The Fundamentals of Operating Systems

Hung-Wei Tseng



Outline

- Operating systems: virtualizing computers
- Process: the most important abstraction in modern OSs
- Restricted operations: kernel and user modes

The goal of an OS



Operating System







Operating systems: virtualizing computers

The idea: virtualization

- The operating system presents an illusion of a virtual machine to each running program and maintains architectural states of a von Neumann machine
 - Processor
 - Memory
 - I/O
- Each virtualized environment accesses architectural facilities through some sort of application programming interface (API)
- Dynamically map those virtualized resources into physical resources



The idea of an OS: virtualization







Latency v.s. Throughput

- A 4K movie clip using H.265 coding takes **70GB** in storage
- If you want to transfer a total of 2 Peta-Byte video clips (roughly 29959 movies) from UCSD
 - 100 miles from UCR
 - Assume that you have a **100Gbps** ethernet
 - Throughput: 100 Gbits per second
 - 2 Peta-byte (16 Peta-bits) over 167772 seconds = 1.94 Days
 - Latency: first 70GB (first movie) in 6 seconds





	Toyota Prius	
	100 miles from UCSD 75 MPH on highway! 50 MPG Max load: 374 kg = 2,770 hard drives (2TB per drive) = 5.6 PB	
Throughput/ bandwidth	450GB/sec	
latency	3.5 hours	2 Peta-by
response time	You see nothing in the first 3.5 hours	You can sta





yte over 167772 seconds = 1.94 Days

rt watching the first movie as soon as you get a frame!

Process: the most important abstraction in modern operating systems

The idea of an OS: virtualization









- The most important abstraction in modern operating systems.
- A process abstracts the underlying computer.
- A process is a **running program** a dynamic entity of a program.
 - Program is a static file/combination of instructions
 - Process = program + states
 - The states evolves over time
- A process may be dynamically switched out/back during the execution

Virtualization

- The operating system presents an illusion of a virtual machine to each running program — process
 - Each virtual machine contains architectural states of a von Neumann machine
 - Processor
 - Memory
 - I/O
- Each virtualized environment accesses architectural facilities through some sort of application programming interface (API)
- Dynamically map those virtualized resources into physical resources — policies, mechanisms

system calls

Demo: Virtualization

 We use the getcpu system call to identify the processor and node on which the calling process (can be a thread as well) is currently running and writes them into the integers pointed to by the cpu and node arguments.



Demo: Virtualization

- Some processes may use the same processor
- Each process has the same address for variable a, but different values.
- You may see the content of a compiled program using objdump



What happens when creating a process



Dynamic allocated data: malloc()

Local variables, arguments

> Linux contains a .bss section for uninitialized global variables







The illusion provided by processes



only a few of them are physically executing/using the installed DRAM.





What the OS must track for a process?

- Which of the following information does the OS need to track for each process?
 - A. Stack pointer
 - B. Program counter
 - C. Process state
 - D. Registers

E. All of the above

 You also need to keep other process information like an unique process id, process states, I/O status, and etc...



Process control block

- OS has a PCB for each process
- Sometimes called Task Controlling Block, Task Struct, or Switchframe
- The data structure in the operating system kernel containing the information needed to manage a particular process.
- The PCB is the manifestation of a process in an operating system



Example: struct task_struct in Linux

Process state struct task_struct { volatile long state; /* -1 unrunnable, 0 runnable, >0 stopped */ vold *stack; atomic_t usage; unsigned int flags; /* per process flags, defined below */ unsigned int ptrace; int on_rq; int prio, static_prio, normal_prio; const struct sched_class *sched_class; struct sched_entity se; struct sched_rt_entity rt; unsigned int policy; int nr_cpus_allowed; Process ID cpumask_t cpus_allowed; pid t pid; struct task_struct __rcu *real_parent; struct task_struct __rcu *parent; struct list_head children; struct list head sibling; struct list_nead tasks; Virtual memory pointers struct mm_struct *mm, *active_mm; Low-level architectural states /* CPU-specific state of this task */ struct thread_struct thread; You may find this struct in /usr/src/linux-headers-x.x.x-xx/include/linux/sched.h } 33



Memory pointers in struct mm_struct



Processor states in struct thread_struct

_				_				
	struct	struct thread_struct {						
		struct desc_struct	ruct desc_struct tls_array[GDT_ENTRY_TLS_ENT					
		unsigned long	sp0;					
		unsigned long	sp;					
	#ifdef	CONFIG_X86_32						
		unsigned long	sysenter_cs;					
	#else							
		unsigned short	es;					
		unsigned short	ds;					
		unsigned short	fsindex;					
		unsigned short	gsindex;					
	#endif							
	#ifdef	CONFIG_X86_32						
		unsigned long	ip;	Pro				
	#endif							
	#ifdef CONFIG_X86_64							
		unsigned long	fs;					
	#endif							
		unsigned long struct perf_event	gs; *ptrace_bps[HBP_NUM];					
		unsigned long	debugreg6;					
		unsigned long unsigned long	ptrace_dr7; cr2;					
unsigned long		unsigned long	trap_nr;					
		unsigned long	error_code;					
#ifdef CONFIG_VM86								
struct vm86		struct vm86	*vm86;					
#en	dif							
unsigned long		unsigned long	<pre>*io_bitmap_ptr;</pre>					
		unsigned long	iopl;					
unsigned		unsigned	io_bitmap_max;					
		struct fpu	fpu;					

Some x86 Register values

rogram counter

Restricted operations: kernel and user modes

Restricted operations

- Most operations can directly execute on the processor without OS's intervention
- The OS only takes care of protected resources, change running processes or anything that the user program cannot handle properly
- Divide operations into two modes
 - User mode
 - Restricted operations
 - User processes
 - Kernel mode
 - Can perform privileged operations
 - The operating system kernel
- Requires architectural/hardware supports



How applications can use privileged operations?

- Through the API: System calls
- Implemented in "trap" instructions
 - Raise an exception in the processor
 - The processor saves the exception PC and jumps to the corresponding exception handler in the OS kernel



user mode

kernel/privileged mode

Architectural support: privileged instructions

- The processor provides
 normal instructions and privileged
 instructions
 - Normal instructions: ADD, SUB, MUL, and etc ...
 - Privileged instructions: HLT, CLTS, LIDT, LMSW, SIDT, ARPL, and etc...
- The processor provides different modes
 - User processes can use normal instructions
 - Privileged instruction can only be used if the processor is in proper mode

Least privileged

Most privileged

Ring 3

Ring 2

Ring 1

Ring 0

Kernel

Device Drivers

Device Drivers

Applications

user mode

Latency Numbers Every Programmer Should Know

Operations	Latency (ns)	Latency (us)	Latency (ms)	
L1 cache reference	0.5 ns			~ 1 CPU cycle
Branch mispredict	5 ns			
L2 cache reference	7 ns			14x L1 cache
Mutex lock/unlock	25 ns			
Main memory reference	100 ns			20x L2 cache, 200x L1 cache
Compress 1K bytes with Zippy	3,000 ns	3 us		
Send 1K bytes over 1 Gbps network	10,000 ns	10 us		
Read 4K randomly from SSD*	150,000 ns	150 us		~1GB/sec SSD
Read 1 MB sequentially from memory	250,000 ns	250 us		
Round trip within same datacenter	500,000 ns	500 us		
Read 1 MB sequentially from SSD*	1,000,000 ns	1,000 us	1 ms	~1GB/sec SSD, 4X memory
Disk seek	10,000,000 ns	10,000 us	10 ms	20x datacenter roundtrip
Read 1 MB sequentially from disk	20,000,000 ns	20,000 us	20 ms	80x memory, 20X SSD
Send packet CA-Netherlands-CA	150,000,000 ns	150,000 us	150 ms	

Demo: Kernel Switch Overhead

 Measure kernel switch overhead using Imbench http:// www.bitmover.com/lmbench/

