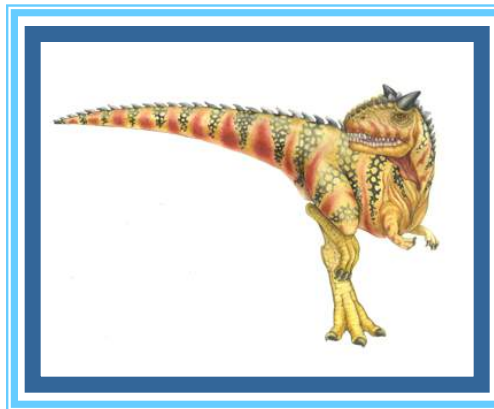# Chapter 4:  Threads

# Chapter 4: Threads

- Overview
- Multicore Programming
- Multithreading Models
- Thread Libraries
- Implicit Threading
- Threading Issues
- Operating System Examples

# Objectives

- To introduce the notion of a thread—a fundamental unit of CPU utilization that forms the basis of multithreaded computer systems

- To discuss the APIs for the Pthreads, Windows, and Java thread libraries

- To explore several strategies that provide implicit threading

- To examine issues related to multithreaded programming

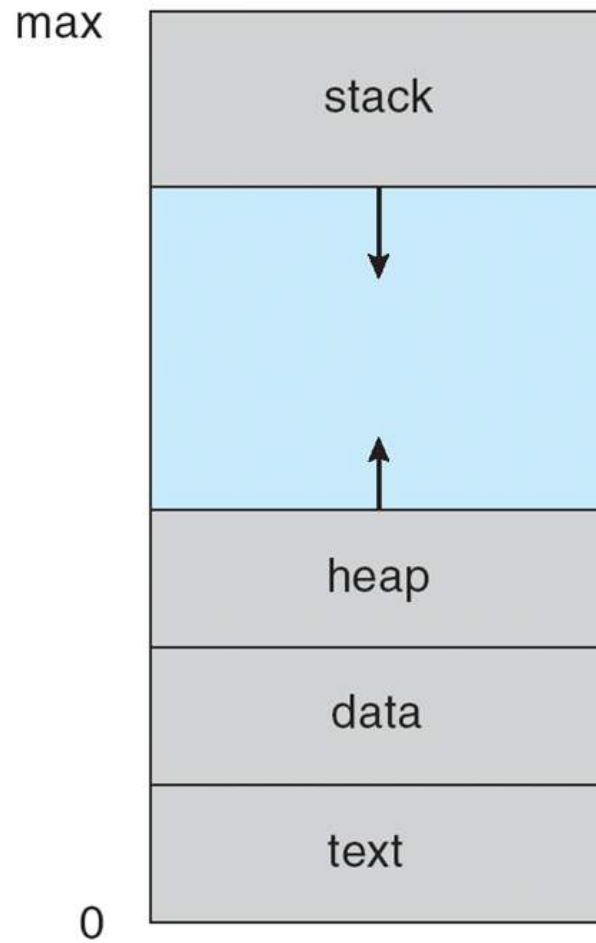- To cover operating system support for threads in Windows and Linux

# Motivation

- Most modern applications are multithreaded

- Threads run within application

- Multiple tasks with the application can be implemented by separate threads
  - Update display
  - Fetch data
  - Spell checking
  - Answer a network request

- Process creation is heavy-weight while thread creation is light-weight

- Can simplify code, increase efficiency

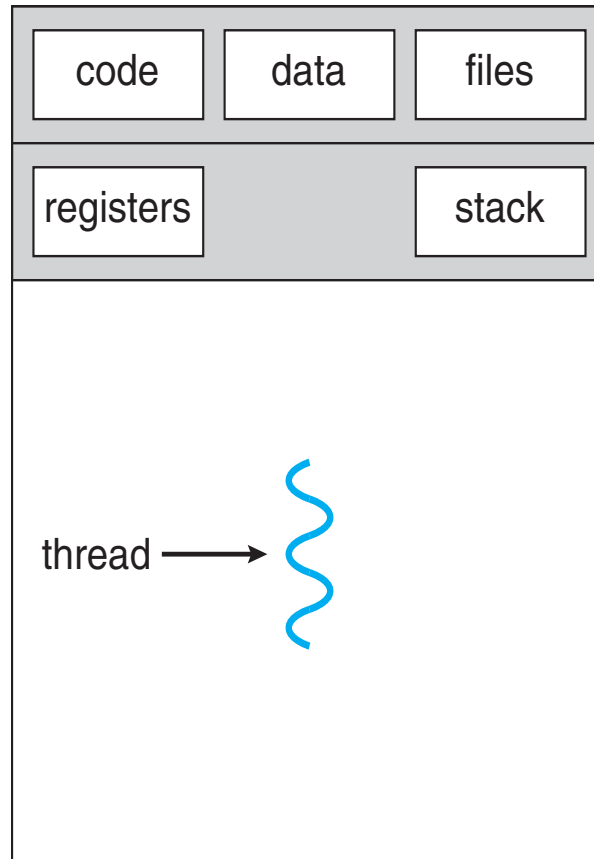- Kernels are generally multithreaded

# Thread

- Unit of CPU Utilization
- Threads have
  - Thread ID
  - program counter
  - register set
  - stack
- Shares with other OS resources such as
  - open files
  - signals
  - code
  - data section
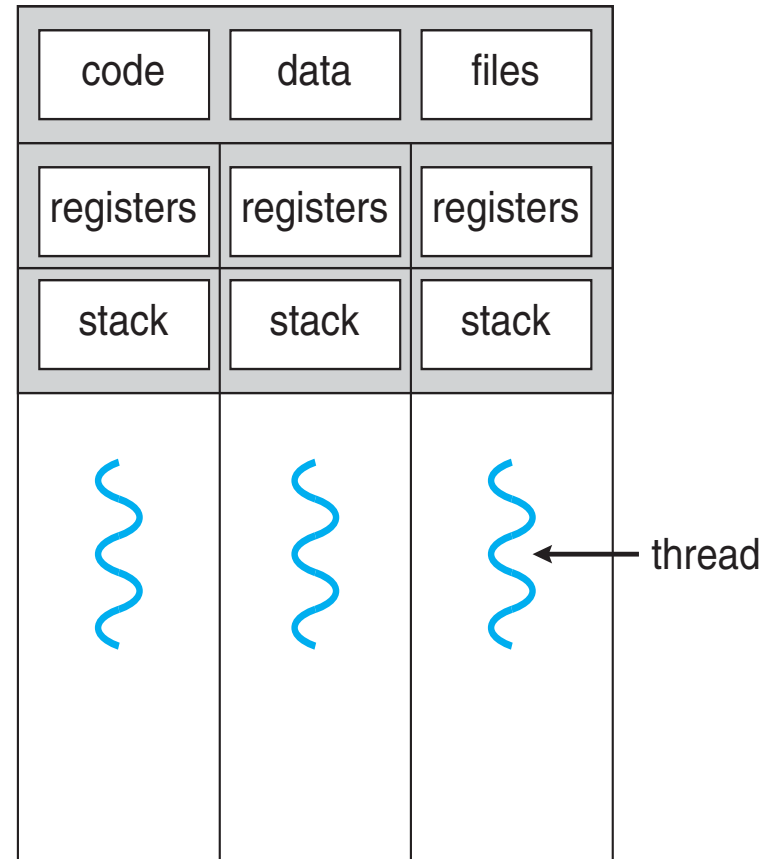
# Virtual Memory

# Single and Multithreaded Processes



single-threaded process                       multithreaded process
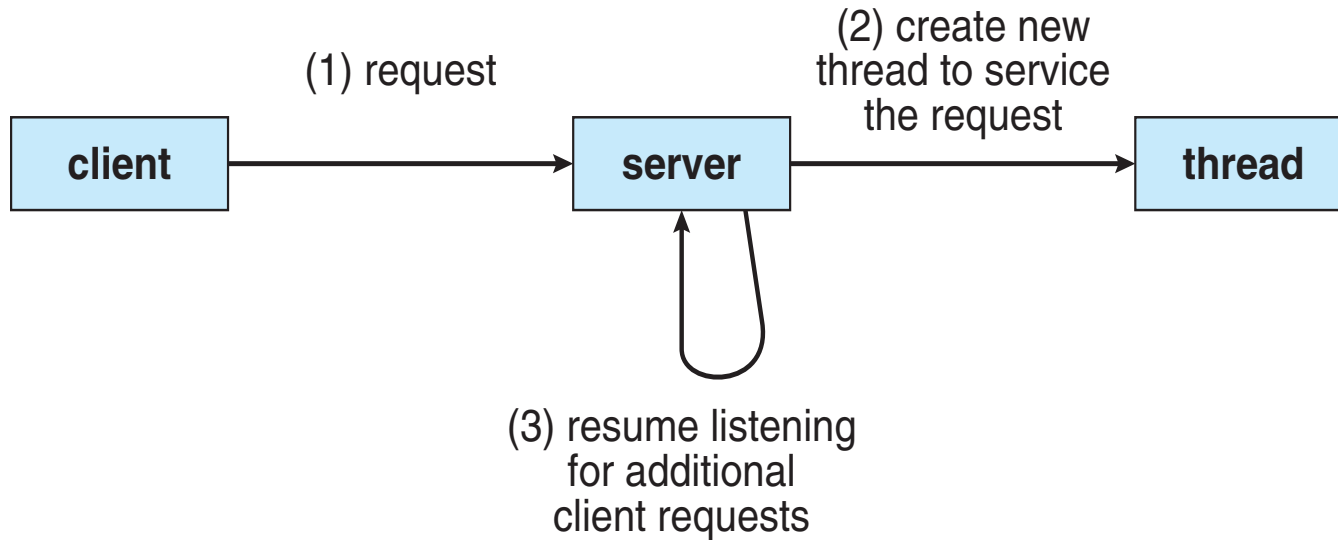
# Hyper Threading

- Intel and Sun technology

- Puts the threading concept in hardware

- Two logical processors per core,

- Each has its own processor state.

# Hyper Threading

■ Each logical processor can be halted, interrupted or directed to execute a specified thread, independently from the other logical processor sharing the same physical core.

■ Hyper-threading works by duplicating certain sections of the processor—those that store the architectural state, but not duplicating the main execution resources (CPU).
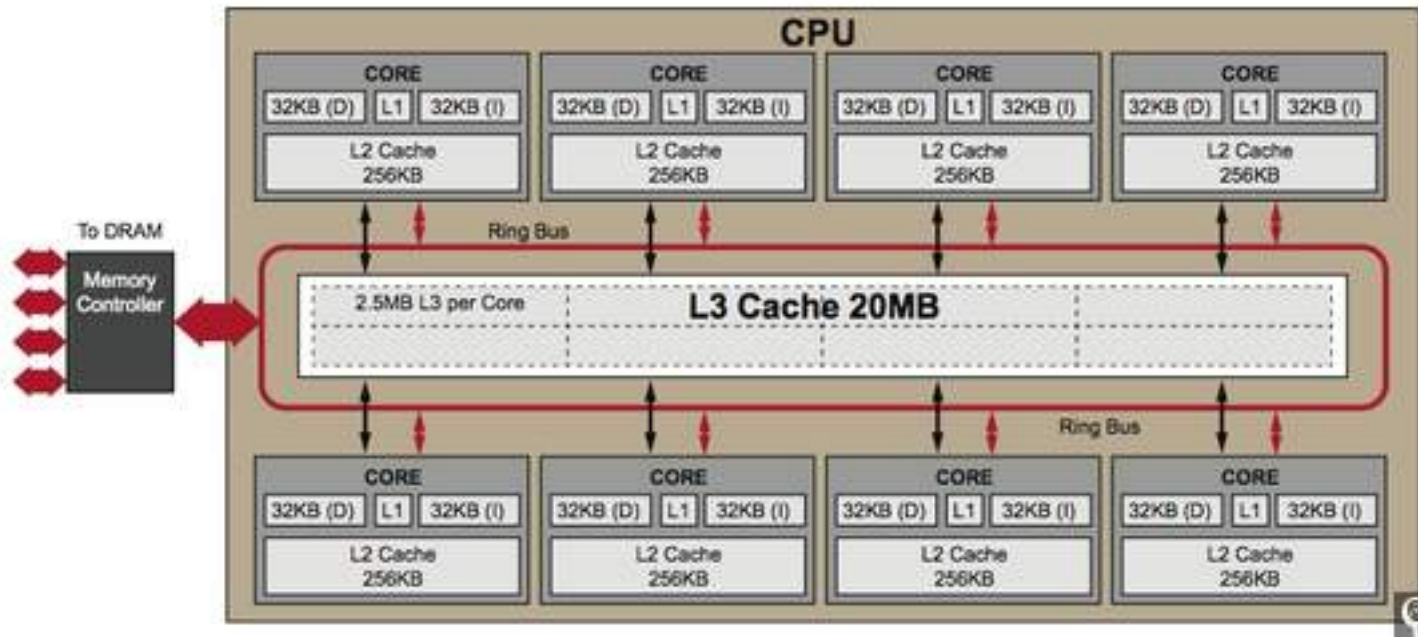
  – Source: www.wikipedia.org

# Multithreaded Server Architecture

# Benefits

- **Responsiveness –** may allow continued execution if part of process is blocked, especially important for user interfaces

- **Resource Sharing –** threads share resources of process, easier than shared memory or message passing

- **Economy –** cheaper than process creation, thread switching lower overhead than context switching

- **Scalability –** process can take advantage of multiprocessor architectures. More specifically, multi core.

# Multi Core Processor

# Multicore Programming

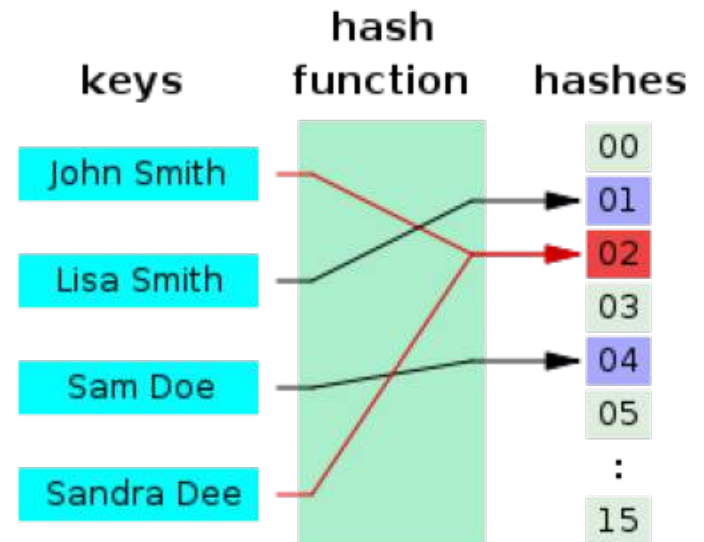- **Multicore** or **multiprocessor** systems put pressure on programmers, challenges include:
    - **Dividing activities**
    - **Balance**
    - **Data splitting**
    - **Data dependency**
    - **Testing and debugging**
- *Parallelism* implies a system can perform more than one task simultaneously
- *Concurrency* supports more than one task making progress
    - Single processor / core, scheduler providing concurrency

# Multicore Programming (Cont.)

- Types of parallelism
  - **Data parallelism** – distributes subsets of the same data across multiple cores, same operation on each
  - **Task parallelism** – distributing threads across cores, each thread performing unique operation
- As # of threads grows, so does architectural support for threading
  - CPUs have cores as well as ***hardware threads***
  - Consider Oracle SPARC T4 with 8 cores, and 8 hardware threads per core
- http://www.zdnet.com/article/understanding-task-and-data-parallelism-3039289129/
- http://www.intel.com/content/www/us/en/embedded/products/xeon/overview.html

# Hashing

- One way encryption. What does that mean?

- Data goes into hash function, mixed and has a unique number representing the hash.

- Because of the hash function, it is difficult to reconstruct the original input.

- Easier and more efficient to store the hash.

- Passwords, images, digital signatures.

- Since the number are big, collisions are rare.

# John The Ripper

- Password cracker

- People tend to use common passwords

- http://www.passwordrandom.com/most-popular-passwords

- Passwords are hashed.

- People tend to use common passwords.

- If you have a list of passwords, pre hash them, and save the hashes, solving them might be practical.

- Rainbow tables contain a lot of pre hashed words.

- John The Ripper takes a hash and compares it to its list of words.

- Back to threads, if we were able to use threading, this would go faster.
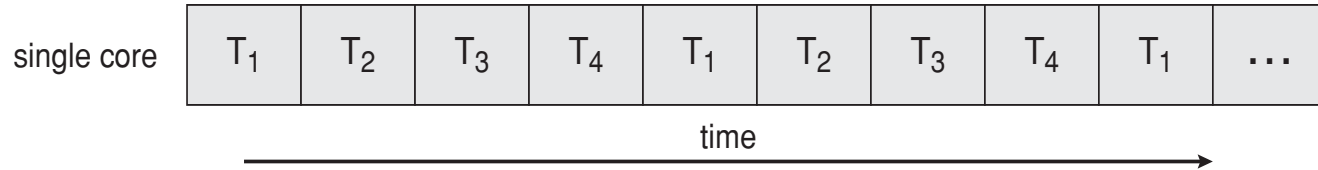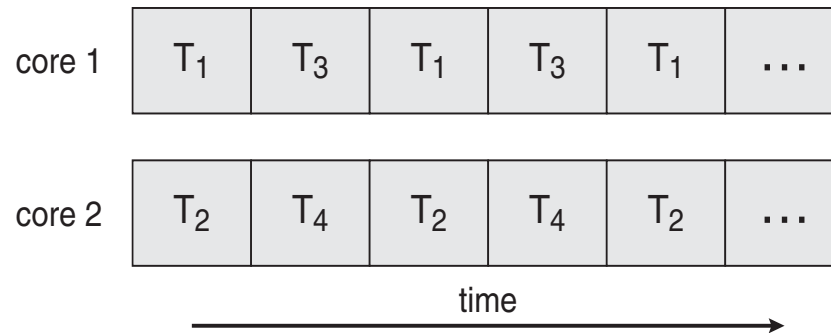
# GPU

- [http://www.nvidia.com/object/what-is-gpu-computing.html](http://www.nvidia.com/object/what-is-gpu-computing.html)
- [https://www.youtube.com/watch?v=ZrJeYFxpUyQ](https://www.youtube.com/watch?v=ZrJeYFxpUyQ)

# Concurrency vs. Parallelism

■ **Concurrent execution on single-core system:**

| single core | $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_1$ | ... |

time →

■ **Parallelism on a multi-core system:**

| core 1 | $T_1$ | $T_3$ | $T_1$ | $T_3$ | $T_1$ | ... |

| core 2 | $T_2$ | $T_4$ | $T_2$ | $T_4$ | $T_2$ | ... |

time →

# Single and Multithreaded Processes



single-threaded process

multithreaded process

# Gene Amdahl

# Amdahl's Law

- Identifies performance gains from adding additional cores to an application that has both serial and parallel components

- $S$ is serial portion
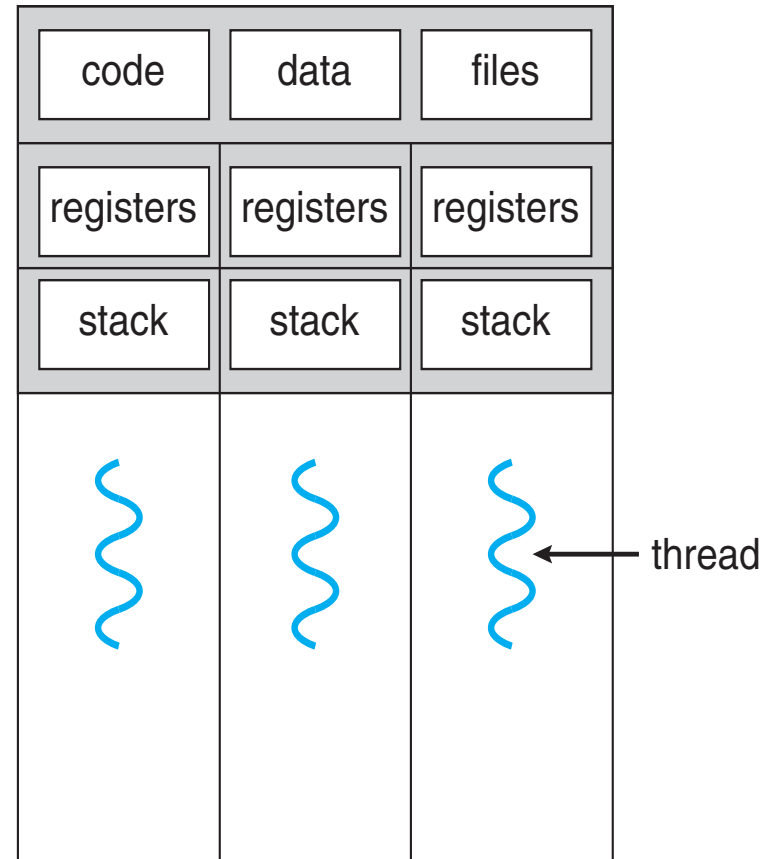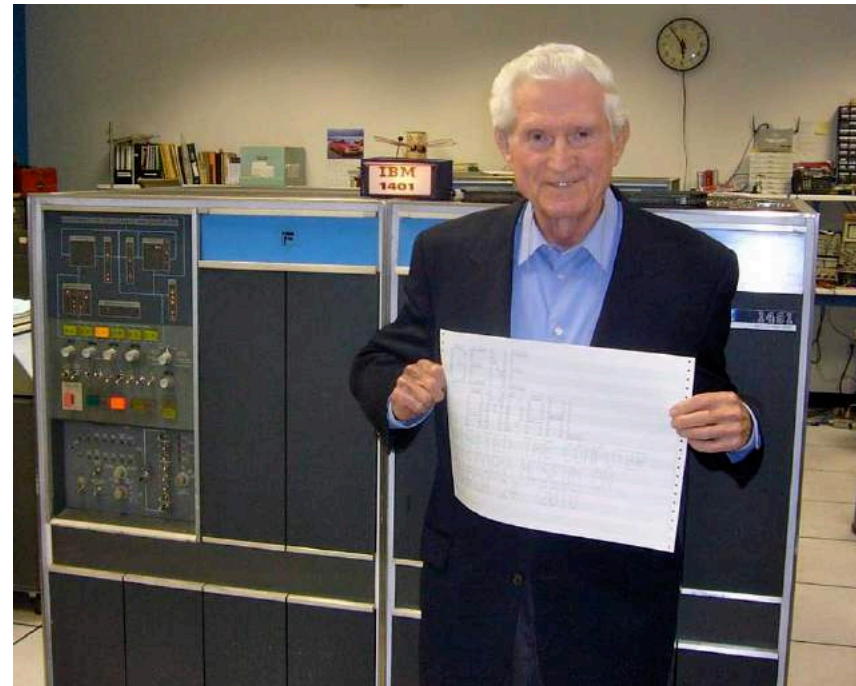
- $N$ processing cores

$$speedup \leq \frac{1}{S + \frac{(1-S)}{N}}$$

- That is, if application is 75% parallel / 25% serial, moving from 1 to 2 cores results in speedup of 1.6 times

- As $N$ approaches infinity, speedup approaches $1 / S$

**Serial portion of an application has disproportionate effect on performance gained by adding additional cores**

- But does the law take into account contemporary multicore systems?

# User Threads and Kernel Threads

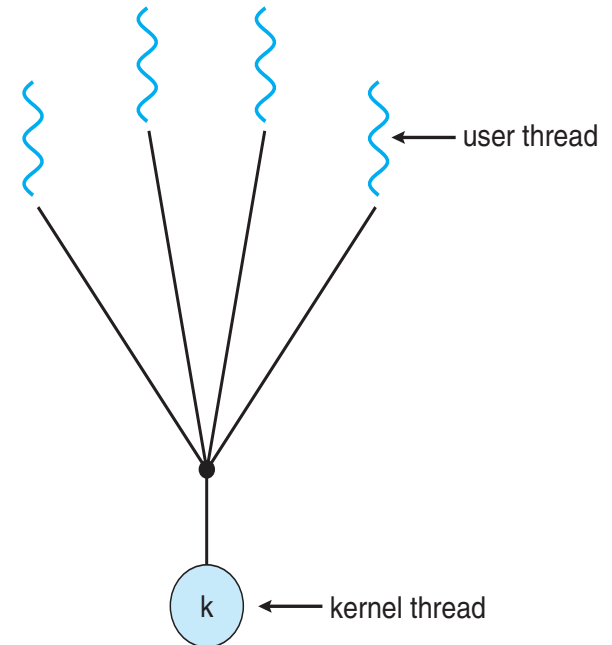- **User threads** - management done by user-level threads library
- Three primary thread libraries:
    - POSIX **Pthreads**
    - Windows threads
    - Java threads
- **Kernel threads** - Supported by the Kernel
- Examples – virtually all general purpose operating systems, including:
    - Windows
    - Solaris
    - Linux
    - Tru64 UNIX
    - Mac OS X

# Multithreading Models

- Many-to-One

- One-to-One

- Many-to-Many

# Many-to-One

- Many user-level threads mapped to single kernel thread

- One thread blocking causes all to block

- Multiple threads **may not run in parallel** on muticore system because only one may be in kernel at a time

- Few systems currently use this model

- Examples:
  - **Solaris Green Threads**
  - **GNU Portable Threads**

← user thread

k ← kernel thread

# Traffic Jam

# One-to-One

- Each user-level thread maps to kernel thread

- Creating a user-level thread creates a kernel thread

- More concurrency than many-to-one

- Number of threads per process sometimes restricted due to overhead

- Examples
  - Windows
  - Linux
  - Solaris 9 and later

← user thread

k    k    k    k    ← kernel thread

# Many-to-Many Model

- Allows many user level threads to be mapped to many kernel threads

- Allows the operating system to create a sufficient number of kernel threads

- Solaris prior to version 9

- Windows with the *ThreadFiber* package

user thread

kernel thread

k    k    k

# Two-level Model

- Similar to M:M, except that it allows a user thread to be **bound** to kernel thread

- Examples
  - IRIX
  - HP-UX
  - Tru64 UNIX
  - Solaris 8 and earlier



user thread

kernel thread

- https://www.youtube.com/watch?v=dpT-8C-QeCY

# Thread Libraries

- **Thread library** provides programmer with API for creating and managing threads

- Two primary ways of implementing
  - Library entirely in user space
  - Kernel-level library supported by the OS

# Pthreads

- May be provided either as user-level or kernel-level

- A POSIX standard (IEEE 1003.1c) API for thread creation and synchronization

- *Specification*, not *implementation*

- API specifies behavior of the thread library, implementation is up to development of the library

- Common in UNIX operating systems (Solaris, Linux, Mac OS X)



A POSIX Standard for Better Multiprocessing

Pthreads

Programming

O'REILLY®

Bradford Nichols, Dick Buttlar & Jacqueline Proulx Farrell

# Pthreads Example

```c
#include <pthread.h>
#include <stdio.h>

int sum; /* this data is shared by the thread(s) */
void *runner(void *param); /* threads call this function */

int main(int argc, char *argv[])
{
  pthread_t tid; /* the thread identifier */
  pthread_attr_t attr; /* set of thread attributes */

  if (argc != 2) {
    fprintf(stderr,"usage: a.out <integer value>\n");
    return -1;
  }
  if (atoi(argv[1]) < 0) {
    fprintf(stderr,"%d must be >= 0\n",atoi(argv[1]));
    return -1;
  }
```

# Pthreads Example (Cont.)

```c
    /* get the default attributes */
    pthread_attr_init(&attr);
    /* create the thread */
    pthread_create(&tid,&attr,runner,argv[1]);
    /* wait for the thread to exit */
    pthread_join(tid,NULL);

    printf("sum = %d\n",sum);
}

/* The thread will begin control in this function */
void *runner(void *param)
{
    int i, upper = atoi(param);
    sum = 0;

    for (i = 1; i <= upper; i++)
        sum += i;

    pthread_exit(0);
}
```

# Pthreads Code for Joining 10 Threads

```
#define NUM_THREADS 10

/* an array of threads to be joined upon */
pthread_t workers[NUM_THREADS];

for (int i = 0; i < NUM_THREADS; i++)
   pthread_join(workers[i], NULL);
```

# Implicit Threading

- Growing in popularity as numbers of threads increase, program correctness more difficult with explicit threads

- Creation and management of threads done by compilers and run-time libraries rather than programmers

- Three methods explored

  - Thread Pools

  - OpenMP

  - Grand Central Dispatch /* Apple

- Other methods include Microsoft Threading Building Blocks (TBB), `java.util.concurrent` package

# Thread Pools

- Create a number of threads in a pool where they await work

- Advantages:

  - Usually slightly faster to service a request with an existing thread than create a new thread

  - Allows the number of threads in the application(s) to be bound to the size of the pool

  - Separating task to be performed from mechanics of creating task allows different strategies for running task

    - i.e.Tasks could be scheduled to run periodically

- Windows API supports thread pools:

```
DWORD WINAPI PoolFunction(AVOID Param) {
    /*
     * this function runs as a separate thread.
     */
}
```

# OpenMP

- Set of compiler directives and an API for C, C++, FORTRAN

- Provides support for parallel programming in shared-memory environments

- Identifies **parallel regions** – blocks of code that can run in parallel

```
#pragma omp parallel
```

Create as many threads as there are cores

```
#pragma omp parallel for
   for(i=0;i<N;i++) {

     c[i] = a[i] + b[i];

}
```

Run for loop in parallel

```
#include <omp.h>
#include <stdio.h>

int main(int argc, char *argv[])
{
  /* sequential code */

  #pragma omp parallel
  {
    printf("I am a parallel region.");
  }

  /* sequential code */

  return 0;
}
```

# Grand Central Dispatch

- Apple technology for Mac OS X and iOS operating systems

- Extensions to C, C++ languages, API, and run-time library

- Allows identification of parallel sections

- Manages most of the details of threading

- Block is in "^{ }" - `^{ printf("I am a block"); }`

- Blocks placed in dispatch queue

  - Assigned to available thread in thread pool when removed from queue

# Threading Issues

■ The programmer must account for the following:

- Semantics of **fork()** and **exec()** system calls

- Signal handling

  ▸ Synchronous and asynchronous

- Thread cancellation of target thread

  ▸ Asynchronous or deferred

- Thread-local storage

- Scheduler Activations

# Semantics of fork() and exec()

- Does `fork()` duplicate only the calling thread or all threads?

  - Some UNIXes have two versions of fork

- `exec()` usually works as normal – replace the running process including all threads

# Signal Handling

■ **Signals** are used in UNIX systems to notify a process that a particular event has occurred.

■ A **signal handler** is used to process signals

  1. Signal is generated by particular event
  2. Signal is delivered to a process
  3. Signal is handled by one of two signal handlers:
     1. default
     2. user-defined

■ Every signal has **default handler** that kernel runs when handling signal

  ● **User-defined signal handler** can override default
  ● For single-threaded, signal delivered to process

# Signal Handling (Cont.)

- Where should a signal be delivered for multi-threaded?
  - Deliver the signal to the thread to which the signal applies
  - Deliver the signal to every thread in the process
  - Deliver the signal to certain threads in the process
  - Assign a specific thread to receive all signals for the process

# Thread Cancellation

- Terminating a thread before it has finished

- Thread to be canceled is **target thread**

- Two general approaches:

  - **Asynchronous cancellation** terminates the target thread immediately

  - **Deferred cancellation** allows the target thread to periodically check if it should be cancelled

- Pthread code to create and cancel a thread:

```
pthread_t tid;

/* create the thread */
pthread_create(&tid, 0, worker, NULL);

    . . .

/* cancel the thread */
pthread_cancel(tid);
```

# Thread Cancellation (Cont.)

■ Invoking thread cancellation requests cancellation, but actual cancellation depends on thread state

| Mode | State | Type |
|------|-------|------|
| Off | Disabled | – |
| Deferred | Enabled | Deferred |
| Asynchronous | Enabled | Asynchronous |

■ If thread has cancellation disabled, cancellation remains pending until thread enables it

■ Default type is deferred

  ● Cancellation only occurs when thread reaches **cancellation point**

    ‣ I.e. `pthread_testcancel()`

    ‣ Then **cleanup handler** is invoked

■ On Linux systems, thread cancellation is handled through signals

# End of Chapter 4