010c:

400110:

400114:

400118:

40011c:

400120:

400124:

400128:

40012c:

400130:

passed on the stack?

a. v0

d. a0

#include <stdio.h> #include <unistd.h>

{ return a; } int arg2(int a, int b) return a + b; int arg3(int a, int b, int c) { return a + b + c; } int arg4(int a, int b, int c, int d)

return a + b + c + d; } int arg5(int a, int b, int c, int d, int e)

The following code is the disassembled code that is generated by the C compiler (with certain optimisations turned

ra

v0,a0

v0,a0,a1

v0,a0,a1

v0, v0, a2

a0,a0,a1

v0,a0,a2

v0, v0, a3

a0,a0,a1

a2,a0,a2

a3,a2,a3 v0,16(sp)

v0,a3,v0

a0,a0,a1

a2,a0,a2

a3,a2,a3

a0,a3,v0

v0,a0,v0

v0, zero

b. There are no local variables inside the function, so the compiler does not need space on the stack

c. It jumps (changes the program counter) to the address in the ra register. The ra register is set by

arriving at the destination of the jump. Thus, logically, the move instruction is executed before

a jal instruction to the address of the instruction after jal. Thus function calls can be

e. The instruction after a jump (i.e. the instruction in the branch delay slot is executed prior to

f. Up to 4 arguments can be passed to a function in registers. Arguments beyond the fourth are

3. The following code provides an example to illustrate stack management by the C compiler. Firstly, examine the C

The following code is the disassembled code that is generated by the C compiler (with certain optimisations turned

sp, sp, -24

ra,20(sp)

s0,16(sp)

v0,400124 <reverse\_print+0x34>

sp, sp, -24

ra,20(sp)

s0,16(sp)

v0,0(a0)

s0,a0

a0,a0,1

4002bc <write>

a0,1

a2,1

ra,20(sp)

s0,16(sp)

sp, sp, 24

ra

4000f0 <reverse\_print>

4000f0 <reverse print>

v0,0(a0)

s0,a0

a1,s0

a2,1

ra

a0,a0,1 a0,1

ra,20(sp)

s0,16(sp)

sp,sp,24

addiu

Allocate 24 bytes on the stack, 16 for a0-a3 (unused) and 8 for ra

function calls other functions, which means the ra register will be

Recall the 's' registers must be preserved when we return from this function. We only use s0, so save it on the stack so we can use the

Save the return address for the function on the stack. This

register in this function, but restore it before returning.

sw

lb Load a character from the pointer passed as the first argument.

nop

Test is the character is zero, if so, jump forward to 400124

This is on the delay slot, save the pointer in s0

print start on the next character in the string.

move

jal

addiu

This is in the delay slot, add 1 to the pointer to have reverse

li

move Remember s0 is preserved across function calls above, so s0 still contains the original pointer passed into the function. Pass the

jal

Another delay slot, load the number of bytes write should output (1

Restore the return address of this function in prep for return from

li

lw

lw

Restore s0 to whatever it was before this function was called.

jr

b. The stack of each invocation of reverse print is 24 bytes, but the function is recursive. The allocation is 24 bytes times the length of the string, and thus if the string is unbounded, so is the

addiu

4002bc <write>

ra

v0,20(sp)

v0,16(sp)

jr

jr

addu

lw

lw

jr

jr

a. arg1 (and functions in general) returns its return value in what register?

d. Which register contains the first argument to the function? e. Why is the move instruction in arg1 after the jr instruction.

implemented with jal and jr ra instructions.

code in the provided example to understand how the recursive function works.

char teststr[] = "\nThe quick brown fox jumps of the lazy dog.\n";

b. What is the maximum depth the stack can grow to when this function is called?

addiu

sw

sw

lb

nop

beqz

move

jal

li

addiu

move

jal

li

lw

lw

jr

27bdffe8

afbf0014

afb00010

80820000

0000000

10400007

00808021

24840001

24040001

02002821

0c1000af

24060001

8fbf0014

8fb00010

03e00008

27bd0018

recursion, and thus stack growth is also unbounded.

4. Why is recursion or large arrays of local variables avoided by kernel programmers?

time to each thread, even when a thread is uncooperative or malicious.

The kernel has no knowledge of the user-level threads.

system call variants to emulate the blocking calls.

Switching between each process (in-kernel thread) is performed by the function

and sets the stack pointer to the stack pointer stored in destination tcb.

stack, and returns to the destination thread to continue where is left off.

They are restored via a rfe instruction (restore from exception).

12. Why must kernel programmers be especially careful when implementing system calls?

program crashing, or compromising the operating system.

switch thread(cur tcb, dst tcb). What does this function do?

scheduler is the normal in-kernel scheduler.

and exit required).

affecting the others.

Kernel Entry and Exit

interrupted.

restored?

9. What is the EPC register? What is it used for?

The value of ExcCode is 8.

R3000 from the lecture slides.

s0-s8, ra).

find them.

appears in the program.

notionally responsible.

Timer Device

CPU

CPU

CPU

Who

Kernel Interrupt Handler

Kernel Interrupt Handler

Kernel Interrupt Handler

Timer Interrupt Handler

Timer Interrupt Handler

Scheduler

Scheduler

Kernel

Kernel Kernel

Kernel

Kernel

Kernel

wrapper function does very little.

have the same name? If not, which name is important?

function?

points to the previous (branch) instruction.

kernel stack assocaited with each process.

return to the caller) onto the current stack.

In the branch delay slot, deallocate the stack.

The kernel stack is usually a limited resource. A stack overflow crashes the entire machine.

Cooperative multithreading is where the running thread must explicitly yield() the CPU so that the dispatcher can select a ready thread to run next. Preemptive multithreading is where an external event (e.g. a regular timer interrupt) causes the dispatcher to be invoked and thus *preempt* the running thread,

Cooperative multithreading relies on the cooperation of the threads to ensure each thread receives

6. Describe *user-level threads* and *kernel-level threads*. What are the advantages or disadvantages of each approach?

User-level threads are implemented in the application (usually in a "thread library"). The thread management structures (Thread Control Blocks) and scheduler are contained withing the application.

Kernel threads are implemented in the kernel. The TCBs are managed by the kernel, the thread

• User threads don't take advantage of parallelism available on multi-CPU machines.

• User-level threads can be implemented on OSes without support for kernel threads.

level threads are being used, this action will block the entire process, destroying the value of

8. Assume a multi-process operating system with single-threaded applications. The OS manages the concurrent

7. A web server is constructed such that it is multithreaded. If the only way to read from a file is a normal blocking read system call, do you think user-level threads or kernel-level threads are being used for the web server? Why?

user-level threads use alarms or timeouts to provide a tick for preemption).

• User threads are generally faster to create, destroy, manage, block and activate (no kernel entry

• If a single user-level thread blocks in the kernel, the whole process is blocked. However, some libraries (e.g., most UNIX pthreads libraries) avoid blocking in the kernel by using non-blocking

• Kernel threads are usually preemptive, user-level threads are usually cooperative (Note: some

A worker thread within the web server will block when it has to read a Web page from the disk. If user-

multithreading. Thus it is essential that kernel threads are used to permit some threads to block without

application requests by having a thread of control within the kernel for each process. Such a OS would have an in-

The function saves the registers required to preserve the compiler calling convention (and registers to

The function saves the resulting stack pointer into the thread control block associated with cur\_tcb,

The function then restores the registers that were stored previously on the destination (now current)

This is a 32-bit register containing the 32-bit address of the return point for the last exception. The instruction causing (or suffering) the exception is at EPC, unless BD is set in Cause, in which case EPC

It is used by the exception handler to restart execution at the at the point where execution was

10. What happens to the KUc and IEc bits in the STATUS register when an exception occurs? Why? How are they

11. What is the value of ExcCode in the Cause register immediately after a system call exception occurs?

System calls with poor argument checking or implementation can result in a malicious or buggy

13. The following questions are focused on the case study of the system call convention used by OS/161 on the MIPS

3. At minimum, what additional information is required beyond that passed to the system-call wrapper

1. How does the 'C' function calling convention relate to the system call interface between the application and

2. What does the most work to preserve the compiler calling convention, the system call wrapper, or the OS/161

a. The 'C' function calling convention must always appear to be adhered to after any system-call wrapper function completes. This invovles saving and restoring of the *preserved* registers (e.g.,

b. The OS/161 kernel code does the saving and restoring of preserved registers. The system call

usually added by setting an agreed-to register to the value of the system call number.

14. In the example given in lectures, the library function *read* invoked the read system call. Is it essential that both

Note: The kernel programmer may can the code in the operating system that implements read a

15. To a programmer, a system call looks like any other call to a library function. Is it important that a programmer

As far as program logic is concerned it does not matter whether a call to a library function results in a system call. But if performance is an issue, if a task can be accomplished without a system call the program will run faster. Every system call involves overhead in switching from the user context to the kernel context. Furthermore, on a multiuser system the operating system may schedule another process to run when a system call completes, further slowing the progress in real time of a calling process.

Note: In the table below, almost everything that is not the timer device or CPU is actually just code executing within the kernel. The distinction of "who" is there to clarify which kernel subsystem is

Begins executing the Kernel Interrupt Handler

Acknowledges the interrupt to the timer device

Calls the *Kernel* to switch to the new process

Restores the user registers from the stack

Sets the processor back to user mode

Jumps to the new user process's PC

Saves the current process's in-kernel context to the stack

Reads the new process's in-kernel context off the stack

Switches to the new process's kernel stack by changing sp

What

Changes the sp to the kernel stack for the interrupted process

Determines the source of the interrupt to be the timer device

similar name, e.g., sys read. This is purely a convention for the sake of code clarity.

16. Describe a plausible sequence of activities that occur when a timer interrupt results in a context switch.

Generates an interrupt exception

Switches from user to kernel mode

*Kernel Interrupt Handler* Saves the user registers for the interrupted process onto the stack

Calls the *Timer Interrupt Handler* 

Calls the *Scheduler* (dispatcher)

Chooses a new process to run

Raises interrupt line

know which library function result in system calls? Under what circumstances and why?

System calls do not really have names, other than in a documentation sense. When the library function read() traps in to the kernel, it puts the number of the system call into a register or on the stack. This number is used to index into a jump table (e.g. a C switch statement). There is really no name used anywhere. On the other hand, the name of the library function is very important, since that is what

c. The interface between the system-call wrapper function and the kernel can be defined to provide additional information beyond that passed to the wrapper function. At minimum, the wrapper function must add the system call number to the arguments passed to the wrapper function. It's

The system call convention also uses the calling convention of the C-compiler to pass arguments to OS/161. Having the same covention as the compiler means the system call wrapper can avoid moving arguments around and the compiler has already placed them where the OS expects to

The 'c' (current) bits are shifted into the corresponding 'p' (previous) bits, after which KUC = 0, IEC = 0 (kernel mode with interrupts disabled). They are shifted in order to preserve the current state at the point of the exception in order to restore that exact state when returning from the exception.

regular CPU time. Preemptive multithreading enforces a regular (at least systematic) allocation of CPU

Load the file descriptor for write (1).

addiu

move

nop

addu

addu

lw

jr

jr

jr

move

## **Questions and Answers**

## **Memory Hierarchy and Caching**

1. Describe the memory hierarchy. What types of memory appear in it? What are the characteristics of the memory as one moves through the hierarchy? How can do memory hierarchies provide both fast access times and large capacity?

The memory hierarchy is a hierarchy of memory types composed such that if data is not accessible at the top of the hierarchy, lower levels of the hierarchy are accessed until the data is found, upon which a

Registers, cache, main memory, magnetic disk, CDROM, tape are all types of memory that can be composed to form a memory hierarchy.

copy (usually) of the data is moved up the hierarchy for access.

In going from the top of the hierarchy to the bottom, the memory types feature decreasing cost per bit, increasing capacity, but also increasing access time.

As we move down the hierarchy, data is accessed less frequently, i.e. frequently accessed data is at the top of the hierarchy. The phenomenon is called "locality" of access, most accesses are to a small subset of all data.

2. Given that disks can stream data quite fast (1 block in tens of microseconds), why are average access times for a block in milliseconds?

rotational latency (1/2 rotation) is in milliseconds (e.g. 2 milliseconds for 15,000rpm disk).

Seek times are in milliseconds (e.g. .5 millisecond track to track, 8 millisecond inside to outside), and

3. You have a choice of buying a 3 Ghz processor with 512K cache, and a 2 GHz processor (of the same type) with a 3 MB cache for the same price. Assuming memory is the same speed in both machines and is much less than 2GHz (say 400MHz). Which would you purchase and why? Hint: You should consider what applications you expect to run on the machine.

If you are only running an small application (or a large one, that accesses only a small subset), then the 3GHz processor will be much faster. If you are running a large application access a larger amount of memory than 512K but generally less than 3MB, the 2GHz processor should be faster as the 3 GHz processor will be limited by memory speed.

### 4. Consider a file currently consisting of 100 records of 400 bytes. The filesystem uses *fixed blocking*, i.e. one 400

Files and file systems

- byte record is stored per 512 byte block. Assume that the file control block (and the index block, in the case of indexed allocation) is already in memory. Calculate how many disk I/O operations are required for contiguous, linked, and indexed (single-level) allocation strategies, if, for one record, the following conditions hold. In the contiguous-allocation case, assume that there is no room to grow at the beginning, but there is room to grow at the end of the file. Assume that the record information to be added is stored in memory.
  - c. The record is added at the end.

a. The record is added at the beginning.b. The record is added in the middle.

- d. The record is removed from the beginning.e. The record is removed from the middle.
- f. The record is removed from the end.

0r/1w

0r/1w

- Contiguous Linked Indexed
  - b. 50r/51w 50r/2w 0r/1w

100r/101w

external fragmentation?

- c. 0r/1w 100r/2w 0r/1w
   d. 0r/0w 1r/0w 0r/0w
- e. 49r/49w 50r/1w 0r/0w 51r/1w
- f. 0r/0w 99r/1w 0r/0w
- Internal fragmentation

6. Old versions of UNIX allowed you to write to directories. Newer ones do not even allow the superuser to write to

5. In the previous example, only 400 bytes is stored in each 512 byte block. Is this wasted space due to internal or

To prevent total corruption of the fs. eg cat /dev/zero > /

them? Why? Note that many unices allow you read directories.

7. Given a file which varies in size from 4KiB to 4MiB, which of the three allocation schemes (*contiguous*, *linked-list*, *or i-node based*) would be suitable to store such a file? If the file is access randomly, how would that influence the suitability of the three schemes?

Contiguous is not really suitable for a variable size file as it would require 4MiB to be pre-allocated, which would waste a lot of space if the file is generally mush smaller. Either linked-list of i-node-

based allocation would be preferred. Adding random access to the situation (supported well by contiguous or i-node based), which further motivate i-node-based allocation to be the most appropriate.

8. Why is there VFS Layer in Unix?

- each file system to be aware of other file system types.Provides transparent access to all supported file systems including network file systems (e.g.
  - NFS, CODA)
    It provides a clean interface between the file system independent kernel code and the file system specific kernel code.

• It provides a framework to support multiple file system types concurrently without requiring

- Provide support for *special* file system types like /proc.
- 9. How does choice of block size affect file system performance. You should consider both sequential and random access.
  - The larger the block size, the fewer I/O operations required and the more contiguous the disk accesses. Compare loading a single 16K block with loading 32 512-byte blocks.

Sequential Access

Random Access
 The larger the block size, the more unrelated data loaded. Spatial locality of access can improve the situation.

10. Is the open() system call in UNIX essential? What would be the consequence of not having it?

It is not absolutely essential. The read and write system calls would have to be modified such that:

• The filename is passed in on each call to identify the file to operate on

- With a file descriptor to identify the open session that is returned by open, the sycalls would also need to specify the offset into the file that the syscall would need to use.

• Effectively opening and closing the file on each read or write would reduce performance.

- 11. Some operating system provide a *rename* system call to give a file a new name. What would be different compared to the approach of simply copying the file to a new name and then deleting the original file?
- The rename system call would just change the string of characters stored in the directory entry. A copy operation would result in a new directory entry, and (more importantly) much more I/O as each block of the original file is copied into a newly allocated block in the new file. Additionally, the original file

blocks need de-allocating after the copy finishes, and the original name removed from the directory. A rename is much less work, and thus way more efficient than the copy approach.

12. In both UNIX and Windows, random file access is performed by having a special system call that moves the

current position in the file so the subsequent read or write is performed from the new position. What would be

the consequence of not having such a call. How could random access be supported by alternative means?

Without being able to move the file pointer, random access is either extremely inefficient as one would have to read sequentially from the start each time until the appropriate offset is arrived at, or the an extra argument would need to be added to read or write to specify the offset for each operation.

## **Questions and Answers**

## Files and file systems

1. Why does Linux pre-allocate up to 8 blocks on a write to a file.

Pre-allocating provides better locality when many writes to independent files are interleaved.

2. Linux uses a *buffer cache* to improve performance. What is the drawback of such a cache? In what scenario is it problematic? What alternative would be more appropriate where a buffer cache is inappropriate?

The buffering writes in the buffer cache provides the opportunity for data to be lost if the system stops prior to the cache being flushed.

Removable storage devices are particular problematic if users don't "unmount" them first.

Robustness can be improved by using a write-through cache at the expense of poor write performance.

- 3. What is the structure of the contents of a directory? Does it contain attributes such as creation times of files? If not, where might this information be stored?
  - See lecture slides.
  - No, directories only have a name-to-inode mapping
  - Attributes of the file are stored in the inode itself.
- 4. The Unix inode structure contains a reference count. What is the reference count for? Why can't we just remove the inode without checking the reference count when a file is deleted?

Inodes contain a reference count due to hard links. The reference count is equal to the number of directory entries that reference the inode. For hard-linked files, multiple directory entries reference a

single inode. The inode must not be removed until no directory entries are left (ie, the reference count is 0) to ensure that the filesystem remains consistent.

- 5. Inode-based filesystems typically divide a file system partition into *block groups*. Each block group consists of a number of contiguous physical disk blocks. Inodes for a given block group are stored in the same physical location as the block groups. What are the advantages of this scheme? Are they any disadvantages?
  - Each group contains a redundant superblock. This make the file system more robust to disk block failures.
     Block groups keep the inodes physically closer to the files they refer to then they would be
  - Block groups keep the inodes physically closer to the files they refer to than they would be (on average) on a system without block groups. Since accessing and updating files also involves accessing or updating its inode, having the inode and the file's block close together reduces disk seek time, and thus improves performance. The OS must take care that all blocks remain within the block group of their inode.
- 6. Assume an inode with 10 direct blocks, as well as single, double and triple indirect block pointers. Taking into account creation and accounting of the indirect blocks themselves, what is the largest possible number of block reads and writes in order to:
  - a. Read 1 byte
  - b. Write 1 byte

Assume the inode is cached in memory.

- a. To write 1 byte, in the worst case:
  - 4 writes: create single indirect block, create double indirect block, create triple indirect block, write data block.
  - 3 reads, 2 writes: read single indirect, read double indirect, read triple indirect, write triple indirect, write data block
  - Other combinations are possible
- b. To read 1 byte, in the worst case:
  - 4 reads: read single indirect, read double indirect, read triple indirect, read data block
- 7. Assume you have an inode-based filesystem. The filesystem has 512 byte blocks. Each inode has 10 direct, 1 single indirect, 1 double indirect, and 1 triple indirect block pointer. Block pointers are 4 bytes each. Assume the inode and any block free list is always in memory. Blocks are not cached.
  - a. What is the maximum file size that can be stored before1. the single indirect pointer is needed?
    - 2. the double indirect pointer is needed?
  - 3. the triple indirect pointer is needed? b. What is the maximum file size supported?
  - c. What is the number of disk block reads required to read 1 byte from a file
  - in the best case?
     in the worst case?
  - d. What is the number of disk block reads and writes required to write 1 byte to a file
    - in the best case?
       in the worst case?
  - a. 1. 5K 2. 69K
    - 3. 8261K
  - b. 1056837K c. 1. 1
  - 2. 4d. What is the number of disk block reads and writes required to write 1 byte to a file
  - 1. 1w 2. 4r/1w
    - 2. 4r/1

both? Should not blocks = size / block size?

8. A typical UNIX inode stores both the file's size and the number of blocks currently used to store the file. Why store

• The blocks used to store a file includes and indirect blocks used by the filesystem to keep track

Blocks used to store the file are only indirectly related to file size.

- of the file data blocks themselves.

   File systems only store blocks that actually contain file data. Sparsely populated files can have
- large regions that are unused within a file.
- 9. How can deleting a file leave a inode-based file system (like ext2fs in Linux) inconsistent in the presence of a power failure.

Mark disk blocks as free.

Deleting a file consists of three separate modifications to the disk:

- Remove the directory entry.Mark the i-node as free.
- If the system only completes a subset of the operations (due to power failures or the like), the file

system is no longer consistent. See lecture slide for example of things that can go wrong.

10. How does adding journalling to a file system avoid corruption in the presence of unexpected power failures.

Simply speaking, adding a journal addresses the issue by grouping file system updates into transactions that should either completely fail or succeed. These transactions are logged prior to manipulating the file system. In the presence of failure the transaction can be completed by replaying the updates remaining in the log.

## **Questions and Answers Memory Management**

# 1. Describe internal and external fragmentation.?

- External Fragmentation: total memory space exists to satisfy a request, but it is not contiguous.
  - Internal Fragmentation: allocated memory may be slightly larger than requested memory; this size
- difference is memory internal to a partition, but is not being used.
- 2. What are the problems with multiprogrammed systems with fixed-partitioning?
  - Inability to run processes greater in size than a partition, but smaller then memory.

Advantages

• Internal fragmentation.

4. A program is to run on a multiprogrammed machine. Describe at which points in time during program development

physical memory, the loader can bind the addresses to the correct location within physical memory as the program is loaded (copied) into memory. This process slows loading (increases

The compiler/linker generated addresses are bound to a logical address space (an abstract memory layout). The program executes using the logical addresses independent of where it is loaded into physical memory. At run-time the logical addresses are translated to the appropriate

VM MMU), but it incurs the cost of translation on every memory reference, which can be

physical addresses by specialised hardware on each memory reference.

5. Describe four algorithms for allocating regions of contiguous memory, and comment on their properties.

Breaks up large block at end of memory with any reduction in searching.

Swapping is where a process is brought into main memory in its entirety, run for a while, and then put

address space (virtual address space) is also divided up into matching-sized chunks (pages). Memory is

Swapping allows the OS to run more programs than what would fit in memory if all programs

Swapping is slow as it has to copy the entire program's in-memory image out to disk and back.

Tends to skip over potentially many regions at the start of list.

Pick the closest free region in the entire list.

Find the worst fit in the entire list

startup latency), and increases executable file size (to hold annotations, minor point).

- 3. Assume a system protected with base-limit registers. What are the advantages and problems with such a protected system (compared to either a unprotected system or a paged VM system)?
  - Disadvantages
  - Partitions must be contiguous external fragmentation.
    - Entire process must be in memory. • Cannot practically share memory with other processes.
      - Memory for applications can be allocated dynamically and hardware translates the application (logical) addresses to allocated addresses. Multiple concurrent executions of the same application is possible.

      - Compaction is also possible.

Applications are protected from each other

- to execution time where addresses within the program can be bound to the actual physical memory it uses for
- execution? What are the implication of using each of the three binding times?

  - Compile/Link time binding. The executable itself contains the actual physical addresses it will use during execution. It can only run at one location, and only a single copy can run at a time, unless the executable is

    - recompiled/relinked to a new location in physical memory prior to each execution.
  - Load time binding
  - Addresses within the executable are annotated such that when the program is loaded into
  - Run-time binding
  - Run-time binding is the most flexible (depending of the translation hardware available, e.g. page

significant.

Scan memory region list from start for first fit.

1. First-Fit

Scan memory region list from point of last allocation to next fit.

2. Next-Fit

Tends to leave small unusable regions, and slower due to requirement of search the entire list.

4. Worst-Fit

3. Best-Fit

Slower as it searches the entire list, fragmentation still an issue.

6. What is compaction? Why would it be used?

- Moving all the allocated regions of memory next to each other (e.g. to the bottom of memory) to free up larger contiguous free regions.
- 7. What is swapping? What benefit might it provide? What is the main limitation of swapping?

completely back on disk.

remained resident in memory.

transfered to and from disk in units of pages.

Virtual Memory

- 8. What is Paging? In brief: Paging is where main memory is divided into equal-sized chunks (frames) and the programs
- The lower bits of a virtual address is not translated and passed through the MMU to form a physical

- 12. Given a two-level page table (in physical memory), what is the average number of physical memory accesses per virtual memory access in the case where the TLB has a 100% miss ratio, and the case of a 95% hit ratio

hardware.

1 \* .95 + .05 \* (1 + 2) = 1.1

3

- 13. What are the two broad categories of events causing page faults? What other event might cause page faults?
  - The virtual addresses are 0x00028123, 0x0008a7eb, 0x0005cfff,0x0001c642, 0x0005b888,  $0 \times 00034000$ **TLB**

0x00028200 0x0063f400 0x00034200 0x001fc600

**EntryLo** 

EntryHi

- 0x0005b200 0x002af200 0x0008a100 0x00145600 0x0005c100 0x006a8700 0x0001c200 0x00a97600
  - VAddr PhysAddr 0x00028123 0x0008a7eb
- 0x0005cfff 0x006a8fff Global bit, R/W  $0x0001c642\ 0x00a97642\ R/W$

18. What is temporal and spatial locality?

together in time.

- 0x0005b888 0x002af888 Read-only  $0x00034000\ 0x001fc000\ R/W$
- 15. Describe an inverted page table and how it is used to translate a virtual address into a physical address.
- See page 6-2 of the MIPS R3000 Hardware Guide for revising TLB operation.

**Results** 

Invalid

Access

Invalid (ASID mismatch)

16. Describe a hashed page table and how it is used to translate a virtual address into a physical address.

that are sparsely populated (e.g. many single pages scattered through memory)?

17. Of the three page table types covered in lectures, which ones are most appropriate for large virtual address spaces

Temporal locality: states that recently accessed items are likely to be accessed in the near future.

Spatial locality: says that items whose addresses are near one another tend to be referenced close

best as it is searched via a hash, and not based on the structure of the virtual address space.

The 2-level suffers from internal fragmentation of page table nodes themselves. The IPT and HPT is

- See lecture slides in the virtual memory lecture.

See lecture slides in the virtual memory lecture.

See lecture slides in the virtual memory lecture.

Page faults may be used to set reference and dirty bits on architectures that do not support them in

14. Translate the following virtual addresses to Physical Addresses using the TLB. The system is a R3000. Indicate if

Illegal memory references and access to non-resident pages.

the page is mapped, and if so if its read-only or read/write.

The EntryHi register currently contains 0x00000200.

9. Why do all virtual memory system page sizes have to be a power of 2? Draw a picture.

- address.
- 10. What is a TLB? What is its function? A translation lookaside buffer is an associative cache of page table entries used to speed up the translation of virtual addresses to physical addresses.
- 11. Describe a two-level page table and how it is used to translate a virtual address into a physical address.

1. What is the difference between the different MIPS address space segments? What is the use of each segment?

This question aims to have you refresh your understanding of the MIPS virtual address space layout. See the last 5 slides of the "Virtual Memory" lecture, or the user-manual on the class web site.

2. What functions exist to help you manage the TLB? Describe their use. (Hint: look in kern/arch/mips/include/tlb.h)

```
* MIPS-specific TLB access functions.
*
     tlb_random: write the TLB entry specified by ENTRYHI and ENTRYLO
          into a "random" TLB slot chosen by the processor.
          IMPORTANT NOTE: never write more than one TLB entry with the
          same virtual page field.
     tlb write: same as tlb random, but you choose the slot.
     tlb read: read a TLB entry out of the TLB into ENTRYHI and ENTRYLO.
          INDEX specifies which one to get.
     tlb_probe: look for an entry matching the virtual page in ENTRYHI.
          Returns the index, or a negative number if no matching entry
          was found. ENTRYLO is not actually used, but must be set; 0
          should be passed.
          IMPORTANT NOTE: An entry may be matching even if the valid bit
          is not set. To completely invalidate the TLB, load it with
          translations for addresses in one of the unmapped address
          ranges - these will never be matched.
*/
void tlb_random(uint32_t entryhi, uint32_t entrylo);
void tlb_write(uint32_t entryhi, uint32_t entrylo, uint32_t index);
void tlb_read(uint32_t *entryhi, uint32_t *entrylo, uint32_t index);
int tlb_probe(uint32_t entryhi, uint32_t entrylo);
3. What macros are used to convert from a physical address to a kernel virtual address?
```

\* The first 512 megs of physical space can be addressed in both kseg0 and \* kseq1. We use kseq0 for the kernel. This macro returns the kernel virtual

not using systems with more physical space than that anyway.)

This question aims to refresh your understanding of the relationship between KSEG0 and physical memory.

\* address of a given physical address within that range. (We assume we're

\* The top of user space. (Actually, the address immediately above the

```
* N.B. If you, say, call a function that returns a paddr or 0 on error,
* check the paddr for being 0 *before* you use this macro. While paddr 0 * is not legal for memory allocation or memory management (it holds
 * exception handler code) when converted to a vaddr it's *not* NULL, *is*
* a valid address, and will make a *huge* mess if you scribble on it.
 */
#define PADDR_TO_KVADDR(paddr) paddr_to_kvaddr(paddr)
static inline vaddr_t
paddr_to_kvaddr(paddr_t paddr){
    return ((paddr) + MIPS_KSEG0);
#define KVADDR_TO_PADDR(vaddr) kvaddr_to_paddr(vaddr)
static inline paddr_t
kvaddr_to_paddr(vaddr_t vaddr){
    return ((vaddr) - MIPS_KSEG0);
}
4. What address should the initial user stack pointer be?
```

\* The starting value for the stack pointer at user level. Because

\* last valid user address.)

\* TLB entry fields.

is attempted.

Functions in addrspace.c:

as copy

#define USERSPACETOP MIPS\_KSEG0

```
* the stack is subtract-then-store, this can start as the next
    address after the stack area.
  * We put the stack at the very top of user virtual memory because it
    grows downwards.
 #define USERSTACK
                           USERSPACETOP
  5. What are the entryli and entrylo co-processor registers? Describe their contents.
See the start of the "Virtual Memory II" lecture for an overview of how these registers are used to load the TLB, and the content of various
bit-fields within the register. The follow macros are defined in OS161 for setting and masking values written to or read from the Entry*
registers.
```

Note that the MIPS has support for a 6-bit address space ID. In the interests of simplicity, we don't use it. The fields related to it

ever set by the processor. If you set it, writes are permitted. If you don't set it, you'll get a "TLB Modify" exception when a write

(TLBLO GLOBAL and TLBHI\_PID) can be left always zero, as can the bits that aren't assigned a meaning. The TLBLO\_DIRTY bit is actually a write privilege bit - it is not

```
There is probably no reason in the course of CS161 to use TLBLO_NOCACHE.
 /* Fields in the high-order word */
 #define TLBHI_VPAGE 0xfffff000
          TLBHI_PID
                          0x00000fc0 */
 /* Fields in the low-order word */
 #define TLBLO_PPAGE
                         0xfffff000
 #define TLBLO_NOCACHE 0x00000800
 #define TLBLO_DIRTY 0x00000400
 #define TLBLO_VALID
                           0 \times 00000200
          TLBLO_GLOBAL 0x00000100 */
  6. What do the as_* functions do? Why do we need as_prepare_load() and as_complete_load()?
These functions are used by OS/161 internally to manage the address spaces of processes. These functions are dependent on the data
structures (e.g. region lists and page tables) used to book-keep virtual memory, and thus you have to implement them. A summary of them
is available in the header file.
as_prepare_load() is required as normally code segments are mapped read-only. as_prepare_load() enable writing to the code
segment while the OS loads the code associated with the process. as_complete_load() then removes write permission to the code
segment to revert it back to read-only.
```

as\_create - create a new empty address space. You need to make sure this gets called in all the right places. You may find you want to change the argument list. May

- create a new address space that is an exact copy of

an old one. Probably calls as\_create to get a new empty address space and fill it in, but that's up to

return NULL on out-of-memory error.

```
you.
        as_activate - make curproc's address space the one currently
                       "seen" by the processor.
        as_deactivate - unload curproc's address space so it isn't
                       currently "seen" by the processor. This is used to avoid potentially "seeing" it while it's being
                       destroyed.
        as_destroy - dispose of an address space. You may need to change the way this works if implementing user-level threads.
        as define region - set up a region of memory within the address
        as_prepare_load - this is called before actually loading from an
                       executable into the address space.
        as_complete_load - this is called when loading from an executable
                       is complete.
        as_define_stack - set up the stack region in the address space.
                       (Normally called *after* as_complete_load().) Hands
                       back the initial stack pointer for the new process.
    Note that when using dumbvm, addrspace.c is not used and these
    functions are found in dumbvm.c.
  7. What does vm_fault() do? When is it called? vm_fault() is called for any virtual memory related faults, i.e. TLB misses, writes to read-
     only pages, and accesses to invalid pages. It is responsible for resolving the fault by either returning an error, or loading an
     appropriate TLB entry for the application to continue
  8. Assuming a 2-level hierarchical page table (4k pages), show for the following virtual addresses:
       1. The page number and offset;
       2. the translated address (after any page allocation); and
       3. the contents of the page table after the TLB miss.
The page table is initially empty, with no L2 pages. You may assume that the allocator returns frames in order, so that the first frame
allocated is frame 0, then frames 1, 2, 3, etc.
```

Answers: Start out by determining which bits of the virtual address are used for each component of the virtual address. The level 1 index is the 10

most significant bits (i.e. vaddr >> 22), the level 2 index is the next 10 bits for the top, i.e. (vaddr >> 12) & 0x3ff, the offset bits are

The components of all the virtual addresses are as follows.

800x0

0x0f0

TLB\_VALID

TLB VALID

Virtual Address Level 1 Index Level 2 Index Offset 0x100008 800x0 0x00x100

0x101

0x100

the last 12 bits of the address, i.e. vaddr & 0xfff.

0x0

0x0

0x101008

0x1000f0

0x41000	0x0	0x41	0x000	
0x41b00	0x0	0x41	0xb00	
0x410000	0x1	0x10	0x000	
If we examine the assume we allow			there are onl	y two indexes used, i.e. there are only two pointers to level-2 page tables. Let's
	_	O_KVADDR(1 < O_KVADDR(2 <	, ,	
				of the start of the frame by shifting. The level-2 page tables are walked by the lresses in that range. If you get the types correct, you can then index off the

```
pointers in the level-1 page tables.
 (level1[0x0])[0x100] = (3 << 12)
                                              TLB VALID
                                                              TLB DIRTY
 (level1[0x0])[0x101] = (4 << 12)
                                              TLB_VALID
                                                              TLB_DIRTY
 (level1[0x0])[0x41] = (5 \ll 12)
(level1[0x1])[0x10] = (6 \ll 12)
```

TLB\_DIRTY

TLB DIRTY

/irtual Address	Physical Address
0x100008	0x3008
0x101008	0x4008
0x1000f0	0x30f0
0x41000	0x5000
0x41b00	0x5b00
0x410000	0x6000

9. What C expressions would you use to set or reset the valid bit in a page table entry?

```
pte = pte | TLBLO_VALID /* set */
pte = pte & ~TLBLO VALID /* reset */
```

### Answer:

kprintf() uses a lock to serialise access to the console. If the lock blocks, it context switches, which will call as\_activate and flush the TLB, ejecting any newly inserted entry. Hence, the system ends up in an infinite loop replacing the TLB entry, and then ejecting it via kprintf()

10. What C expression would you use to test if the valid bit is set? if (pte & TLBLO VALID)

Answer:

see Q8

11. How would you extract the 12-bit offset (bits 0 - 11) from the virtual address? offset = vaddr & ~ TLBHI\_VPAGE

Answer:

see Q8

12. How would you convert the 10 most significant bits (22-31) of a virtual address into an index? index = vaddr >> 22

### Answer:

alloc\_kpages() and free\_kpages()

13. How would you convert the next 10 most significant bits (12-21) into an index? index = (vaddr << 10) >> 22

#### Answer:

Using a bump pointer.

14. How would you round down a virtual address to the base of the page? vaddr = vaddr & TLBHI\_VPAGE

### Answer:

free\_kpage() can't sensibly return free memory to a bump pointer allocator (except in the rare case that the last memory allocated is the first free'd), so it simply loses memory that is passed to it.

allowing the process to continue. Large pages have to wait for a longer loading time per fault. • Increases swapping I/O throughput. Given a certain amount of memory, larger pages have higher I/O throughput as the content of a page is stored contiguously on disk, giving sequential access.

• Increases the working set size. Work set is defined as the size of memory associates with all pages accessed within a window of time. With large the pages, the more potential for pages to include significant amounts of unused data, thus working set size generally increases with page

Pre-paging requires predicting the future (which is hard) and the penalty for making a mistake may be

• Only good as a theoretic reference point: The closer a practical algorithm gets to optimal,

• Implementation requires a time stamp to be kept for each page, updated on every reference

size.

o Optimal

• FIFO

• LRU

- - memory, reducing miss rate.

  - Increases page fault latency as a page fault must load the whole page into memory before

2. Why is demand paging generally more prevalent than pre-paging?

3. Describe four replacement policies and compare them.

the better

Easy to implement

■ Impossible to implement

expensive (may page out a needed page for an unneeded page).

• First-in, first-out: Toss the oldest page

Impossible to implement efficiently

■ Toss the first page with a zero reference bit.

- - table entries, thus smaller (and potentially faster to lookup) page table. • Increases TLB coverage as the TLB has a fixed number of entries, thus larger entries cover more
- Decreases number of pages. For a given amount of memory, larger pages result in fewer page
- issue, and round up an application to the next page boundary usually results in a relative small amount wasted compared to the application size itself.
- for un-used memory when allocation is rounded up. Practically speaking, this is not usually an

- Increases internal fragmentation, and the unit of allocation is now greater, hence greater potential
- 1. What effect does increasing the page size have?
- Virtual Memory

- **Questions and Answers**

 Toss the least recently used page Assumes that page that has not been referenced for a long time is unlikely to be referenced in the near future Will work if locality holds

Age of a page is isn't necessarily related to usage

■ Toss the page that won't be used for the longest time

- Clock • Employs a usage or reference bit in the frame table. • Set to one when page is used • While scanning for a victim, reset all the reference bits
- What is thrashing? How can it be detected? What can be done to combat it?
  - Thrashing is where the sum of the working set sizes of all processes exceeds available physical memory and the computer spends more and more time waiting for pages to transfer in and out. It can be detected by monitoring page fault frequency.
- Some processes can be suspended and swapped out to relieve pressure on memory. **Multi-processors**
- 5. What are the advantages and disadvantages of using a global scheduling queue over per-CPU queues? Under which circumstances would you use the one or the other? What features of a system would influence this decision? Global queue is simple and provides automatic load balancing, but it can suffer from contention on the
  - global scheduling queue in the presence of many scheduling events or many CPUs or both. Another disadvantage is that all CPUs are treated equally, so it does not take advantage of hot caches present on the CPU the process was last scheduled on. Per-CPU queues provide CPU affinity, avoid contention on the scheduling queue, however they require
- 6. When does spinning on a lock (busy waiting, as opposed to blocking on the lock, and being woken up when it's
  - free) make sense in a multiprocessor environment? Spinning makes sense when the average spin time spent spinning is less than the time spent blocking

more complex schemes to implement load balancing.

- the lock requester, switching to another process, switching back, and unblocking the lock requester. 7. Why is preemption an issue with spinlocks?
  - Spinning wastes CPU time and indirectly consumes bus bandwidth. When acquiring a lock, an overall system design should minimise time spent spinning, which implies minimising the time a lock holder holds the lock. Preemption of the lock holder extends lock holding time across potentially many time slices, increasing the spin time of lock aquirers.
  - See the lecture notes for code. It improves scalability as the spinning is on a read instruction (not a test-and-set instruction), which allows the spinning to occur on a local-to-the-CPU cached copy of the lock's memory location. Only when the lock changes state is cache coherency traffic required across

8. How does a read-before-test-and-set lock work and why does it improve scalability?

instruction which require exclusive access to the memory across all CPUs.

the bus as a result of invalidating the local copy of the lock's memory location, and the test-and-set

I/O bound processes typically have many short CPU bursts. If such a process is scheduled, it will typically not use up it's time slice. Priorities are recomputed taking into account the consumed CPU, and hence an I/O-bound process will end up having a higher priority than a process that started at the

(reduced) over time, and hence CPU-bound processes also increase in priority is they are not

12. Why would a hypothetical OS always schedule a thread in the same address space over a thread in a different

CPU is consumed when switching from one task to another. This switching does not contribute to application progress. If done very frequently (a very short time slice), a significant portion of available

Programmed I/O is normally less efficient as the CPU is busy waiting for input when input is

15. A device driver routine (e.g. read\_block() from disk) is invoked by the file system code. The data for the

16. Describe how I/O buffering can be formulated as a bounded-buffer producer-consumer problem.

until a free slot is available in the bounded-size buffer cache.

buffer cache entry as clean (consumed), freeing it up for further writes.

Programmed I/O can be more efficient in circumstances where input it frequent enough such that the overhead of interrupts (getting into and out of the interrupt handler) is more than the average time spent busy waiting. A realistic example is fast networking where packets are nearly always available.

filesystem is requested from the disk, but is not yet available. What do device drivers generally do in this scenario?

They block on a synch primitive that is later woken up by the disk interrupt handler when the block is

Take the a file system buffer cache as an example. File system writes are *produced* by application

requests. The OS must select an available entry in the buffer cache (or re-use an existing one), or block

The interrupt handler can be considered a consumer as after the write to disk completes, it marks the

The interrupt handler must not block waiting for a resource, as it indirectly blocks the application that

If the application has the resource the interrupt handler requires (e.g. memory buffers), the system is

17. An example operating system runs its interrupt handlers on the kernel stack of the currently running application.

What restriction does this place on the interrupt handler code? Why is the restriction required?

Context switch is faster. Better locality. If done too often (always) it starves other tasks.

- **Scheduling** 9. What do the terms I/O bound and CPU bound mean when used to describe a process (or thread)?
  - The time to completion of a *CPU-bound* process is largely determined by the amount of CPU time it receives.
- The time to completion of a *I/O-bound* process is largely determined by the time taken to service its I/O requests. CPU time plays little part in the completion time of I/O-bound processes.
- 10. What is the difference between cooperative and pre-emptive multitasking?
- Cooperative the thread specifically release the CPU, pre-emptive the thread has no choice. 11. Consider the multilevel feedback queue scheduling algorithm used in traditional Unix systems. It is designed to
- favour IO bound over CPU bound processes. How is this achieved? How does it make sure that low priority, CPU bound background jobs do not suffer starvation? Note: Unix uses low values to denote high priority and vice versa, 'high' and 'low' in the above text does not refer to the Unix priority value.
- same priority level and which used more CPU cycles. Even a process with low priority will be scheduled eventually, since the priority of processes that are continually scheduled eventually also receive a low or lower priority. Also note that the recorded amount of CPU consumed (that is used to calculate priority) is aged
- 13. Why would a round robin scheduler NOT use a very short time slice to provide good responsive application behaviour?

I/O

scheduled.

address space? Is this a good idea?

CPU will be wasted on scheduler overhead.

reading the status register wainting for further input.

- 14. Describe programmed I/O and interrupt-driven I/O in the case of receiving input (e.g. from a serial port). Which technique normally leads to more efficient use of the CPU? Describe a scenario where the alternative technique is more efficient. For programmed I/O, the CPU waits for input by continuously reading a status register until it indicates input data is ready, afterwhich, the CPU can read the incoming data, and then return to
- For interrupt-driven I/O, the serial port device sends an interrupt to the CPU when data is ready to be read. The interrupt handler is then invoked, which acknowledges receiving the interrupt, reads the data from the device, and returns from the interrupt. The CPU is free for other activities while not interrupt processing.

unavailable.

available.

was running.

deadlocked.