**From Processes to Threads**

---

**Processes, Threads and Processors**

- Hardware can interpret N instruction streams at once
  - ➢ Uniprocessor, N==1
  - ➢ Dual-core, N==2
  - ➢ Sun's Niagara T2 (2007) N == 64, but 8 groups of 8
- An OS can run 1 process on each processor at the same time
  - ➢ Concurrent execution increases perforamnce
- An OS can run 1 thread on each processor at the same time

## Processes and Threads

- Process abstraction combines two concepts
  - Concurrency
    - Each process is a sequential execution stream of instructions
  - Protection
    - Each process defines an address space
    - Address space identifies all addresses that can be touched by the program
- Threads
  - Key idea: separate the concepts of concurrency from protection
  - A thread is a sequential execution stream of instructions
  - A process defines the address space that may be shared by multiple threads
  - Threads can execute on different cores on a multicore CPU (parallelism for performance) and can communicate with other threads by updating memory

## The Case for Threads

Consider the following code fragment
for(k = 0; k < n; k++)
    a[k] = b[k] * c[k] + d[k] * e[k];

Is there a missed opportunity here? On a Uni-processor?
On a Multi-processor?

## The Case for Threads

Consider a Web server
  get network message (URL) from client
  get URL data from disk
  compose response
  send response

How well does this web server perform?

## Programmer's View

```
void fn1(int arg0, int arg1, …) {…}

main() {
    …
    tid = CreateThread(fn1, arg0, arg1, …);
    …
}
```

At the point CreateThread is called, execution continues in parent thread in main function, and execution starts at fn1 in the child thread, *both in parallel  (concurrently)*

## Introducing Threads

- A thread represents an abstract entity that executes a sequence of instructions
  - It has its own set of CPU registers
  - It has its own stack
  - There is no thread-specific heap or data segment (unlike process)

- Threads are lightweight
  - Creating a thread more efficient than creating a process.
  - Communication between threads easier than btw. processes.
  - Context switching between threads requires fewer CPU cycles and memory references than switching processes.
  - Threads only track a subset of process state (share list of open files, pid, …)

- Examples:
  - OS-supported: Windows' threads, Sun's LWP, POSIX threads
  - Language-supported: Modula-3, Java
    - These are possibly going the way of the Dodo

## Context switch time for which entity is greater?

1. Process
2. Thread

## How Can it Help?

- How can this code take advantage of 2 threads?

  ```
  for(k = 0; k < n; k++)
      a[k] = b[k] * c[k] + d[k] * e[k];
  ```

- Rewrite this code fragment as:

  ```
  do_mult(l, m) {
      for(k = l; k < m; k++)
          a[k] = b[k] * c[k] + d[k] * e[k];
  }
  main() {
      CreateThread(do_mult, 0, n/2);
      CreateThread(do_mult, n/2, n);
  ```
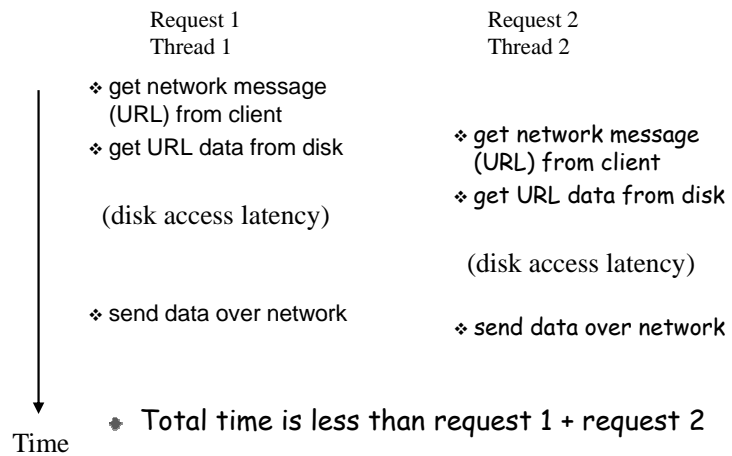
- What did we gain?

## How Can it Help?

- Consider a Web server

  Create a number of threads, and for each thread do
  - ❖ get network message from client
  - ❖ get URL data from disk
  - ❖ send data over network

- What did we gain?

## Overlapping Requests (Concurrency)

Request 1
Thread 1

Request 2
Thread 2

❖ get network message
  (URL) from client
❖ get URL data from disk

   (disk access latency)

❖ get network message
  (URL) from client
❖ get URL data from disk

   (disk access latency)

❖ send data over network

❖ send data over network

Time

* Total time is less than request 1 + request 2

---

**Threads have their own...?**

1. CPU
2. Address space
3. PCB
4. Stack
5. Registers

## Threads vs. Processes

Threads

- A thread has no data segment or heap
- A thread cannot live on its own, it must live within a process
- There can be more than one thread in a process, the first thread calls main & has the process's stack
- If a thread dies, its stack is reclaimed
- Inter-thread communication via memory.
- Each thread can run on a different physical processor
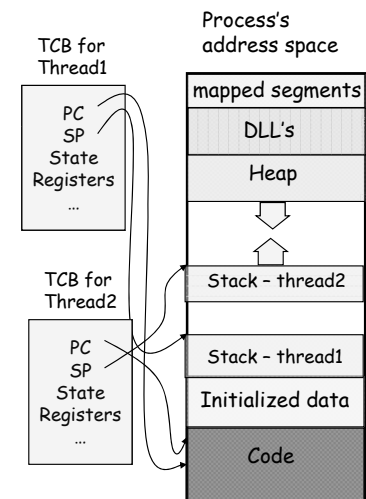- Inexpensive creation and context switch

Processes

- A process has code/data/heap & other segments
- There must be at least one thread in a process
- Threads within a process share code/data/heap, share I/O, but each has its own stack & registers
- If a process dies, its resources are reclaimed & all threads die
- Inter-process communication via OS and data copying.
- Each process can run on a different physical processor
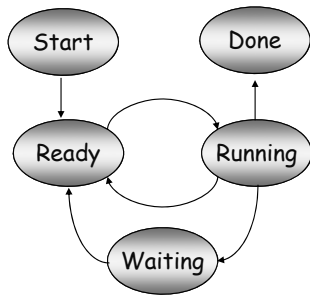- Expensive creation and context switch

13

## Implementing Threads

- Processes define an address space; threads share the address space

- Process Control Block (PCB) contains process-specific information
  - Owner, PID, heap pointer, priority, active thread, and pointers to thread information

- Thread Control Block (TCB) contains thread-specific information
  - Stack pointer, PC, thread state (running, …), register values, a pointer to PCB, …

TCB for Thread1

PC
SP
State
Registers
…

TCB for Thread2

PC
SP
State
Registers
…

Process's address space

mapped segments

DLL's

Heap

Stack – thread2

Stack – thread1

Initialized data

Code

14

## Threads' Life Cycle

- Threads (just like processes) go through a sequence of *start*, *ready*, *running*, *waiting*, and *done* states

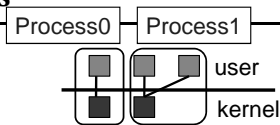## Threads have the same scheduling states as processes

1. True
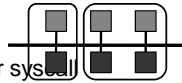2. False

## User-level vs. Kernel-level threads

Process0 — Process1

- ● User-level threads (M to 1 model)
  - ➢ + Fast to create and switch
  - ➢ + Natural fit for language-level threads
  - ➢ - All user-level threads in process block on OS calls
    - ❖ E.g., read from file can block all threads
  - ➢ -User-level scheduler can fight with kernel-level scheduler

- ● Kernel-level threads (1 to 1 model)
  - ➢ + Kernel-level threads do not block process for syscall
  - ➢ + Only one scheduler (and kernel has global view)
  - ➢ - Can be difficult to make efficient (create & switch)

## Languages vs. Systems

- ● Kernel-level threads have won for systems
  - ➢ Linux, Solaris 10, Windows
  - ➢ pthreads tends to be kernel-level threads
- ● User-level threads still used for languages (Java)
  - ➢ User tells JVM how many underlying system threads
    - ❖ Default: 1 system thread
  - ➢ Java runtime intercepts blocking calls, makes them non-blocking
  - ➢ JNI code that makes blocking syscalls can block JVM
  - ➢ JVMs are phasing this out because kernel threads are efficient enough and intercepting system calls is complicated
- ● Kernel-level thread vs. process
  - ➢ Each process requires its own page table & hardware state (significant on the x86)

## Latency and Throughput

- Latency: time to complete an operation
- Throughput: work completed per unit time
- Multiplying vector example: reduced latency
- Web server example: increased throughput
- Consider plumbing
  - Low latency: turn on faucet and water comes out
  - High bandwidth: lots of water (e.g., to fill a pool)
- What is "High speed Internet?"
  - Low latency: needed to interactive gaming
  - High bandwidth: needed for downloading large files
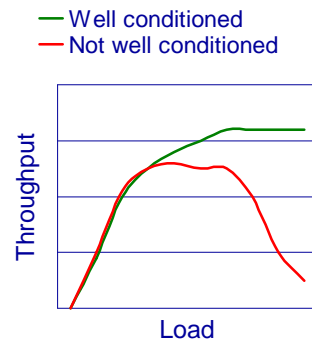  - Marketing departments like to conflate latency and bandwidth…

## Relationship between Latency and Throughput

- Latency and bandwidth only loosely coupled
  - Henry Ford: assembly lines increase bandwidth without reducing latency
- My factory takes 1 day to make a Model-T ford.
  - But I can start building a new car every 10 minutes
  - At 24 hrs/day, I can make 24 * 6 = 144 cars per day
  - A special order for 1 green car, still takes 1 day
  - Throughput is increased, but latency is not.
- Latency reduction is difficult
- Often, one can buy bandwidth
  - E.g., more memory chips, more disks, more computers
  - Big server farms (e.g., google) are high bandwidth

## Thread or Process Pool

- Creating a thread or process for each unit of work (e.g., user request) is dangerous
  - High overhead to create & delete thread/process
  - Can exhaust CPU & memory resource
- Thread/process pool controls resource use
  - Allows service to be well conditioned.

— Well conditioned
— Not well conditioned

Throughput

Load

---

**When a user level thread does I/O it blocks the entire process.**

1. True
2. False