**ADRIAN PERRIG & TORSTEN HOEFLER**

# Networks and Operating Systems (252-0062-00)
# Chapter 4: Synchronization

IT TOOK A LOT OF WORK, BUT THIS LATEST LINUX PATCH ENABLES SUPPORT FOR MACHINES WITH 4,096 CPUs, UP FROM THE OLD LIMIT OF 1,024.

DO YOU HAVE SUPPORT FOR SMOOTH FULL-SCREEN FLASH VIDEO YET?

NO, BUT WHO USES *THAT*?

Source: xkcd

---

## Real Time

2015-02-26

# Real-time scheduling

- **Problem: giving real time-based guarantees to tasks**
  - Tasks can appear at any time
  - Tasks can have deadlines
  - Execution time is generally known
  - Tasks can be periodic or aperiodic

- **Must be possible to reject tasks which are unschedulable, or which would result in no feasible schedule**

3

# Example: multimedia scheduling

Starting moment for A1, B1, C1 — Deadline for A1 — Deadline for B1 — Deadline for C1

A: A1, A2, A3, A4, A5
B: B1, B2, B3, B4
C: C1, C2, C3

Time (msec) 0 10 20 30 40 50 60 70 80 90 100 110 120 130 140

4

2

# Rate-monotonic scheduling

- **Schedule periodic tasks by always running task with shortest period first.**
  - Static (offline) scheduling algorithm

- **Suppose:**
  - m tasks
  - $C_i$ is the execution time of i'th task
  - $P_i$ is the period of i'th task

- **Then RMS will find a feasible schedule if:**

$$\sum_{i=1}^{m} \frac{C_i}{P_i} \le m(2^{1/m} - 1)$$

- **(Proof is beyond scope of this course)**

5

# Earliest Deadline First

- **Schedule task with earliest deadline first (duh..)**
  - Dynamic, online.
  - Tasks don't *actually* have to be periodic…
  - More complex - O(n) – for scheduling decisions

- **EDF will find a feasible schedule if:**

$$\sum_{i=1}^{m} \frac{C_i}{P_i} \le 1$$

- **Which is very handy.  Assuming zero context switch time…**

6

**ETH**zürich

## Guaranteeing processor rate

- **E.g. you can use EDF to guarantee a rate of progress for a long-running task**
  - Break task into periodic jobs, period $p$ and time $s$.
  - A task arrives at start of a period
  - Deadline is the end of the period

- **Provides a *reservation* scheduler which:**
  - Ensures task gets $s$ seconds of time every $p$ seconds
  - Approximates weighted fair queuing

- **Algorithm is regularly rediscovered…**

7

**ETH**zürich

## Multiprocessor Scheduling

8

**ETH**zürich

## Challenge 1: sequential programs on multiprocessors

- **Queuing theory ⇒ straightforward, although:**
  - More complex than uniprocessor scheduling
  - Harder to analyze

Task queue

Core 0
Core 1
Core 2
Core 3

But…

9

**ETH**zürich

## It's much harder

- **Overhead of locking and sharing queue**
  - Classic case of scaling bottleneck in OS design

- **Solution: per-processor scheduling queues**

Core 0
Core 1
Core 2
Core 3

In practice, each
is more complex
e.g. MFQ

10

**ETH**zürich

# It's much harder

- **Threads allocated arbitrarily to cores**
  - ⇒ tend to move between cores
  - ⇒ tend to move between caches
  - ⇒ really bad locality and hence performance

- **Solution: affinity scheduling**
  - Keep each thread on a core most of the time
  - Periodically rebalance across cores
  - Note: this is *non-work-conserving!*

- **Alternative: hierarchical scheduling (Linux)**

11

**ETH**zürich

# Challenge 2: parallel applications

- **Global barriers in parallel applications ⇒
  One slow thread has huge effect on performance**
  - Corollary of *Amdahl's Law*

- **Multiple threads would benefit from cache sharing**

- **Different applications pollute each others' caches**

- **Leads to concept of "co-scheduling"**
  - Try to schedule all threads of an application together

- **Critically dependent on *synchronization* concepts**

12

**ETH**zürich

# Multicore scheduling

- **Multiprocessor scheduling is two-dimensional**
  - When to schedule a task?
  - Where (which core) to schedule on?

- **General problem is NP hard** ☹

- **But it's worse than that:**
  - Don't want a process holding a lock to sleep
    ⇒ Might be other running tasks spinning on it
  - Not all cores are equal

- **In general, this is a wide-open research problem**

13

**ETH**zürich

# Little's Law

- **Assume, in a train station:**
  - 100 people arrive per minute
  - Each person spends 15 minutes in the station
  - How big does the station have to be (house how many people)
- **Little's law: "*The average number of active tasks in a system is equal to the average arrival rate multiplied by the average time a task spends in a system*"**

14

## Our Small Quiz

- **True or false (raise hand)**
  - Throughput is an important goal for batch schedulers
  - Response time is an important goal for batch schedulers
  - Realtime schedulers schedule jobs faster than batch schedulers
  - Realtime schedulers have higher throughput than batch schedulers
  - The scheduler has to be invoked by an application
  - FCFS scheduling has low average waiting times
  - Starvation can occur in FCFS scheduling
  - Starvation can occur in SJF scheduling
  - Preemption can be used to improve interactivity
  - Round Robin scheduling is fair
  - Multilevel Feedback Queues in Linux prevent starvation
  - Simple Unix scheduling fairly allocates the time to each user
  - RMS scheduling achieves full CPU utilization
  - Multiprocessor scheduling is NP hard

15

## Last time: Scheduling

- **Basics:**
  - Workloads, tradeoffs, definitions
- **Batch-oriented scheduling**
  - FCFS, Convoys, SJF, Preemption: SRTF
- **Interactive workloads**
  - RR, Priority, Multilevel Feedback Queues, Linux, Resource containers
- **Realtime**
  - RMS, EDF
- **Multiprocessors**

- **This time: OSPP Section 5 (not including IPC)**
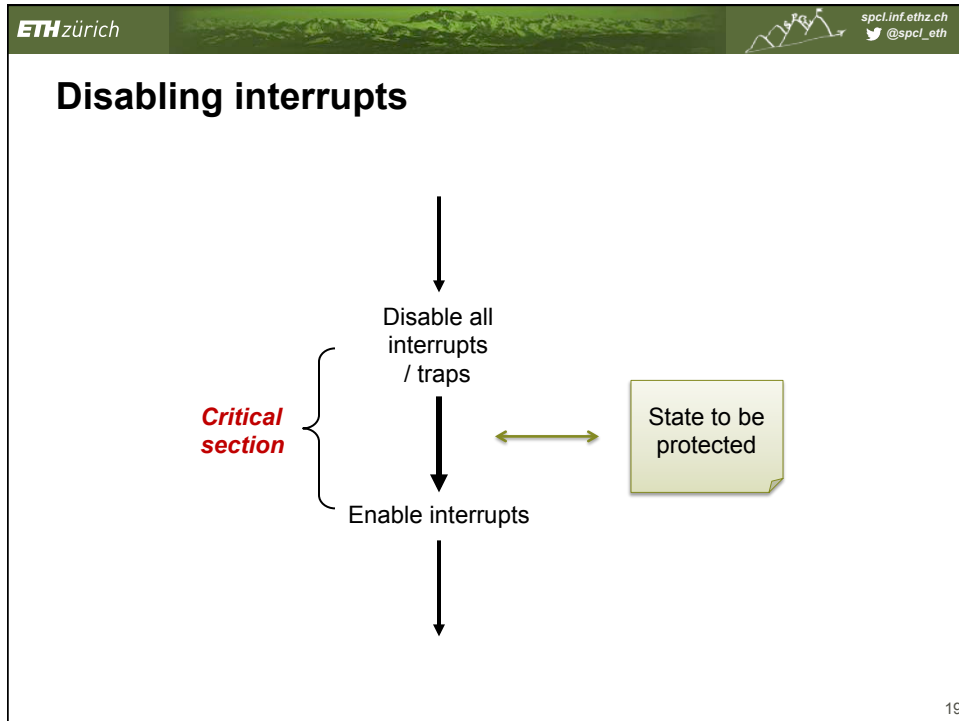
16

## Goals today

- **Overview of inter-process communication systems**
  - Hardware support
  - With shared memory
  - Without shared memory
  - Upcalls

- **Generally: very broad field**
  - Quite competitive… especially with microkernels

17

**Recap: Hardware support for synchronization**

18

**ETH**zürich

## Disabling interrupts

Disable all
interrupts
/ traps

***Critical
section***

Enable interrupts

State to be
protected

19

---

**ETH**zürich

## Disabling interrupts

- **Nice and simple**
- **Can't be rescheduled inside critical section**
  **⇒ data can't be altered by anything else**
- **Except…**
- **Another processor!**
  - Hmm….

- **Very efficient if in kernel on a *uniprocessor*.**

20

**ETH**zürich

## Test-And-Set instruction

- **Atomically:**
  - Read the value of a memory location
  - Set the location to 1

- **Available on some hardware (e.g., PA-RISC)**
  - (actually, more a RAC – Read-And-Clear)

21

**ETH**zürich

## Compare-And-Swap (CAS)

```
word cas(word *flag, word oldval, word newval) {
        atomically {
                if (*flag == oldval) {
                        *flag = newval;
                        return oldval;
                } else {
                        return *flag;
                }
        }
}
```

- **Available on e.g., x86, IBM/370, SPARC, ARM,…**
- **Theoretically, *slightly* more powerful than TAS**
  - Why?
  - Other variants e.g., CAS2, etc.

22

ETH*zürich*

# Load-Link, Store-Conditional

**Factors CAS, etc. into two instructions:**

1. **LL: load from a location and mark as "owned"**
2. **SC: Atomically:**
   1. Store *only* if already marked by this processor
   2. Clear any marks set by other processors
   3. Return whether it worked.

**Available on PPC, Alpha, MIPS, etc…**

23

---

ETH*zürich*

# Back to TAS…



old = TAS(flag)
if (old == True)

Spin forever waiting?

Critical section

flag ← False

24

# Spinning

- **On a uniprocessor:**
  - Not much point in spinning at all.  What's going to happen?
  - Possibly an interrupt

- **On a multiprocessor:**
  - Can't spin forever
  - Another spin is always cheap
  - Blocking thread and rescheduling is expensive
  - Spinning only works if lock holder is running on another core

25

# Competitive spinning

- **How long to spin for?**

- **"Competitive spinning":**
  - Within a factor of 2 of optimal, offline (i.e., impossible!) algorithm

- **Good approach: spin for the context switch time**
  - Best case: avoid context switch entirely
  - Worst case: twice as bad as simply rescheduling

26

# IPC with shared memory

27

## Techniques you already know ☺

- **Semaphores**
  - P, V operations

- **Mutexes**
  - Acquire, Release

- **Condition Variables**
  - Wait, Signal (Notify), Broadcast (NotifyAll)

- **Monitors**
  - Enter, Exit

28

**ETH**zürich

# Focus here: interaction with scheduling

- **Most OSes provide some form of these**

- **Key issue not yet covered: interaction between scheduling and synchronization**

- **Example: Priority inversion**
  - Assuming a priority scheduler, e.g., Unix, Windows

29

---

**ETH**zürich

# Priority Inversion

High priority

Acquire
lock

Low priority

Time

30

---

---

## Priority Inheritance

- **Process holding lock *inherits* priority of highest priority process that is waiting for the lock.**
  - Releasing lock $\Rightarrow$ priority returns to previous value
  - Ensures forward progress

- **Alternative: *Priority Ceiling***
  - Process holding lock acquires priority of highest-priority process that *can ever* hold lock
  - Requires static analysis, used in embedded RT systems

Priority Inheritance



IPC without shared memory

2015-02-26

# Asynchronous (buffered) IPC

Process 1 → RECV ⊗ →

Receiver blocks waiting for msg

Sender does not block

Process 2 → SEND →

Time

43

# Synchronous (unbuffered) IPC

Process 1 → RECV →

Process 2 → SEND ⊗ →

Time

44

Synchronous (unbuffered) IPC



Duality of messages and shared-memory

- **Famous claim by Lauer and Needham (1978):**

  *Any shared-memory system (e.g., one based on monitors and condition variables) is equivalent to a non-shared-memory system (based on messages)*

- **Exercise: pick your favourite example of one, and show how to build the dual.**

## Unix Pipes

- **Basic (first) Unix IPC mechanism**

- **Unidirectional, buffered communication channel between two processes**

- **Creation:**

  ```
  int pipe(int pipefd[2])
  ```

- **Q. How to set up pipe between two processes?**

- **A. Don't!  Create the pipe first, then fork…**

47

## Pipe idiom (man 2 pipe)

```
int
main(int argc, char *argv[])
{
    int pipefd[2];
    pid_t cpid;
    char buf;

    assert(argc == 2);

    if (pipe(pipefd) == -1) {
        perror("pipe");
        exit(EXIT_FAILURE);
    }

    cpid = fork();
    if (cpid == -1) {
        perror("fork");
        exit(EXIT_FAILURE);
    }

    if (cpid == 0) {    /* Child reads from pipe */
        close(pipefd[1]);          /* Close unused write end */

        while (read(pipefd[0], &buf, 1) > 0)
            write(STDOUT_FILENO, &buf, 1);

        write(STDOUT_FILENO, "\n", 1);
        close(pipefd[0]);
        _exit(EXIT_SUCCESS);

    } else {            /* Parent writes argv[1] to pipe */
        close(pipefd[0]);          /* Close unused read end */
        write(pipefd[1], argv[1], strlen(argv[1]));
        close(pipefd[1]);          /* Reader will see EOF */
        wait(NULL);                /* Wait for child */
        exit(EXIT_SUCCESS);
    }
}
```

Create a pipe

48

24

## Pipe idiom (man 2 pipe)

```
int
main(int argc, char *argv[])
{
    int pipefd[2];
    pid_t cpid;
    char buf;

    assert(argc == 2);

    if (pipe(pipefd) == -1) {
        perror("pipe");
        exit(EXIT_FAILURE);
    }

    cpid = fork();
    if (cpid == -1) {
        perror("fork");
        exit(EXIT_FAILURE);
    }

    if (cpid == 0) {     /* Child reads from pipe */
        close(pipefd[1]);          /* Close unused write end */

        while (read(pipefd[0], &buf, 1) > 0)
            write(STDOUT_FILENO, &buf, 1);

        write(STDOUT_FILENO, "\n", 1);
        close(pipefd[0]);
        _exit(EXIT_SUCCESS);

    } else {            /* Parent writes argv[1] to pipe */
        close(pipefd[0]);          /* Close unused read end */
        write(pipefd[1], argv[1], strlen(argv[1]));
        close(pipefd[1]);          /* Reader will see EOF */
        wait(NULL);                /* Wait for child */
        exit(EXIT_SUCCESS);
    }
}
```

Read from pipe and write to standard output until EOF

51

## Pipe idiom (man 2 pipe)

```
int
main(int argc, char *argv[])
{
    int pipefd[2];
    pid_t cpid;
    char buf;

    assert(argc == 2);

    if (pipe(pipefd) == -1) {
        perror("pipe");
        exit(EXIT_FAILURE);
    }

    cpid = fork();
    if (cpid == -1) {
        perror("fork");
        exit(EXIT_FAILURE);
    }

    if (cpid == 0) {     /* Child reads from pipe */
        close(pipefd[1]);          /* Close unused write end */

        while (read(pipefd[0], &buf, 1) > 0)
            write(STDOUT_FILENO, &buf, 1);

        write(STDOUT_FILENO, "\n", 1);
        close(pipefd[0]);
        _exit(EXIT_SUCCESS);

    } else {            /* Parent writes argv[1] to pipe */
        close(pipefd[0]);          /* Close unused read end */
        write(pipefd[1], argv[1], strlen(argv[1]));
        close(pipefd[1]);          /* Reader will see EOF */
        wait(NULL);                /* Wait for child */
        exit(EXIT_SUCCESS);
    }
}
```

In parent: close read end and write argv[1] to pipe

52

26

**ETH**zürich

# Unix shell pipes

- **E.g.:**

```
curl --silent http://spcl.inf.ethz.ch/Teaching/2014-osnet/ | sed
  's/[^A-Za-z]/\n/g' | sort -fu | egrep -v '^\s*$' | wc -l
```

- **Shell forks each element of the pipeline**
  - Each process connected via pipes
  - Stdout of process $n$ → stdin of process $n+1$
  - Each process then exec's the appropriate command
  - Exercise: write it! (hint: 'man dup2'…)

53

**ETH**zürich

# Messaging systems

- **A good textbook will examine options:**
  - End-points may or may not know each others' names
  - Messages might need to be sent to more than one destination
  - Multiple arriving messages might need to be demultiplexed
  - Can't wait forever for one particular message

- **BUT: you'll see most of this somewhere else!**
  - In networking
  - Many parallels between message-passing operating systems and networks

54

**ETH**zürich

## Example

- **The concept of a "port" allows:**
  - Naming of different end-points within a process
  - Demultiplexing of messages
  - Waiting selectively for different kinds of messages

- **Analogous to "socket" and "TCP port" in IPv4**
  - In Unix, "Unix domain sockets" do exactly this.
  - **`int s = socket(AF_UNIX, type, 0);`**

55

**ETH**zürich

## Naming pipes

- **Pipes so far are only named by their descriptors**
  - Namespace is *local* to the process
  - Copied on **`fork()`**

- **How to put a pipe in the global namespace?**
  - Make it a "named pipe"
  - Special file of type "pipe" (also known as a FIFO)

56

**ETH**zürich

## Named pipes

```
htor@lenny:
htor@lenny:~$ mkfifo /tmp/fifo
htor@lenny:~$ echo "Hello" > /tmp/fifo
htor@lenny:~$
```

```
htor@lenny: ~
htor@lenny:~$ cat /tmp/fifo
Hello
htor@lenny:~$
```

57

---

**ETH**zürich

## Local Remote Procedure Call

- **Can use RPC locally:**
    - Define procedural interface in an IDL
    - Compile / link stubs
    - Transparent procedure calls over messages

- **Naïve implementation is slow**
    - Lots of things (like copying) don't matter with a network, but do matter between local processes
    - Can be made very fast: more in the AOS course…

58

29

**ETH** zürich

# Unix signals

- *Asynchronous* notification from the kernel

- **Receiver doesn't wait: signal just happens**

- **Interrupt process, and:**
  - Kill it
  - Stop (freeze) it
  - Do "something else" (see later)

59

---

**ETH** zürich

# Signal types (some of them)

| Name | Description / meaning | Default action |
|------|----------------------|----------------|
| SIGHUP | Hangup / death of controlling process | Terminate process |
| SIGINT | Interrupt character typed (CTRL-C) | Terminate process |
| SIGQUIT | Quit character typed (CTRL-\) | Core dump |
| SIGKILL | `kill -9 <process id>` | **Terminate process** |
| SIGSEGV | Segfault (invalid memory reference) | Core dump |
| SIGPIPE | Write on pipe with no reader | Terminate process |
| SIGALRM | `alarm()` goes off | Terminate process |
| SIGCHLD | Child process stopped or terminated | Ignored |
| SIGSTOP | Stop process | Stop |
| SIGCONT | Continue process | |
| SIGUSR1,2 | User-defined signals | Terminate process |

"Hanging up" the phone (terminal)

Can't be disabled!

E.g., after other side of pipe has closed it

Used by debuggers (e.g., `gdb`) and shell (`CTRL-Z`)

Etc. – see `man 7 signal` for the full list

60

30

**ETH**zürich

## Where do signals come from?

- **Memory management subsystem:**
  - **SIGSEGV**, etc.

- **IPC system**
  - **SIGPIPE**

- **Other user processes**
  - **SIGUSR1,2**, **SIGKILL**, **SIGSTOP**, **SIGCONT**

- **Kernel trap handlers**
  - **SIGFPE**

- **The "TTY Subsystem"**
  - **SIGINT**, **SIGQUIT**, **SIGHUP**

61

**ETH**zürich

## Sending a signal to a process

- **From the Unix shell:**
  ```
  $ kill –HUP 4234
  ```

- **From C:**
  ```
  #include <signal.h>
  int kill(pid_t pid, int signo);
  ```

- **"Kill" is a rather unfortunate name** ☹

62

**ETH**zürich

## Unix signal handlers

- **Change what happens when a signal is delivered:**
  - Default action
  - Ignore signal
  - Call a user-defined function in the process
    → the *signal handler*

- **Allows signals to be used like "user-space traps"**

63

---

**ETH**zürich

## Oldskool: `signal()`

- **Test your C parsing skills:**

  ```
  #include <signal.h>

  void (*signal(int sig, void (*handler)(int))) (int);
  ```

- **What does this mean?**

64

## Oldskool: `signal()`

```
void (*signal(int sig, void (*handler)(int))) (int);
```

- **Unpacking this:**
  - A handler looks like
    ***void my_handler(int);***
  - Signal takes two arguments…
    *An integer (the signal type, e.g. SIGPIPE)*
    *A pointer to a handler function*
  - … and returns a pointer to a handler function
    *The previous handler,*

- **"Special" `handler` arguments:**
  - `SIG_IGN` (ignore), `SIG_DFL` (default), `SIG_ERR` (error code)

65

## Unix signal handlers

- **Signal handler can be called at *any time!***

- **Executes on the current user stack**
  - If process is in kernel, may need to retry current system call
  - Can also be set to run on a different (alternate) stack

⇒ **User process is in *undefined* state when signal delivered**

66

# Implications

- **There is very little you can safely do in a signal handler!**
  - Can't safely access program global or static variables
  - Some system calls are *re-entrant*, and can be called
  - Many C library calls cannot (including `_r` variants!)
  - Can sometimes execute a `longjmp` if you are careful
  - With `signal`, cannot safely change signal handlers…

- **What happens if another signal arrives?**

# Multiple signals

- **If multiple signals of the *same* type are to be delivered, Unix will *discard all but one.***

- **If signals of *different* types are to be delivered, Unix will deliver them *in any order*.**

- **Serious concurrency problem:
  How to make sense of this?**

**ETH**zürich

# A better `signal()` POSIX `sigaction()`

New action for signal `signo`

```
#include <signal.h>

int sigaction(int signo,
              const struct sigaction *act,
              struct sigaction *oldact);

struct sigaction {
    void (*sa_handler)(int);
    sigset_t    sa_mask;
    int         sa_flags;
    void (*sa_sigaction)(int, siginfo_t *, void *);
};
```

Previous action is returned

Signal handler

Signals to be blocked in this handler (cf., `fd_set`)

More sophisticated signal handler (depending on flags)

69

---

**ETH**zürich

# Signals as upcalls

- **Particularly specialized (and complex) form of *Upcall***
  - Kernel RPC to user process

- **Other OSes use upcalls much more heavily**
  - Including Barrelfish
  - "Scheduler Activations": dispatch every process using an upcall instead of return

- ***Very* important structuring concept for systems!**

70