

- Chapter 4

◦ process creation is heavy-weight while thread creation is light-weight

◦ Threads can simplify code and increase efficiency.

◦ Examples for tasks can be implemented by separate threads

↳ Update display, Fetch data, spell checking, Answer a network request

◦ Client ⁽¹⁾ request → Server ⁽²⁾ → Thread ⁽³⁾ (2) create new thread to service the request
(Asynchronous) (3) resume listening for additional client requests

◦ Synchronous threading: the parent waits for the forked children to finish

↳ lots of data sharing between threads

ex: merge sort

◦ Asynchronous threading: the parent does not wait, all run concurrently

↳ lots of duplication of data, but no communication.

◦ Threading benefits:

↳ 1) Responsiveness: allow continued execution if part of process is blocked

2) Resource Sharing: easier than shared memory or message passing

3) Economy: 1) cheaper than process creation

2) thread switching lower overhead than context switch

4) Scalability: run threads on different cores, in parallel

◦ Parallelism: a system can perform more than one task simultaneously

↳ not possible on single core

◦ Concurrency: support more than one task making progress

↳ single core → scheduler providing concurrency

o types of parallelism:

↳ 1) data parallelism: distributes subsets of the same data across multiple cores with same operation on each

2) Task parallelism: distributing threads across cores, each thread perform unique operation.

of threads grows, so does architectural support for threading / CPU has cores and hardware

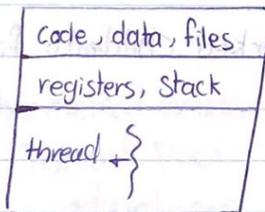
o Multicore programming challenges: 1) Task identification (which part to be parallel/concur)

2) Balance (equally valued tasks)

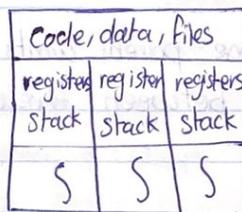
3) Data splitting

4) Data dependency

5) Testing and debugging



Single-thread
Process



Multi thread
Process

o Amdahl's law: $Speed\ up \leq \frac{1}{S + \frac{(1-S)}{N}}$ S: Serial portion, N: processing cores

o as $N \rightarrow \infty$, $Speed\ up \rightarrow \frac{1}{S}$

↳ Serial portion has disproportionate effect on the performance

o user threads: management done by user-level threads library

↳ like pthreads in POSIX, windows threads, Java threads

o kernel threads: Supported by the kernel

↳ treated as virtual processor

- o Mapping user threads to kernel threads
 - ↳ Many-to-One, One-to-One, Many-to-Many

- o Many-to-One: many user-level threads mapped to single kernel thread

- ↳ 1) One thread blocking causes all to block
- 2) Multiple threads may not run in parallel because only one may be in the kernel at a time (in multicore system)

Like Solaris Green, GNU Portable threads

- o One to One: Each user-level thread maps to kernel thread

- ↳ 1) Creating a user thread creates a kernel thread
- 2) More concurrency than many to one
- 3) # of threads/process sometimes restricted due to overhead

Like Windows, Linux, Solaris 9 and later

- o Many to Many: many user threads mapped to many kernel threads

- ↳ Allows the OS to create sufficient # of kernel threads

Like Windows with the ThreadFiber package, Solaris prior to version 9

- o Two-level model: it allows a user thread to be bound to kernel thread

Like IRIX, HP-UX, Tru64 UNIX, Solaris 8 and earlier

- o Thread library: provides programmer with API for creating/managing threads

- ↳ Implemented way: 1) Library entirely in user space
- 2) Kernel-level library supported by the OS

- o Pthread: provided either as user or kernel level

- ↳ POSIX Standard (IEEE 1003.1C) defining an API for thread creation and synchronization
- a specification for thread behavior not implementation

- o Java thread: managed by the JVM
 - ↳ may be created by: 1) Extending thread class
 - 2) Implement the runnable interface

o Program correctness more difficult with explicit threads

why implicit ~~threading~~ ^{threading}?
 why implicit ~~threading~~ ^{threading}?
 why implicit ~~threading~~ ^{threading}?

o Creation and management of threads done by Compilers and run-time Libraries rather than programmers, This is named as implicit threading

o three methods on Implicit threading: 1) Thread Pool

2) OpenMP

3) Grand Central Dispatch

Other method include Microsoft Threading Building Block (TBB), java.util.concurrent package

o Thread Pool: Create number of threads in a pool where they await work

↳ advantages: 1) faster to service a request with an existing

thread than create a new one

2) Numbers of threads in the application are

bound to the size of the Pool.

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

task (Ex: a task scheduled to run periodically)

o OpenMP: Set compiler directives and API for C, C++, FORTRAN

↳ 1) provide support for parallel programming in shared memory environments

2) Identifies parallel regions - blocks of code that can run in parallel.

[#pragma omp parallel] creates as many threads as there are processing cores

o Grand Central dispatch: Apple tech for Mac OS X and IOS

↳ Extensions of C, API, run-time library that allows developers to identify sections of codes to run in parallel (OpenMP)

→

- Chapter 4 - cont...

- Block: is a self-contained unit of work $\{ \text{Some code} \}$;
 - ↳ are scheduled by being placed in dispatch queue.
 - , when it is removed from the queue, it is assigned to available thread from the thread pool.

◦ dispatch queue

- ↳ 1) Serial (FIFO), each process has its own serial queue known as (main queue)
- 2) Concurrent (FIFO), but several blocks may be removed at once.
- with priority: low, default, high

- Signals in UNIX are used to notify a process that a particular event has occurred.

◦ Signal handler is used to process signals

- ↳ signals are: 1) generated by particular event
- 2) delivered to a process
- 3) handled by one of two signal handlers (default, user-defined)

- Every signal has default handler runs by the kernel

↳ can be override by user-defined signal handler

→ For ~~any~~ single-thread, signal delivered to process.

◦ For multithreaded processes:

- ↳ 1) signals ~~is~~ delivered to the thread which it applies
- 2) signal is delivered to every thread in the process.
- 3) signal is delivered to a certain thread in the process.
- 4) Assign a specific thread to receive ~~at~~ all signals for the process.

o thread to be canceled is target thread.

o Asynchronous cancellation: the target is terminated immediately.

o Deferred cancellation: allows the target to check periodically if it should be cancelled.

```
pthread_create(&tid, 0, worker, NULL); / pthread_cancel(tid);  
pthread_t tid;
```

o thread cancellation requests cancellation
↳ actual cancellation depends on thread state.

o If thread has cancellation disabled, cancellation remains pending until thread enables it.

o Default type is deferred.

→ Cancellation only occurs when thread reaches cancellation point

- pthread_testcancel()

- then cleanup handler is invoked.

o Thread cancellation is handled through signals in Linux.

o Thread Local Storage (TLS): each thread have its own copy of data

↳ useful when you do not have control over thread creation

(ex: when using thread pool)

2) Different from local variables → TLS visible across func invocations

local variables visible only during

single func invocation

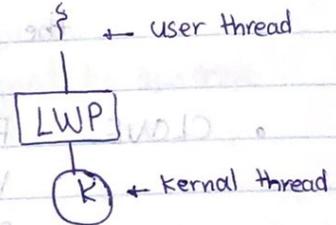
3) Similar to static data (TLS is unique to each thread)

(Scheduler activation)

o an intermediate data structure between user and kernel threads
is used (lightweight process (LWP)) in both M:M and two-level models
when requiring communication to maintain the appropriate number of kernel threads allocated to the application

→ it appears to be virtual process on which process can schedule user thread to run.

+ Each LWP is attached to kernel thread.



o Scheduler activations provide upcalls

→ upcalls: a communication mechanism from the kernel to the upcall handler in the thread library.

→ it allows an application to maintain the correct number kernel threads.

o Windows threads

↳ 1) implements the 1:1 mapping kernel level

2) Each thread contains: - thread id

Context of the thread ← [- Register set
- Separate user and kernel stack
- Private data storage area used by run-time and dynamic link libraries

o Primary data structures in windows threads include:

↳ 1) ETHREAD (executive thread block) ~~includes pointer to KTHREAD~~

2) KTHREAD (kernel thread block)

3) TEB (thread environment block) - thread id, user-mode stack, TLS
in user space.

→ Scheduling and synchronization info, kernel-mode stack, pointer to TEB
in kernel space.

→ includes pointer to process to which thread belongs to KTHREAD
in kernel space.

Linux threads (or tasks)

↳ clone() for thread creation

↳ allows a child task to share the address space of the parent task (process)

CLONE

FS: File-Sys information is shared

VM: the same memory space are shared

SIGHAND: signal handlers are shared

FILES: set of open files is shared

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