

- $T_{cpu} = (\text{instruções} * cpi) / f(\text{Hz})$
- $T_{cpu} = \text{instruções} * CPI * T$
- $T_{cpu} = N^{\circ} \text{ciclos} * T$
- $T_{cpu} = N^{\circ} \text{ciclos} / f$
- $Ciclos = \text{instruções} * cpi = T_{cpu} * f$
- $IPS = \text{instruções} / T_{cpu}$
- $\text{Instruções} = \text{ciclos} / CPI$
- $Ciclos = \text{instruções} * CPI = T_{cpu} * f$
- $f = (\text{instruções} * cpi) / T_{cpu}$
- **global CPI** = $\sum \% * CPI \text{ class}$
- **Mips p** = $\text{instruções p} / (T_{cpu} p * 10^6)$
- $CPI = T_{cpu} / (\text{instruções} * T)$
- $T_{cpu} = \text{instruções} * CPI * T$
- **Speedup_{x/y}** = $\text{desempenho}_x / \text{desempenho}_y = T_{cpu} y / T_{cpu} x = n$
- **Speedup depois/antes** = $\text{desempenho depois} / \text{desempenho antes} = T \text{ antes} / T \text{ depois} = T \text{ antes} / (\text{Tafetado} / \text{melhoria}) + T \text{ não afetado}$
- **Tempo depois da melhoria** = $(\text{Tempo afetado pela melhoria} / \text{Valor da melhoria}) + \text{Tempo não afetado pela melhoria}$
- KMGTPPEZY 1024 bytes
- **Picossegundos(Ps)** -> $Ps * 10^{-12} \text{ s}$
- **GHz** -> $\text{GHz} * 10^9 \text{ Hz}$
- **ns** -> $ns * 10^{-9} \text{ s}$

Partes usadas pela situação Type-R:

- PC, PC+4, Instructor Memory, Register, ALU, Mux (ALUSrc), Mux (MemToReg)

Em Lw:

- PC, PC+4, Instructor Memory, Register, ImmGen, Mux (ALUSrc), ALU, Mux (MemToReg).

Em beq:

- PC, PC+4, Instructor Memory, Register, ImmGen, ALU, Mux (Branch).
- O desempenho associa-se diretamente ao tempo de execução, precisa-se do CPI
- CPI em Risc-v é 1.
- O CPI depende da implementação do processador e do programa

Métricas que permitem concluir desempenho de um processador -> CPI

- IF -> Instruction fetch -> Memória de instruções
- ID -> Instruction decode -> Banco de registos (leitura)
- EX -> Execute -> ALU
- MEM -> Memory access -> Memória de dados
- WB -> Write back -> Banco de registos (escrita)

Tempo entre 2 instruções pipelined = Tempo entre 2

Não Instruções/Número de andares do pipeline

Instruções não pipelined vs Número de andares no

pipeline: Duração do ciclo de relógio diminui | Tempo para executar uma instrução não diminui | Tende a aumentar, sobretudo se os andares do pipeline não são perfeitamente equilibrados | Todas as instruções demoram o mesmo tempo | Aumenta o número de instruções executadas por unidade de tempo (throughput)

- Num processador pipelined, o período do relógio não pode ser inferir à latência do andar que tem maior latência.

- Num processador não pipelined, o período não pode ser inferior ao máximo dos somatórios das latências das operações que são efetuadas por alguma instrução. **T_{não pipelined}** = \sum
- $CPI = \text{ciclos} / \text{instruções} = ((\text{andares} - 1) + \text{instruções}) / \text{instruções}$
- **Pipeline** - 4 andares
- **T_{lw pipelined}** = $\text{andares} * T_{\text{pipelined}}$
- **T_{lw não pipelined}** = $T_{\text{não pipelined}} \text{ lw} > T_{\text{não pipelined}}$
- $\text{ciclos} = (\text{andares} - 1) + \text{instruções}$
- $IPC = \text{instruções} / \text{ciclos}$
- **speedup** = $\text{Texecução não pipelined} / N^{\circ} \text{ andares}$

Context switch:

- O programa em execução pelo processador muda periodicamente
- A execução é interrompida para se dar a mudança
- A mudança de contexto consiste em guardar o contexto do programa em execução e repor o contexto do programa que vai passar a ser executado.
- A mudança é sempre controlada pelo sistema operativo.

Exceções:

- Assinalam uma circunstancia excepcional, ocorrida durante o processamento, que requer atenção com urgência | Constituem um mecanismo para passar o controlo ao sistema operativo (SO) | Uma exceção pode ter origem interna ou externa ao processador | As exceções externas ao processador são também conhecidas como interrupções (interrupts)
- Uma pipeline é uma forma de ILP. Várias instruções estão em execução em cada ciclo de relógio. Cada instrução está numa diferente da execução.
- **Pipeline scheduling dinâmico** -> Apesar de a ordem de execução das instruções poder não corresponder à ordem das instruções no programa, os efeitos visíveis das instruções tem de corresponder aos da sua execuções puramente sequencial.

-
- **Hit Rate** = $n^{\circ} \text{ de hits} / n^{\circ} \text{ acessos}$

- **Miss Rate** = $n^{\circ} \text{ misses} / n^{\circ} \text{ acessos} = 1 - \text{hit rate}$

- **n^oacessos** = $n^{\circ} \text{ bits} + n^{\circ} \text{ acessos}$

- **Localidade Temporal** é uma posição de memória acedida tenderá a ser acedida outra vez em breve

- **Localidade espacial** as posições de memória perto de uma posição de memória acedida tenderão a ser acedidas em breve

- **Palavra** = endereço/bytes por palavra

- **Bloco** = Palavra/palavra por bloco = endereço/bytes por bloco

- **Índice ou conjunto** = $\text{bloco} \% n^{\circ} \text{ conjuntos}$

- **tag** = $\text{bloco} / n^{\circ} \text{ conjuntos}$

- **Tempo de CPU** = $\text{ciclos execução} \times \text{duração de 1 ciclo}$

Contando com os acessos à memória:

- **Tempo de CPU** = $(\text{ciclos execução} + \text{ciclos memory-stall}) \times \text{duração de 1 ciclo}$

- **ciclos memory-stall** = $\text{ciclos read-stall} + \text{ciclos write-stall}$

- **ciclos read-stall** = $n^{\circ} \text{ reads} \times \text{read miss-rate} \times \text{read miss-penalty}$

- **ciclos write-stall** = $n^{\circ} \text{ writes} \times \text{write miss-rate} \times \text{write miss-penalty} + \text{write-buffer stalls}$

Tempo medio de acesso à memória (Average memory access time):

-AMAT = hit time + miss rate × miss penalty

Com caches com vários níveis:

-Miss rate global= \prod miss rate

Acesso a um disco magnético

-tempo de acesso=seek time+latência rotacional+tempo de transferência + overhead do controlador

- latência rotacional= $0.5 \times (60/\text{velocidade de rotação})$

Desafio do paralelismo

-Tempo depois da paralelização=Tempo antes da paralelização/ Número de processadores

-Speedup=(Tempo antes da paralelização/Tempo depois da paralelização)=Número de processadores

Lei de Amdhal (no contexto da paralelização)

-Tempo depois da paralelização=(Tempo afetado pela paralelização/Número de processadores)+Tempo não afetado pela paralelização

Speedup=Tempo antes/ (Tempo afetado/Numero de processadores)+tempo não afetado)

(Tempo antes da paralelização = Tempo afectado + Tempo não afectado)

Acessos atômicos à memória em RISC-V

-lr.w rd, (rs1) (instrução load-reserved)

-> Funciona como lw rd, 0(rs1)

-> Associa uma marca ao endereço acedido

-sc.w rd, rs2, (rs1) (instrução store-conditional)

Se for executada pela mesma thread, no mesmo processador em que foi executado o lr.w, sobre o mesmo endereço, e a marca ainda lhe está associada

-> Escreve o conteúdo de rs2 no endereço pretendido

-> Remove a marca

-> Põe o valor 0 em rd

Caso contrário

-> Não escreve nada na memória

-> Remove qualquer marca pertencente à thread

-> Põe um valor diferente de 0 em rd

Organização da memória em sistemas SMP

Uniform memory access (UMA)

-> Todos os processadores acedem igualmente a toda a memória

-> Competição no acesso à memória limita número de processadores

Non-uniform memory access (NUMA)

-> Cada processador acede a uma zona da memória com menor latência do que às outras

-> Diminui o tráfego no bus de acesso à memória

-> Permite mais processadores

-> Programas devem ter em conta diferenças nos acessos à memória