

Serial vs Parallel

Serial

En una ejecución serial, el tiempo de ejecución viene dado por:

$$N = \text{no instr}$$

$$T = N/F$$

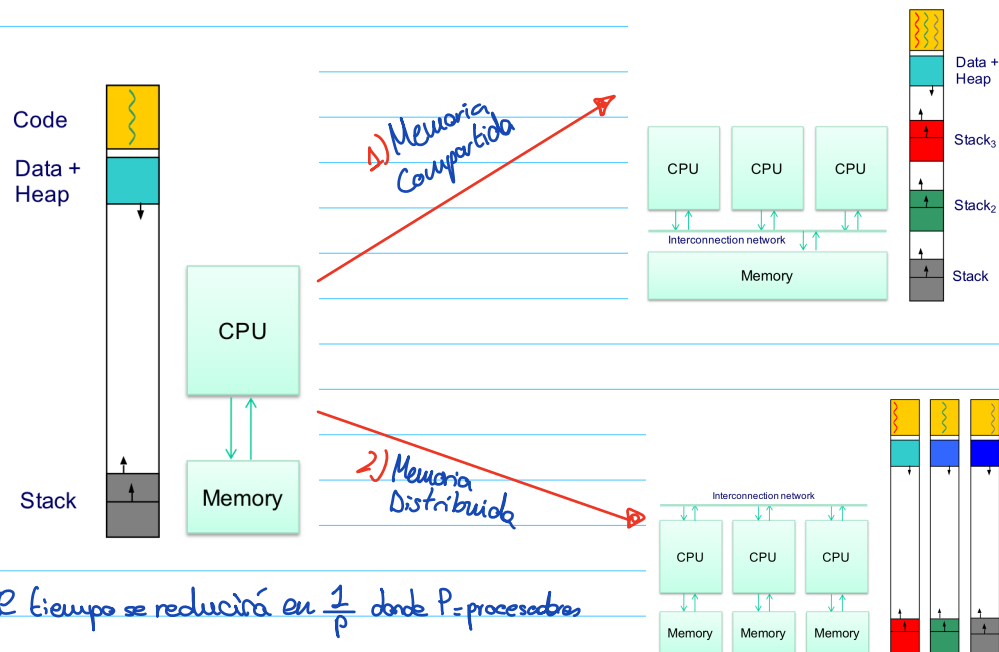
$$F = \text{IPC} \cdot f$$

↳ Inst. por Ciclo

→ Para reducir tiempo: ↑ F

Parallel

Otra manera de reducir el tiempo es mediante **arquitecturas paralelas**, pueden ser:



El tiempo se reducirá en $\frac{1}{P}$ donde P = procesadores

Nos queda entonces que:

nos centraremos en esta

$$T = (N/P)/F$$

N = no instr

P = no procesadores

F = Frecuencia

⚠ Habrá que coordinar las tareas (trozos de la aplicación) para un correcto acceso

Throughput computing

También podemos explotar la arquitectura paralela corriendo diferentes aplicaciones en diferentes procesadores a la vez [definición throughput computing]

K programas y P procesadores \longrightarrow cada programa recibe $\frac{P}{K}$ procesadores

El SO multiplexa su ejecución

Concurrencia

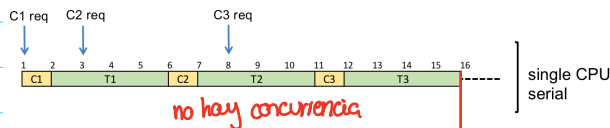
Ejecución concurrente

Consiste en romper un problema en partes (tareas) para asegurar su ejecución simultánea

Ejemplos

1) Serial \rightarrow 1 CPU

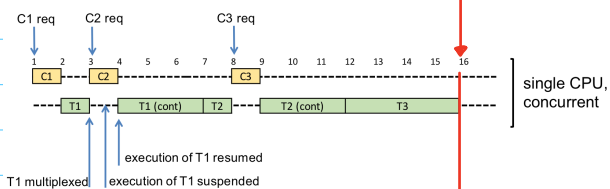
Task Ck: receives client requests (duration 1 time unit)
Task Tk: executes a single bank transaction (e.g. withdraw/deposit some money in bank account, with a duration of 4 time units)



single CPU serial

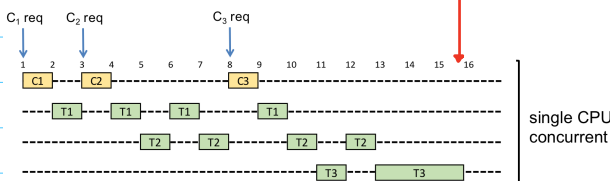
! No tienen por que ser en orden las tareas

2) Concurrente \rightarrow 1 CPU



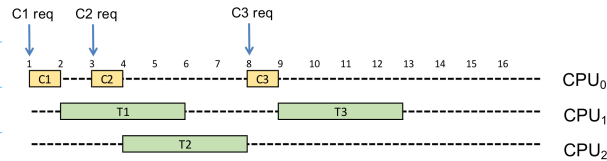
single CPU, concurrent

1 CPU, mismo tiempo siempre



single CPU concurrent

3) Paralelo → 3 CPU



Ejecución paralela: Reducir el tiempo de ejecución (**response**) de un programa.

parallel computing

Se usan múltiples CPU para ejecutar tareas identificadas para concurrencia.

1 programa → P procesadores

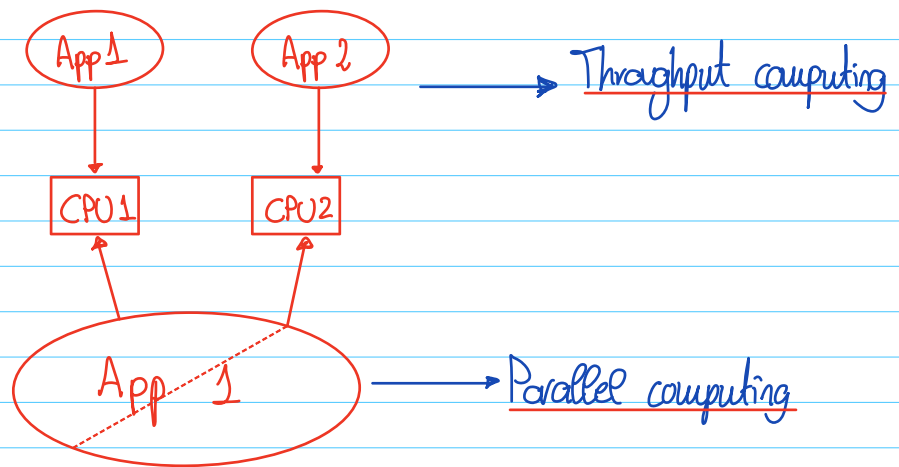
Cada CPU recibe $1/P$ del programa, se reduce el tiempo en P

Throughput computing: Se usan múltiples procesadores para incrementar el número de programas por unidad de tiempo.

n programas en P procesadores

→ $n > p$: cada programa recibe P/n procesadores

$n \leq p$: cada programa recibe 1 procesador



Problemas Concurrencia

Existen diferentes problemas con la concurrencia:

- Database condition

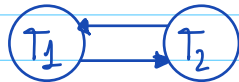
recursos compartidos con dos o más tareas que modifican (leen y escriben)

→ $sum = sum + 1$

- Si tenemos 3 tareas que acceden a la vez, $sum = 5$ no es!!
se solapan

- Deadlock

Das o más tareas que se bloquean por que se espera a que otra haga algo y esta otra también y así sucesivamente. Algunos



- Livelock

Das más tareas causan continuamente de estado en respuesta a cambios de otras tareas sin hacer nada útil

- Starvation

Una tarea no puede ganar acceso a un recurso compartido porque está esperando indefinidamente

- Performance

Problemas

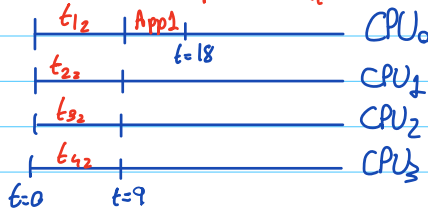
1. Assume we want to execute two different applications in our parallel machine with 4 processors: application *App1* is sequential; *App2* is parallelised defining 4 tasks, each task executing one fourth of the total application. The sequential time for the applications is 8 and 40 time units, respectively. Assuming: that *App1* starts its execution at time 4 and *App2* starts at time 0, draw a time line showing how they will be executed if the operating system:

- Does not allow multiprogramming, i.e. only one application can be executed at the same time in the system. \rightarrow entendiendo que solo un CPU no concurrente
- Allows multiprogramming so that the system tries to have both applications running concurrently, each application making use of the number of processors is able to use.
- The same as in the second case, but now *App2* is parallelised defining 3 tasks, each task executing one third of the total application.

4 CPU

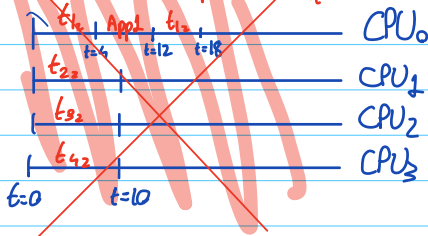
App1 \rightarrow Secuencial \rightarrow 8 u.t \rightarrow empieza en $t=4$
 App2 \rightarrow t_1, t_2, t_3, t_4 \rightarrow 40 u.t \rightarrow empieza en $t=0$
 0 10 20 30 40 u.t
 \downarrow
 Cada tarea 10 u.t

a) Asumiendo no tiempo de request

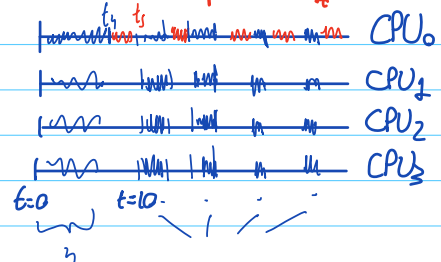


el SO permite que se ejecute cualquier 1 u.t \rightarrow b superamos puede ser mejor

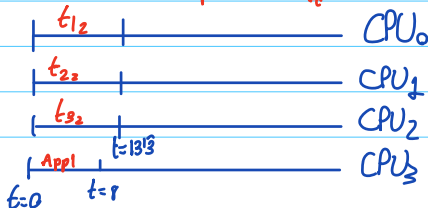
b) Asumiendo no tiempo de request



b) Asumiendo no tiempo de request



c) Asumiendo no tiempo de request



el throughput computing mejor

2. Assume we want to execute two different applications (*app1* and *app2*) in our parallel machine with p processors. Both applications can be (ideally) parallelised defining up to p tasks, each task executing one over p of the total application. The sequential time for *app1* and *app2* is 1200 and 2000 time units, respectively. Assuming that the operating system allows the multiprogrammed execution of parallel applications, for $p = 8$ decide the best allocation of processors to both applications in order to minimise the time the user has to wait for them to finish and giving the programmer the impression that both applications are running from the beginning.

App 1 y App2 $\rightarrow p$ tareas $\rightarrow \frac{2}{p}$ del total
 p procesadores

$A_1: 1200 \text{ u.t}$
 $A_2: 2000 \text{ u.t}$
 $p = 8$

buscar que corran desde el principio \rightarrow $\left\{ \begin{array}{l} A_1: 8 \text{ tareas de } 150 \text{ u.t} \\ A_2: 8 \text{ tareas de } 250 \text{ u.t} \end{array} \right.$

No como requests
 may que han sido?

