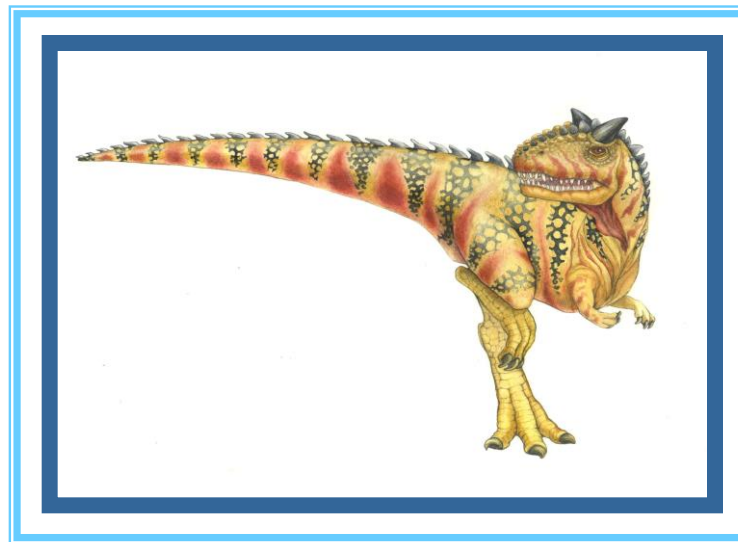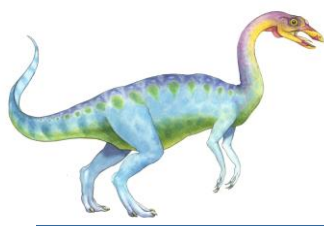# Chapter 4:  Multithreaded Programming

# Chapter 4: Multithreaded Programming

- Overview
- Multicore Programming
- Multithreading Models
- Thread Libraries
- Implicit Threading
- Threading Issues

# 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

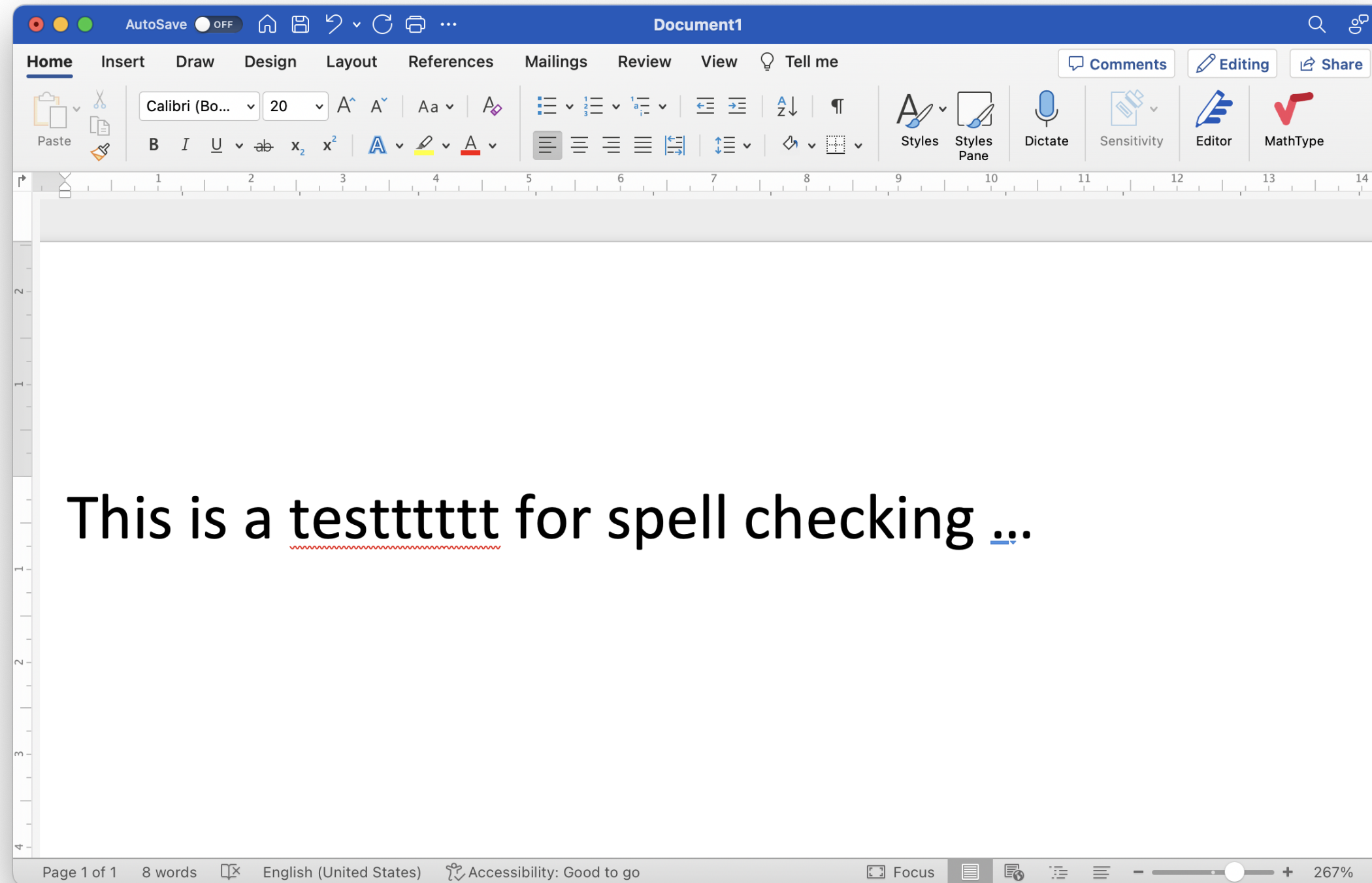- To cover operating system support for threads in Windows and Linux

# 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

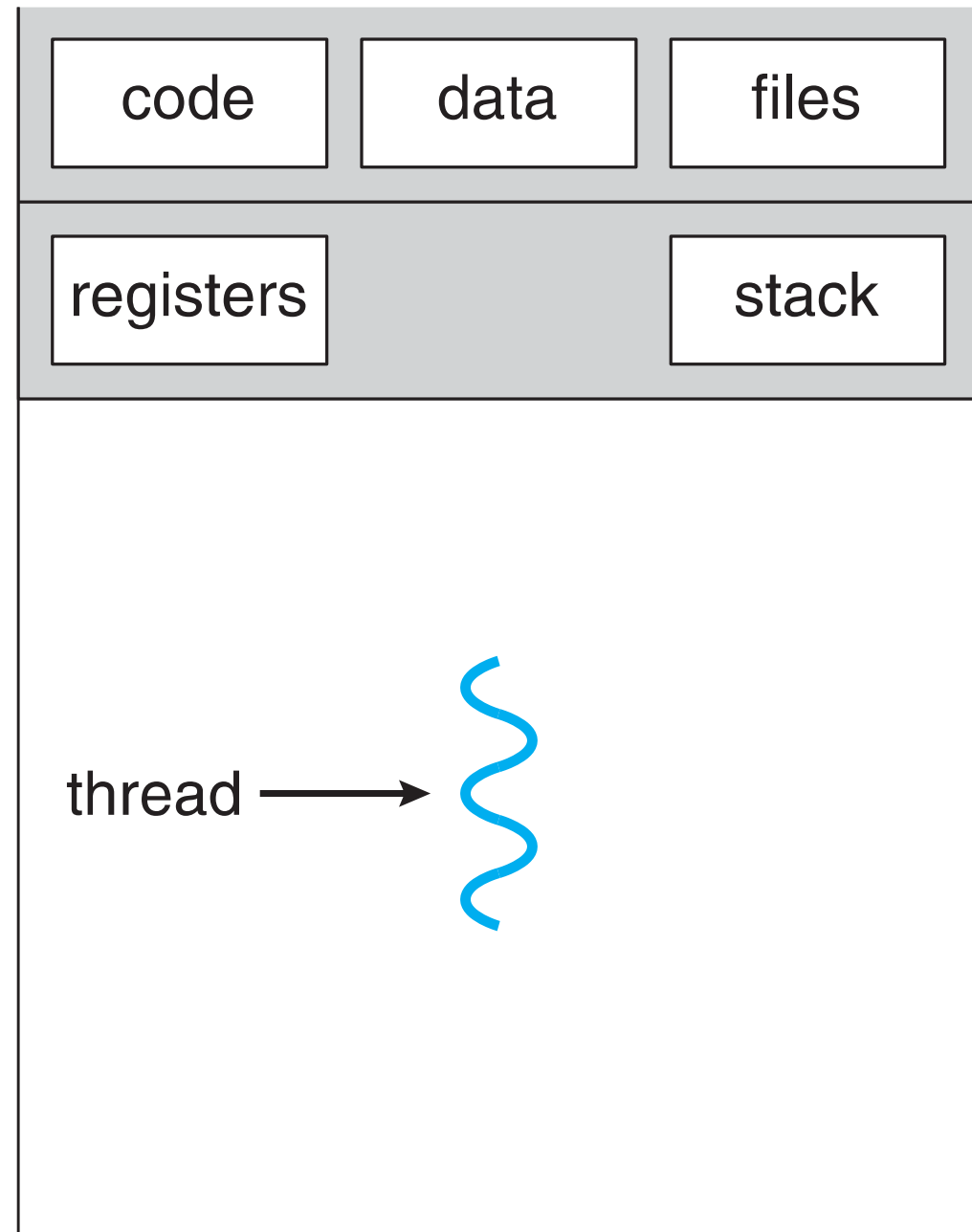- Kernels are generally multithreaded
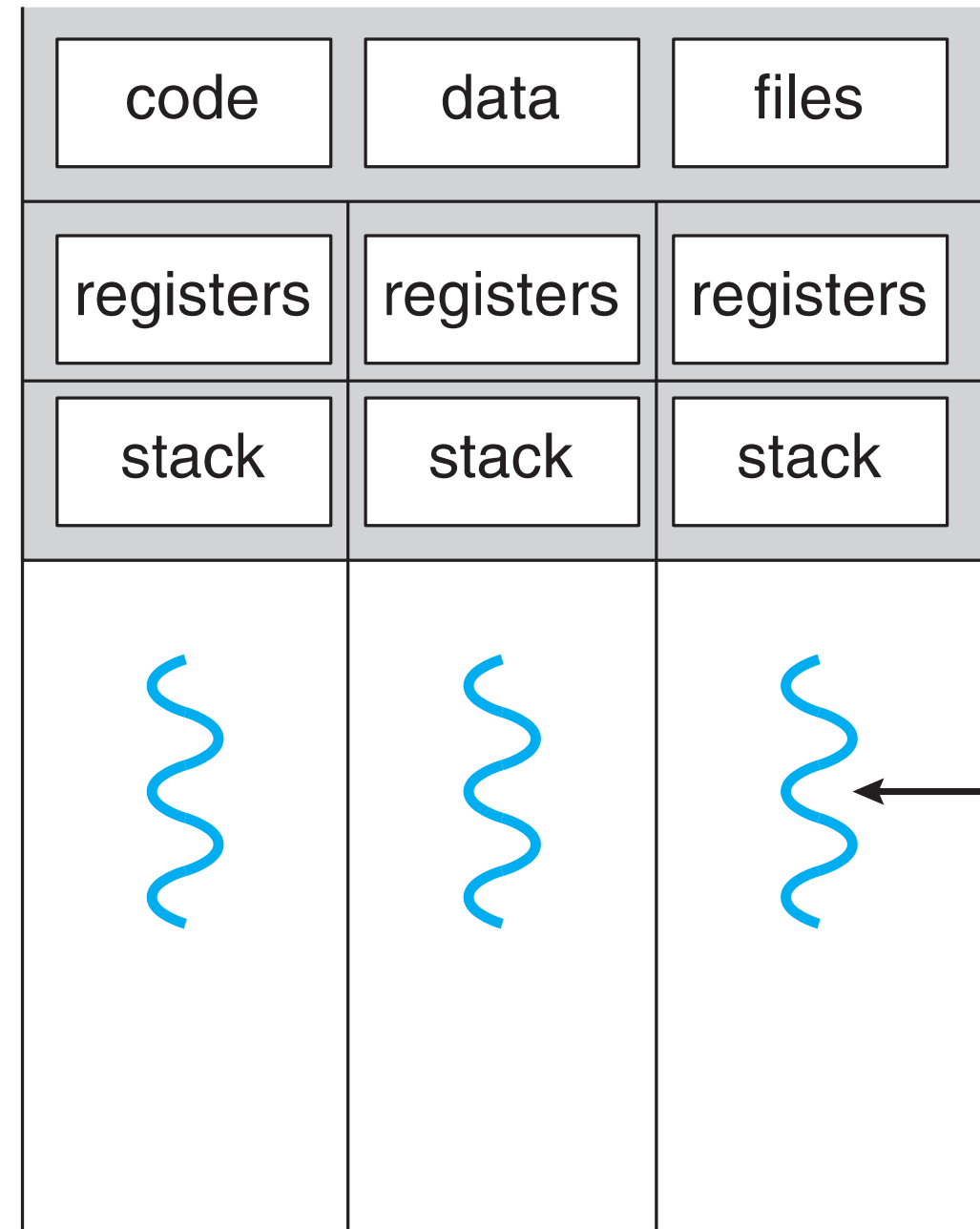
# Microsoft Word Spell Checking

# Single and Multithreaded Processes

| code | data | files |
|------|------|-------|

| registers | | stack |
|-----------|---|-------|

thread ⟶ 〰

single-threaded process

| code | data | files |
|------|------|-------|

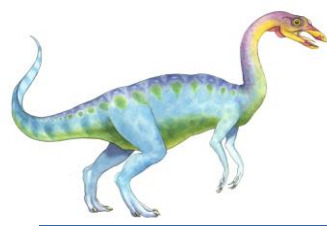| registers | registers | registers |
|-----------|-----------|-----------|
| stack | stack | stack |

〰 〰 〰 ⟵ thread

multithreaded process

```
Type as fast as you can, press q to quit
The pressed key is █
```

Count number of key pressed in a second

```
key_hit_count = 0

last_second_count = 0

t0 = current timestamp

While not stop:

        Check keyboard or timeout in 9,5 seconds

        If there is a key hit:

                increase key_hit_count by one

        t1 = current timestamp

        if t1 - t0 >= 1 second:

                key_hit_last_second = key_hit_count - last_second_count

                rate = key_hit_last_second / (t1 - t0)

                print rate

                last_second_count = key_hit_count

                t0 = t1
```

# Keyhit Rating – Two Threads

**Key_count_thread:**

    While not stop:

        Wait for keyhit

        increase key_hit_count by one
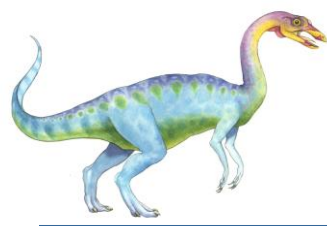
**Rating_thread:**

    While not stop:

        sleep for 1 second

        key_hit_last_second = key_hit_count – last_second_count

        rate = key_hit_last_second / (t1 – t0)

        print rate

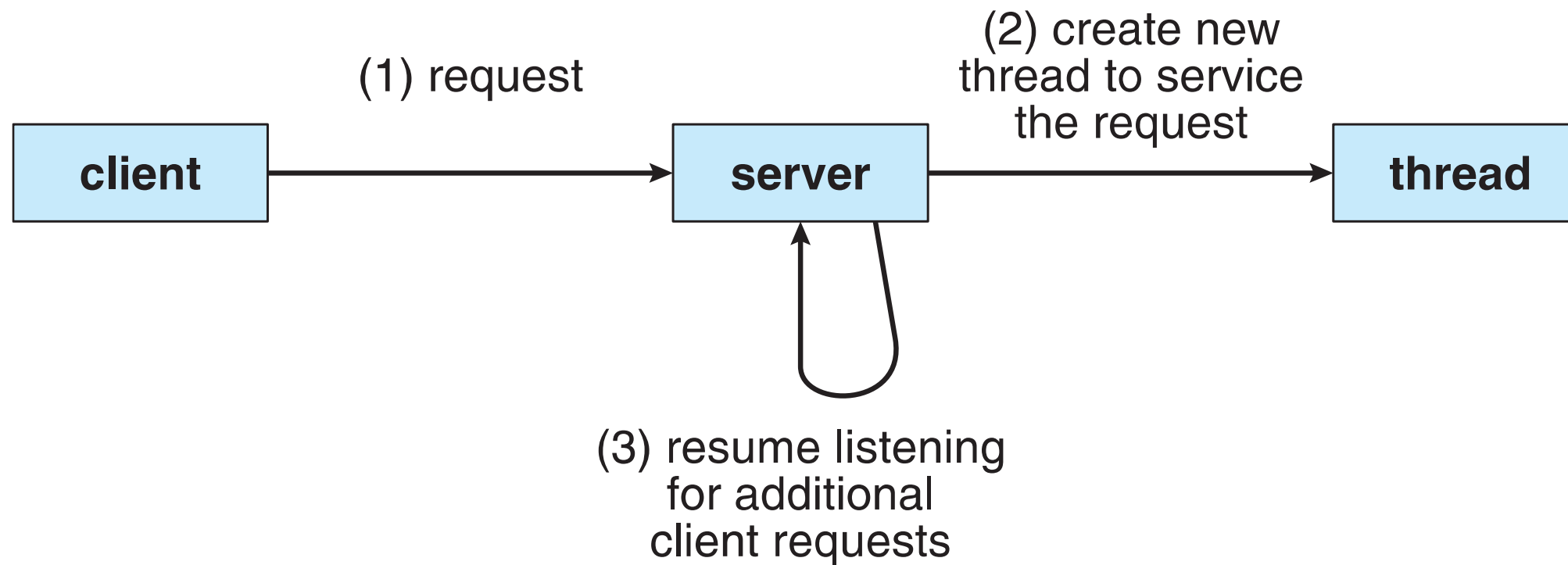        last_second_count = key_hit_count

**Main:**

    key_hit_count = 0

    last_second_count = 0

    Start Rating_thread
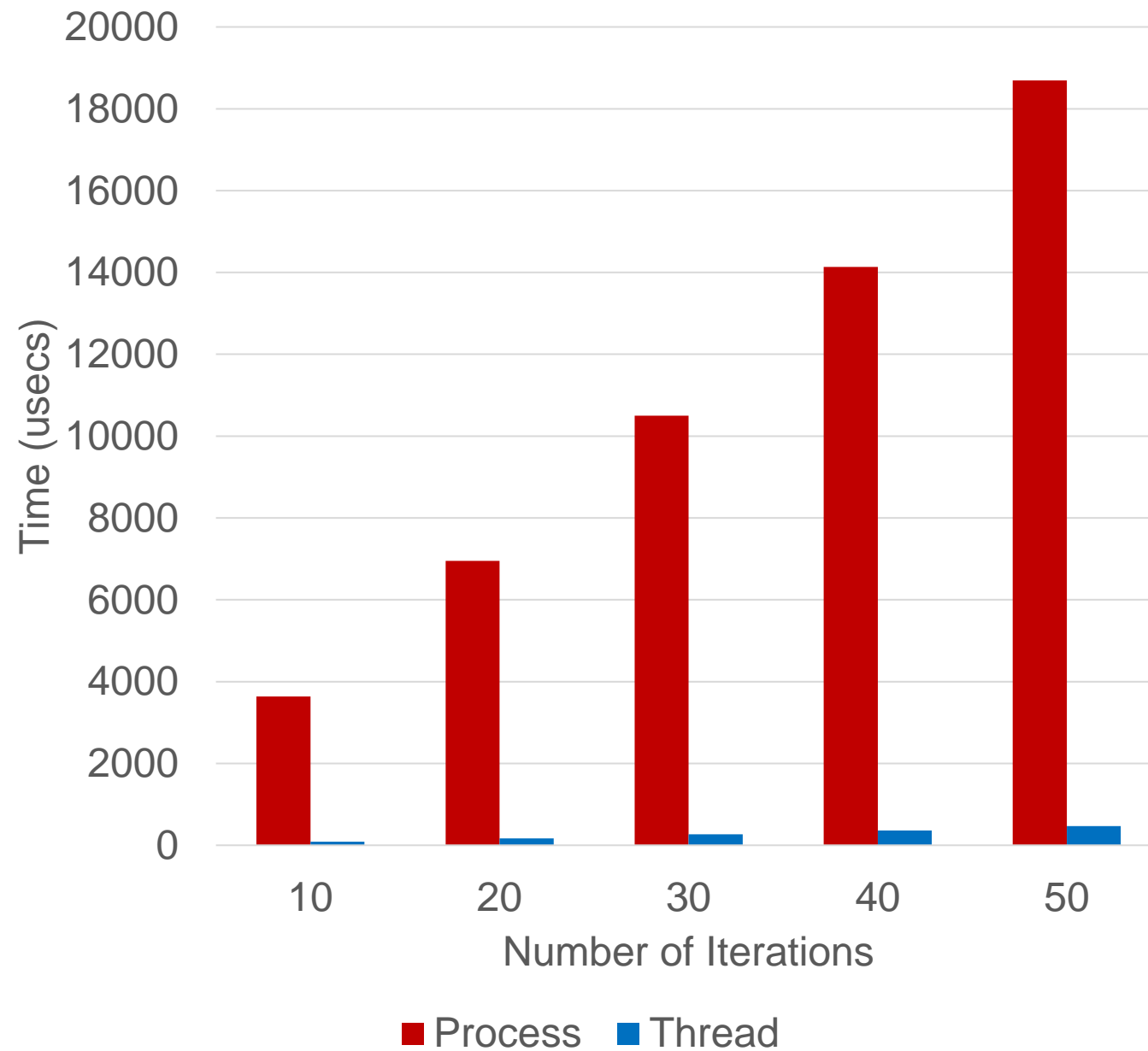
    Start Key_count_thread

# Multithreaded Server Architecture



(1) request

(2) create new
thread to service
the request

client → server → thread

(3) resume listening
for additional
client requests

# Process/Thread Creation Overheads



| Iterations | Process | Thread |
|:---:|:---:|:---:|
| 10 | 3633.80 | 87.80 |
| 20 | 6947.80 | 168.20 |
| 30 | 10501.00 | 269.80 |
| 40 | 14135.00 | 359.40 |
| 50 | 18693.00 | 470.60 |

# 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
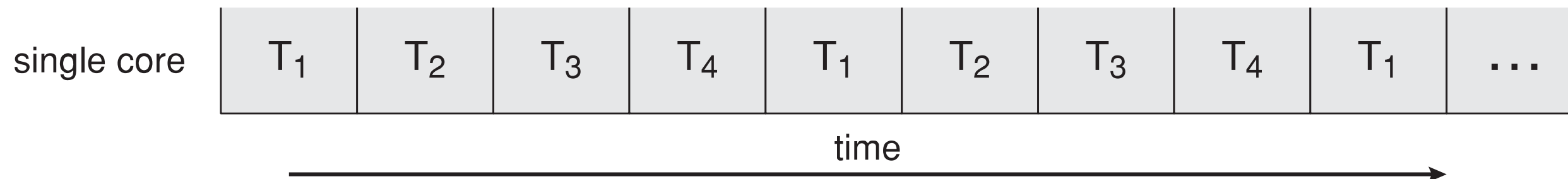
# Multicore Programming

- **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
- 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**
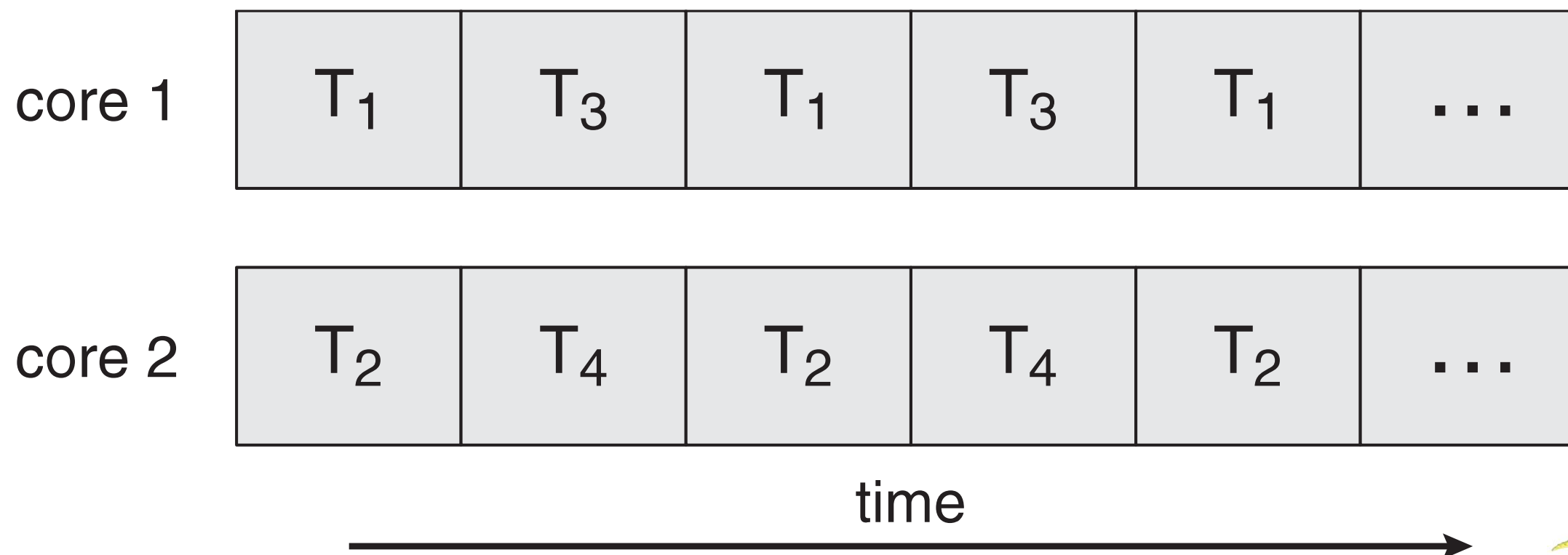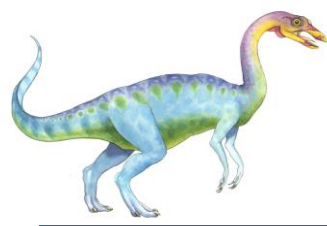  - Consider Oracle SPARC T4 with 8 cores, and 8 hardware threads per core

# Concurrency vs. Parallelism
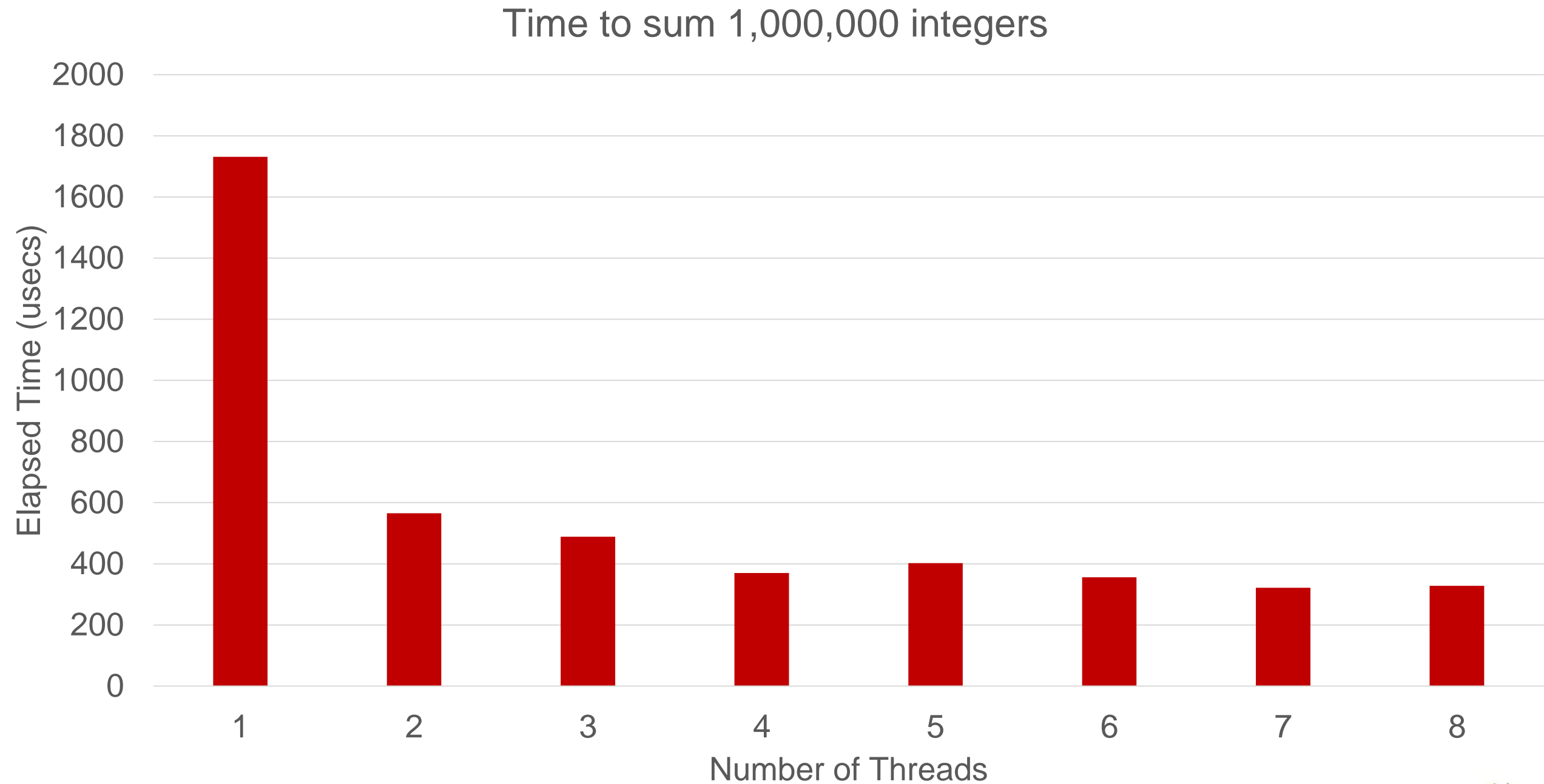
- **Concurrent execution on single-core system:**

| 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 ⟶

# Multi-Threading Performance
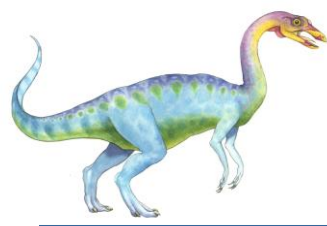
Time to sum 1,000,000 integers

# User Threads and Kernel Threads

- **User threads** - management done by user-level threads library

- Three primary thread libraries:
  - POSIX **Pthreads**
  - Win32 threads
  - Java threads

- **Kernel threads** - Supported by the Kernel

- Examples – virtually all general purpose operating systems, including:
  - Windows
  - Solaris
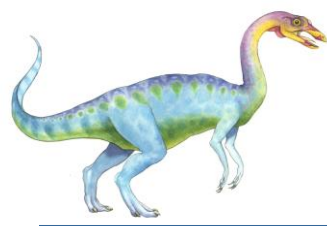  - Linux
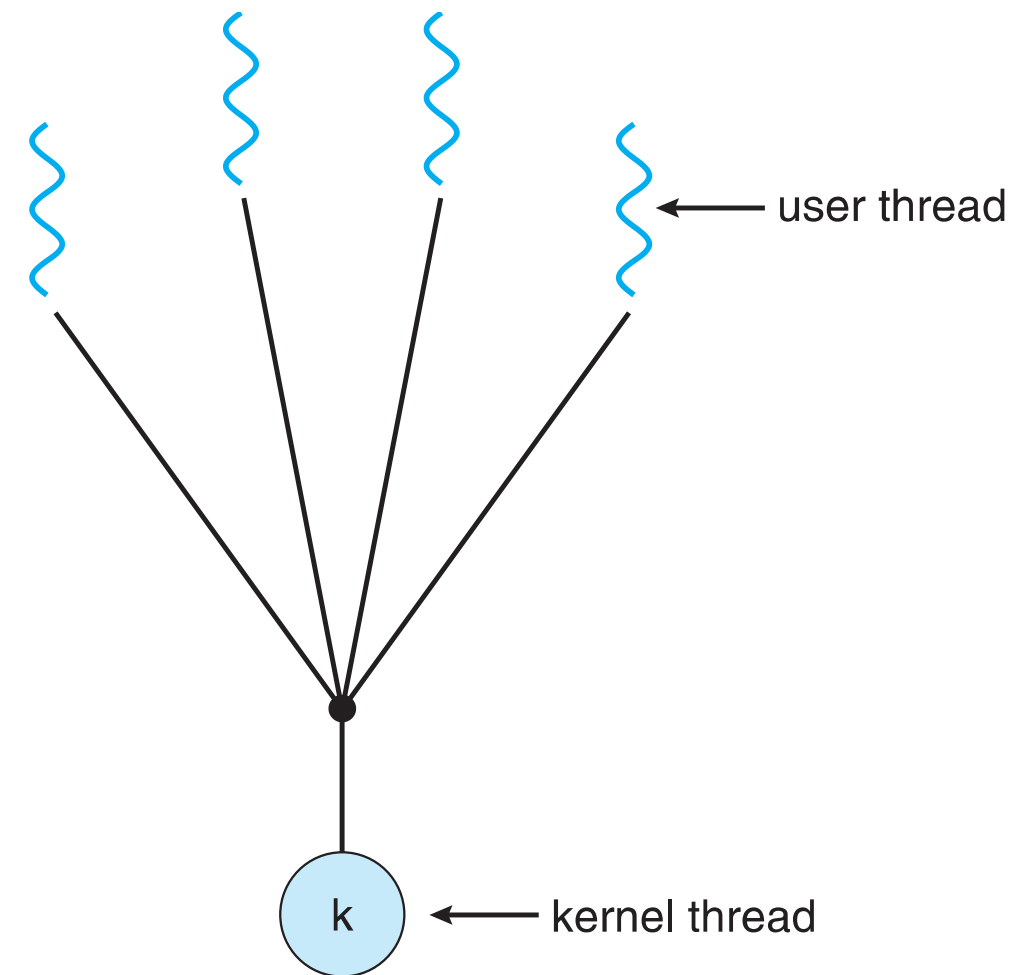  - Tru64 UNIX
  - Mac OS X

# 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**
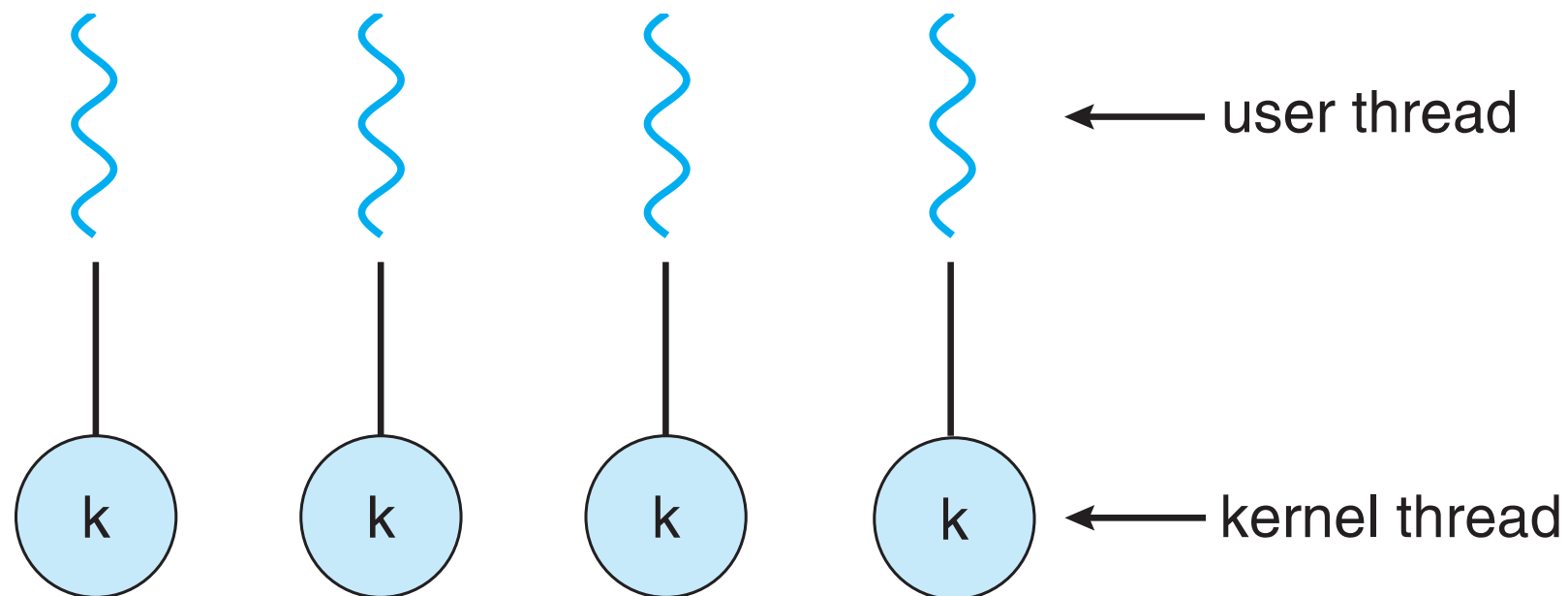
user thread

k ← kernel thread

# 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

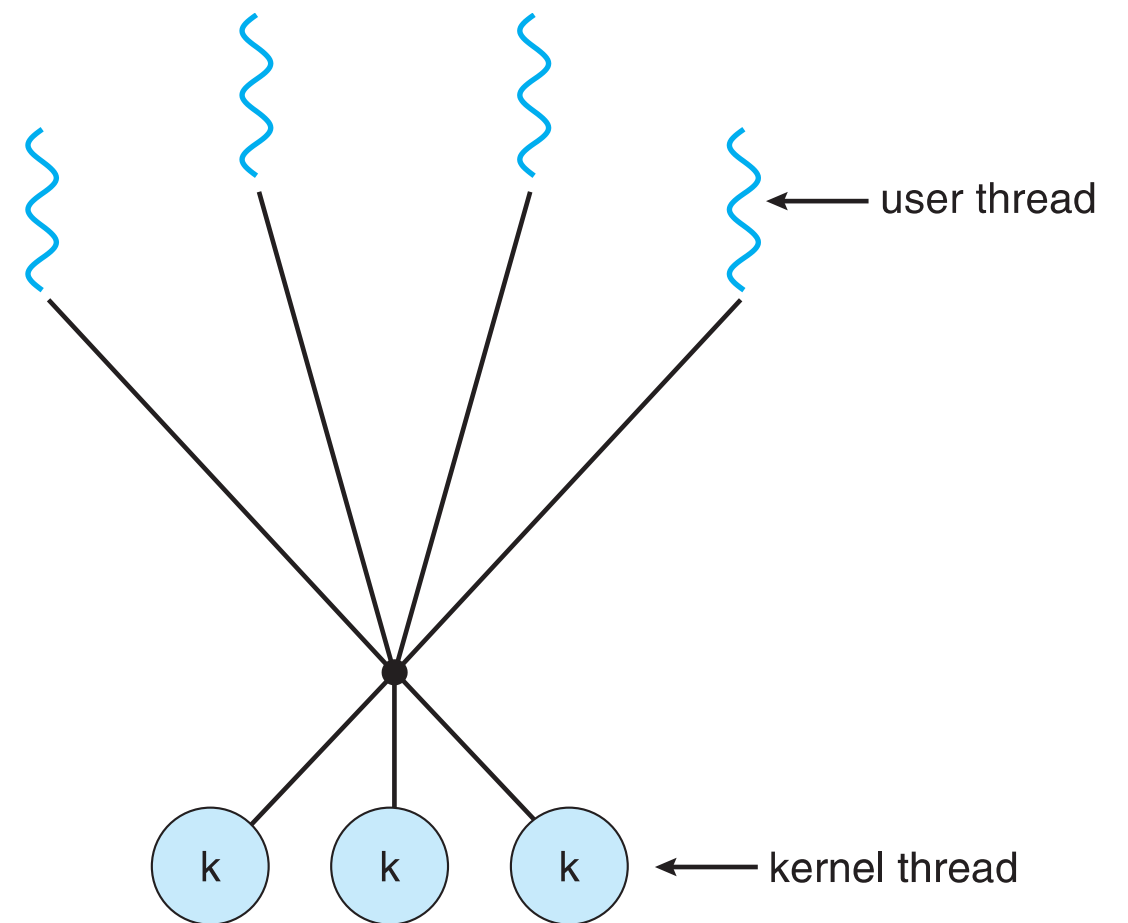- Number of threads per process sometimes restricted due to overhead

- Examples
  - Windows NT/XP/2000
  - Linux
  - Solaris 9 and later



user thread

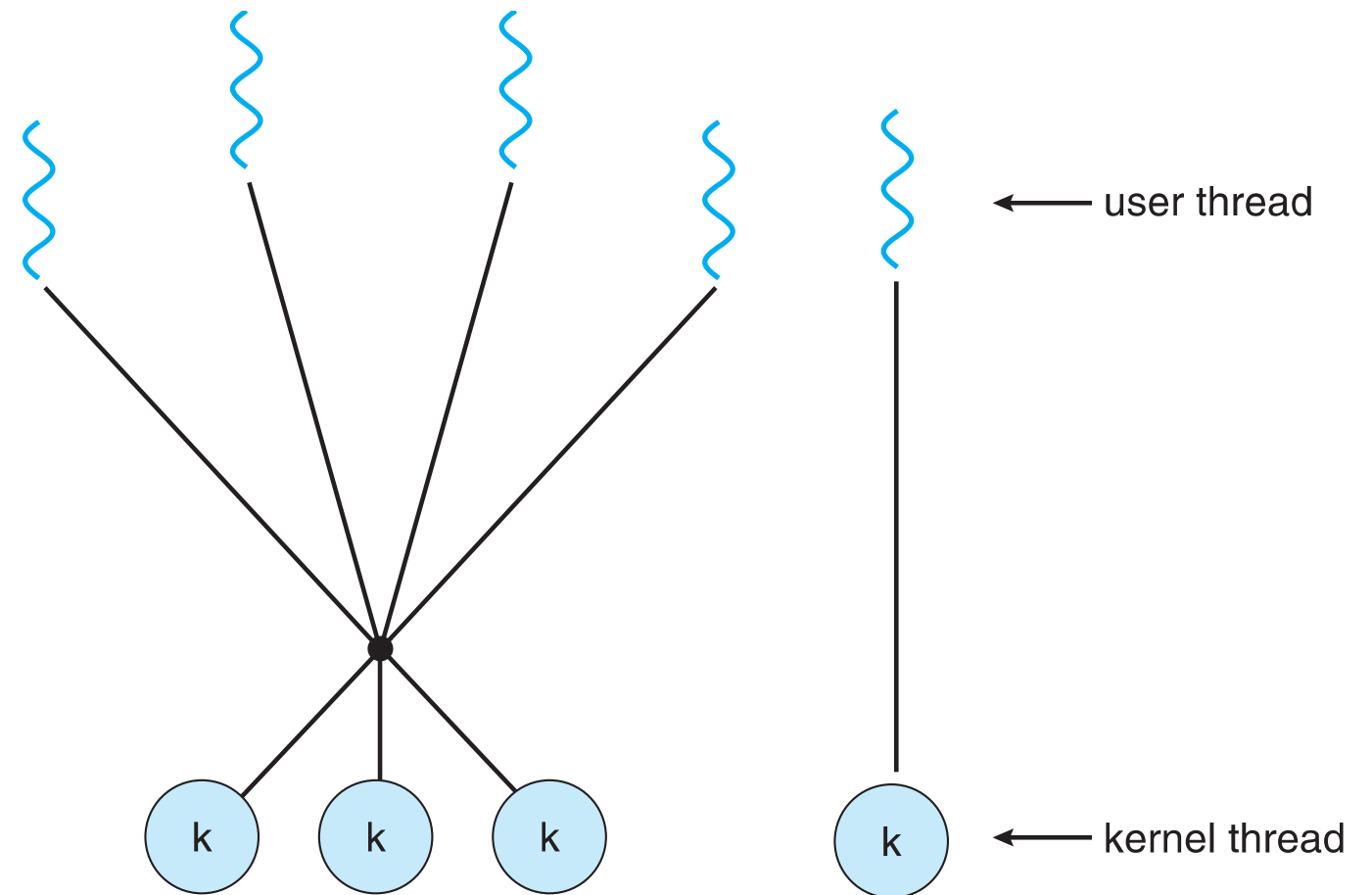kernel thread

# 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 NT/2000 with the ThreadFiber package

← user thread

k    k    k    ← kernel thread

# 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

# Thread Libraries

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

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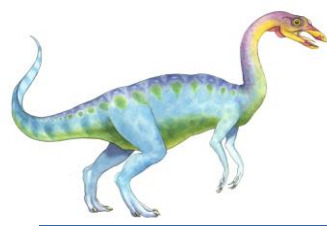  - Kernel-level library supported by the OS

# 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

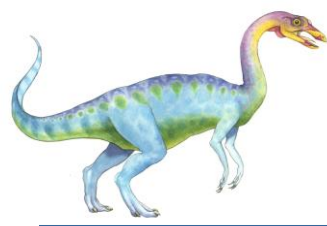- Common in UNIX operating systems (Solaris, Linux, Mac OS X)

# 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

- 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.
     */
}
```

# Threading Issues

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

- Signal handling
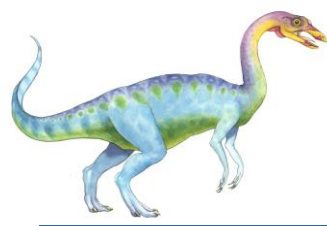  - Synchronous and asynchronous

# 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

- 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);
```