

AC 21/22 - Tema 1: Arquitecturas paralelas, clasificación y prestaciones

1. Introducción

Tanto el diseño de nuevos procesadores como el desarrollo de programas eficientes requieren de evaluar las prestaciones de un sistema. En este capítulo vamos a ver formas de obtener esas métricas, que ponen de manifiesto las relaciones existentes entre el tiempo de la CPU y parámetros como el número de instrucciones, capacidad de la microarquitectura o número de instrucciones por ciclo. También estudiaremos los cuellos de botella y el uso de la Ley de Amdal para razonar acerca de éstos.

2. Tiempo de CPU

Tiempo de respuesta y de CPU

Tiempo de respuesta: tiempo transcurrido entre el inicio de la ejecución de un programa y la obtención de los resultados

- Tiempo de CPU: contenido en el tiempo de respuesta, es el tiempo que el procesador pasa ejecutando instrucciones del programa.

El resto del tiempo de respuesta, no incluido en el T_{cpu} es, por ejemplo, tiempo de E/S o tiempo que el procesador dedica a otros programas distintos durante la ejecución.

Es decir: el tiempo de respuesta puede ser un indicador de las prestaciones de un ordenador, pero puede verse afectado por operaciones de E/S o por la carga del procesador en el momento de la ejecución. Nos centraremos en el T_{cpu} , considerando siempre T_e/s despreciable.

$$T_{cpu} = \text{Ciclos del programa} \cdot T_{ciclo} = \text{Ciclos del programa} / \text{frecuencia}$$

Es decir, T_{ciclo} es el periodo y frecuencia es frecuencia y; como sabemos, $f = 1/T$

- Ciclos de programa: número de ciclos de reloj del procesador que tarda el programa en ejecutarse.

$$\text{Ciclos de programa} = \text{Num instr} \cdot \text{Ciclos por instrucción}$$

- Frecuencia: frecuencia del reloj

$$f = 1/T_{ciclo}$$

Por tanto, otra expresión para el T_{cpu} sería:

$$T_{cpu} = CPI \cdot NI \cdot T_{ciclo} = (NI \cdot CPI) / f$$

Especificación

En una máquina en la que todas las instrucciones requieren el mismo número de ciclos, su CPI medio puede calcularse de la siguiente manera:

$$\text{CPI} = (\text{sumatorio from } i:0 \text{ to } W:\text{num instr distintas } [NI \cdot \text{CPI}]) / NI$$

Normalmente las instrucciones del programa se codifican en una sola instrucción máquina. Pero para repertorios que aprovechan el paralelismo de datos esto no es necesariamente cierto, así que podemos calcular su NI de la siguiente manera:

$$NI = \text{Num operaciones del programa} / \text{Num operaciones que puede codificar cada instr}$$

$$NI = \text{Noper} / \text{OPI}$$

El CPI suele depender de las características de la arquitectura y su disposición física. Pero hay ocasiones (procesador segmentado NO superescalar) en las que el compilador es decisivo para el aprovechamiento de los recursos, en cuyo caso nos interesaría expresar el CPI a partir de las características de la microarquitectura.

Podemos usar el num medio de ciclos que tienen que transcurrir desde que se emiten una o varias instr, hasta que se puede realizar otra emisión (llamado CPE); y el número medio de instr que se emiten (IPE), de forma que:

$$\text{CPI} = \text{num medio ciclos entre emisión de una instr hasta sig emisión} / \text{num medio instr emitidas}$$

$$\text{CPI} = \text{CPE} / \text{IPE}$$

- En un procesador no segmentado (las instr se ejecutan de una en una), CPE es igual al número medio de ciclos de reloj que tarda en procesarse la instr; y IPE es 1 porque las instr van de una en una.
- En un procesador segmentado, el valor máximo de CPE es 1 porque cada ciclo podrían empezar a emitirse instrucciones. Si el procesador segmentado sólo puede empezar a ejecutar una instrucción por ciclo, se tiene:

$$\text{CPI} = 1/\text{IPE} = 1/1 = 1$$

Lo cual sería un caso ideal, dado que en los procesadores segmentados normalmente no es posible emitir una instrucción por ciclo (por problemas con dependencias, por ejemplo). Así que, normalmente, $\text{CPE} > 1$ y por tanto $\text{CPI} > 1$. Cuanto más cercano a 1 sea el CPI, mejor estamos aprovechando el cauce del procesador.

En el caso de procesadores superescalares o VLIW, donde se pueden emitir varias instr por ciclo, el CPI puede ser menor que 1.

Conclusión final tras estas aclaraciones:

$$\text{Tcpu} = NI \cdot \text{CPI} \cdot \text{Tciclo} = (\text{Noper}/\text{OPI}) \cdot (\text{CPE}/\text{IPE}) \cdot \text{Tciclo}$$

3. Medidas de velocidad de procesamiento y benchmarks

Referidas a la velocidad de ejecución de instrucciones

MIPS: Millones de instrucciones por segundo

$$\text{MIPS} = \text{NI} / (\text{Tcpu} * 10^6) = \text{NI} / (\text{NI} * \text{CPI} * \text{Tciclo} * 10^6) = 1 / (\text{CPI} * \text{Tciclo} * 10^6) = F / \text{CPI} * 10^6$$

Con la tecnología actual, no es raro sustituir el 10^6 por 10^9 para obtener GIPS, y así no tener que manejar MIPS tan grandes.

Cuando trabajamos con MIPS hay que tener en cuenta que es un cálculo específico al repertorio: una máquina de repertorio RISC necesita muchas más instr para ejecutar el mismo programa que una CISC así que, incluso si el Tcpu de las dos máquinas fuera el mismo para un mismo programa (ambas máquinas son igualmente eficientes), el MIPS de la máquina RISC sería mayor y caeríamos en error si dijéramos que por ello es mejor. Los MIPS miden la velocidad con la que cada procesador ejecuta las instrucciones de su repertorio - por ello muchas veces se refiere a ellos como la velocidad pico del procesador.

MFLOPS: Millones de operaciones float por segundo

$$\text{MFLOPS} = \text{operaciones float} / \text{Tcpu} * 10^6$$

De nuevo, es común sustituir el 10^6 por 10^9 y obtener GFLOPS; por 10^{12} , TFLOPS; por 10^{15} , PFLOPS; por 10^{18} , EFLOPS...

Benchmarks

Tanto cuando evaluamos MIPS como FLOPS, si nos referimos a los valores pico de esas magnitudes hay que emplear un programa o cjto de programas sobre los que hacer medidas de tiempo y NI o NIfloat.

Para tener un marco común de programas de pruebas para evaluar prestaciones, tenemos benchmarks. Estos varían según el fin de la evaluación o el tipo de recurso que se va a evaluar. Algunos ejemplos conocidos son SPEC, TPC, Linpacks.

4. Estimación de tiempo de CPU

Hay veces que es necesario estimar el Tcpu que puede necesitar un programa dadas las características del computador en el que se ejecuta, el cjto de operaciones, el patrón de acceso a memoria del programa...

Las expresiones dadas hasta ahora asumen que los datos a los que se accede están en el nivel 1 de la caché - por eso el CPI suele ser pequeño. Pero si los datos están en memoria o más profundos en la caché el tiempo puede aumentar considerablemente. Podemos obtener una expresión del Tcpu mínimo con la expresión del Tcpu: ese es el tiempo mínimo que va a necesitar el programa, sin fallos de caché y con el procesador a pleno rendimiento.

Pero si tenemos información de las características de la jerarquía de memoria de la máquina, podemos tener una mejor estimación del tiempo mínimo de CPU, porque puede existir un solapamiento casi total entre el tiempo de ejecución de instrucciones en el procesador, y el tiempo de acceso a memoria principal (para load o store, para los que se producen los fallos de caché): es decir, el tiempo mínimo será el mayor de estos dos tiempos:

$$\text{Tejec_min} = \max(\text{Tcpu}, \text{Tmemoria})$$

$$\text{Tmemoria} = \text{Num accesos a memoria} \cdot \text{tiempo medio de acceso a memoria}$$

$$T_{\text{memoria}} = N_{\text{acc}} \cdot t_{\text{mem}}$$

El tiempo de acceso a memoria en una caché look-through (mira toda la caché y si no está ahí lo que quiere, va a memoria) se calcula:

$$t_{\text{mem}} = a_1 \cdot t_1 + (1 - a_1) \cdot (t_1 + t_m)$$

- a_1 - tasa de aciertos de la caché, la probabilidad de que un acceso a memoria esté en caché
 - Para estimar la tasa de aciertos se determina el número de accesos a memoria y en cuántos de ellos falla
- t_1 - tiempo de acceso a caché L1
- t_m - tiempo de acceso a memoria
 - Determinado por la arquitectura

Si además del acceso tenemos en cuenta el tiempo necesario para reemplazar una línea de caché con el contenido que traemos desde memoria, la fórmula queda:

$$t_{\text{mem}} = a_1 \cdot t_1 + (1 - a_1) \cdot (t_1 + t_m + \text{preemplazo} \cdot t_{\text{linea}})$$

- preemplazo - probabilidad de que, cuando se produzca una falta, haya que cambiar una línea de la caché actualizando en memoria la línea reemplazada
 - A partir del patrón de accesos de memoria y las características de la caché se puede determinar también la probabilidad de reemplazo
- t_{linea} - tiempo necesario para actualizar la línea
 - Determinado por la arquitectura

Para mejorar las prestaciones de la jerarquía de memoria se puede:

- Reducir los tiempos de acceso a los distintos niveles de caché
- Reducir tasa de fallos a niveles de caché
- Reducir la penalización cuando se producen fallos
 - Técnicas basadas en inclusión de recursos (cachés de víctimas o pseudo-asociativas)
 - Procedimientos a nivel de programación como precaptación u optimización de código (mezcla de arrays, fusión de bucles, operaciones con submatrices (blocking))

5. Ganancia de velocidad y ley de Amdahl

Para medir el resultado de una mejora se utiliza la ganancia de velocidad, que compara la velocidad de una máquina antes y después de los cambios. La expresión a ganancia de velocidad es:

$$S_p = V_p / V_1 = (W/T_p) / (W/T_1) = T_1 / T_p$$

- V_1 - velocidad previa a la mejora
- V_p - velocidad posterior a la mejora
- W - carga de trabajo que se aplica, es la misma post y pre mejora
- T_1 - T_{cpu} previo a la mejora
- T_p - T_{cpu} previo a la mejora

La ley de Amdahl establece una cota superior a la ganancia que se puede conseguir mejorando alguno de los recursos en un factor p y según la frecuencia con la que se utiliza ese recurso:

LEY DE AMDAHL

$$S_p \leq 1 / (1 + f \cdot (p-1))$$

- S_p - ganancia
- f - fracción del tiempo de ejecución ANTES DE APLICAR LA MEJORA en un recurso QUE NO EMPLEA LA MEJORA
 - O sea, si hay una mejora en el 25% del código, $f=0.75$
 - Así, si $f=1$ es que no se aprovecha la mejora en ningún punto del código y S_p será menor o igual a 1, así que no hay mejora de velocidad
 - Si $f=0$ es que la mejora se aprovecha en todo el código y S_p puede alcanzar un valor igual al factor p