

並行計算

EP05: 執行線 (Thread)

在學習時你想你了解了，
到要寫下來時更加確定，
要教他人時就會更肯定，
但只有在寫程式時才能真正體認自己充分了解。

Alan J. Perlis

本集內容

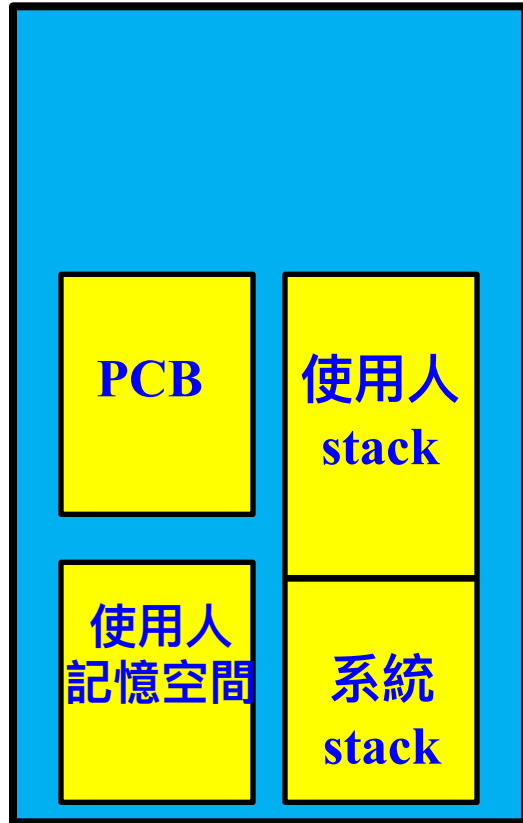
- 什麼是執行線 (thread) ?
 - 使用Thread的好處
 - 使用人 (User) Thread和系統 (Kernel) Thread
 - 複線作業 (Multithreading) 模型
 - 多核心程式寫作基本觀念
 - 取消Thread執行 (Thread Cancellation)
 - Thread專用和Thread安全
 - 協同函數 (Coroutine) 和纖維線 (Fiber)
 - 簡單的歷史回顧
 - 再說本質上的順序性 (Inherently Sequential)
- 請看影片的說明，那兒有取得投影片和程式的連結

什麼是Thread（執行線）

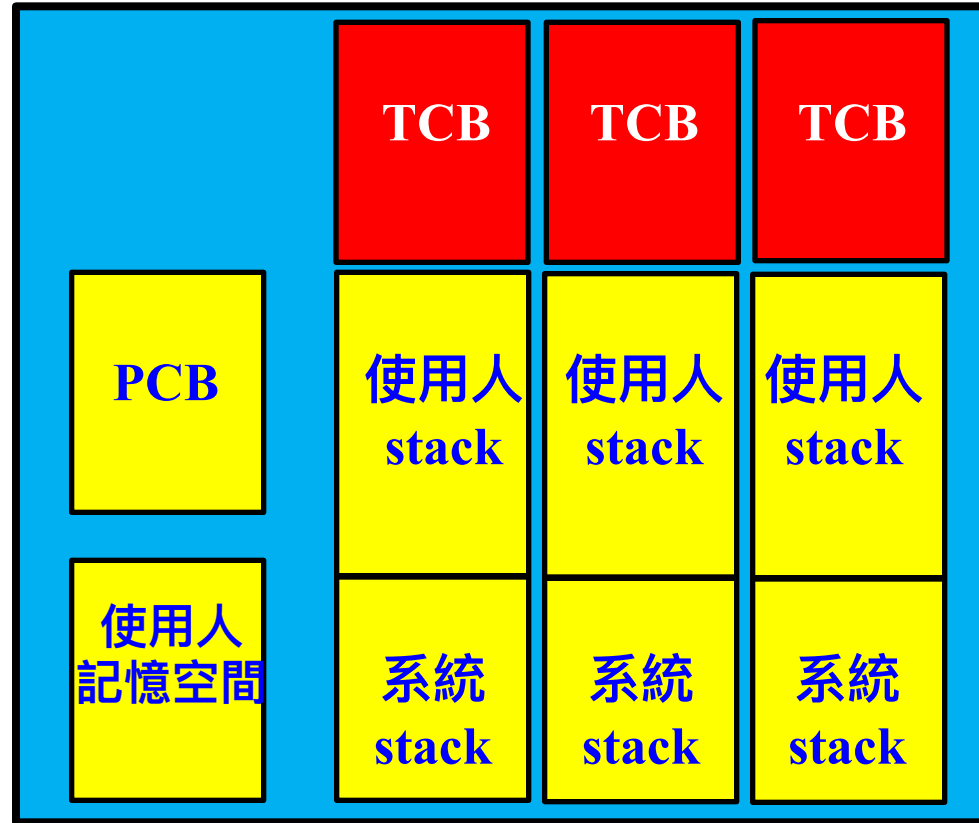
- **Thread（執行線）** 是一個基本的可以在CPU/Core中執行、由某個調度程式（Scheduler）管理的執行單位。這個調度程式可能是作業系統中的CPU調度程式（CPU Scheduler）。
- **Thread**由**process**產生。
- 一個**Thread**擁有**thread ID（TID）**、**program counter**、一組暫存器和**stack**。這一點和**process**相同。
- 但是，因為**thread**是由**process**產生，所以由某一個**process**產生的所有**thread**共享該**process**擁有的**code**和**data**區域，以及其它作業系統的資源（譬如記憶體、檔案等等）。
- 因為**thread**需要的系統資源比較低，**thread**有時會稱為**輕量process（Light Weight Process、LWP）**，而**process**則叫做**重量process（Heavy Weight Process、HWP）**。

單線和複線Process

單線 (Single-Threaded) Process



複線 (Multithreaded) Process



使用Thread的好處

- **反應能力**：如果一條**thread**（基於任何理由）不能執行了，包含它的**process**的其它**thread**仍然可能可以執行。
- **共享資源**：在約定的情況下，一個**process**產生的所有**thread**會**共用**許多系統資源（譬如檔案、記憶體等等）。
- **經濟效益**：因為一個**process**使用很多系統資源，產生或摧毀一個**process**、為它配置記憶體和資源、甚至於**切換執行環境**（**context switch**）都相當費時；但**thread**使用資源少，因此從事這些作業也相對省時間。
- **使用多處理機架構**：在多處理機系統中，若干個處理機可以同時執行同一個**process**中的若干**thread**（包含該**process**本身）。更重要的是，**不需要修改程式**就可以做到。

使用人Thread和系統Thread: 1/3

■ 使用人Thread

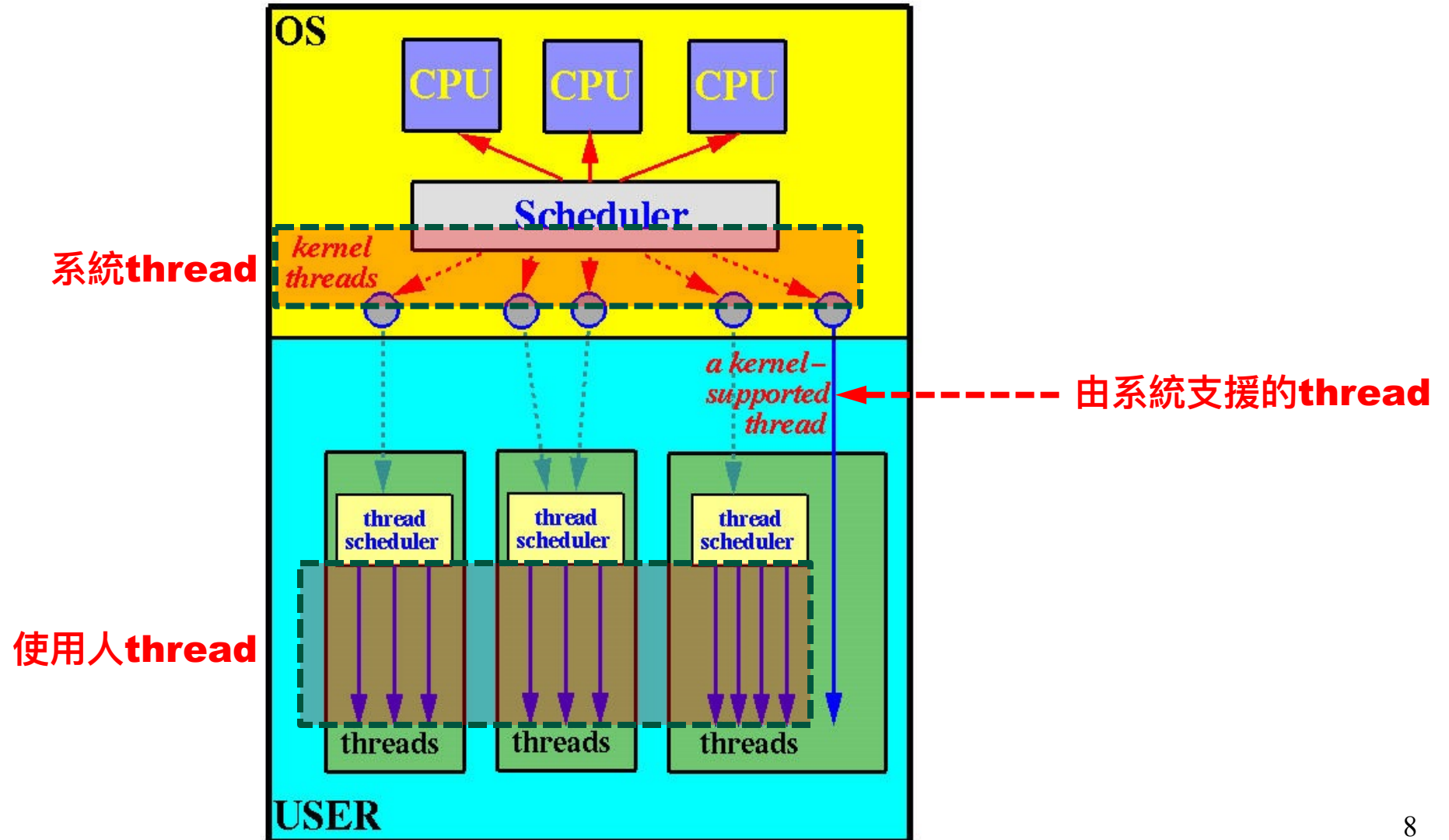
- ✓ 在使用人的層面支援使用人**thread** (**user thread**) ，系統不知道使用人**thread**的存在。
- ✓ 一般而言，使用人**thread**是由程式庫支援，可以產生、停止一個**thread**；可以和另一個**thread**匯合；甚至於調度 (**schedule**) 使用人**thread**。
- ✓ 因為使用人**thread**不需要系統介入 (不需要透過叫用系統) 處理，使用人**thread**效率會比較高。
- ✓ 然而，因為系統只認得包含了產生這些使用人**thread**的**process**，一旦一個**process**被暫停執行，它產生的使用人**thread**也都不能執行。換言之，當產生這些使用人**thread**的**process**失去了CPU使用權，它產生的使用人**thread**也就無CPU可用。

使用人Thread和系統Thread: 2/3

■ 系統Thread

- ✓ 作業系統支援**系統thread**。系統負責產生、停止一個**thread**、和另一個**thread**匯合、也負責調度 (schedule) 系統**thread**。
- ✓ 從事上述工作時都得叫用系統、在**OS**內完成，因為這些額外的系統負擔，**系統thread**的效率 (和**使用人thread**相比) 是較低的。但因為**thread**所需資源較少，所以**系統thread**的效率仍然比**process**高。
- ✓ 另一方面，當一個**系統thread**由於某種原因不能執行時，同一個**process**的其它**系統thread**仍然有機會執行，CPU調度程式 (scheduler) 可以挑另一個**系統thread**執行。
- ✓ 在多處理機環境下，作業系統可能會把這些**系統thread**安排到不同的處理機上執行。

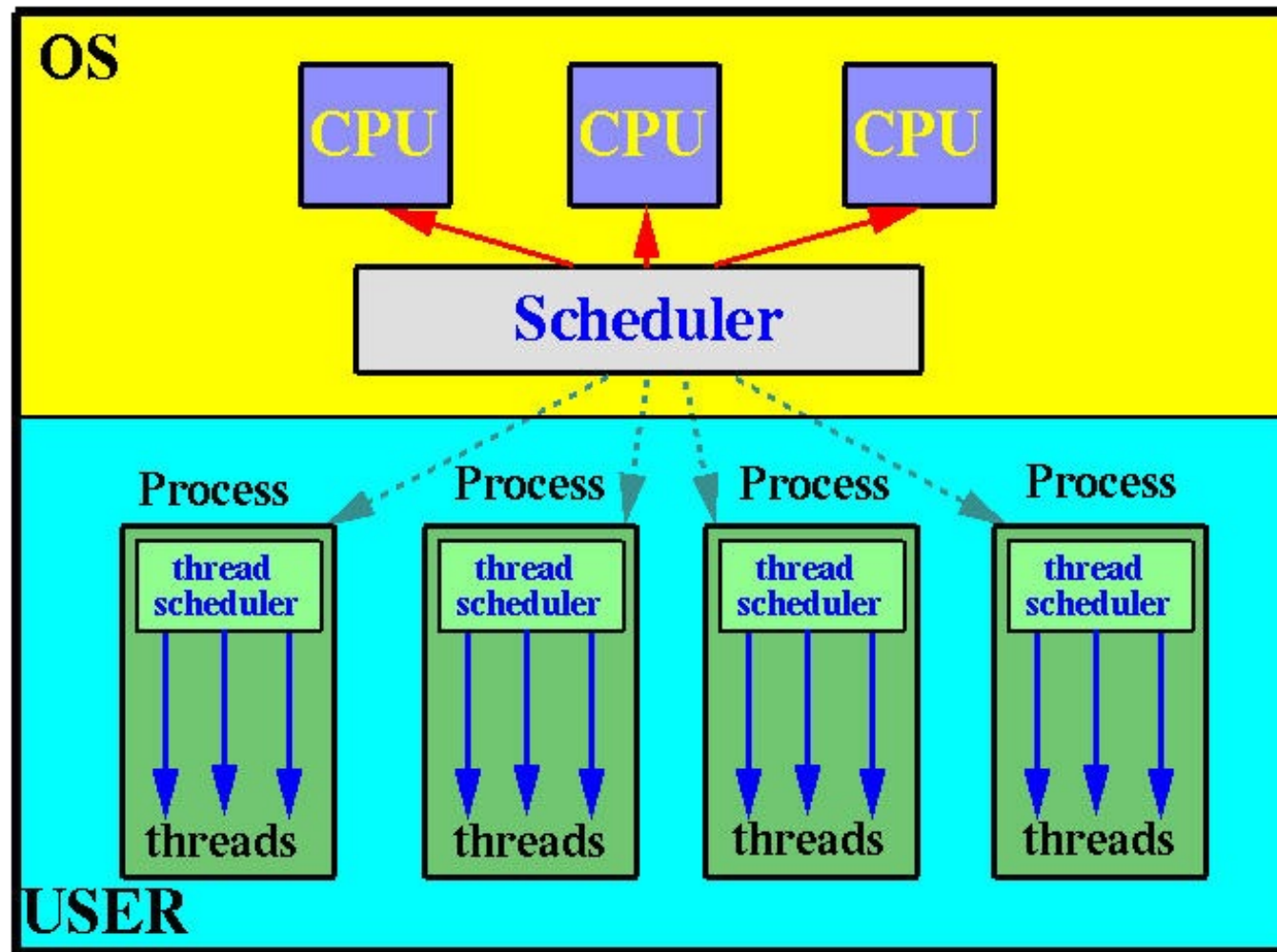
使用人Thread和系統Thread: 3/3



複線作業 (Multithreading) 模型

- 不同的系統可能會用不同的方式支援**thread**，下面是三種常見的方式：
 1. **多-對-1 (Many-to-One)**: 一個**系統thread** (或**process**) 給若干個**thread**使用，所以這是個**使用人thread**模型。
 2. **1-對-1 (One-to-One)**: 每一條**使用人thread**都對應著一條**系統thread**；一些**UNIX/Linux**系統和**Windows**用此方式。
 3. **多-對-多 (Many-to-Many)**: 若干條**使用人thread**對應著若干條**系統thread**。

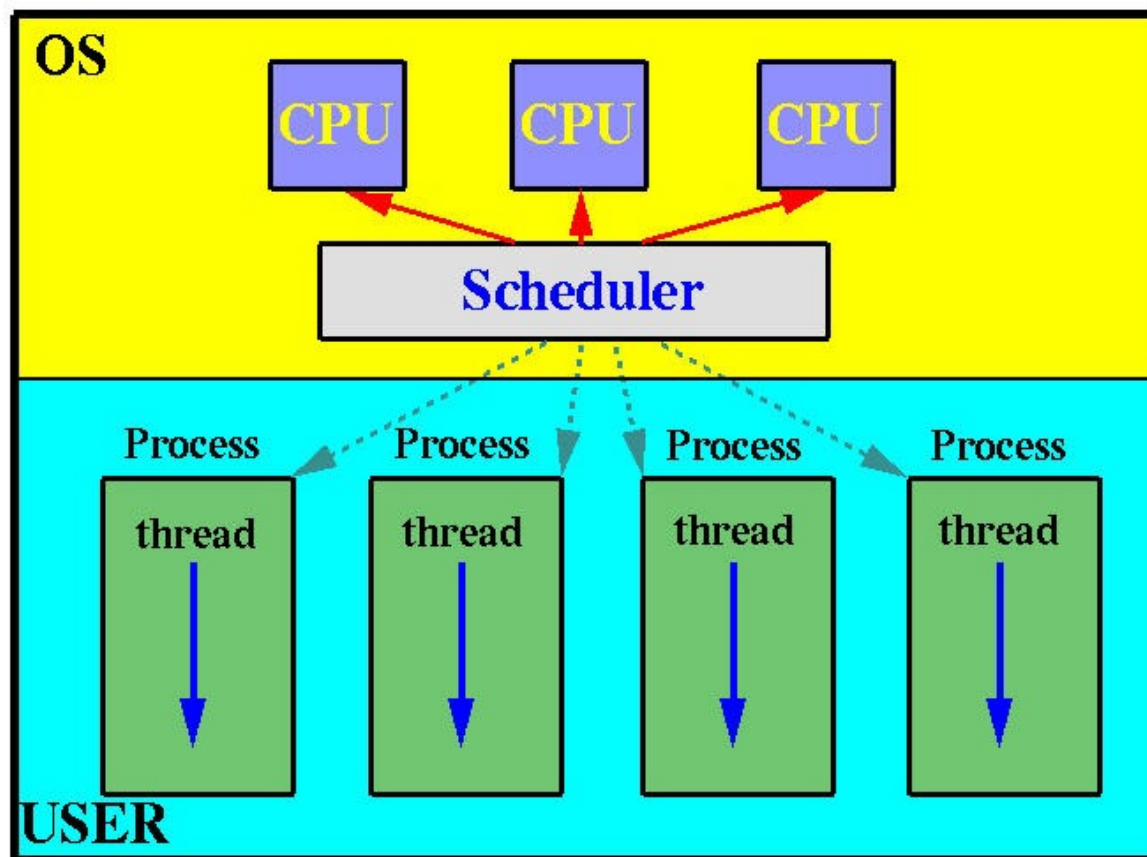
多-對-1 (Many-to-One) 模型



每一個**process**都有若干使用人**thread**，它們都對應到一條（給**process**用的）系統**thread**。
若**process**被暫停（失去**CPU**使用權），在它之下的所有使用人**thread**也都不能執行。

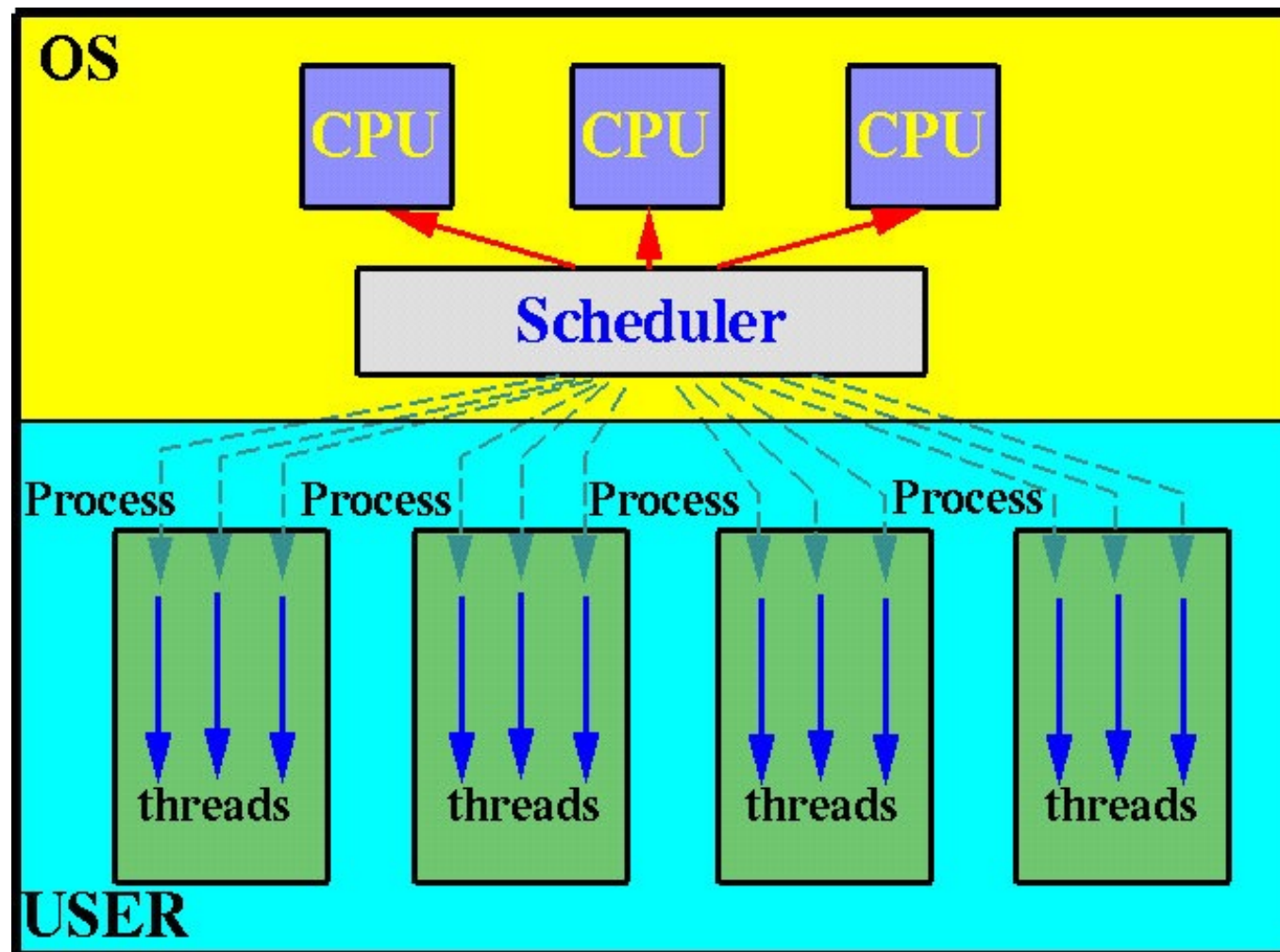
1-對-1 (One-to-One) 模型: 1/2

一個極端的個案: 傳統Unix



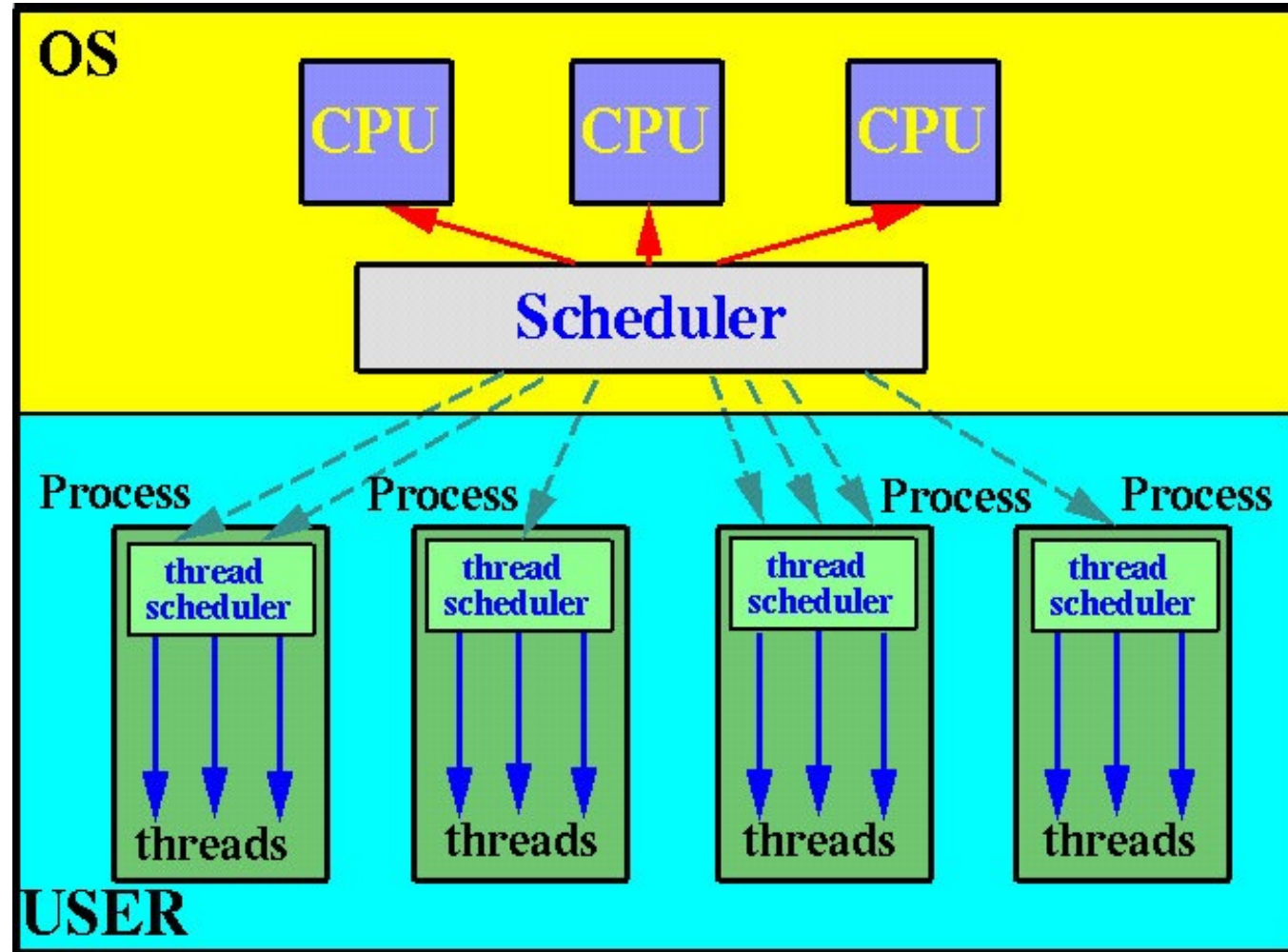
每一個process都只有一條使用者thread，它對應到一條系統thread。

1-對-1 (One-to-One) 模型: 2/2



每一個**process**都有若干條使用人**thread**，每一條使用人**thread**都對應著一條系統**thread**。
若一條系統**thread**不能執行，它對應的使用人**thread**就不能執行。

多-對-多 (Many-to-Many) 模型



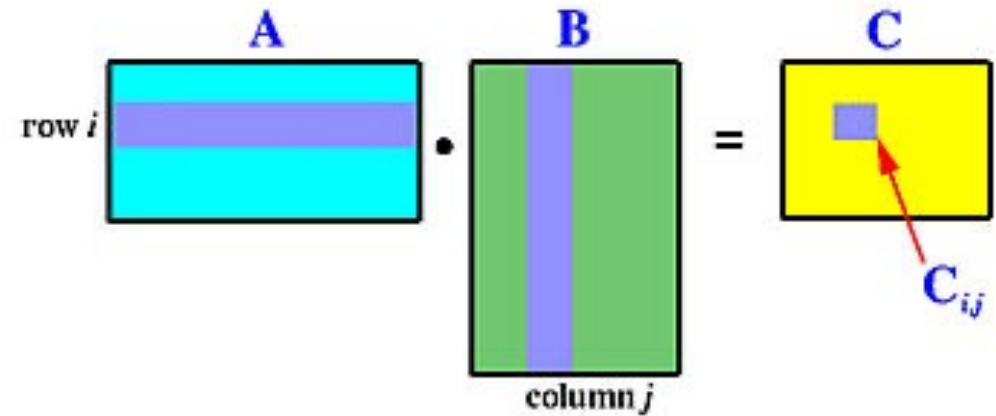
每一個**process**都有若干條**thread**，它們對應到若干條系統**thread**（數目未必相同）。若一條系統**thread**不能執行，在當時對應著該系統**thread**的使用人**thread**就不能執行。

多核心程式寫作基本觀念: 1/6

- 在單核心 (**core**) 的CPU之下，CPU調度程式每次只能選一條 **thread** 在唯一的CPU中執行。
- 若一個CPU 有若干個核心 (**core**) ， CPU 調度程式可以選若干條 **thread** 、每一條 **thread** 都佔用一個核心執行。
- 支援多核心的系統要比想像中複雜得多，寫作程式時也如此。
- 在多核心CPU 系統下工作的程式必須注意到五項課題：**工作分配、工作量均衡、資料劃分、資料獨立性、測試與偵錯。**

多核心程式寫作基本觀念: 2/6

- 工作分配：因為每一條**thread**可以在不同的核心下執行，我們得分析問題，把程式的工作分拆出來、使得可以分散到不同核心上並行地執行。
- 矩陣相乘是一個極好的例子。
- 不幸的是，有一些問題**本質上就是順序**（**inherently sequential**）的，很難被轉換成平行（或並行）。



$$C_{i,j} = \sum_{k=1}^n A_{i,k} \times B_{k,j}$$

對每一個 C_{ij} 都可以產生一條 **thread**

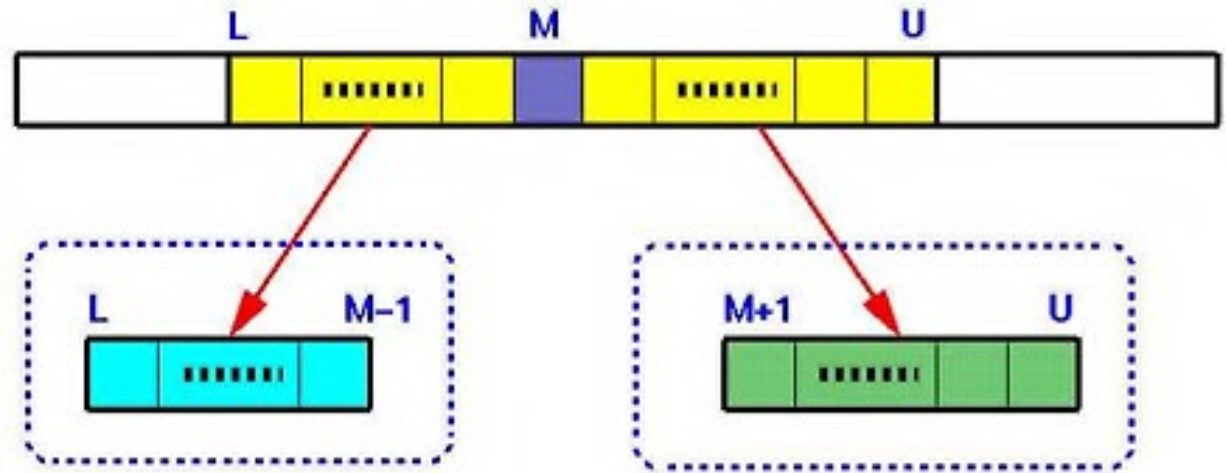
- ✓ 每一個 C_{ij} 都需要 n 個乘法
- ✓ C 有 $n \times n = n^2$ 個元素，總共需要 $O(n^3)$ 個乘法，所以總乘法數目是 $O(n^3)$ 。
- ✓ 如果我們用 n^2 條 **thread**（每一個 C 中的元素一條），每一條 **thread** 只要個 n 乘法算出，所以用複線做法比不用時快 $O(n^2)$ 倍，雖然我們需要 n^2 條 **thread**。

多核心程式寫作基本觀念: 3/6

- **工作量均衡**：儘可能讓分配給每一條**thread**的工作量佔總工作量的貢獻比率大致相同。
- 若一條工作量不顯著的**thread**常常被執行，它就常常佔用一個核心，於是其它從事重要計算工作的**thread**能被執行的機會就降低，整個程式的效率也就跟著降低。

多核心程式寫作基本觀念: 4/6

- **資料劃分**：需要處理的資料可以分成若干段落、每一個段落都可以獨立處理。
- 矩陣相乘是個好例子。
- 快速排序 (Quicksort) 是另一個好例子。當我們用某個中間值把輸入陣列切割後，中間值左 (或右) 邊的值都比中間值小 (或大) ；換言之，中間值已經在它最終的位置。於是，左邊和右邊就可以分開並且獨立地排序。
- 程式習題2中的**並行合併排序**則是第三個好例子。



把 $a[L..U]$ 分割成 $a[L..M-1]$ 和 $a[M+1..U]$ 之後，我們可以產生兩條 **thread**，一條處理左段、另一條處理右段。每一條 **thread** 都只排交給它的那一段、而不會碰到另一段。因此，產生 **thread** 的方式就是個二元樹狀結構 (binary tree) ；
程式習題2的合併排序也完全相同。

多核心程式寫作基本觀念: 5/6

- 資料獨立性：要注意到被若干**thread**共用的資料；譬如說，有沒有若干**thread**同時更新一個共用的變數等等。事實上，如果若干**thread**會用到同一筆資料，這筆資料在各**thread**之間就不再有獨立性，而是某些**thread**會依賴另一些**thread**更新該筆資料後的新值。
- 若這種情況發生，就很可能會得到意想不到的結果。因此，各**thread**之間就得**同步**（**synchronization**）、使得在任何時刻都只會有一條**thread**會更新一個共用變數的值。
- 寫作複線作業程式時，**同步**（**synchronization**）是一項有相當難度、也不容易做得好的課題。

多核心程式寫作基本觀念：6/6

- 測試與偵錯：複線作業程式的行為是**動態**（**dynamic**）的，特別是有共享的資源時就更是如此。在這次測試中出現的錯誤在下一
次測試中未必會發生；或者是這次測試沒有出現錯誤、但在下一
次測試中就跑了出來；或者是每次都出錯、但錯誤卻都不同。
- 有些錯誤在整個程式的生命期都不會出現，也有可能完全出乎
意料的時刻出現。
- 有一些偵錯的課題（譬如**競爭狀況** — 同時更新一個共享的資料
，或**系統死鎖**）沒有很效率的解法，我們能夠做的就是設計
和寫作程式時不要把這些錯誤放進去。
- 所以，測試和偵錯是一項藝術，而且寫程式時就需要很仔細地規
劃和設計。往後我們會陸續討論這些課題。

取消Thread執行 (Thread Cancellation) : 1/2

- **取消thread執行**指的是中止正在執行的thread；也就是說，在一條thread完成執行之前就中止它。這條被取消的thread叫做**目標執行線 (target thread)**。
- 有兩種取消的方式：
 - ✓ **非同步取消 (asynchronous cancellation)**：目標thread立即停止。
 - ✓ **延遲取消 (deferred cancellation)**：目標thread固然會停止執行，但會容許它週期性地檢查自己是否應該進入已結束 (**terminated**) 狀態，讓系統收回給它的資源。這樣就給目標thread有序地結束的機會。
 - ✓ **目標thread**可以讓自己進入結束狀態的地方叫做**取消點 (cancellation point)**。

取消Thread執行 (Thread Cancellation) : 2/2

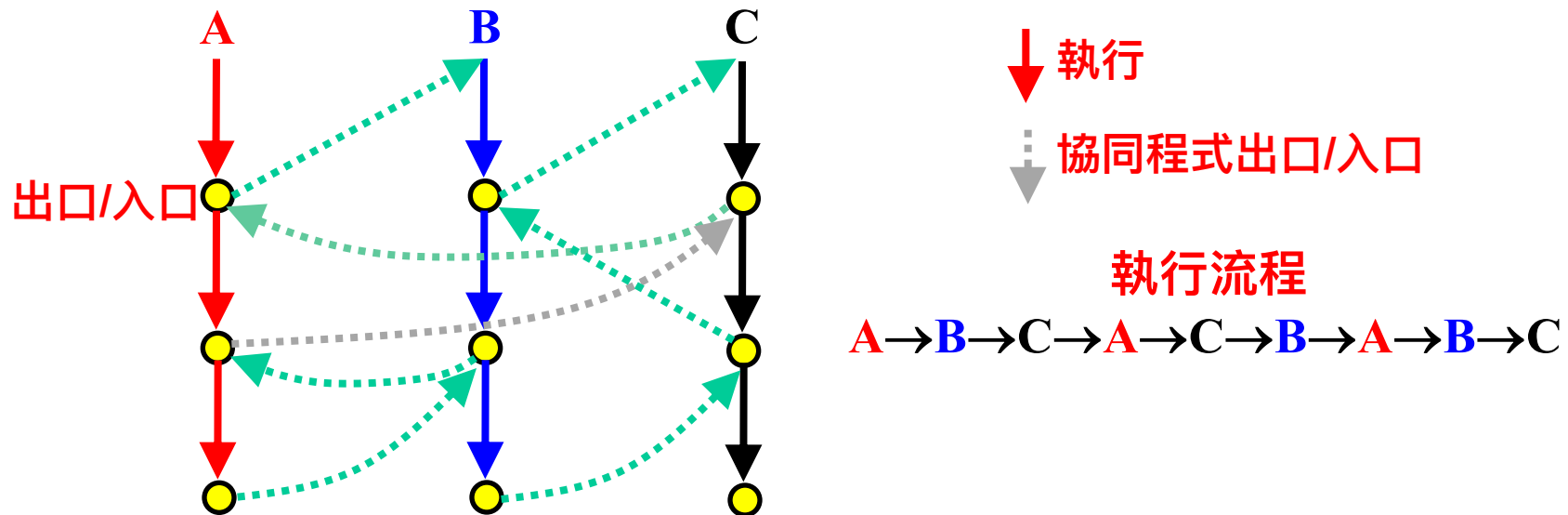
- 在**非同步取消**下，如果一條**目標thread**擁有一些系統中的共用資源，**目標thread**被取消時系統不一定能收回這些資源，因為可能有其它**thread**正在使用。
- 在**延遲取消**下，**目標thread**可以自己決定何時**終止**執行，因此收回這些資源不是個問題。
- 很多系統對**process** 和**thread**使用**非同步取消**，譬如**Unix**中用**kill**命令或**kill()** 叫用系統對**process**執行**非同步取消**。
- **POSIX Threads** (Pthreads) 支援**延遲取消**。

Thread專用和Thread安全

- **Thread專用** (**Thread-specific**) 資訊指的是那些**thread**用來做自己作業的資訊。
- 對**thread專用**資料支援不足就可能發生問題；譬如說，雖然每一條**thread** 都有自己的stack，但它們都共用產生它們的**process**的堆積 (**heap**)，請回顧第三集中有關一個**process** 記憶空間劃分的說明 (**第3講第16頁**)。
- 如果兩條**thread**同時叫用`malloc()` 在**heap**中配置記憶體 (**第3講第6-7頁**) ？或者是，兩條**thread**同時使用`printf()` 顯示資訊 (**第三講第32-34頁**) ？
- 一個程式庫的函數若同時被若干**thread**叫用時都能正常工作，這個程式庫 (或函數) 就是**thread安全**的。

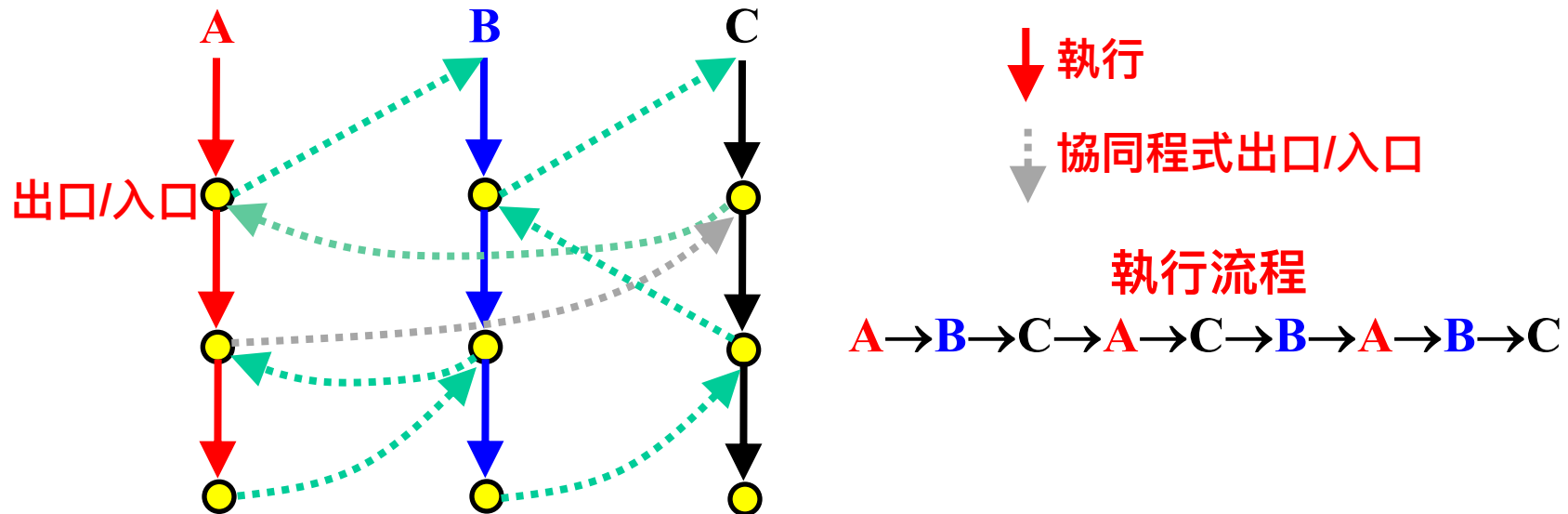
協同函數和纖維線: 1/3

- 一般叫用函數時，函數總是從該函數的第一條指令（入口）開始執行，雖然離開該函數時可能有若干出口。
- 一個協同函數（coroutine）有若干個入口、也有若干個出口。於是「叫用」一個協同函數時，會從上一次離開該函數的出口的下一條指令開始執行。



協同函數和纖維線: 2/3

- 從上一頁的**出口/入口**的方式看來，它是否有一點像**調度程式 (scheduler)** 呢？
- 沒錯！一個出口相當於被**調度程式暫停**，而入口則是被調度程式選中繼續執行，當然這就是從被暫停時的下一道指令開始。
- **協同程式幾乎就是調度程式進行環境切換 (context switch)** 的翻版。

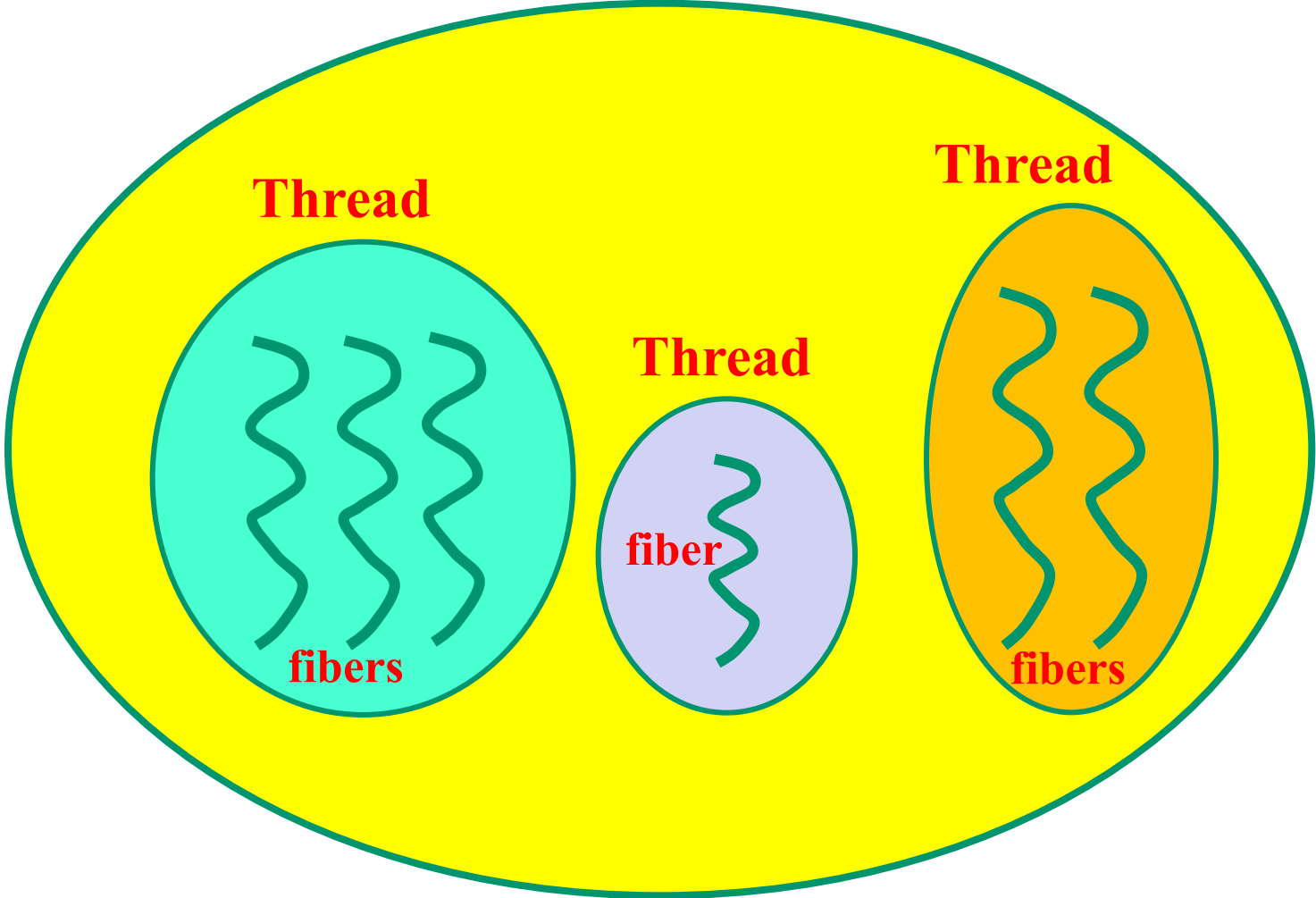


協同函數和纖維線: 3/3

- 執行纖維線 (簡稱纖維線 — **fiber**) 就是輕量的**thread** , 這就像**thread**是一條輕量的**process** ; 但是 , 纖維線必須人工地被應用程式調度 , 也就是沒有調度纖維線的調度程式。
- 纖維線是由**thread**產生 , 它和其他被該**thread**產生的纖維線共享該**thread**的資源。
- 一條纖維線在產生時擁有一個**stack**、一組暫存器、它的資料區域 (或是它的局部記憶體) 。
- 纖維線透過一種合作 (**co-operative**) 方式調度。
- 合作調度 (**co-operative scheduling**) 是使用某種「自願」的方式放棄CPU而由下一條纖維線使用。一些有纖維線的語言使用**YIELD**或類似的函數 , 讓叫用者讓出CPU使用權。
- 纖維線比執行線 (**thread**) 簡單 , 而且用類似協同程式的方式作業。

Processes \Rightarrow Threads \Rightarrow Fibers

Process



簡單的歷史回顧（Process和Thread）：1/4

- **Process** 這個詞最遲在**1960**年代中期之前就已經被廣泛使用。
- **Multiprogramming**的出現也差不多是**1960**年代前後，**Edsger W. Dijkstra**在荷蘭**Technische Hogeschool Eindhoven**（**Eindhoven University of Technology**，簡稱**THE**）發表了**THE**作業系統（**1965**年），這可能是**Multiprogramming**這個字正式浮上檯面的時刻，不過學術論文要到**1968**年才刊出。
- 另一方面，**IBM**在**1964**年**4**月發表**IBM System/360**，開始了電腦系統設計和相容性的全新視野。
- **System/360**的作業系統是**IBM System/360 Operating System**（**OS/360**），它在**1966**年**3**月**31**日交給客戶使用。

簡單的歷史回顧 (Process和Thread) : 2/4

- 事實上，**OS/360**有四套不同版本：
 - ✓ **PCP** (**Primary Control Program**) : 這是一套一次只能執行一個程式的**OS**，給那些小機型用的，很快就停用。
 - ✓ **Multiprogramming with a Fixed number of Tasks** (**MFT**) : 這是給中階機型用的系統，最多只能支援同時執行**15**個程式。
 - ✓ **Multiprogramming with a Variable number of Tasks** (**MVT**) : 這是給大型機用的，它去掉了最多只能執行**15**個程式的限制。
 - ✓ **Disk Operating System/360** (**DOS/360**) : 給小機型用的系統。
- 在**IBM**術語中，上述的**task**相當於**process**。
- 但是，**OS/360**和**DOS/360**中也支援**thread**，不過也叫做**task**。

簡單的歷史回顧（Process和Thread）：3/4

- 在**OS/360**和**DOS/360**之下，有不少**集體指令**（macro）可以產生**task**（也就是**thread**）並且進行**同步**（synchronization）；這些一直沿用至今：
 - ✓ 使用人可以用ATTACH產生**task**、用DETACH把已經終結的**task**刪除。
 - ✓ **IBM**用*serially reusable*（**序列性可以再使用**）資源表示任何時刻只容許一個**task**可以使用的資源，用今天的話來說就是**互斥**地使用。
 - ✓ 對這些**互斥**的資源而言，程式員可以產生一個**等待線**（queue），在使用一個共同資源前用enqueue（ENQ）集體指令等待使用該資源；使用完後用dequeue（DEQ）離開該等待線、讓下一個在等待線中的**task**可以使用該共用資源。
 - ✓ 使用人可以定義**事件**（event），用WAIT（事件名稱）等待該事件發生（有沒有像叫用系統wait()呢）；若某**task**工作完成，該**task**可以用POST宣告某事件已經發生。
 - ✓ **Hercules**是一個很好用的**IBM System/370**，**System/390**和**z/Architecture**的模擬程式（<http://www.hercules-390.org>）。目前**System/370**的幾個作業系統（包含**VM/370**）和編譯程式都可以在網上找得到，您不妨試試在個人電腦上玩玩這些幾十年前的經典系統。

簡單的歷史回顧 (Process和Thread) : 4/4

- **OS/360**並沒有用目前教科書方式進行**同步** (**synchronization**) , 但是也用另一方式 (ENQ/DEQ 和 WAIT/POST) 完成我們期望的結果。
- **IBM PL/I F (1966)** 中有相當完整的**複線** (**multithreading**) 執行功能, 但編譯程式卻並沒有公開給客戶使用, 因此使用人似乎只能使用集體指令解決這個問題。目前的**IBM z/OS**繼承了幾十年前的傳統, 這些功能都有了、而且也相容。
- 這些被產生的task會和其它task**共同被系統調度**。所以, **OS/360**對thread的支援應該是前述的**系統thread**。Thread這個字是Jerome Haward Saltzer在他的1966年博士論文中提到的 (第20頁腳註)、V. Vyssotsky建議的字, 但那是指process而不是今天熟知的執行線。

再說 **Inherently Sequential** 問題: 1/4

- 在理論計算機科學中，如果輸入個數是 n ，那麼 **P** 是所有我們可以用某個變數為 n 的多項式步驟能夠完成的問題的集合。
- 譬如，搜尋一個已經排序的陣列得用 $O(\log_2(n))$ 個比較（**二分搜尋法**）、用比較排序會有 $O(n\log_2(n))$ 的演算法（**合併法**、**堆積法**、等等）、矩陣相乘 $O(n^3)$ 等等。這些都被視為「**簡單**」的題目。
- 如果我們有足夠的處理機（或核心），而且在理想情況下每一個 **thread** 都可以取得一個處理機執行。
- **NC** 這個集合是所有用 $O(n^p)$ 個處理機、在 $O(\log_2^q(n))$ 步驟下可以完成的問題的集合，此地 p 和 q 是兩個和 n 無關的常數。

再說 **Inherently Sequential** 問題: 2/4

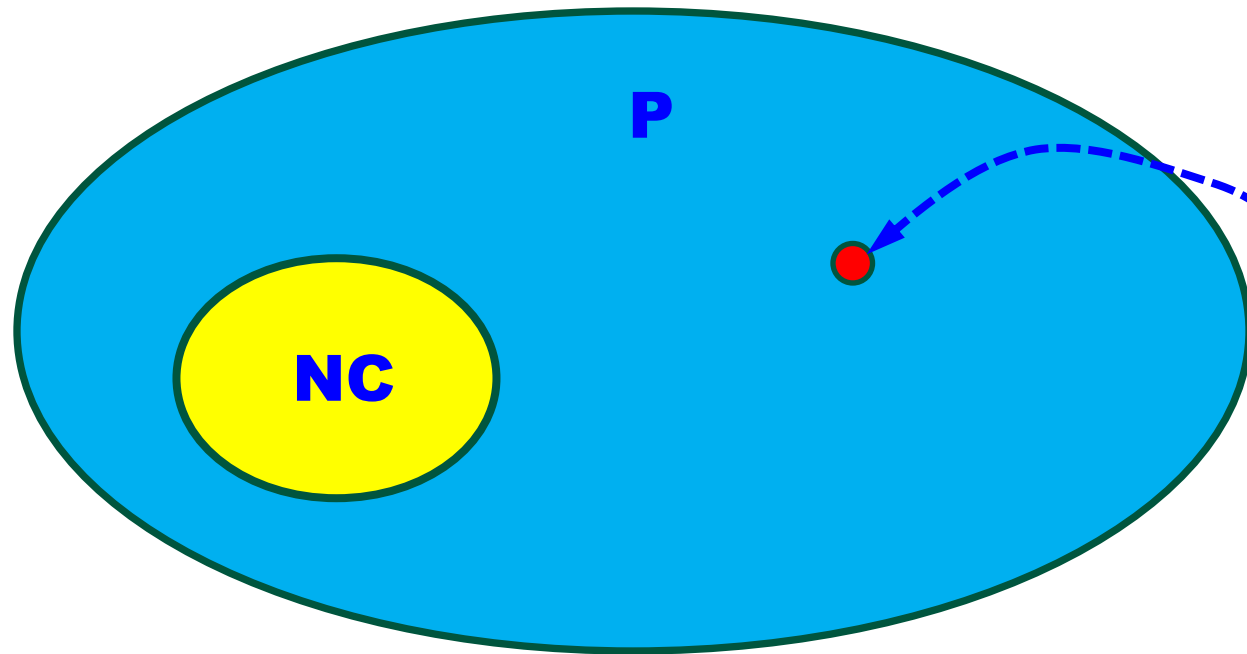
- 每一個 **NC** 中的問題也是 **P** 中的問題，因為我們只要把每一個處理機序列化就行了（譬如，把給第一個處理機的工作做完後再讓給第二個處理機繼續、等等），於是總工作量（在只有一個處理機的情況下）就是 $O(n^p \log_2^q(n))$ 。
- 很明顯地， $O(n^p \log_2^q(n)) < O(n^{p+q})$ ，因為 $\log_2(n) < n$ 。所以，每一個 **NC** 中的問題也都是 **P** 中的問題。
- **NC** 是 **Nick's Class**，這是 **Stephen Cook** 為這一類型問題的命名，此地 **Nick** 是另一位在這個領域有很廣很深研究的學者 **Nick Pippenger** 的名字。

再說 **Inherently Sequential** 問題: 3/4

- **問題是**： **NC** 是 **P** 的部分集合，但 **NC** 是否等於 **P** 卻是個到目前都未曾解答的問題（和 **P** 是否等於 **NP** 屬於同一類）。若 **NC** 不等於 **P**，那麼 **P** 中就至少會有一個問題，它不能用 $O(n^p)$ 個處理機、每一個處理機用 $O(\log_2^q(n))$ 個步驟做出來；換言之，這個問題無法很有效地用平行（當然也是並行）的方式解決。
- 您在演算法複雜度（**complexity**）或理論計算機科學的課程中學到，此地不打算再深入討論。

再說 **Inherently Sequential** 問題: 4/4

- 直覺地說，**本質上為順序性** (**inherently sequential**) 的問題就是難以 (或不可能) 被平行化的問題，它很可能不在 **NC** 中。



若 **P = NC**，則**本質上為順序性**的問題不存在，每一個在 **P** 中的問題都可以**平行化**。

若 **P \neq NC**，則至少有一個在 **P** 中的問題不在 **NC** 中，這是個**難以 (或無法) 被平行化**的問題。

我們學到了什麼？

- 了解**thread**的意義和不同系統下支援**thread**的方式
- 使用人**thread**和系統**thread**
- 複線作業模型
- 複線作業在多處理機（或核心）下的程式寫作基本概念
- 其它和**thread**相關的專有名詞
 - 協同程式（**Coroutine**）和執行纖維（**Fiber**）
 - 歷史回顧
 - 本質上順序性（**Inherently Sequential**）的意義

結束，謝謝收看！
期望您再次觀看下一集

請看影片的說明，那兒有取得投影片和程式的連結