

1 byte = 8 bits

2 GHz = 2×10^9 Hz

- Read after Write (RAW): a instrução li que depende de leitura que foi executada depois dela ultima.
- Write After Read (WAR): a instrução li que depende de li pode ser executada em simultâneo com a ultima li.
- Write After Write (WAW): a instrução li depende de li e é só de resultado de depois dessa ultima li.

- # pragma omp serial / for com que operações únicas thread execute devido ao modo imparidade de threads que fazem o mesmo e esperam a sua vez fazer o trabalho.
- # pragma omp master / com que a thread 0 (thread master) executa bloco enquanto os outros threads continuam executando.
- # pragma omp parallel for paralelo combina variáveis locais calculadas por múltiplos threads em uma única variável global. Permite que cada thread tenha a sua própria variável que seja controlada individualmente especifica (*, max, min, etc.)
- # pragma omp atomic garante que operações simples em memória paralela sejam executadas de forma individual e sem interrupção, permanecendo condicionais. #pragma omp atomic é preciso realmente tornar a execução segura.

$$Fator = CPI \times \#I$$

$$CPI = \frac{\#CPI}{\#I}$$

$$CPI_{processador} = \frac{\#CPI}{\#Op}$$

$$CPI_{processador} = \frac{\sum_{i=0}^{P-1} \#CPI_i}{\sum_{i=0}^{P-1} \#Op_i}$$

$$CPI_{processador} = \frac{\max(\#CPI)}{\sum_{i=0}^{P-1} \#Op_i}$$

Eficiência:

$$Ep = \frac{Sp}{P}$$

P: número de processadores, T: tempo

de execução p: T: tempo de execução com p processadores, Sp: SpeedUp com p processadores

Quando falamos de ciclos:

$\#I = \text{nº de vezes que o ciclo se repete} \times \text{nº de instruções que cada ciclo tem}$

Ciclo não Vetonizado:

- Unmovable loop (for i=0; a[i]=0; i++) for seja executado em paralelo (a variação do i é privada)
- estruturas condicionais grandes (if else)
- invocações de funções dentro do ciclo

[dependência RAW]

Ex: `for (i=1; i < size; i++)`
 $c = a[i-1]*2; a[i] = c > 0 ? c : 1;$
 $d = i - (i-1) = 1$
 $a[1] = a[0]*2 + 1; RAW$
 $a[2] = a[1]*2 + 1$ → vetonizado

Ex: `for (i=0; i < size-1; i++)`
 $c = a[i]*2; a[i] = (c>0 ? c : 1);$
 $d = -(i-(i-1)) = -i+i-1 = -1$
 $a[size-1] = a[size-2]*2 + 1$ WAW
 $a[size-2] = a[size-3]*2 + 1$ Vetonizado

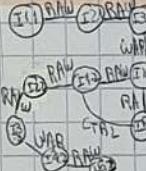
Combina variáveis locais calculadas por múltiplos threads em uma única variável global. Permite que cada thread tenha a sua própria variável que seja controlada individualmente especifica (*, max, min, etc.)

pragma omp atomic garante que operações simples em memória paralela sejam executadas de forma individual e sem interrupção, permanecendo condicionais. #pragma omp atomic é preciso realmente tornar a execução segura.

• schedule (static): divide o loop de forma fixa e paralelo; mas quando todos os iterações terminam o mesmo tempo

• schedule (dynamic): divide o loop sobre demanda, quando os threads têm tempos variáveis para obter um balanceamento de carga dinâmico.

• considera o programa abaixo:



I1: movl 1(%bx,%edi),%esi

I2: addl \$20,%esi

I3: movl %esi,(%bx,%edi)

I4: dec %ecx

I5: jne I1

(CPI=1, autoorden, todas paralelas como torradas)

loop unrolling = 2:

I1: movl 1(%bx,%edi),%esi

I2: addl \$20,%esi

I3: movl %esi,(%bx,%edi)

I4: dec %edi

I5: addl \$20,%edi

I6: movl %edi,-1(%bx,%edi)

I7: subl \$2,%edi

I8: jne I1

• SupE: CPI (reduzido) $\rightarrow \#I$ (igual) \rightarrow Texec (reduzido)

• loop unrolling: CPI (reduzido) $\rightarrow \#I$ (pode aumentar) \rightarrow Texec (reduzido)

• Vetonizado: CPI (pode aumentar) $\rightarrow \#I$ (muito reduzido) \rightarrow Texec (muito reduzido)

• Paralelismo: CPI (igual) $\rightarrow N$ de Itens (igual ou mais) \rightarrow Texec (reduzido)

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

(I1) RAV (I2) RAV (I3) RAV (I4) RAV (I5) RAV (I6) RAV (I7) RAV (I8) RAV

• O ganho de desempenho obtido com a Vetonização de ciclos relativamente à respectiva versão escalonada deve-se à diminuição do número total de instruções executadas (#I).

• As unidades de processamento gráfico minimizam o impacto das necessidades matriciais no desempenho do programa omitindo rapidamente entre grupos de threads dependentes o tempo de acesso de alguns threads com execução de outras.

• Ciclo "for (i=1; i < size, i+=1) a[i] = b[i-1]/(i*10/123)" não vetoniza devido a uma dependência RAW entre iterações.

• Considera o código: #define size... float a[size]; float aux[size]; main() { int i; for (i=1; i < size; i++)

aux[i] = powf(a[i], 2) + 10.0/a[i] - a[i-1]; }

• Ciclo não vetonizado por ter uma

dependência RAW (a[i-1] == a[i-1])

utilizando funções dentro do loop (powf())

• Considera o código com vetonização por processamento paralelo float a[size], aux[size];

main() { int i; for (i=1; i < size; i++)

aux[i] = aux[i] - a[i-1]; }

• Ciclo aberto explora Thread Level Parallelism (TLPI)

#define W 400000 int a[W], int sum, i, j;

pragma omp parallel for

for (i=0; i < W; i++) sum=0;

for (j=1; j < i; j++) sum+= a[j];

a[i][0] = sum; }

• Erros 2 condições: Este programa produz comportamento

indefinido porque existem raias condicionais causadas pelo acaso e alteradas

imprevisivel a dados partilhados pelos vários threads. E provavelmente facticas com os elementos

de linhas de uma matriz nem sempre independentes em cada linha para paralelizar esta tarefa. No entanto, existe informação estruturalmente prévia da estrutura dos

threads: variáveis globais sum e j. No conteúdo do programa só faz sentido tal abordagem. Existem duas soluções:

- Declarar variáveis dentro de ciclo exterior (int sum, j);

- Em vez de declarar global: int j; antes de ciclo interior ou for(j=1; ...)

openMP (# pragma omp parallel for private (sum, j))

Faz mais sentido utilizar estrutura dinâmica porque não é necessário que a matriz é dividida e

é melhor dividir a

o loop por demanda.

CPI = $\frac{5}{5} = 1$

CPI = $\frac{3}{5} = 0.6$

$\Delta mpc = mpc$	$T_{exec} = \frac{CPI \times \#I}{f} = \frac{\#CC}{f}$	$Tag \quad \text{ret} \quad offset$	$T_{cc} = \max(T_{tag}, T_{offset}) + T_{tag+data}$	$m - \# \text{ de linhas}$	\bullet Mapamento direto: onde $adrs = m \cdot \# \text{ de linhas}$
ganhos da ΔT_{exec}	$\bullet CPI = \frac{\#CC}{\#I}$ $T_{cc} = \frac{1}{3}$	$\bullet Tag = m - b - n$	$\bullet T_{cc} = \max(T_{tag}, T_{offset}) + T_{tag+data}$	\bullet Torna leitura por set ($E=1$); padroniza para	
	$\bullet CPI = CPI_{CPU} + CPI_{MEM}$	$\bullet \text{BlockOffset} = \text{tag}_B \oplus (\text{dado} \text{ dentro de bloco})$	$\bullet T_{tag+data} = \min(\# \text{ de linhas}, \# \text{ de blocos}) \cdot \text{latência da memória}$	\bullet latência = tempo de leitura de linhas, latência de bloco + latência de bloco	
	$\bullet CPI_{MEM} = (\%Mem + \%Mem \times \%Mem) \times mpc$	$\bullet m = \#I \text{ (número de instruções)} / \text{número total de instruções}$	$\bullet S = 4 \quad E = 1 \quad B = 2 \quad m = 5 \rightarrow \text{muito lento}$	\bullet Mapamento completamente associativo: leitura única de bloco ($B=1$), reduz latência de leitura, leitura menor para	
	$\bullet CPI_{mem-I} = mpc \times mpc$		\bullet necessita comparar com todos os tags	\bullet Mapamento m-way set associativo: leitura com E linhas cada; comparação	
	$\bullet CPI_{mem-D} = \%Mem \times mpc$			\bullet m-way de leitura de blocos, leitura menor para	

• Aumento ($\Delta mpc = mpc$) implica frequentemente um aumento do CPI. Tcc diminui, mas o tempo de leitura à memória mantém-se com necessários mais ciclos para a leitura da memória ($T_{cc} \text{ diminui} \Rightarrow \text{CPI aumenta}$)

• Diminuição de $\#I$ (removendo instruções mais complexas (inserção de mais ciclos)) resulta num aumentado CPI.

• Localidade temporal: um elemento de memória medida pelo CPU é a localidade de uso num futuro passado
• Localidade espacial: um elemento é acessado pelo CPU, os elementos com endereços próximos tendem a ser usados num futuro próximo

$[j] = 0; i \leftarrow i + 1; a[i] = b[i * 10^3] + 2, \dots$
→ existe localidade espacial nos acessos $a[i]$
→ existe localidade temporal nos acessos a i
• A técnica de pipeline é altamente suscetível a uma seqüência de uso único, ou seja, o desempenho de um processador pode resultar numa diminuição do desempenho, tanto vez que deve aprender o uso longo quanto a instância demorada de pipeline

Pipeline: resumo

- execução de instruções simultaneamente em diferentes estágios
- permite aumentar a frequência de relações
- Dependência de dados: stallings (injecção de bolhas (Nops); realimentação); elementos penalizações
- Dependência de controle: saltos condicionais implicam execução especulativa (previsão); envolvem muita complexidade de pipeline

discretas	binárias	1 bit por posição
10^3 Kbytes/B	2^{10} Kbytes/KB	10^{-9}
10^6 Megabytes/MB	2^{20} Megabytes/MB	10^{-6}
10^9 Gigabytes/GB	2^{30} Gigabytes/GB	10^{-3}
10^{12} Petabytes/PB	2^{40} Petabytes/PB	
10^{15} Exabytes/EB	2^{50} Exabytes/EB	
$(10^3 \text{ Bytes/B})^2$	2^{100} Zettabytes/ZB	$3 \cdot 10^{-12}$

→ Cache

• re-total de linhas: $S \cdot E$
• Capacidade de cache: $S \cdot E \cdot B$

Fetch	ponto de Mem	Instrução	Int. mem	Int. mem	Add L	mapamento direto	rotacionando associativo	m-way set associativo
					$S=2 \quad E=1 \quad B=2 \quad m=5$	$S=1 \quad E=4 \quad B=2 \quad m=5$	$S=2 \quad E=2 \quad B=2 \quad m=5$	
1: mov \$4,%eax	1	mov \$4,%eax	mov \$4,%eax	x	5(00101)	ret 2/linha 0 miss	ret 0/linha 0 miss	ret 0/linha 0 miss
2: mov \$2,%eax	2	mov \$2,%eax	mov \$2,%eax	x	14(01110)	ret 3/linha 0 miss	ret 0/linha 1 miss	ret 0/linha 0 miss
3: mov(%eax),%eax	3	mov(%eax),%eax	mov(%eax),%eax	x	10(01010)	ret 1/linha 0 miss	ret 0/linha 2 miss	ret 1/linha 1 miss
4: add %eax,%eax	4	add %eax,%eax	add %eax,%eax	x	29(11101)	ret 2/linha 0 miss	ret 0/linha 3 miss	ret 0/linha 1 miss
5: dec %eax	5	dec %eax	dec %eax	x	4(00101)	ret 2/linha 0 miss	ret 0/linha 0 hit //	ret 0/linha 0 hit
6: jmp loop	6	jmp loop	jmp loop	x	21(00101)	ret 2/linha 0 miss	ret 0/linha 0 miss	ret 0/linha 0 miss
I1: mov \$10,%eax	7	mov \$10,%eax	mov \$10,%eax	x	16(10000)	ret 0/linha 0 miss	ret 0/linha 2 miss	ret 0/linha 0 miss
I2: sub \$5,%eax	8	sub \$5,%eax	add %eax,%eax	x	5(00101)	ret 2/linha 0 miss	ret 0/linha 0 hit //	ret 0/linha 0 miss
I3: jne I1	9	jne	x	x	Miss rate	100%	75%	37,5
I4: add \$10,%eax	10	add \$10,%eax	jmp	x				

• Dependência de controle: jumps; prové-se que o salto é sempre dado; A consequência é a perda de吞吐量, feita posteriormente; Se a estrutura estiver errada, existem instruções entrando linha

→ processador com formulário + data forwarding não consegue lidar com saltos, tem que resultar num desvio de ciclo igual ao nº de bloco

1	2	3	4	5	6	7	8	9	10	11
I1 F D E W	I1 F D E W	I2 F D D D E W	I2 F D E W	I3 F F F D E W	I3 F D E W	I4 B ¹ B ² B ³ B ⁴ B ⁵	I4 F D E W	I5 F D E W	I5 F D E W	I6 F D E W
I7 F D D D E W	I7 F D E W	I8 F F F D E W	I8 F D E W	I9 B ¹ B ² B ³ B ⁴ B ⁵	I9 F D E W	I10 B ¹ B ² B ³ B ⁴ B ⁵	I10 F D E W	I11 F D E W	I11 F D E W	I12 F D E W
I12 B ¹ B ² B ³ B ⁴ B ⁵	I12 B ¹ B ² B ³ B ⁴ B ⁵	I13 B ² B ³ B ⁴ B ⁵	I13 B ² B ³ B ⁴ B ⁵	I14 B ³ B ⁴ B ⁵	I14 B ³ B ⁴ B ⁵	I15 B ⁴ B ⁵	I15 B ⁴ B ⁵	I16 B ⁵	I16 B ⁵	I17 B ⁵
I17 F D E W	I17 F D E W	I18 F D E W	I18 F D E W	I19 F D E W	I19 F D E W	I20 F D E W	I20 F D E W	I21 F D E W	I21 F D E W	I22 F D E W
I23 F D E W	I23 F D E W	I24 F D E W	I24 F D E W	I25 F D E W	I25 F D E W	I26 F D E W	I26 F D E W	I27 F D E W	I27 F D E W	I28 F D E W

• Data-forwarding: Para o valor de necessário (ma dependência de dados) da estrutura ordenar estavam erradas as instruções entrando linha

→ processador com formulário + data forwarding não consegue lidar com saltos, tem que resultar num desvio de ciclo igual ao nº de bloco

• Se necessitamos 9 ciclos. A dependência de dados entre I1 e I2 é resolvida por data forwarding, resultado é mal previsto resultando numa penalização de 2 ciclos