

Tipos de miss (3/5)

Compulsory (cold-start): A primeira vez que é acessado um endereço pertencente a um bloco, ele não está em cache. Relacionados com o tamanho dos blocos.

Capacidade: Misses devidos ao bloco ter sido retirado da cache por a cache não ter capacidade para todos os blocos usados pelo programa. Relacionados com o tamanho da cache.

conflito (ou colisão): misses devidos ao bloco ter sido retirado da cache por estar a ocupar

a posição onde deverá passar a estar outro bloco, e que não ocorreriam se a cache fosse fully associativa. Relacionados com a associatividade da cache

bloco está na cache (hit)

bloco não está em cache (miss)

atualizar só a cache: write-back

ler o bloco para a cache, passando a hit: write-allocate

atualizar cache e memória: write-through

não ler o bloco para a cache; é atualizada só a memória: write-allocate

write-through: escritas são imediatamente propagadas para o nível abaixo da memória com ou sem write-allocate. Pode escrever na cache antes de o bloco estar presente.

write-back: escritas só se refletem no nível abaixo da memória quando o bloco é substituído. Usa em geral write-allocate. substituição de um bloco modificado (dirty) só depois de copiado para o nível inferior (ou para um write-back buffer). Pode ser usado um write buffer para guardar as alterações a efetuar.

Acesso à memória e T_{cpu} ignorando o acesso à memória | Hit time = 1 ciclo

contando com os acessos à memória

T_{cpu} = nº ciclos de execução x duração 1 ciclo

T_{cpu} = nº ciclos de execução + nº ciclos memory stall) x duração 1 ciclo

• nº ciclos memory stall = nº ciclos read-stall + nº ciclos write-stall

• nº ciclos read stall = nº reads x read miss-rate x read miss penalty

• nº ciclos write stall = nº writes x write miss rate x write miss penalty + write buffer stalls

latência de acesso = miss penalty

memória virtual: memória organizada em páginas (blocos entre 4Kb e 16Kb)

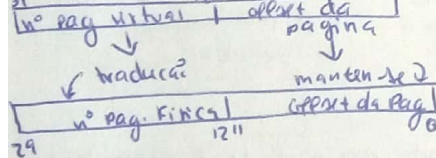
As páginas da memória virtual são mapeadas das páginas da M. física

uma página de M. virtual pode residir em qualquer das páginas da memória física.

MISS = Hit time + miss penalty ciclos

Tradução de endereços: para efetuar um acesso à memória, é necessário traduzir o endereço virtual para o endereço físico correspondente.

End. virtual de 32 bits



endereço físico: 30 bits | Página: 12 bytes = 4KB

implementação da tabela de páginas:

① crescente: uma entrada por cada página virtual

cresce à medida que são acessadas mais páginas

pros: acesso direto | contras: grande parte dos programas acede a pag. no início e no fim do espaço de endereçamento

invertida: uma entrada por cada pag. físico. Acesso através de uma função (hash) aplicada ao nº da pag. virtual

pros: Adaptada ao espaço end. físico | contras: acesso complexo

Multi-nível: partes do nº da pag. virtual são usadas para indexar as tabelas dos vários níveis.

Tabelas do 1º nível: pros: adapta-se à parte utilizada do espaço de endereçamento. | contras: acesso em 2 ou 3 passos

identificam os Pag. físicos. Se existem as correspondentes a ser usadas

Average memory access time

AMAT = hit time + miss rate x miss penalty

caches com vários níveis

cada nível apresenta uma miss rate local

miss rate L₂ = miss rate global: percent.

nº miss L₂ = de acessos que obrigam a

nº acesso à cache = acesso à memória principal

miss rate global = Speedup = tempo L₁

miss rate L₁ = cache 2 níveis: tempo L₂

miss rate L₂ = nº miss L₁

miss rate L₃ = nº acesso à cache L₂

miss rate L₄ = nº miss L₃

miss rate L₅ = nº acesso à cache L₄

miss rate L₆ = nº miss L₅

miss rate L₇ = nº acesso à cache L₆

miss rate L₈ = nº miss L₇

miss rate L₉ = nº acesso à cache L₈

miss rate L₁₀ = nº miss L₉

miss rate L₁₁ = nº acesso à cache L₁₀

miss rate L₁₂ = nº miss L₁₁

miss rate L₁₃ = nº acesso à cache L₁₂

miss rate L₁₄ = nº miss L₁₃

miss rate L₁₅ = nº acesso à cache L₁₄

miss rate L₁₆ = nº miss L₁₅

a memória física comporta-se como uma cache da memória virtual.

TLB: Translation-lookaside buffer é uma cache da tabela de pag que contém traduções de páginas virtuais em físicas.

a tabela só é consultada se a tradução n° pag. virtual → n° pag. física não está presente no TLB.

acesso ao TLB → **TLB hit** o n° da pag física contido no TLB é concatenado com o offset na pag para obter o endereço físico

TLB miss:

é gerada a exceção TLB miss.

Passos de um acesso à memória: 1° Tradução do endereço virtual para físico

2° Acesso à memória

→ Procura na cache:

• hit

• miss: acesso à memória física

acesso a um disco magnetico

velocidade de rotação: velocidade a que o disco gira, medido em rotações por minuto (rpm)

seek time: tempo de edocação, tempo que demora a colocar a cabeça de leitura/escrita na pista a acessar

latência rotacional: tempo que é necessário esperar, em média, até que o setor pretendido esteja sob a cabeça de leitura/escrita, depende da velocidade de rotação e corresponde ao tempo de meia rotação.

taxa de transferência: velocidade a que é possível ler (ou escrever) informação do (ou no) disco, medido em MB/s

controlador: circuito que controla o funcionamento do disco induzindo algum overhead nos acessos.

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

$$\text{latência rotacional} = \frac{1}{2} \times \frac{60}{\text{velocidade de rotação}}$$

exemplo: Qual o tempo necessário para ler um bloco de 512 bytes de um disco com:

seek time médio = 4,0 ms

velocidade de rotação = 15000 rpm

taxa de transferência = 100 MB/s

overhead do controlador = 0,2 ms

Se o período do relógio do processador for de 0,5 ns, quantos ciclos de relógio leva a leitura?

situação ideal:

Tempo depois da paralelização = $\frac{\text{tempo antes da paralelização}}{\text{n° de processadores}}$

desafio do Paralelismo (Lei de Amdahl)

Speed up = $\frac{\text{tempo antes}}{\text{tempo depois}}$

escalabilidade: diz respeito ao modo como evolui o speedup de um programa

escalabilidade forte: evolução do speedup com o n° de processadores

escalabilidade fraca: evolução do speedup quando a dimensão do problema aumenta com o n° de processadores

sistemas de memória partilhada: uniform memory access (UMA) non-uniform ... (NUMA)

GPU: graphics processing unit: processador especializado para criação de elementos gráficos

GPGPU: general programming GPU: multiprocessador que mantém aptidões de GPU's para o processamento de elementos gráficos mas que pode ser usado para progr. genéricos

n-way set associative: conjunto com n posições da cache

fully associative: colocação em qualquer posição da cache

palavra = endereço $\frac{\text{bloco}}{\text{palavra}} = \frac{\text{endereço}}{\text{bytes por bloco}}$

índice / conjunto = $\frac{\text{palavra}}{\text{n° índices do bloco}} \times \frac{\text{n° índices do bloco}}{\text{n° conj.}}$

tag = $\frac{\text{palavra}}{\text{n° índices do bloco}} \times \frac{\text{n° índices do bloco}}{\text{n° conj.}}$

endereço físico = $\text{n° pag física} \times \text{tam. pag} + \text{offset}$

capacidade total cache = $\text{tag} \times (\text{n° conj.} - \text{mdie})$

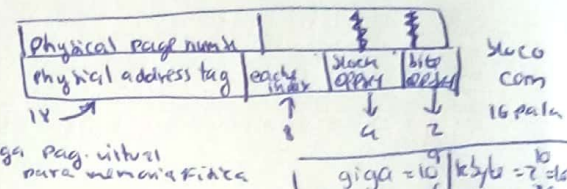
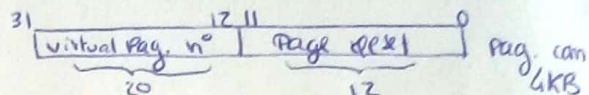
valid tag + 32 x n

• **Relly associative**

• **LRU**

• **write back:** a cada pag. virtual está

associado um dirty bit que indica se o conteúdo da pag foi alterado e se é necessário copiado para memória secundária quando a pag for substituída na memória física.



giga = 10 ⁹	kbyte = 2 ¹⁰
mega = 10 ⁶	megabyte = 2 ²⁰
tera = 10 ¹²	gigabyte = 2 ³⁰
pet = 10 ¹⁵	terabyte = 2 ⁴⁰
exa = 10 ¹⁸	petabyte = 2 ⁵⁰
zeta = 10 ²¹	exabyte = 2 ⁶⁰

temos processamento paralelo de programa paralelo quando o programa utiliza múltiplos processadores em simultâneo.

Speed up = $\frac{\text{tempo antes}}{\text{tempo depois}}$

Speed up = $\frac{\text{tempo antes da paralelização}}{\text{tempo depois da paralelização}} = \text{n° de processadores}$

tempo depois da paralelização = $\frac{\text{Temp antes parali} - \text{Temp n° afetado parali}}{\text{n° de processadores}} + \text{temp n° afetado parali}$

• aumento de CPUs não significa aumento do desempenho.

$$1 \text{ GHz} = 1 \times 10^9 \text{ Hz}$$

6.7
6.0
6.1

quanto menor for o tempo de execução de um programa no computador x, maior será o desempenho deste:

$$\text{desempenho}_x = \frac{1}{\text{Tempo Execução}_x}$$

$$\text{desempenho}_x > \text{desempenho}_y \Rightarrow \text{tempo}_x < \text{tempo}_y$$

- x é mais rápido do que y
- $n > 1$: comp x é n x mais rápido que o comp y

• n é a melhoria de desempenho que se tem usando o comp x, em relação a usar -x o y.
• $n < 1$: slowdown

$$TPI = 250 \text{ ps} = 250 \times 10^{-12} \text{ s}$$

$$P = \frac{1}{TPI} = \frac{1}{250 \times 10^{-12}} = 4 \times 10^9 \text{ Hz} = 4 \text{ GHz}$$

$$t_{cpu} = n^\circ \text{ ciclos} \times T$$

$$t_{cpu} = \frac{n^\circ \text{ ciclos}}{f}$$

$$n^\circ \text{ ciclos} = n^\circ \text{ instruções} \times CPI$$

$$CPI_{global} = \sum_{i=1}^n \frac{1}{f_i} \times CPI_i$$

$$t_{cpu} = n^\circ \text{ instruções} \times CPI \times TPI$$

$$t_{cpu} = \frac{n^\circ \text{ instruções} \times CPI}{f}$$

Instrução
TIPO - R

Lei de Amdahl's: tempo de execução = tempo de execução antes da melhoria + tempo de execução não afetado

Controlo da operação da ALU

maior MIPS não significa maior rapidez

Instruction	opcode	Aluop	Instru. operation	Function Field	Desired ALU action	ALU control input
lw	00		load word	xxxx	add	0010
sw	00		store word	xxxx	add	0010
beq	01		branch equal	xxxx	subtr	0110
R-type	10		add	1000	add	0010
R-type	10		subtract	1000	subtr	0110
R-type	10		and	100100	and	0000
R-type	10		or	100101	or	0001
R-type	10		Set less than	101010	SLT	0111

MIPS 10³ mil milhões

$$MIPS = \frac{n^\circ \text{ de instruções}}{t_{cpu} \times 10^6}$$

$$MIPS = \frac{P}{CPI \times 10^6}$$

opcode	rs	rt	rd	shl
31	26	21	16	0

Instrução
TIPO - R

opcode	rs	rt	imm
(6)	(5)	(5)	(16)

Instrução
TIPO - J

opcode	imm
31 (6) 26 (26)	0

Aluop		Function Field						operation
op1	op2	FS	F4	F3	F2	F1	F0	
0	0	x	x	x	x	x	x	0010
x	1	x	x	x	x	x	x	0110
1	x	x	x	0	0	0	0	0010
1	x	x	x	0	0	1	0	0110
1	x	x	x	0	1	0	0	0000
1	x	x	x	0	1	0	1	0001
1	x	x	x	1	0	1	0	0111

ALU control	Function
0000	And
0001	OR
0010	Add
0110	Substr
0111	SLT
1100	NOR

unidades nativas tipo-R: sign-ext, ALU branch, shift, memória

unidades nativas lw: shift, ALU branch, memória

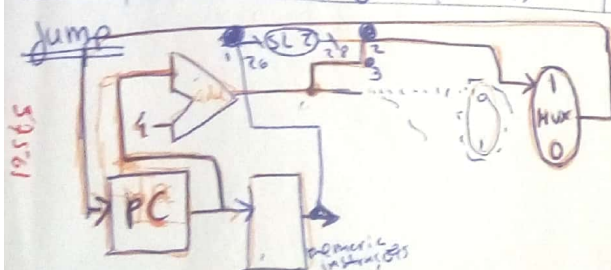
unidade nativa sw: mux (reg), mux (mem), ALU branch, shift

unidade nativa beq: memória, mux (not), ALU branch, shift

unidade nativa jump: nada ativa, excepto o ALU branch, shift

Sinais de controlo

Reg Dst	Reg write	ALUSrc	ALU control	Mem to Reg	Mem Read	Mem Write	Branch	Instruction
1	1	0	Puntia	0	0	0	0	
0	1	1	Puntia	1	1	0	0	R-type
x	0	1	"	x	0	0	0	lw
x	0	0	"	x	0	1	0	sw
							1	beq



1 instrc [25-0] 2 instrc [31-0] 3 PC+4 [31-28]

latência de um circuito: tempo desde que os valores estão presentes nas entradas do circuito até as suas saídas estarem estabilizadas.

Pipeline: IF ID EX MEM WB

1. Todas as inst. tem o mesmo tamanho e podem ser lidas no ciclo no pipeline e decodificadas no 2º.
2. Poucos formatos diferentes de instrução e com os registos sempre nas mesmas posições, pelo que é possível iniciar a malitura em simultâneo com a decodificação.
3. só loads e stores lidam com endereços e é necessário usar a ALU para cálculo de endereços e para outra operação na mesma instrução.
4. valores alinhados em memória permitem que um único acesso seja suficiente para os transferir.
5. Instruções do pipeline produzem um valor que é escrito no último andar.

caminho crítico: caminho que compreende a sequência de operações efetuadas durante a execução da instr. cuja duração é a mais longa. aumento duração inst

6. duração do ciclo de relógio diminui tempo de execução de uma inst. diminui tende a aumentar sobretudo se os andares do pipeline não são perf. equl aumenta o nº inst. executadas por unidade de tempo

7. os registos do pipeline só tem os andares do pipeline e guardam a info necessária sobre a inst. a executar em cada lugar

deca = 10 ¹	octa = 70 ²⁴	quilo byte = 10 ³	giga byte = 10 ⁹
hecto = 10 ²		mega byte = 10 ⁶	tera byte = 10 ¹²
quilo = 10 ³		quilo byte = 2 ¹⁰ = 1024	
mega = 10 ⁶		mega byte = 2 ²⁰	
giga = 10 ⁹		giga byte = 2 ³⁰	
tera = 10 ¹²		tera byte = 2 ⁴⁰	
petta = 10 ¹⁵			
exa = 10 ¹⁸			
zeta = 10 ²¹			

Trata-se do uso de paralelismo na execução das instruções de um programa sequencial num único CPU.

Não se trata da execução de programas paralelos, compostos por várias partes que podem ser executadas em paralelo por múltiplos CPUs.

ILP: paralelismo ao nível da execução das instruções

Formas de ILP:

- Pipelining: várias instruções estão a ser executadas em cada ciclo de relógio. Quanto mais profundo é o pipeline, maior é o ILP. Cada instrução está numa fase diferente da execução.
- Multiple issue: várias instruções estão a começar a ser executadas em cada ciclo de relógio. Instruções recorrem a unidades funcionais duplicadas. O CPI pode tornar-se inferior a 1, recorrendo ao IPC. O número de instruções executadas por ciclo.

Process width do processador é o nº de instr. que podem ser lançadas em simultâneo.

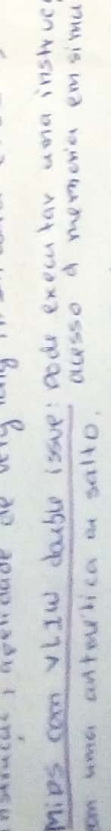
As instruções candidatas a lançamento encontram-se nos issue slots do processador.

As instruções que não efetivamente lançadas substituem issue packet.

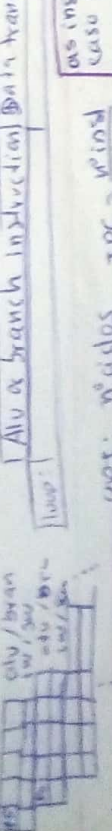
Multiple issue estatico: é o compilador que decide que instr. serão executadas em simultâneo, organizando-as nos issue slots do processador. Entre ao compilador reduzir a garantir que não ocorrem conflitos. As instr. nos issue slots (que vão constituir o issue packet) podem ser encicadas como única grande instrução; a amplitude de very long inst. word (VLIW)

MIPS com VLIW double issue: pode executar uma instrução de acesso à memória em simultâneo em uma antena de salto.

no pipeline fica:



Passando a divisão da seq para o andar 3D, há uma no máximo 1 inst. para cada instrução no pipeline quando é dividido para 4.



CPI: nº ciclos / nº inst. (passando a divisão da seq para o andar 3D, há uma no máximo 1 inst. para cada instrução no pipeline quando é dividido para 4)

Multiple issue dinâmico: processador analisa as inst nos issue slots e decide quantas lançará em simultâneo. O processador trabalha com os compiladores estruturais, de dados e de controle, quantificando a complexidade da execução, um processador super escalas é processador com MI dinâmico. Em alguns casos as instruções podem ser executadas para de ordem no pipeline dynamic pipeline scheduling.

Execução especulativa de instr: tenta-se em decidir que inst que são executadas assumindo que uma inst. anterior terá um desfecho determinado efêto. A execução especulativa permite aumentar o ILP.

Boothrocks: fatores que dificultam o aproveitamento do ILP: dependências longas, acessos à memória, inst que usam muitas operações e saltos irregulares (perda de tempo - erros de especulação)

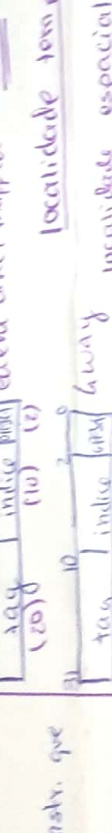
one-way reassociative = direct mapped, codificação numa posição fixa da cache.

B-way set associative -> conjunto de 2 posições de cache

n-way set associative -> conjunto de n posições de cache

fully associative -> colocação em qualquer posição de cache

LRU: least recently used, escolhendo elemento que não é usado há mais tempo.



localidade temporal: posições de memória acessada tendem a ser acessadas novamente em breve.

localidade espacial: posições de memória tendem a ser acessadas em sequência.

cache: nº de conjuntos de cache determina o nº de conjuntos em que o bloco poderá estar.

• nº blocos por conjunto determina em quantas posições um bloco poderá estar.

• dimensão de uma linha de cache, nº de palavras por bloco

• nº de bytes por palavra

• tag identifica o bloco presente em cada posição

previsão do salto: muitas vezes o programa em execução é interrompido para se dar a mudança. Para setar a execução do prog. é necessário saber qual a inst a executar, os valores nos registros (contendo), a mudança de contexto consiste em regular o contexto do prog. em execução e reter o contexto do prog. que o vai substituir. A mudança pode ser especificada em condição de uma ação por preemptiva (em desuso do sistema operacional), assimila uma circunstância operacional, ocorrendo durante o processamento, que requer uma ação urgente. Pode ser implícita (em desuso do sistema operacional) ou explícita (em desuso do sistema operacional).

fixa: o salto não é afetado pelo contexto em execução.

variável: o salto é afetado pelo contexto em execução.

global: o salto é afetado pelo contexto em execução.

local: o salto é afetado pelo contexto em execução.

global: o salto é afetado pelo contexto em execução.

local: o salto é afetado pelo contexto em execução.

global: o salto é afetado pelo contexto em execução.

local: o salto é afetado pelo contexto em execução.

Cache: A info é copiada para a memória do nível superior, que é mais. Essa memória é a cache. Os conteúdos de posições de memória distintas vão ser de ocupação a mesma posição, em momentos diferentes.

Hit: o conteúdo da posição de memória acessada está na cache.

Miss: oposto de hit. miss rate: 1 - hit rate.

Hit time: tempo (CC) que demora o acesso ao conteúdo de uma posição de memória na cache.

Miss penalty: tempo (CC) que demora transferir o conteúdo de uma posição de memória de um nível inferior da memória para o nível acima e se-lo para utilização (podendo substituir o seu conteúdo anterior).

Hit rate: fração das posições de memória acessadas cujo conteúdo foi encontrado na cache.

index: índice de uma posição de memória na cache.

block: unidade de transferência de dados.

repetem-se n posições: índice de uma posição de memória na cache.

1/0 (cache ativa 1): índice de uma posição de memória na cache.

cache: unidade de transferência de dados.

hit: o conteúdo da posição de memória acessada está na cache.

miss: oposto de hit. miss rate: 1 - hit rate.

hit time: tempo (CC) que demora o acesso ao conteúdo de uma posição de memória na cache.

miss penalty: tempo (CC) que demora transferir o conteúdo de uma posição de memória de um nível inferior da memória para o nível acima e se-lo para utilização (podendo substituir o seu conteúdo anterior).

hit rate: fração das posições de memória acessadas cujo conteúdo foi encontrado na cache.

index: índice de uma posição de memória na cache.

block: unidade de transferência de dados.

repetem-se n posições: índice de uma posição de memória na cache.

1/0 (cache ativa 1): índice de uma posição de memória na cache.

cache: unidade de transferência de dados.

hit: o conteúdo da posição de memória acessada está na cache.

miss: oposto de hit. miss rate: 1 - hit rate.

hit time: tempo (CC) que demora o acesso ao conteúdo de uma posição de memória na cache.

miss penalty: tempo (CC) que demora transferir o conteúdo de uma posição de memória de um nível inferior da memória para o nível acima e se-lo para utilização (podendo substituir o seu conteúdo anterior).

hit rate: fração das posições de memória acessadas cujo conteúdo foi encontrado na cache.

index: índice de uma posição de memória na cache.

block: unidade de transferência de dados.

repetem-se n posições: índice de uma posição de memória na cache.

pipeline: • Todas as instruções têm o mesmo tamanho

• Poucos formatos diferentes de instruções e com os registros sempre nas mesmas posições, pelo que é possível iniciar a sua leitura em simultâneo com a decodificação.

• Se loads e stores usam com endereços, não há ALU.

• Valores alinhados em memória permitem que um único acesso seja suficiente para os transferir.

• Instruções só produzem um valor que é escrito no último andar do pipeline.

→ ciclo do relógio diminui, mas o tempo de execução não tende a aumentar se os andares não são perfeitamente equilibrados.

→ todas as inst. levam o mesmo tempo

→ aumenta o n.º de inst. executadas por unidade de tempo (throughput)

• Os registos do pipeline isolam os andares e guardam a informação sobre a inst. a executar em cada andar.

• Os registos dos andares do pipeline também guardam os sinais de controlo da inst. a executar.

• Conflitos estruturais: Existem se não é possível executar uma ou alguma combinação de instruções no mesmo ciclo de relógio.

• Exatidão para o tipo de registos de instruções e de dados.

• Conflitos de dados: Há-os quando uma instrução necessita de um valor produzido por uma instrução anterior, que ainda não está disponível quando a instrução lhe tenta aceder. Por exemplo, se o valor ainda não está no registo no ciclo em que instrução passa pelo andar ID.

• Conflitos de controlo: é necessário saber o destino de um salto condicional para poder executar a próxima instrução.

Sinais de controlo do RISC line

	Reg.DS	ALUop1	ALUop2	Branch	Mem. read	Mem. write	Reg. file	Write back
R-format	1	1	0	0	0	0	1	1
LW	0	0	0	0	1	0	1	1
SW	0	0	0	0	0	1	0	1
BEQ	0	0	1	1	0	0	0	0

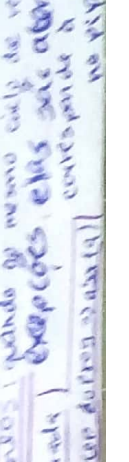
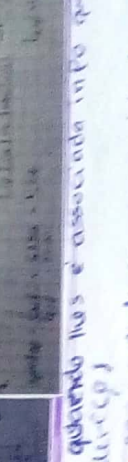
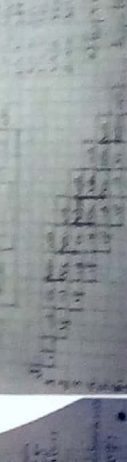
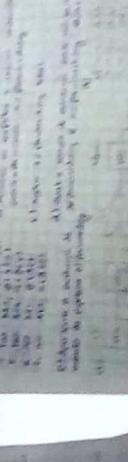
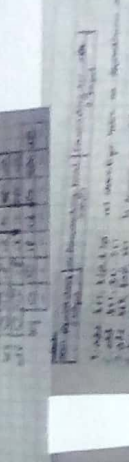
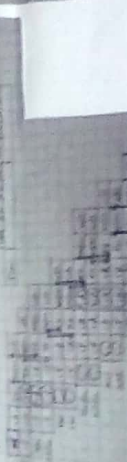
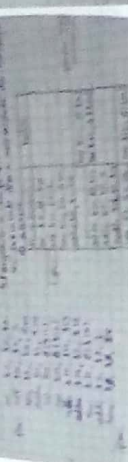
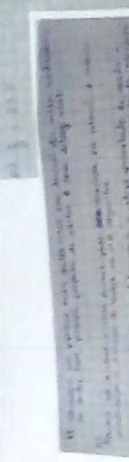
WB → Mem
Mem → EX
WB → EX

Forwarding sem delay slot

• WB → Mem
• Mem → EX
• WB → EX

• WB → Mem
• Mem → EX
• WB → EX

• WB → Mem
• Mem → EX
• WB → EX



ex. opções possíveis, quando lhos é associada info que indica a inst. que lhos dá origem (o seu endereço)

quando o mesmo ciclo de relógio ocorrem múltiplas opções, elas são abeladas, conuando pela que corresponde à instrução mais avançada no pipe line

3 CS: tipos de misses:

compulsory (cold start):

a primeira vez que é acessado um endereço pertencente a um bloco, ele não está na cache. Relacionados com o tamanho dos blocos.

capacidade: misses devidos ao bloco ter sido retirado da cache por a cache não ter capacidade para todos os blocos usados pelo programa. Relacionados com o tamanho da cache.

conflito (colisão):

misses devidos ao bloco ter sido retirado da cache por estar a ocupar a posição onde deveria passar a estar outro bloco e que não é conveniente ter a cache posse fully associative. Relacionados com a associatividade da cache.

memória virtual: a memória organiza-se em páginas, blocos de 4k, 8k, 16k, 32k. as páginas da memória virtual são mapeadas para as páginas da memória física. uma pag. da m.v pode residir em da memória física. A tabela de páginas mantém a correspondência entre pag. virt e pag. físicas. invariável: uma entrada por cada pag. física.

Page Fault: as pag. virtuais que não estão na memória física são guardadas em memória secundária. se a pag. virtual acessada não está em memória física ocorre uma Page Fault. Nesse caso a pag. virtual tem de ser carregada para uma página física.

TLB: é uma cache da tabela de páginas que contém traduções de pag. virtuais em físicas. A tabela de páginas só é consultada se a tradução na TLB.

Acesso a um disco magnético

Tempo de seek time + latência + tempo de transferência

rotacional = $\frac{1}{2} \times \text{velocidade de rotação (rpm)}$

transfer = $\frac{\text{bytes a transferir}}{\text{taxa de transferência}}$

escalabilidade forte: evolução do speedup com o nº de CPUs

escalabilidade fraca: evolução do speedup quando a dimensão do problema aumenta com o nº de CPUs

escalabilidade: modo como evolui o speedup do programa

seja um programa cuja execução num único processador demore 20s, dos quais 1s corresponde à sua parte sequencial

o speedup obtido se o programa for executado em 2 CPUs com o trabalho distribuído igualmente por todos?

se o speedup obtido se um dos 2 CPUs ficar 1/2 atrasado?

Atualizar a cache: write-back

atualizam a cache + memória: write-through

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

se o bloco não está na cache (miss) → ler o bloco da memória e atualizar a cache (miss)

write-through: escritas não imediatamente propagadas para a memória

com a memória atualizada, pode escrever na cache antes do bloco estar presente na memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-through: escritas não imediatamente propagadas para a memória

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

uma write-back substitui o bloco de memória modificado (dirty) só depois de copiado para o nível inferior (ou para um write-back buffer)

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

write-back: escritas só se refletem no nível superior da memória quando o bloco é substituído

Sapendo que o bon CPU com CPI base 1, assumindo que exames na presença de hit na cache mínima sem atrasos e um relógio de $f = 2.0 \text{ GHz}$. Assumindo que a latência de acesso da memória é 100ns. Tem um miss rate por hit 2%.

a) Qual o CPI real, tendo em conta os efeitos de acesso à memória?

b) Quanto mais rápido pode o CPU ser se adicionarmos um cache L2 e é grande o suficiente para reduzir o miss rate global, acesso à memória, para 0.5%?

c) Qual o valor de miss rate da cache L2?

$T = 0.125 \text{ ns}$ Hit = 1 ciclo = 0.125ns

miss global = 100ns

miss = 1 hit + miss global = 0.125 + 100 = 100.125 ns

hit = 0.125 ns CPI miss = 801 ciclos

hit = 100.125

a) $CPI_{global} = \frac{1}{f} \times (CPI_{hit} + \text{miss rate} \times CPI_{miss})$
 $(L1) = \frac{1}{2.0 \times 10^9} \times (1 + 0.02 \times 801) = 9 \text{ ciclos}$

b) miss rate global = 0.5%
 $CPI_{global} = (CPI_{hit} + \text{miss rate} \times CPI_{miss})$
 $9 = (1 + 0.005 \times 801) \times CPI_{hit}$
 $CPI_{hit} = \frac{9}{1.4005} = 6.43$

Speedup = $\frac{CPI_{base}}{CPI_{real}} = \frac{9}{6.43} = 1.4$

c) miss rate L2 = miss rate L1 x miss rate L2

miss rate L2 = $\frac{0.005 \times 801}{0.125} = 0.032$

fully associative: colocar em qualquer posição da cache

way associative: colocar em qualquer posição na cache

37161

EXE 1:

11 \$10, 01 (\$4)

addi \$10, \$10, -2

addi \$17, \$17, 0

sw \$11, 0(\$4)

sw \$17, 0(\$4)

antes da execução das instruções, o endereço guardado no registro \$4 é 10000016, nos 2 processadores e a palavra nesse endereço é 200

Logo os valores nos registros \$10, \$11 do P1 e \$17, \$13 do P2 são os valores de memória em 10000016 de pois de cada instrução.

\$4 = 200

	\$10	\$11	\$12	\$13	mem [10000016]
0	200	200	200	200	200
1	200	200	200	200	200
2	198	200	200	200	200
3	198	200	200	200	200
4	198	200	200	200	200
5	198	200	200	200	200
6	198	200	200	200	200

atomic read/write

read: load aligned: like normal but also saves the address in link register.

write: store conditional: check if address same as in link register. yes? return 1

the idea is that you will load the value stored in memory into a register, modify it. However you have to be careful with the link register. If you modify it, you will lose the address of the next instruction.

if you want to modify the value in memory with your modified one, you need to know the address of the next instruction. If you modify the link register, you will lose it.

if you want to modify the value in memory with your modified one, you need to know the address of the next instruction. If you modify the link register, you will lose it.

if you want to modify the value in memory with your modified one, you need to know the address of the next instruction. If you modify the link register, you will lose it.

if you want to modify the value in memory with your modified one, you need to know the address of the next instruction. If you modify the link register, you will lose it.

if you want to modify the value in memory with your modified one, you need to know the address of the next instruction. If you modify the link register, you will lose it.

if you want to modify the value in memory with your modified one, you need to know the address of the next instruction. If you modify the link register, you will lose it.

if you want to modify the value in memory with your modified one, you need to know the address of the next instruction. If you modify the link register, you will lose it.

Cache: memória de nível superior de pequenas dimensões para onde a informação é copiada.

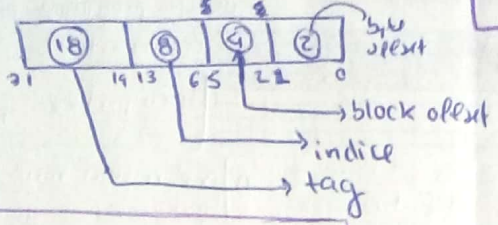
endereço	
palavra	
índice	
tag	
H/M	
bits	

valid	tag	data
0		
1	3	HE3

cache 8 posições

bloco 16 palavras

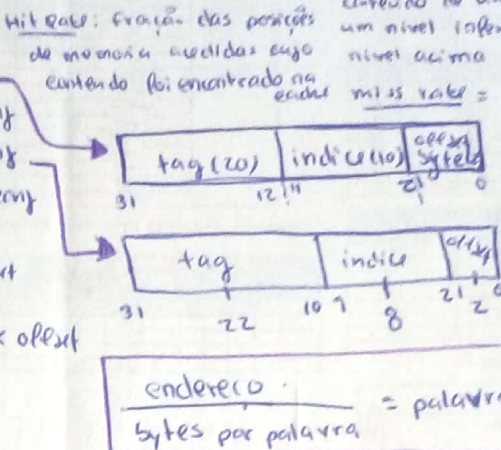
1 way associative \Rightarrow direct mapped
 2 way associative \Rightarrow 2 posições \Rightarrow 4 conj.
 4 way associative \Rightarrow 4 posições \Rightarrow 2 conj.
 8 way associative \Rightarrow 8 posições \Rightarrow 1 conj.



$$\text{bloco} = \frac{\text{palavra}}{\text{palavra por bloco}} = \frac{\text{endereço}}{\text{byte por bloco}}$$

Hit: conteúdo da posição de memória acessada está na cache.
Hit time: tempo (ciclos de relógio) que demora o acesso ao conteúdo de uma posição de memória na cache.

miss: contrário do hit.
miss penalty: tempo que demora transf. o conteúdo de uma posição de memória de um nível inferior da hierarquia para o nível acima e torná-lo disponível para o acesso.
miss rate: $1 - \text{Hit rate}$



LRU: Least recently used
 o escolhido o elemento que não é acessado há mais tempo

$$\text{tag} = \frac{\text{bloco}}{\text{conj.}} = \frac{\text{palavra}}{\text{palavra / índice}}$$

$$\text{índice ou conjunto} = \text{bloco} \% \text{nº de conjuntos} = \text{palavra} \% \text{nº índices}$$

$$\text{off set bits} = \log_2 (\text{nº bits do bloco})$$

$$\text{nº conjuntos / índice} = \frac{\text{nº bits cache}}{\text{nº bit endereço}}$$

$$\text{nº bits índice / conj.} = \log_2 (\text{nº conjuntos})$$

escrita na memória:
 - se o bloco está em cache (Hit) \rightarrow atualizar se a cache (write-back) / atualizar cache + memória (write-through)
 - se o bloco não está em cache (miss) \rightarrow ler o bloco para a cache (passa a Hit) / não ler o bloco para a cache, atualizar se a memória (write-allocate)

write-through: escritas não imediatamente propagadas para o nível abaixo da memória.

write-back: escrito no nível reflete no nível abaixo da memória quando o bloco é substituído (dirty)

cache com 2 níveis

acessos à memória:

$$\text{miss rate } L_1 = \frac{\text{nº misses } L_1}{\text{nº acessos cache } L_1}$$

$$\text{miss rate } L_2 = \frac{\text{nº misses } L_2}{\text{nº acessos cache } L_2}$$

$$\text{miss rate global} = \text{miss rate } L_1 \times \text{miss rate } L_2$$

$$T_{\text{cpu}} = (\text{nº ciclos execução} + \text{nº ciclos memory-stall}) \times T$$

$$T_{\text{cpu}} = \left(\text{nº ciclos read-stall} + \text{nº ciclos write-stall} \right) \times T$$

$$T_{\text{cpu}} = \left(\text{nº reads} \times \text{miss-rate read} \times \text{miss penalty read} + \text{nº writes} \times \text{miss penalty write} \times \text{write-stall} \right) \times T$$

$$A_{\text{mat}} = \text{hit time} + \text{miss rate} \times \text{miss penalty}$$

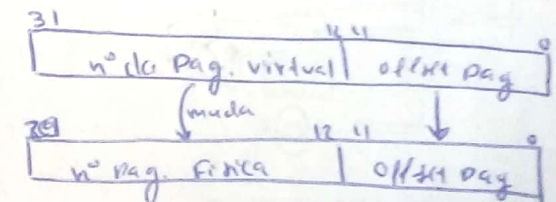
Memória virtual

- os endereços usados por programa referem-se a um espaço de endereçamento virtual.
- memória organizada em páginas (blocos de 4kib até 16kib)
- Pág da memória virtual são mapeadas para pag da memória física
- uma pag. da m.v pode residir em qualquer das pag. m.f

Page Fault: se a pag virtual não está em memória física, ocorre page fault. neste caso, a pag. virtual tem de ser carregada para uma pag. física.

endereços

endereço virtual \Rightarrow 32 bits



endereço físico \Rightarrow 30 bits + 2 offset
 pag 2 bytes = 4kib

$$\text{miss rate} = \frac{\text{nº misses}}{\text{nº acessos}}$$

Tabela Pag			TLB		
P.V	valid	P.F. Addr	Valid	Tag	nº bits

$$\text{tempo acesso} = \text{seek time} + \text{rotational latency} + \text{transfer time} + \text{overhead}$$

endereço	
pag. vir/tag	
TLB (H/M)	
pag. table	
endereço físico	

endereço físico = (pag. física x tamanho pag.) + offset
 tag = endereço // tamanho pag.

ex: tamanho da pag = 4 kib
 $4 \times 2^{10} = 2^7 \times 2^{10} = 2^{17} \text{ offset}$

se TLB for hit, não é necessário ir à tabela pag.

Interactions chosen	Interactions within	Disagrees (ms)	AGD repeats (ms)	Disagrees (ms)	Repeats (ms)	AGD repeats (ms)
Lowest word (100)	200 ms	250 ms	200 ms	200 ms	250 ms	6000 ms
Stems, word (100)	200 ms	250 ms	200 ms	200 ms		5000 ms
Reduced (100, 250, 400, 600, 800)	200 ms	250 ms	200 ms		350 ms	6000 ms
Disagrees (100)	200 ms	250 ms	200 ms			5000 ms

Input/Output	Input/Output	Input/Output	Input/Output	Input/Output	Input/Output
Inputs	Q1	1	1	1	1
	Q2	1	1	1	1
	Q3	1	1	1	1
	Q4	1	1	1	1
	Q5	1	1	1	1
	Q6	1	1	1	1
Outputs	Q7	1	1	1	1
	Q8	1	1	1	1
	Q9	1	1	1	1
	Q10	1	1	1	1
	Q11	1	1	1	1
	Q12	1	1	1	1

quanto menor for o tempo de desempenho
(execução)
de um programa no PCx, melhor CPU e

desempenho $x = \frac{1}{\text{temp. execução}_x}$

Simulation	Overhead address translation stage (cycles/line)					Memory access stage (cycles/line)			Writeback stage (cycles/line)	
	Tag/Line	1/2/3/4/5	1/2/3/4/5	1/2/3/4/5	1/2/3/4/5	Branch	Non-Branch	Non-Branch	Tag Writeback	Value Writeback
Default	1	1	2	3	4	0	0	0	1	0
1st	0	0	2	3	4	0	1	0	1	0
2nd	0	0	2	3	4	0	0	1	0	0
3rd	0	0	3	4	5	0	0	0	0	0

Speedup = $\frac{\text{desempenho}_x}{\text{desempenho}_y} = \frac{t_{\text{exec}}.y}{t_{\text{exec}}.x} = n$ $\begin{matrix} x \text{ é } n \text{ vezes mais rápido do que } y \\ n \geq 1 \quad n \leq 1 \quad \text{slowdown} \end{matrix}$

$$T = 250 \text{ ps} = 250 \times 10^{-12} \text{ s} = 0.25 \times 10^{-9} \text{ s}$$
$$f = \frac{1}{T} = \frac{1}{0.25 \times 10^{-9}} = 4 \times 10^9 \text{ s}^{-1} = 4 \times 10^9 \text{ Hz} = 4 \text{ GHz}$$

WIBS = $\frac{\text{n}^\circ \text{ instructions}}{\text{temp. sec.}} \cdot 10^6$

Adeline:

- todos os instâncias têm o mesmo tamanho;
- poucos formatos diferentes de instâncias/registros ocupam
muitas posições;
- se loads/store indicam o endereço;

1. duração ciclo de relógio diminuir
2. tempo de execução de instruções diminuir
3. tende a aumentar o tempo entre instruções, sobretudo se os endereços não são perfeitamente alinhados.
4. aumenta o nº de instruções executadas por unidade de tempo (Hz).

estruturais: quando não é possível executar uma combinação de instruções no mesmo ciclo.

de dados: quando uma instrução necessita de um valor produzido por uma inst.

do contrato: necessário saber o destino de um salto crível para poder explicar a
 - de \$50. \$10. \$11

--	--	--	--	--

 para inst.

-

norm \rightarrow ex

-

WB \rightarrow Ek

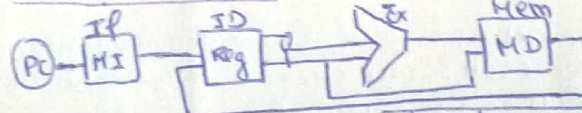
- Forwarding

$$\text{Speed up} = \frac{t_{\text{EA / Parallel}}}{t_{\text{Serial / Parallel}}}$$

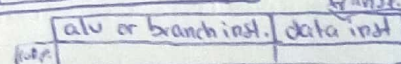
$t_{\text{processor pipeline}}$ = tempo do clock + tempo

$t_{\text{processor } n \text{ pipelined}} = \text{somatonia das } (\%)$

$t_{\text{instrução no pipeline}} = 5 \times \text{tempo do andar com } > \text{tempo}$



desembarcamento de cédigo



exceptions
(transformed to)

- interromper a execução do programa corrente
- executar o código do SO que lida e/ou executa a ocorrência
- retomar a execução do código do programa em abstrato.

$$\text{ciclos acesso} = \frac{t_{\text{acesso}}}{T}$$

$$t_{\text{transf}} = \frac{\text{Bytes a trans}}{\text{trans transf}}$$

$$\text{overhead} = \frac{\text{bytes a frame}}{\text{frame trans} \times 10^6}$$

→ 40050 discs ← $1\text{MB} = 10^6 \text{ Bytes}$