# Chapter 4:  Threads & Concurrency

# Outline

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

# Objectives

- Identify the basic components of a thread, and contrast threads and processes

- Describe the benefits and challenges of designng multithreaded applications

- Illustrate different approaches to implicit threading including thread pools, fork-join, and Grand Central Dispatch

- Describe how the Windows and Linux operating systems represent threads

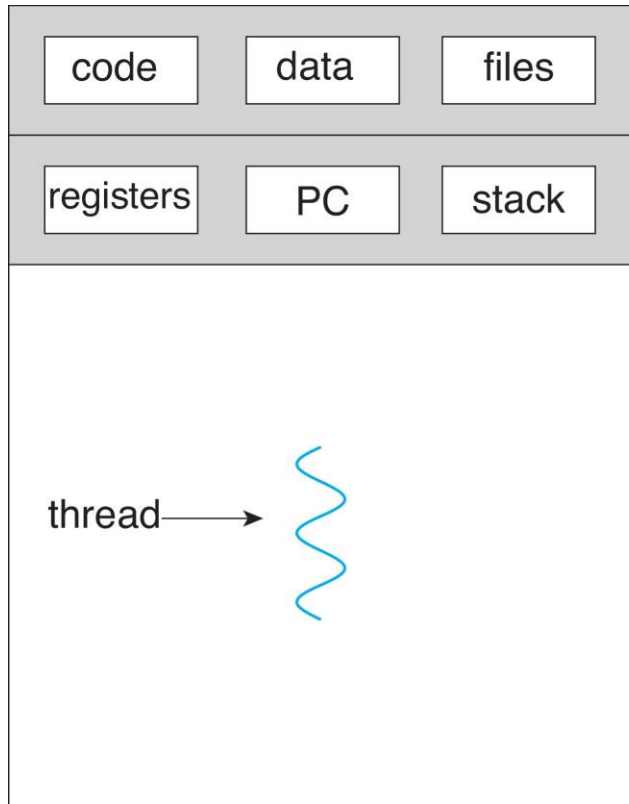- Design multithreaded applications using the Pthreads, Java, and Windows threading APIs

# 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

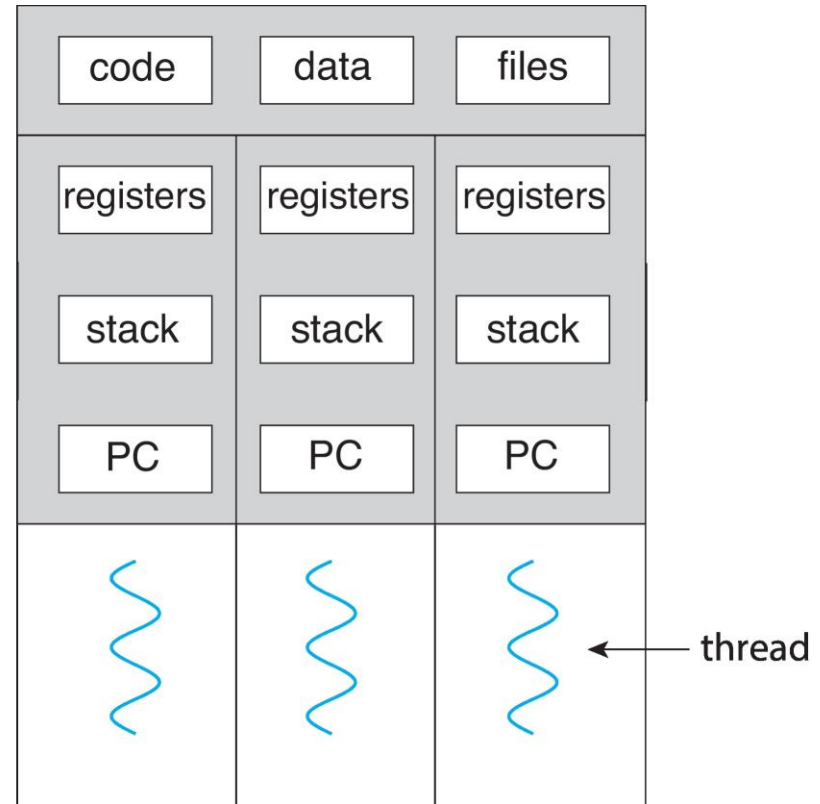- Kernels are generally multithreaded

# Single and Multithreaded Processes
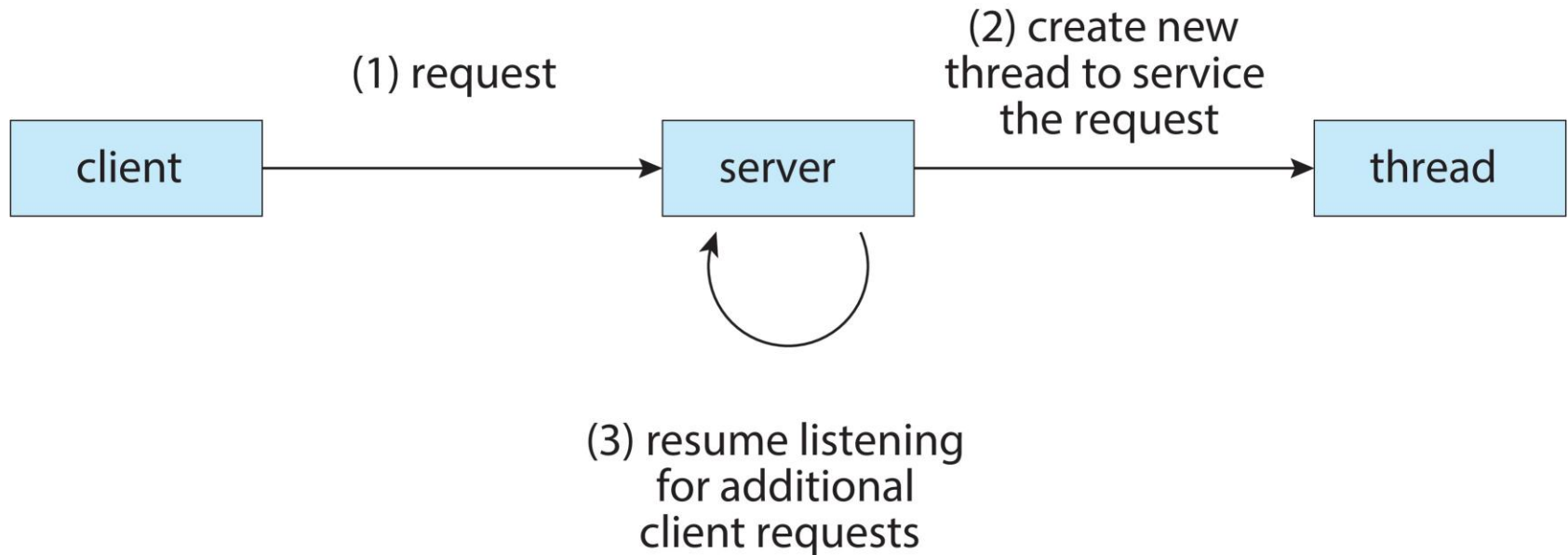


single-threaded process

multithreaded process

# 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 multicore architectures
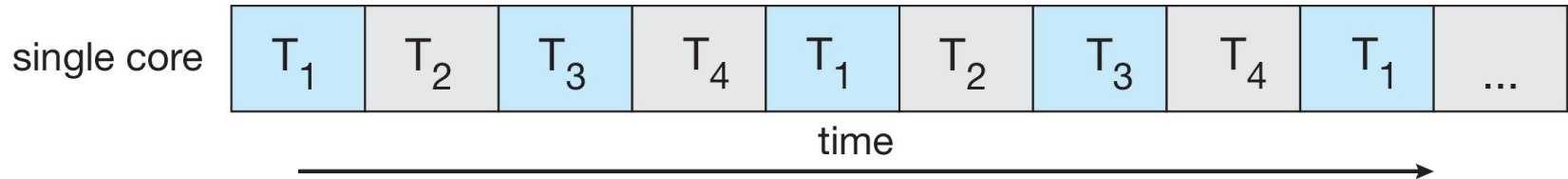
# Multicore Programming

- **Multicore** or **multiprocessor** systems putting 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

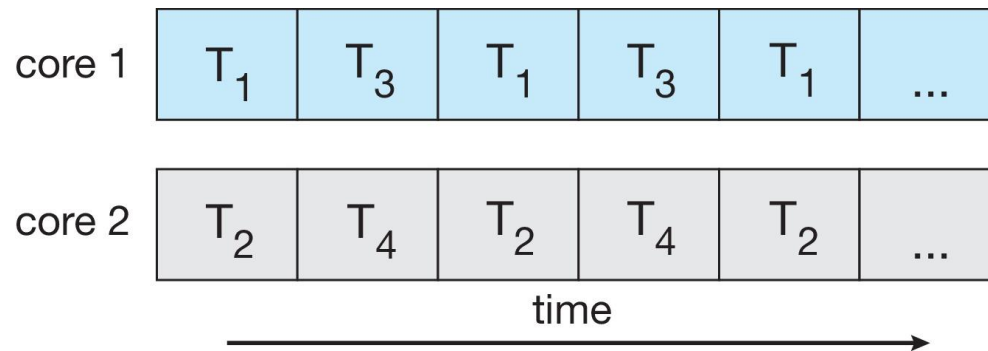    - Single processor / core, scheduler providing concurrency

# 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 | $T_1$ | $T_3$ | $T_1$ | $T_3$ | $T_1$ | ... |
|---|---|---|---|---|---|---|

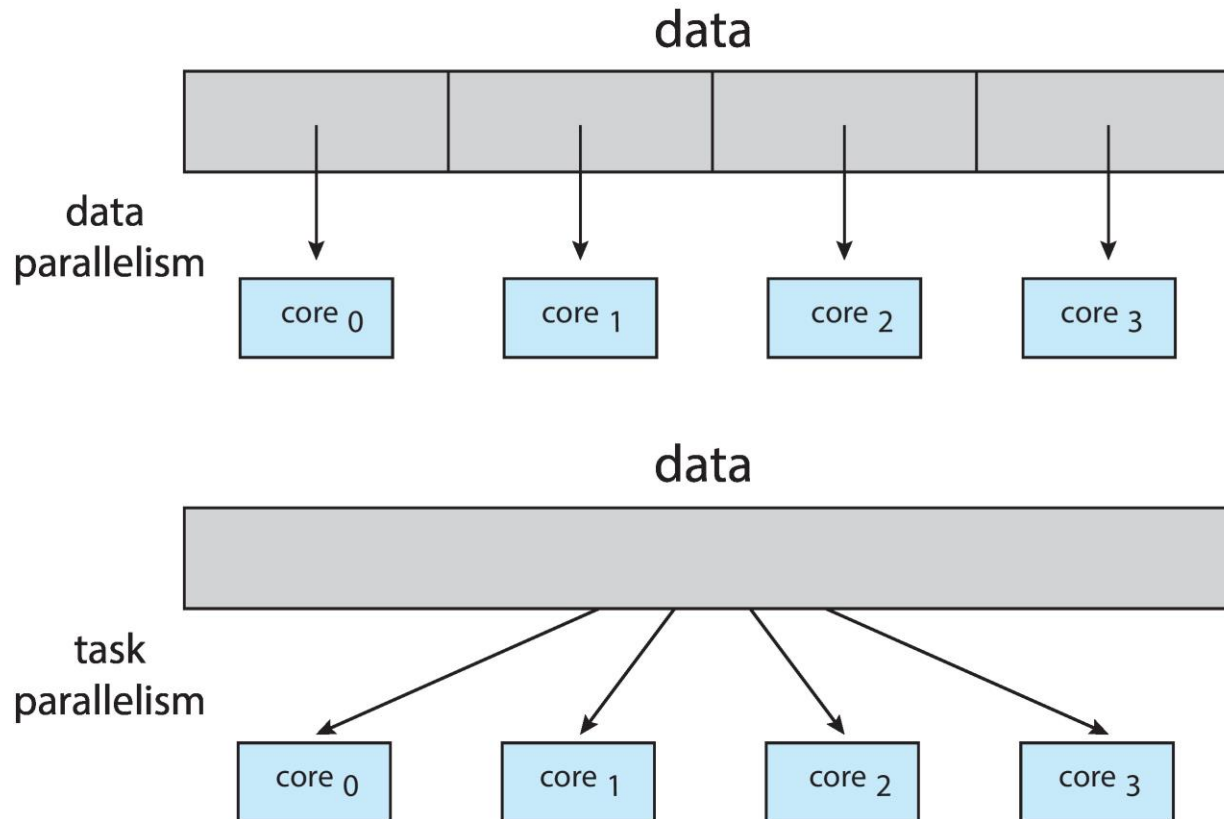| core 2 | $T_2$ | $T_4$ | $T_2$ | $T_4$ | $T_2$ | ... |
|---|---|---|---|---|---|---|

time →

# Multicore Programming

- 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

# Data and Task Parallelism

# 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

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

Yes, it is. The speedup can't be less than zero (if parallelizing your code makes it slower, just run the serial version) or more than 100% (because saving more than 100% of the runtime would make the runtime negative).

Multicore processors are essentially just a matter of terminology changing over time. If you think you have one processor with four cores, Amdahl thinks you have four processors.

Ref: https://bit.ly/3i44GC2

- That is, if application is 75% parallel / 25 ~~cores~~ results in speedup of 1.6 times

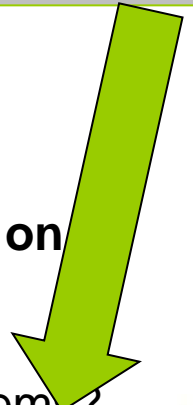- As *N* approaches infinity, speedup approaches 1 / *S*

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

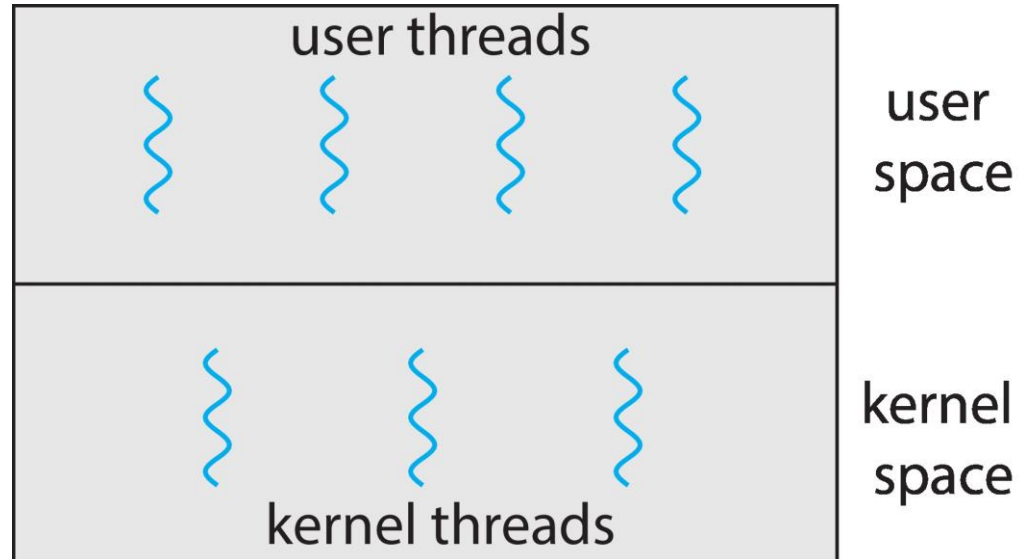- But does the law take into account contemporary multicore systems?

# 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

  - Linux

  - Mac OS X

  - iOS

  - Android
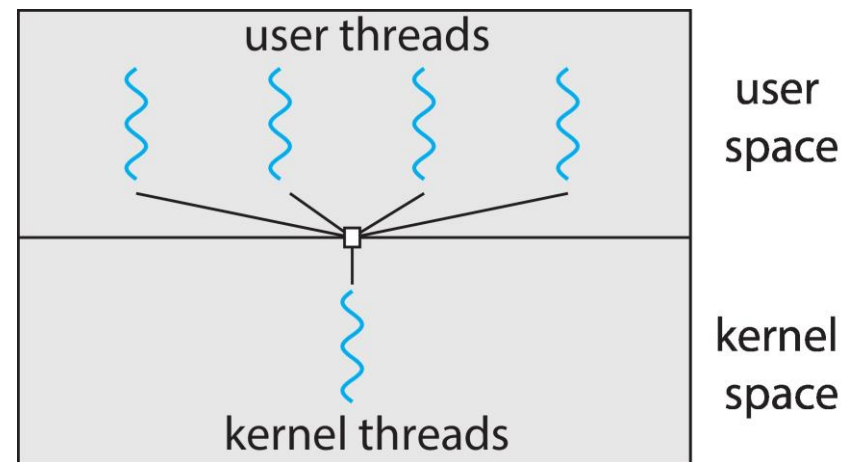
# User and Kernel Threads

# Multithreading Models

- Many-to-One

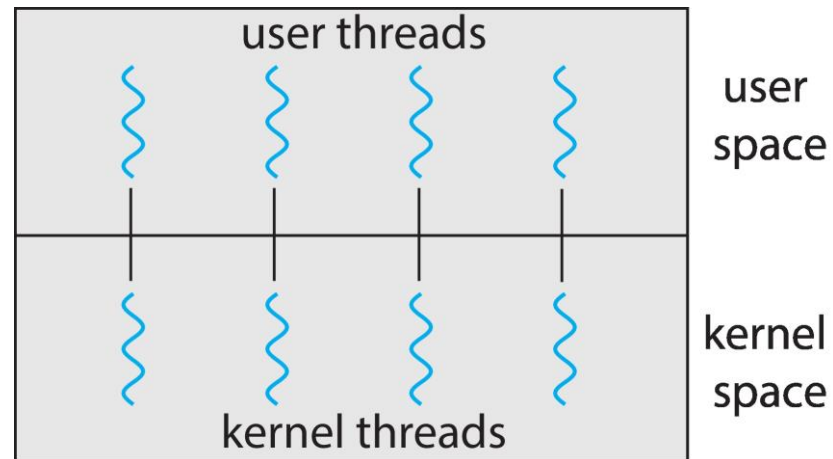- One-to-One

- 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**

# 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

# 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

- Windows with the *ThreadFiber* package

- Otherwise not very common

- Multiplexes

# Two-level Model

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

# 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 (Linux & Mac OS X)

# Java Threads

- Java threads are managed by the JVM

- Typically implemented using the threads model provided by underlying OS

- Java threads may be created by:

  - Extending Thread class

    Create a new class derive from the thread Thread class and override it's run() method.

  - Implementing the Runnable interface

    The Runnable interface in Java is the core element when you are working with Threads. Any Java class intending to execute threads must implement the Runnable interface.

  -
    ```
    public interface Runnable
    {
        public abstract void run();
    }
    ```

  - Standard practice is to implement Runnable interface

# Java Threads

**Implementing Runnable interface:**

```java
class Task implements Runnable
{
   public void run() {
      System.out.println("I am a thread.");
   }
}
```

**Creating a thread:**

```java
Thread worker = new Thread(new Task());
worker.start();
```

**Waiting on a thread:**

```java
try {
   worker.join();
}
catch (InterruptedException ie) { }
```

# Java Executor Framework

- Rather than explicitly creating threads, Java also allows thread creation around the Executor interface:

```
public interface Executor
{
    void execute(Runnable command);
}
```

- The Executor is used as follows:

```
Executor service = new Executor;
service.execute(new Task());
```

# 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

- Five methods explored
  - Thread Pools
  - Fork-Join
  - OpenMP
  - Grand Central Dispatch
  - Intel Threading Building Blocks

# 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

  - Tasks could be scheduled to run periodically

- Windows API supports thread pools:
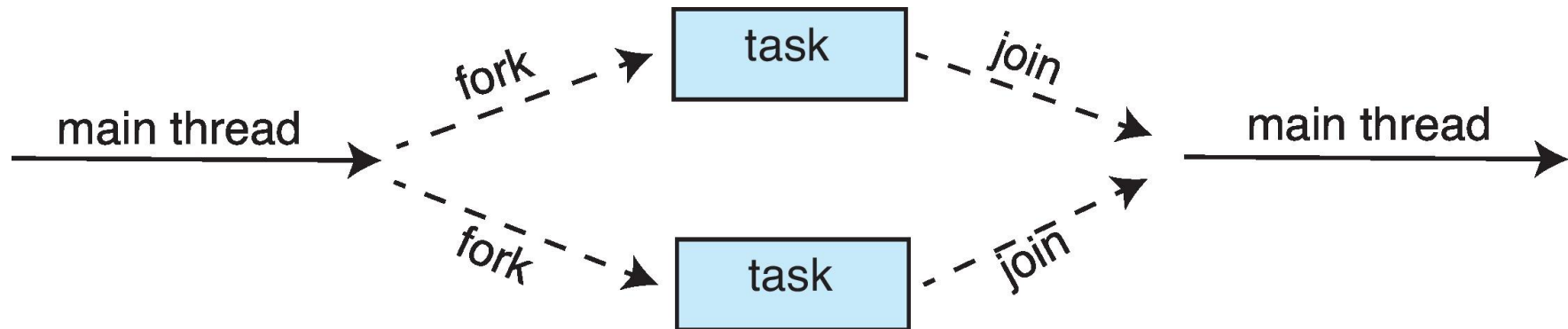
```
DWORD WINAPI PoolFunction(AVOID Param) {
    /*
    * this function runs as a separate thread.
    */
}
```

# Fork-Join Parallelism

- Multiple threads (tasks) are **forked**, and then **joined**.

# 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

```c
#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 macOS and iOS operating systems

- Identifies a language extensions for C, C++ and Objective-C known as blocks

- Allows identification of parallel sections

- Block is in "^{ }" :

```
^{ printf("I am a block"); }
```

- Blocks placed in dispatch queue
    - Assigned to available thread in thread pool when removed from queue

# Intel Threading Building Blocks (TBB)

- Template library for designing parallel C++ programs

- A serial version of a simple for loop

```
for (int i = 0; i < n; i++) {
    apply(v[i]);
}
```

- The same for loop written using TBB with **parallel_for** statement:

```
parallel_for (size_t(0), n, [=](size_t i) {apply(v[i]);});
```

# Threading Issues

- 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

# 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);

/* wait for the thread to terminate */
pthread_join(tid,NULL);
```

# Thread-Local Storage
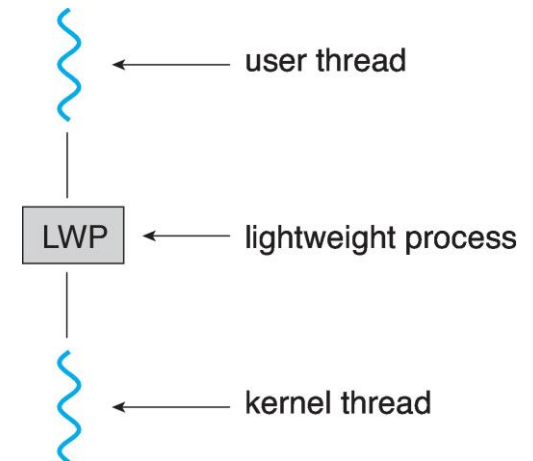
- **Thread-local storage** (**TLS**) allows each thread to have its own copy of data

- Useful when you do not have control over the thread creation process (i.e., when using a thread pool)

- Different from local variables

    - Local variables visible only during single function invocation

    - TLS visible across function invocations

- Similar to `static` data

    - TLS is unique to each thread

# Scheduler Activations

- Schedular activation is the communication scheme between user-thread library and kernel

- For any application (both from M:M and two-level), it is required to main communication with number of allocated kernel level threads

- This done by using an intermediate data structure between user and kernel threads – which is called **lightweight process** (**LWP**)

  - Appears to be a virtual processor on which process can schedule user thread to run

  - Each LWP attached to kernel thread

  - How many LWPs to create?

    ▸ Depends on the application type

- Kernel informs application regarding certain events using **upcalls** procedure

- Thread library handles this **upcalls** using **upcall handler**

user thread

LWP ← lightweight process

kernel thread

# Operating System Examples

- Windows Threads

- Linux Threads

# Windows Threads

- Each thread contains

    - A thread id

    - Register set representing state of processor

    - Separate user and kernel stacks for thread runing in user mode or kernel mode

    - Private data storage area used by run-time libraries and dynamic link libraries (DLLs)

- The register set, stacks, and private storage area are known as the **context** of the thread

# Linux Threads

- Linux refers to them as *tasks* rather than *threads*

- Thread creation is done through `clone()` system call

- `clone()` allows a child task to share the address space of the parent task (process)

- `struct task_struct` points to process data structures (shared or unique)

# End of Chapter 4