Arquitetura II

Cl086 – Tópicos em Arquitetura de Computadores Cl702 – Arquitetura de Computadores

Roberto Hexsel roberto@inf.ufpr.br www.inf.ufpr.br/roberto/CI086.html

LIEPR Dinf BCC

Arquitetura II — introdução

2007-1

Por que estudar Arquitetura?

HEPR Dinf RCC

Arquitetura II — introdução

2007-1

Por que estudar Arquitetura?

- ser um programador competente: usar caches e memória virtual de forma vantajosa;
- aprender algoritmos e técnicas para alta velocidade: tirar proveito da hierarquia de memória;
- usar bem recursos da máquina: idem com relação a E/S, segmentação;
- poder trabalhar com eletrônica embarcada:
 mercado novo e com poucos programadores competentes;
- entender artigos das revistas da ACM, IEEE;
- e, principalmente, ser feliz a vida é de vocês...

Bibliografia

texto:

Computer Architecture: A Quantitative Approach, H&P-QA J L Hennessy e D A Patterson, 3a Ed, Morgan Kaufmann, 2003. cuja tradução ruim é

Arquitetura de Computadores – uma abordagem quantitativa, J L Hennessy e D A Patterson, 1a Ed, Campus, 2003,

textos auxiliares:

Computer Organization and Design, P&H-COD D A Patterson e J L Hennessy, 3a Ed, Morgan Kaufmann, 2005 existe tradução da segunda edição

HEPR Dinf RCC

Arquitetura II — introdução 2007

Programa

- ★ Tendências tecnológicas (1) texto em www.*://roberto/CI702.html
- ★ avaliação de desempenho (1)
- ★ conjuntos de instruções MIPS64 (2) TODOS DEVEM ESTAR COM O LIVRO!
- ★ revisão pipelines, processadores super-escalares (3)
- ★ paralelismo no nível de instrução, hw e sw (4)
- * hierarquia de memória, caches, mem virtual (4)
- ★ 03mai prova
- * multiprocessadores, mem logicamente compartilhada (4)
- ★ sistemas de E/S, discos (2)
- ★ multicomputadores e clusters, mem logicamente distribuída (2)
- * redes diretas (2)
- ★ 21jun prova e 04jul entrega do trabalho
- ★ 03jul final

HEPR DInf RCC

Arquitetura II — introdução

Material adicional para o curso

2007-1

http://www.inf.ufpr.br/roberto/CI702.html http://www.inf.ufpr.br/roberto/CI086.html /public/soft/linux/simplescalar ~roberto/ci702

lista de e-mail em roberto-tac@inf.ufpr.br

não guardem este endereço em máquinas com uindous

Organização de Computadores

Modelo de Von Newman

- computador com programa armazenado (1945)
- memória é um vetor de bits
- * interpretação dos bits definida pelo arquiteto e programador/compilador
- parte da memória contém instruções
- parte da memória contém dados

HEDR DIM RCC

Arquitetura II — introdução 2007-

Definição de Arquitetura de Computadores

 arquitetura do conjunto de instruções (CdI): conjunto de instruções e registradores visíveis ao programador Instruction Set Arquitecture = ISA

2. organização:

blocos como sistema de memória, barramentos, CPU mais de uma implementação de mesmo conjunto de instruções (AMD e Intel, 80{,1,2,3,4,5,6}86)

3. hardware:

tecnologia de implementação, circuitos integrados (CMOS vs NMOS), pipelining vs multi-ciclo

Arquitetura engloba todos os três aspectos

HEDR DIAFRCC

Requisitos a serem atendidos pelo arquiteto I

ÁREA DE APLICAÇÃO:

Arquitetura II — introdução

usos do computador

desktop desempenho balanceado para diversas tarefas; sistemas interativos com gráficos, vídeo e áudio;

servidor de computação científica alto desempenho com operações de ponto flutuante e gráficos;

servidor comercial alto desempenho com bancos de dados e sistemas transacionais; confiabilidade, disponibilidade e escalabilidade;

embutido suporte a aplicações específicas;

baixo custo e pequeno consumo de energia

IIFPR Dinf RCC 9

Requisitos a serem atendidos pelo arquiteto II

COMPATIBILIDADE DE SOFTWARE:

software pré-existente

nível de linguagem de programação maior flexibilidade; novo compilador?

compatibilidade nos binários ISA completamente definido; pouca flexibilidade mas sem investimentos em software novo.

HEPR Dinf RCC 10

Arquitetura II — introdução

2007-1

Requisitos a serem atendidos pelo arquiteto III

SUPORTE A SISTEMA OPERACIONAL:

tamanho do espaço de endereçamento importantíssimo!! pode limitar aplicabilidade;

ightarrow EdE cresce $\geq 1/2$ bit aa

gerenciamento de memória necessário para SOs modernos; paginado ou segmentado;

proteção usos diferentes por SO e usuários; proteção à páginas ou a segmentos.

HEPR Dinf RCC 1

Arquitetura II — introdução 2007-

Requisitos a serem atendidos pelo arquiteto IV

Padrões: requeridos pelo mercado

ponto flutuante formatos e aritmética: IEEE754;
 barramentos dispositivos de E/S: PCI, SCSI;
 sistema operacional Unix, Windows, PalmOS;
 redes suporte a tecnologias distintas (Ethernet, Infiniband);
 linguagens de programação linguagem de alto nível afeta projeto do conjunto de instruções.

IIFPR Dinf RCC 12

Tecnologias Fundamentais

- 1. tecnologia de semicondutores: microprocessadores CMOS
- 2. tecnologia de semicondutores: memória dinâmica (DRAM)
- 3. tecnologia de armazenamento: discos magnéticos
- 4. tecnologia de interconexão: redes locais

HEPR Disf RCC 11

Arquitetura II — introdução 2007-1

Tendências da Tecnologia I

Tecnologias de implementação fundamentais 1- Circuitos Integrados para CPUs

núm de transistores cresce $\approx 35\%$ aa 4X em 4 anos

tamanho do Cl aprox 10 a 20% aa lado do retângulo

efeito combinado núm transistores/Cl cresce ≈55% aa

Ver fabricação de circuitos integrados em www.intel.com→ museum e em P&H-COD

HEPR Dist RCC 14

Arquitetura II — introdução 2007-1

Fotografia de um CI do Pentium IV

Wafer de 8" com 564 MIPS R20K

HEPR DInf RCC 16

Arquitetura II — desempenho Revisão

Arquitetura é:

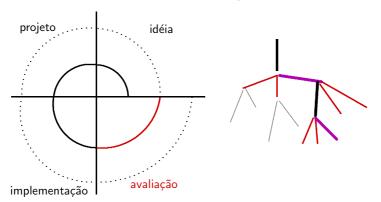
```
conjunto de instruções + organização + hardware especif abstrata + especif refinada + concreta
```

Tendências da Tecnologia			
	capacidade	velocidade	
lógica	2x em 3 anos	2x em 3 anos	
DRAM	4x em 3 anos	2x em 10 anos	
disco	4x em 3 anos	2x em 10 anos	
processador	_	2x em 1.5 anos	

HEPRING RCC 11

Arquitetura II — desempenho 2007-1

Processo de Projetar



Processo iterativo: idéia \Longrightarrow implementação \Longrightarrow avaliação \Longrightarrow re-projeto ...

avaliação geralmente implica em descartar idéias ruins

LIFPR Dinf RCC

Ferramentas para Avaliação de Desempenho

- sistema: programas de teste, rastros (traces)
- hardware: custo, atrasos, área, potência
- simulação nos níveis: Cdl, register transfer, portas, circuitos
- Teoria de Filas
- Regras de projeto
- princípios fundamentais
- Leis (de Moore, de Amdahl, et alli)

HEPR Dist RCC 10

Arquitetura II — desempenho

2007-1

Avaliação de Desempenho

QUAL É O MELHOR AVIÃO?

avião	capacidade	alcance	veloc	produção
	[p]	[km]	[km/h]	[p*km/h]
B-737	101	1008	957	96657
B-747	470	6640	976	458720
Concorde	132	6400	2160	285120
DC-8	146	13952	870	127020

O que importa são DESEMPENHO e CUSTO

LIEDE DIN ERC 20

Arquitetura II — desempenho 2007-1

Métricas para Avaliação de Desempenho

aplicação respostas por hora; operações por segundo

linguagem produtividade do programador (?)

compilador qualidade do código gerado

ISA instruções por segundo (MIPS ou MFLOPS)

circ de dados megaBytes por segundo

unid. funcionais ciclos/operações por segundo

portas, fios ciclos por segundo

Desempenho (i)

$$\mathsf{desempenho}_X \stackrel{\triangle}{=} \frac{1}{\mathsf{tempo}\,\,\mathsf{exec}_X}$$

 $X \notin \boldsymbol{\mathcal{M}} \text{ vezes mais rápido que } Y \text{ se} \\ \text{o tempo de execução de } Y \notin \boldsymbol{\mathcal{M}} \text{ vezes mais longo que o de } X$

$$\frac{\mathsf{desempenho}_{X}}{\mathsf{desempenho}_{Y}} = \frac{\mathsf{tempo}\ \mathsf{exec}_{Y}}{\mathsf{tempo}\ \mathsf{exec}_{X}} = \mathcal{M}$$

Desempenho tem unidade de "coisas" por segundo; maior desempenho é melhor!

HEPR Dinf RCC 22

Arquitetura II — desempenho 2007-

Desempenho (ii)

$$\frac{\mathsf{desempenho}_X}{\mathsf{desempenho}_Y} = \frac{\mathsf{tempo}\ \mathsf{exec}_Y}{\mathsf{tempo}\ \mathsf{exec}_X} = \mathcal{M}$$

Exemplo:

Máquina A executa programa em 10 segundos; máquina B executa mesmo programa em 15 segundos. Quanto A é melhor que B ?

$$\frac{\mathsf{desempenho}_A}{\mathsf{desempenho}_B} = \frac{\mathsf{tempo}\ \mathsf{exec}_B}{\mathsf{tempo}\ \mathsf{exec}_A} = \frac{15}{10} \Rightarrow 50\%$$

HEPR Dist RCC 23

Arquitetura II — desempenho

2007-1

Avaliação de Desempenho

- tempo de resposta, tempo decorrido [s/tarefa]
- tempo de CPU (usuário + sistema) [s]
- ciclo de relógio [ns] frequência do relógio [GHz]
- vazão/produção (throughput) [tarefas/s]

desempenho do sistema = tempo decorrido (sem carga) desempenho de CPU = tempo de CPU dedicado ao usuário

desempenho de pico NUNCA é atingido na prática desempenho de pico \gg desempenho sustentado

Equação do Desempenho (i)

tempo de CPU
$$=$$
 ciclos da CPU \times ciclo de relógio
$$= \frac{\text{ciclos da CPU}}{\text{freq de relógio}}$$

ciclos da CPU = núm de instr
$$\times$$
 núm de ciclos por instr
$$CPI \stackrel{\triangle}{=} ciclos \ por \ instrução$$

HEPR Dinf RCC 2

Arquitetura II — desempenho

Equação do Desempenho (ii)

tempo de CPU = núm de instr
$$\times$$
 CPI \times ciclo de relógio = $\frac{\text{núm de instr} \times \text{CPI}}{\text{freq de relógio}}$

esta é a equação mais importante do semestre!

HEPR Dinf RCC 2

Arquitetura II — desempenho 2007-

Equação do Desempenho (iii)

Exemplo:

Programa executa na máquina A em 10s com relógio de 100MHz. Queremos máquina B que execute mesmo programa em 6s. Por causa da mudança no projeto da CPU, máquina B vai usar 1.2 vezes mais ciclos de relógio que A. Qual o relógio de B?

tempo de
$$\mathsf{CPU}_A = \frac{\mathsf{ciclos}\;\mathsf{da}\;\mathsf{CPU}_A}{\mathsf{freq}\;\mathsf{de}\;\mathsf{relógio}_A}$$

$$10s = \frac{\mathsf{ciclos}\;\mathsf{da}\;\mathsf{CPU}_A}{100\times10^6}$$
 $\Rightarrow \;\#\mathsf{ciclos}\;\mathsf{da}\;\mathsf{CPU}_A = 1000\times10^6$

$$\begin{array}{ll} \text{tempo de CPU}_B &=& \frac{1.2*\text{ciclos da CPU}_A}{\text{freq de relógio}_B} \\ &\Rightarrow& \text{freq de relógio}_B = \frac{1.2*1000\times10^6}{6s} \\ &\Rightarrow& \text{freq de relógio}_B = 200\times10^6 \\ \\ \text{ganho de } \frac{10}{6} = 0.67 &\Rightarrow& \text{freq de relógio } 100\% \text{ maior} \end{array}$$

Exemplo: Programa executa na máquina A em 10s com relógio de 100MHz. Queremos máquina B que execute mesmo programa em 6s. Por causa da mudança no projeto da CPU, máquina B vai usar 1.2 vezes mais ciclos de relógio que A. Qual o relógio de B?

HEPR Dist RCC 23

Arquitetura II — desempenho 2007-

Equação do Desempenho (iv)

tempo de CPU
$$= \mathcal{N} \times \text{CPI} \times \mathcal{R}$$
 núm instr $\times \text{CPI} \times \text{período do relógio}$

$$\mathsf{CPI} \ = \ \sum_{j}^{n} \mathsf{CPI}_{j} imes \mathsf{F}_{j}$$
 onde $\mathsf{F}_{j} = \mathcal{I}_{j}/\mathcal{N}$ F_{j} é a freqüência da instrução \mathcal{I}

HEPR Dist RCC 20

Equação do Desempenho (v)

$$\mathsf{CPI} = \sum_{j=0}^{n} \mathsf{CPI}_{j} \times \mathsf{F}_{j}$$

instr	freq[%]	ciclos	CPI_j
ALU	40	1	.40
load	25	3	.75
store	10	3	.30
desvios	25	2	.50
		total	1.95

Fatores Determinantes

tempo de CPU = núm de instr \times CPI \times ciclo de relógio

tempo de CPU	desempenho do sistema
núm de instruções	compilador & processador
CPI	organização e arquitetura da CPU
freqüência de relógio	tecnologia de CIs & arquitetura

HEPR Dist RCC

Arquitetura II — desempenho

Administrativamente falando...

Todos já estão com o livro?

Mestrandos já estão com os artigos?

Ler:

How not to lie with statistics: the correct way to summarize benchmark results, Fleming & Wallace, Comm ACM 29(3), Mar 1986

Characterizing computer performance with a single number, J E Smith, Comm ACM 31(10), Oct 1988

HEPR Dinf RCC 2'

Arquitetura II — desempenho 2007-1

Medidas de desempenho

MIPS = milhões de instruções por segundo

MIPS =
$$\frac{\text{num de instr}}{\text{tempo decorrido} \times 10^6}$$
$$= \frac{\text{freq de relógio}}{\text{CPI} \times 10^6}$$

Problemas

- independente do conjunto de instruções (RISC/CISC)
- varia para programas na mesma máquina (int×PF)
- pode variar na proporção inversa ao desempenho

Medidas de desempenho

MFLOPS = milhões de instr de ponto flutuante p/s

MFLOPS =
$$\frac{\text{núm de instr de pto flutuante}}{\text{tempo decorrido} \times 10^6}$$

Problemas

- independente do conjunto de instruções (Cray/68882)
- varia para programas na mesma máquina (soma×div)
- média ponderada de custo de instruções: $soma(a,b) \propto 1$ $seno(x) \propto 8$

HEPR Dinf RCC

Arquitetura II — desempenho

2007-1

Medir desempenho com programas de teste

programas simples: quicksort, números primos simples de implementar, fora da realidade

programas sintéticos: Dhrystone, Whetstone simples de implementar, não são código usável

núcleos de programas: SPEC, Livermore Loops

fáceis de medir, não testam sistema de forma realista

programas de verdade: gcc, LaTeX, Spice

mistura deve refletir uso "normal" (browser?)

HEPR Dinf RCC 34

Arquitetura II — desempenho

2007-1

Medindo desempenho - SPEC

- primeiro conjunto em 1989
- ★ 10 programas produzem um só número: SPECmarks
- segundo conjunto em 1992
- * SPECint92 com 6 programas com inteiros
- \star SPECfp92 com 14 programas com ponto flutuante
- terceiro conjunto em 1995
- ★ SPECint95 com 8 programas com inteiros
- ★ SPECfp95 com 10 programas com ponto flutuante
- * conjunto caduca em três anos
- * versão base com mesmas flags de compilação (todos programs)
- quarto conjunto em 2000
- ★ CINT2000 com 11 programas com inteiros (C e C++)
- \star CFP2000 com 14 programas com ponto flutuante (fortran{77,90}, C)

Comparação de resultados

	máq ${\cal A}$	máq <i>B</i>
prog 1 [s]	1	10
prog 2 [s]	1000	100
total [s]	1001	110

com programa 1, \mathcal{A} é 10 vezes mais rápido que \mathcal{B} com programa 2, \mathcal{B} é 10 vezes mais rápido que \mathcal{A} erm...

Média Aritmética =
$$\frac{1}{n}\sum_{i=1}^n \mathsf{Tempo}_i$$

 ${\cal B}$ é 1001/110=9.1 vezes mais rápido que ${\cal A}$

HEPR Dist RCC 37

Arquitetura II — desempenho 2007-1

Resumindo desempenho com um só Número

• média aritmética (ponderada) segue tempo de execução:

$$\sum T_i/n$$
 ou $\sum (p_i \cdot T_i)$

• média harmônica (ponderada) de taxas (MIPS/MFLOPS) segue tempo de execução:

$$n/\sum (1/R_i)$$
 ou $n/\sum (p_i\cdot R_i)$

- tempo de execução normalizado é útil para comparações de escalabilidade (máquina-A é X vezes mais rápida que modelo-base)
- para tempo de execução normalizado, deve usar média geométrica:

$$(\prod T_i/N_i)^{1/n}$$

HEPR Dinf RCC 3

Arquitetura II — desempenho 2007-1

Lei de Amdahl (i)

Tempo de execução após melhoria = (tempo de execução afetado / quanto melhorou) + tempo de execução não-afetado

esta também é importante

Exemplo:

programa executa em 100s, multiplicações consomem 80% do tempo total. Quanto devo melhorar o circuito multiplicador se quero tempo total em 20s ?

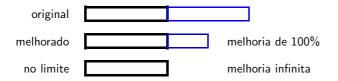
$$20 = 80/n + 20$$

Idem se quero tempo total em 40s ?

$$40 = 80/n + 20$$

Lei de Amdahl (ii)

$$\begin{array}{ll} \mathsf{Ganho}_{\mathsf{total}} & = & \frac{\mathsf{Tempo}_{\mathsf{orig}}}{\mathsf{Tempo}_{\mathsf{melhor}}} \\ \\ & = & \frac{1}{(1 - \mathsf{Frac}_{\mathsf{melhor}}) + (\mathsf{Frac}_{\mathsf{melhor}} \slash \mathsf{Ganho}_{\mathsf{melhor}})} \end{array}$$



LIEPR Dinf RCC

Arquitetura II — desempenho 2007-1

Lei de Amdahl (iii)

Exemplo:

programa executa em 100s, multiplicações consomem 80% do tempo total. Melhorando o circuito multiplicador em 50%, qual é o ganho total de velocidade ?

$$\begin{array}{ll} \mathsf{Ganho_{total}} &=& \frac{1}{(1-0.80)+\frac{0.80}{1/0.50}} \\ &=& \frac{1}{0.20+0.80*0.50} \\ &=& \frac{1}{0.60} \approx 67\% \end{array}$$

HEPR Dist RCC

Arquitetura II — conj de instruções 2007-1

Resumo

Equação do Desempenho: o que vale é o tempo de execução!

tempo de CPU =
$$\frac{\text{instruções}}{\text{programa}} \times \frac{\text{ciclos}}{\text{instrução}} \times \frac{\text{tempo}}{\text{ciclo}}$$

Lei de Amdahl: o que interessa é o desempenho global!

$$\mathsf{Ganho}_{\mathsf{total}} = \frac{1}{(1 - \mathsf{Frac}_{\mathsf{melhor}}) + (\mathsf{Frac}_{\mathsf{melhor}} \, / \, \, \mathsf{Ganho}_{\mathsf{melhor}})}$$

revisão: Definição de Arquitetura

- Conjunto de Instruções: visão do programador, compilador
- Organização: visão do projetista do computador
- Implementação I: visão do projetista do processador
- Implementação II: visão do projetista dos circuitos integrados

HEPR Dist RCC 42

Arquitetura II - conj de instruções

2007-1

revisão: Equação do Desempenho

tempo de CPU =
$$\frac{\text{instruções}}{\text{programa}} \times \frac{\text{ciclos}}{\text{instrução}} \times \frac{\text{tempo}}{\text{ciclo}}$$

$$|\text{código}| \times \text{CPI} \times |\text{ciclo}|$$

- tamanho do código depende de algoritmo e compilador
- CPI depende do Conj de Instruções e organização da CPU
- ciclo depende da organização e implementação da CPU
- objetivo: minimizar tempo e não termos isolados

HEPR Dinf RCC

Arquitetura II — conj de instruções 2007-1

Arquitetura ao longo das décadas

DÉCADA	ARQUITETO ENVOLVE-SE COM	
1960	projetar circuitos aritméticos eficientes	
1970	projeto de conjunto de instruções [1]	
1980	RISC vs CISC, super-escalaridade, compi	ladores [2]
1990	espaço de endereçamento dobra (32 $ ightarrow$ 6	64 bits)
	otimização de desvios através de execuçã	o condicional
	suporte a multi-mídia; operações de pont	o flutuante
2000	instruções largas (VLIW)	
maior suporte a especulação (execução condicional)		ondicional)
	emulação de 80×86	(Transmeta)

[1] Maior complexidade para aumentar produtividade dos programadores
[2] If a compiler cannot generate it, desktop and server programs
generally won't use it. H&P QA pg 91

Conjunto de Instruções

Recorte:

- operações
- operandos
- endereçamento de operandos
- codificação
- implementação

Instruction Set Architecture \longrightarrow ISA = parte visível ao programador e ao compilador \longleftrightarrow interface

HEPR Dist RCC

Arquitetura II — conj de instruções

2007-1

Conjunto de Instruções ≡ interface hw-sw

CDI especifica a funcionalidade do processador

- quais operações ele suporta
- quais mecanismos de armazenamento suporta e sua utilização
- como programador/compilador define programas a executar

Estudo de CDIs é a parte de arquitetura deste curso; restante é micro-arquitetura

HEPR Dist RCC 47

Arquitetura II — conj de instruções

2007-1

Conjunto de Instruções ≡ interface hw-sw

Interface bem projetada:

- sobrevive a muitas implementações (portabilidade, compatibilidade)
- é usada de muitas maneiras (generalidade)
- provê funcionalidade conveniente aos níveis acima
- permite implementação eficiente nos níveis abaixo

Receita para um "bom" Conjunto de Instruções

implementabilidade

permite implementações numa faixa de preço/desempenho

DEVE permitir implementações de alto desempenho

programabilidade

deve facilitar expressão de programas pelo programador/compilador

compatibilidade para o futuro e com o passado implementabilidade e programabilidade através de gerações: deve garantir que todo sw existente execute x86: 8086,286,386,486,Pentium{,Pro,II,III,4}

Estas são as 3 "-idades" de um CdI

HEPR Dist RCC

Arquitetura II — conj de instruções 2007-1

ISA - operações

tipo	exemplo
aritmética e lógica	add, sub, and, sll
transferência de dados	move, load
controle	branch, jump, call
sistema	syscall, trap
ponto flutuante	fadd, fmul, fdiv, sqrt
decimal	addd, convert
string	move, comp
multimidia 2D, 3D	MMX/SSE

LIEPR DIM RCC St

Arquitetura II — conj de instruções 2007-

ISA - operações

10 instruções do x86 mais populares em 5 programas SPECint95:

ord	instrução	% total	
1	load	22	
2	desv cond	20	1/5
3	compare	16	
4	store	12	Id+st o 1/3
5	add	8	
6	and	6	
7	sub	5	
8	move reg-reg	4	por que?
9,10	call+ret	1 + 1	
	total	96	

Amdahl legisla: estas devem ser otimizadas

HEPR Dinf RCC 5:

ISA - operandos

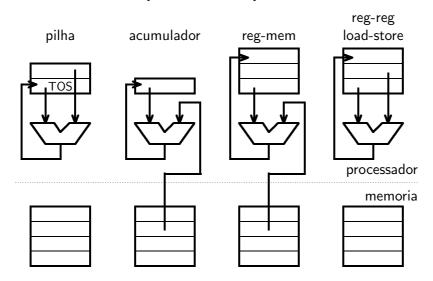
- Localização dos operandos
- * no processador
- * em memória
- Quantos bits são necessários para identificar um operando?
- * no processador
- * em memória
- Quanto tempo é necessário para acessar um operando?
- * no processador
- * em memória

HEPR Dist RCC 55

Arquitetura II — conj de instruções

2007-1

ISA - operandos no processador



HEPR Dinf RCC

Arquitetura II — conj de instruções

2007-1

ISA - operandos no processador

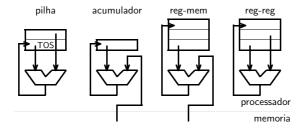
Operandos no processador: acesso rápido & identificação compacta

- Acumulador
- * código compacto (endereço implícito)
- * hardware simples
- * tráfego elevado de/para memória
- Pilha
- * código compacto (endereço do topo da pilha implícito)
- Registradores
- * registradores devem ser identificados
- * espaço de nomes adicional (variáveis em memória e registradores)
- * menos acessos à memória
- * maior velocidade porque regs. são mais rápidos que memória

ISA - operandos no processador

código para computar C = A + B

PILHA	ACUMUL	REG-MEM	REG-REG
push A	load A	load R1,A	load R1,A
push B	add B	add R3,R1,B	load R2,B
add	store C	store R3,C	add R3,R1,R2
pop C			store R3,C



LIEPR Dinf BCC

Arquitetura II - conj de instruçõe

2007-1

ISA - operandos para operações de ULA

- Instruções da ULA combinam operandos
- número de operandos é implícito:
- * dois: r1 ← r1 OP r2 código denso destrutivo porque sobre-escreve destino assimétrico (A-B, A/B)
- * três: r1 ← r2 OP r3
- operandos em registradores ou em memória
- * registrador + memória (IBM, x86)

NÃO ORTOGONAL

* reg + reg (RISCs, Cray)

ortogonal com load/store

ortogonal,

* qualquer combinação (VAX)

instr com tamanho variável

HEPR Dinf RCC

Arquitetura II — conj de instruções

2007-1

ISA - de volta a Lilliput

Ordem dos bytes nas palavras

• big endian: byte com endereço xxxx002 na posição mais significativa de uma palavra de 32 bits Motorola (powerPC), IBM (r6000), SPARC

end byte	MS			ms
0	0	1	2	3
4	4	5	6	7

• little endian: byte com endereço xxxx002 na posição menos significativa de uma palavra de 32 bits DEC (alpha), Intel (x86)

end byte	MS			ms
0	3	2	1	0
4	7	6	5	4

ambos – MIPS (e powerPC)

ISA - alinhamento de operandos

• Palavras tem "tamanho natural":

byte, half-word, word, doubleword float, double, quadword

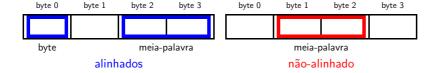
• Operando alinhado pelo tamanho natural se

endereço MOD tamanho = 0

 \star palavra alinhada: ender mod 4 = 0 1024, 1028, 102c, 1030,... $\mod 4$ porque palavra tem tamanho de 4 bytes

★ palavra desalinhada ender mod $4 \neq 0$

1025 ou 1026 ou 1027



TIEPR DINFRCC 5

Arquitetura II — conj de instruções

2007-1

ISA - alinhamento de operandos

- sem restrições
- ★ software mais simples
- $\star\,$ hardware deve detectar e fazer ≥ 2 referências à memória
- ★ lógica complexa; tempo extra em todas as instruções (pior caso)
- * dificulta adiantamento de loads no pipeline
- com restrições
- ★ software deve garantir alinhamento
- * hardware verifica e provoca excessão se for o caso
- \star adiantamento de loads mais simples
- mais ou menos
- * mais de uma instrução para acessos desalinhados (load-half)
- * compilador ajuda no alinhamento
- ★ hardware verifica e provoca excessão se for o caso

HEPR Dinf RCC

Arquitetura II — conj de instruções

2007-1

ISA - modos de endereçamento

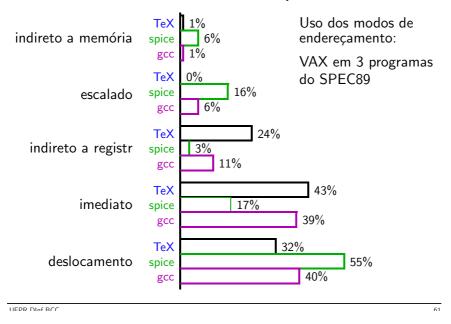
MODO	ENDER EFETIVO	EXEMPLO
a registrador	R	add r4,r3
imediato	imed	add r4,#8
deslocamento	M[R+imed]	add r4,100(r1)
indireto a registrador	M[R]	add r4,(r1)
indexado	$M[R_i {+} R_j]$	add $r3$, $(r1+r2)$
absoluto	M[imed]	add r1,(1001)
indireto a memória	M[M[R]]	add r1,@(r3)
auto-incremento	M[R]; R += d	add $r1,(r2)+$
auto-decremento	R -= d; M[R]	add r1,–(r2)
escalado	$M[R_i + R_j*d]$	add r1,100(r2)[r3]

Apenas absoluto e indireto a registrador são fundamentais; outros podem ser derivados destes

Arquitetura II - conj de ins

2007-1

ISA - modos de endereçamento



ISA - modos de endereçamento

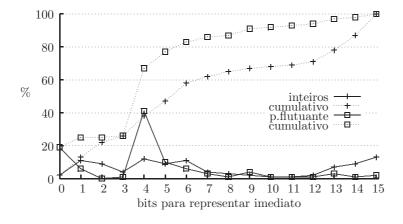
- Distância dos deslocamentos (LD e ST)
- * 16 bits cobre 99% dos casos (dist $\leq \pm$ 32K bytes/pals) Alpha (SPEC CINT2000 e CFP2000)
- Uso de imediatos (tamanho das constantes)
- *~1/5 a 1/4 das instruções de ULA e movimentação de dados usam constantes
- $*\approx 90\%$ das constantes são representadas em ≤ 16 bits instruções suportam constantes de 16bits (como representa maiores?)
- * 20 a 30% das constantes são negativas

HEPR DInf RCC

Arquitetura II — conj de instruções 2007-1

ISA - modos de endereçamento

Uso de imediatos: Alpha SPEC CINT2000 e CFP2000 1/5-1/4 das instr de ULA e de movim de dados usam constantes



Administrativamente falando...

- Avaliação:
- * BCC
 - \star 2 provas 35% + 35%
 - ★ trabalho 20%
 - * participação nas aulas de exercícios 10%
- * Mestrado
 - \star 2 provas 30% + 30%
 - ★ trabalho 20% (entrega na semana das finais)
 - * participação nas aulas de exercícios 20%
- 5 aulas de exercícios
- ★ listas distribuídas na aula

HEPR Dinf RCC

Arquitetura II — conj de instruções

2007-1

Administrativamente falando...

Todos devem ler o artigo sobre o Cray-1 para a próxima aula The CRAY-1 Computer System, R M Russel, CACM 21(1), Jan 78

instruções vetoriais:

```
• addv v1,v2,v3 # for (i=0;i<64;i++) v1[i]=v2[i]+v3[i];
```

• encadeamento (chaining):

"saída" do ldv passa direto para muliv, para addv e para stv

HEPR Dinf RCC

Arquitetura II — conj de instruções

2007-1

ISA - controle de fluxo de execução

- Cada instrução é auto-contida (≠ if-then-else, do-until)
- Seqüência de execução de instruções é determinada pelo PC
- * execução seqüencial: PC++
- * instruções mudam conteúdo do PC explicitamente: PC ← ENDER
- Aspectos importantes:
- * desvios condicionais (toma desvio ou segue seqüencialmente)
- * determinação do endereço de destino
- * ligação para endereço de retorno de funções
- * salvamento de estado
- * interface com sistema operacional

Tipos de instruções de controle:

desvios condicionais (desvia ou prossegue)
 saltos incondicionais
 chamada de função
 geralmente incondicionais
 necessita ligação para endereço de retorno
 pode salvar estado (registradores)
 retorno de função

* geralmente incondicionais

* retorno indireto através do endereço de ligação

* pode salvar estado (registradores)

IIEDR Dief RCC 67

Arquitetura II — conj de instruções 2007-1

ISA - controle de fluxo de execução

Tipos de instruções de controle:

• chamadas de sistema syscall, trap

* semelhante à chamada de função

* endereço de destino definido no opcode

* atravessa fronteira de proteção

• retorno de chamada de função reti

* semelhante ao retorno de função

* opcode diferente da chamada

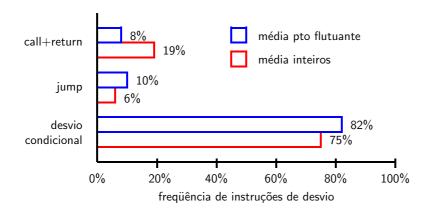
* atravessa fronteira de proteção

Arquitetura II — conj de instruções

HEDRING RCC 69

ISA - controle de fluxo de execução

Tipos de desvios: Alpha SPEC CINT2000 e CFP2000



Tipos de mudança no fluxo:

- condicional ou incondicional
- salva conteúdo do PC
- o endereço de destino endereço fixo (imediato ou PC+imediato) definido em tempo de execução (registrador)

```
HEPR Ninf RCC 70
```

Arquitetura II - conj de instruções

2007-1

2007-1

ISA - controle de fluxo de execução

- Mecanismos de decisão de desvios:
- * registrador de status (condition code)

 $PC \Leftarrow (R1>R2 ? PC+desloc : PC)$

operação de ALU no caminho crítico da resolução do

desvio

Arquitetura II - conj de instruções

* registrador de condição (CC = condition code register)

 $R1 \leftarrow (R2 - R3) \text{ mudaCC}$

 $PC \leftarrow (CC ? PC + desloc : PC)$

necessita estado adicional (CC)

dificulta reordenação de código (pelo efeito colateral)

* registrador de uso geral

 $R1 \leftarrow (R2 > R3)$

 $PC \leftarrow (R1>0 ? PC+desloc : PC)$

pouco estado adicional

usa registrador inteiro para poucos bits

HEPR Dist RCC

ISA - controle de fluxo de execução

Mecanismos para especificar endereço de destino:

• endereço arbitrário como load/store

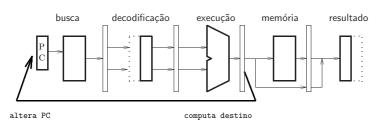
```
PC \Leftarrow PC + (R1 + desloc)
```

versátil e ortogonal

mais demorado para decodificar

executa desvio e computa destino em estágios separados

do pipeline (muda PC na busca e computa ender em exec)



Mecanismos para especificar endereço de destino:

• relativo ao PC (deslocamento em imediato)

PC ← (PC+desloc)

código independente de posição facilita ligação

destino pode ser computado no estágio de busca/decodificação imediato com 12 bits cobre 99% dos casos (MIPS)

 $(\mathsf{dist} \leq \pm \ \mathsf{2K} \ \mathsf{pals})$

imediato com 16 bits cobre 99,5% dos casos (ALPHA)

(dist $\leq \pm$ 32K pals)

MAS endereço de destino não pode ser muito distante

HEPR DInf RCC

Arquitetura II — conj de instruções

2007-1

ISA - controle de fluxo de execução

Mecanismos para especificar endereço de destino:

• indireto a registrador (endereço em registrador)

 $PC \Leftarrow R1$

código compacto

destino a qualquer distância

destino pode ser computado dinamicamente (retorno de função)

instrução adicional para carregar registrador

muda PC e "calcula" destino em estágios diferentes

endereço implícito no opcode (vetor de syscalls)
 necessário para interface com sistema operacional

LIEPR Dief RCC 74

Arquitetura II — conj de instruções

2007-1

ISA - controle de fluxo de execução

Mecanismos para especificar endereço de ligação/retorno:

- necessário para retorno de função e chamada de sistema
- registrador implícito (R31 no MIPS)
 registrador deve ser salvo antes da próxima chamada
- registrador explícito

pouquinho mais de flexibilidade

 pilha de chamada de funções
 SE pilha é parte da arquitetura complexidade adicional:

instruções tipo push e pop

hardware deve detectar over/under-flow na pilha

Salvar e recompor estado nos saltos/desvios

- que estado salvar
- * chamadas de função ⇒ registradores (CISCs)
- * chamadas de sistema \rightharpoonup registradores, PC, flags, PSW
- ISA/hardware não necessita salvar registradores
- * função que chama salva regs em uso

caller save

* função chamada salva regs que vai usar

callee save

- ISA/hardware pode salvar registradores
- * VAX call
- * pressupõe convenção de chamada/retorno pelo compilador
- processadores recentes n\u00e3o fazem salvamento mas SPARC tem janelas de registradores

HEPR Dinf RCC 7

Arquitetura II — conj de instruções

2007-1

ISA - codificação das instruções

Tipos de codificação: comprimento variável, fixo e híbrido



Variavel: VAX, Intel x86



Fixo: Alpha, ARM, MIPS, PowerPC, SPARC

operação	modo1	ender1	
	modo1	ender1	ender2
operação	modo1	modo2	ender

Hibrido: IBM 360/370, MIPS16, TI TMS5320c54x

HEPR Dist RCC 7

Arquitetura II — MIPS 2007-1

ISA - resumo

Instruction Set Architecture = ISA = parte visível ao programador e ao compilador

operações
 operandos
 popularidade vs implementação
 popularidade vs implementação

• endereçamento de operandos código compacto vs flexibilidade

codificação
 código compacto vs decodificação
 facilidade de implementação em pipeline

revisão: Desempenho

Equação do desempenho: o que interessa é o tempo total!

tempo de CPU =
$$\frac{\text{instruções}}{\text{programa}} \times \frac{\text{ciclos}}{\text{instrução}} \times \frac{\text{tempo}}{\text{ciclo}}$$

$$|\text{código}| \times \text{CPI} \times |\text{ciclo}|$$

Lei de Amdahl: o que interessa é o desempenho global!

$$\mathsf{Ganho}_{\mathsf{total}} = \frac{1}{(1 - \mathsf{Frac}_{\mathsf{melhor}}) + (\mathsf{Frac}_{\mathsf{melhor}} \, / \, \, \mathsf{Ganho}_{\mathsf{melhor}})}$$

HEPR Dist RCC 70

Arquitetura II — MIPS 2007

revisão: Conjunto de Instruções

- Instruction Set Architecture = ISA
 = parte visível ao programador e ao compilador
- contrato entre hardware e software
- simplifica compilador se
- * ortogonal
- * regular
- * facilitar composições
- codificação simples e regular simplifica hardware
- ⋆ operações popularidade vs implementação
- ★ operandos popularidade vs implementação
- ★ endereçamento de operandos código compacto vs flexibilidade
- * codificação código compacto vs decodificação facilidade de implementação em pipeline

HEPR Dinf RCC

Arquitetura II — MIPS 2007-1

Compiladores 0

Função do compilador:

- todos os programas corretos executam corretamente
- maioria dos programas compilados executa rapidamente
- compilação rápida
- suporte a depuração

bloco básico = trecho de código entre desvios

Compiladores 0

estágio	depende de	função
front end	linguagem	produz representação intermediária
otimizador de	linguagem	transformação de loops
alto nível	(pouco)	integração de funções
otimizador	máquina	otimização local e global
global	(pouco)	alocação de registradores
gerador de	máquina	seleção de instruções
código	(muito)	otimização dependente de máquina

 IIEPP Dief RC
 8

 Arquitetura II — MIPS
 2007

Compiladores 0 - estilos de transformação

- otimização no alto nível
- * efetuada na linguagem de alto nível para uso pelos outros estágios
- otimização local
- * otimização somente em blocos básicos
- otimização global
- * otimização que atravessa blocos básicos
- * transformação de loops
- alocação de registradores
- * associa registradores com operandos
- otimização dependente do processador
- * depende da arquitetura

LIEPR DIM RCC 8:

Arquitetura II — MIPS 2007-

Compiladores 0 - operações nos estágios

- tradução para representação intermediária
- expansão de funções inlining
- otimização de loops desenrolar; código invariante
- eliminação da variável de indução cálculo de índices
- eliminação de sub-expressões comuns cálculos repetidos
- otimizações de saltos
- propagação de constantes
 variável que é constante
- alocação de registradores
- strength reduction mult vs soma+desloca
- escalonamento do pipeline reordenação; branch delay slots
- tradução para linguagem de máquina

Arquitetura II - MIPS

Compiladores 0 - do que eles gostam?

- quem escreve um compilador deseja
- * regularidade simplifica análise de casos
- * ortogonalidade suporta todas as combinações * composabilidade primitivas ao invés de soluções
- * as três permitem operações simples combinadas em operações complexas
- compiladores efetuam análise de casos gigantesca
- * opções demais dificultam escolhas
- conjuntos de instruções ortogonais quanto a
- * operações tipos de dados modos de endereçamento
- * completude

condições de desvio → eq lt uma, ou cond de desvio → eq ne lt gt le ge todas as soluções, mas não só algumas escolhas idiossincráticas

LIEPR Dinf BCC

Arquitetura II - MIPS 2007-1

Conjunto de Instruções do MIPS

Projeto de RISCs na década de 80 para obter implementação de pipelines num único CI

> Reduced Instruction Set Computers reduced == simples, e não pequeno

- ênfase em
- * decodificação rápida
- * instruções com tamanho fixo
- * codificação regular
- compilador poderia escalonar instruções para execução
- código grande/esparso

LIEPR Dinf BCC

Arquitetura II - MIPS

MIPS64

- endereços alinhados de 32 bits
- modo de endereçamento é deslocamento load/store
- tipos de dados simples: byte half word doubleword, single double
- registradores
- * 32 regs de uso geral, de 64 bits (R0 = 0) regs
- * 32 regs de ponto flutuante de 64 bits regsPF
- * registrador de status para ponto flutuante
- * sem registrador de status para inteiros
- três formatos de instrução com mesmo tamanho

LIEPR Dinf RCC

MIPS - modos de endereçamento

MODO	ENDER EFETIVO	EXEMPLO
a registrador	R	add r4,r3,r2
imediato	imed	add r4,#8
deslocamento	M[R+imed]	add r4,100(r1)
indireto a registrador	[R]	jr r4
absoluto	ender	j ender

HEPR Dinf RCC Arquitetura II — MIPS

2007-1

MIPS

• transferência de dados

- * load/store byte/half/word/doubleword
- * load/store PF single/double
- * move de-para regs e regsPF
- ULA
- * add/sub/mult/div
- * and/or/xor
- * sll/srl (lógicos), sra (aritmético) deslocamentos
- * loadHigh (usado para constantes de 32 bits)

HEPR Dinf RCC

Arquitetura II — MIPS

MIPS

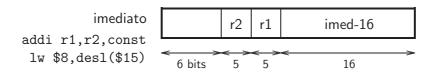
- ponto flutuante
- * add/sub/mult/div single/double
- * conversões de-para inteiros
- * desvios (liga/desliga bits para desvios)
- controle
- $\neq 0$ bits_PF * desvios condicionais: =0

* jump/jr jump-register * jal jump-and-link-register * trap/rte return from exception

MIPS

Formatos das instruções

- formato-I
- * instruções de ALU com imediatos
- * load e store
- * desvios condicionais
- * jump-register



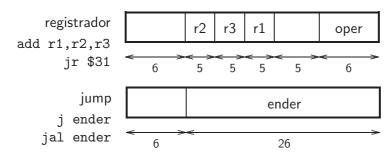
TIEDR Dist RCC

Arquitetura II — MIPS 2007-1

MIPS

Formatos das instruções

- formato-R
- * instruções de ALU com três operandos
- formato-J
- * saltos incondicionais

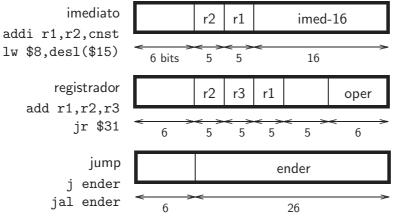


HEPR Dinf RCC

Arquitetura II — MIPS 2007-1

MIPS

- ullet instruções regulares facilitam decodificação \longrightarrow hw rápido
- instruções simples facilitam construção do compilador e geração de código → sw rápido



Administrativamente falando...

ou

filmes imperdíveis para computeiros:

Enigma → ponto de vista dos europeus que quebraram o código U571 → marinheiros/cowboys capturam máquina de codificação assistam aos dois

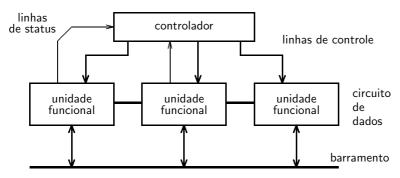
Quem conseguir, leia *The enigma of inteligence* que é a biografia de Alan Turing, tido como o pai da IA além de muitas outras coisas

HEPR Dist RCC 04

Arquitetura II — MIPS 2007-1

MIPS - microarquitetura

Microarquitetura = implementação do conjunto de instruções



Estrutura ↔como componentes são interligados Seqüenciamento ↔movimento de dados entre componentes

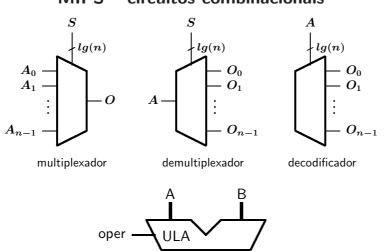
estrutura = datapath estática següenciador = controle dinâmico

HEPR Dinf RCC

MIPS - circuitos combinacionais

2007-1

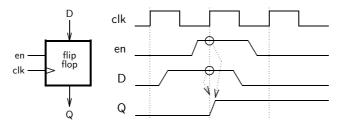
Arquitetura II — MIPS

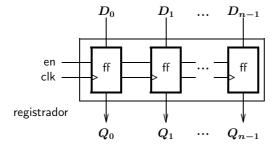


HEPR Dinf RCC 96

stat Res

MIPS – circuitos seqüenciais

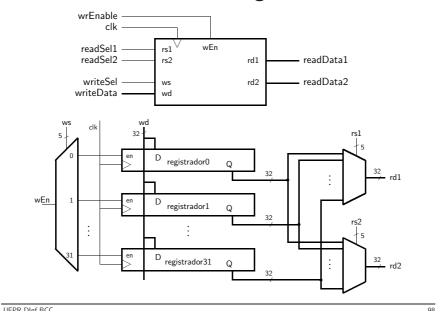




TIEPR Dist RCC 07

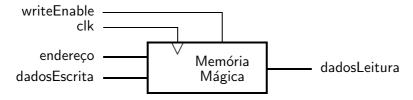
Arquitetura II — MIPS 2007-1

MIPS - bloco de registradores



Arquitetura II — MIPS 2007-1

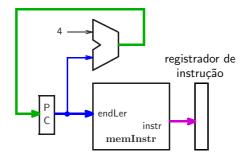
MIPS - memória idealizada



Modelo simplificado da memória

- leituras e escritas completam em 1 ciclo
- * leitura pode ocorrer a qualquer instante (combinacional)
- * escrita efetuada na borda ascendente do relógio escrita habilitada
 - → dados e endereços estáveis na borda do relógio
- modelo mais realista será discutido adiante

MIPS - busca de instruções



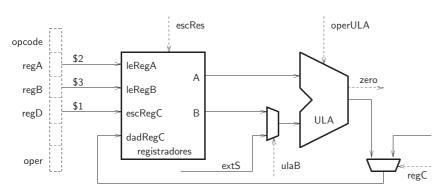
- conteúdo do PC indexa memória
- ao final do ciclo
- * nova instrução no Registrador de Instrução
- * PC é incrementado

HEPR Dist RCC 100

Arquitetura II — MIPS 2007-1

MIPS - instruções de ALU (i)

add \$1, \$2, \$3 # \$1 <- \$2 + \$3

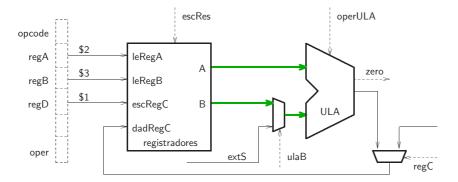


HEPR Dinf RCC 101

Arquitetura II — MIPS 2007-1

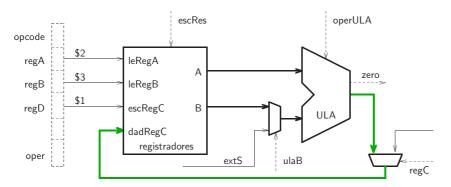
MIPS - instruções de ALU (ii)

add \$1, \$2, \$3 # \$1 <- \$2 + \$3



MIPS - instruções de ALU (iii)

add \$1, \$2, \$3 # \$1 <- \$2 + \$3

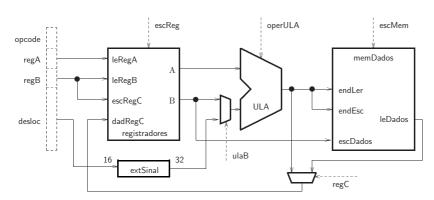


TIEPR DISF RCC 103

Arquitetura II — MIPS 2007-1

MIPS - acesso à memória (i)

lw \$8, desloc(\$15) # \$8 <- M[desloc + \$15]
sw \$8, desloc(\$15) # M[desloc + \$15] <- \$8</pre>

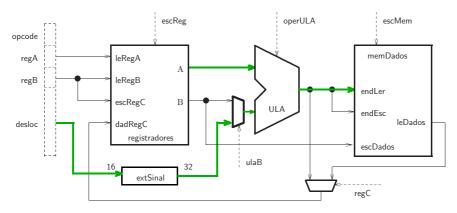


HEPR Dinf RCC 10A

Arquitetura II — MIPS 2007-1

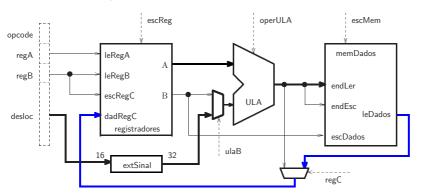
MIPS - acesso à memória (ii)

lw \$8, desloc(\$15) # \$8 <- M[desloc + \$15]</pre>



MIPS - acesso à memória (iii)

lw \$8, desloc(\$15) # \$8 <- M[desloc + \$15]</pre>

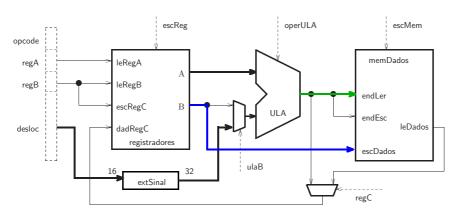


TIEPR DISF RCC 106

Arquitetura II — MIPS 2007-1

MIPS - acesso à memória (iv)

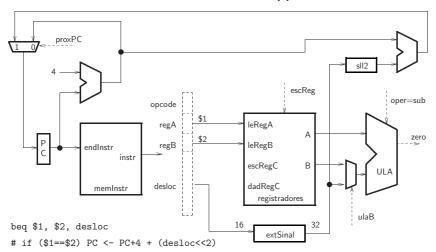
sw \$8, desloc(\$15) # M[desloc + \$15] <- \$8



HEPRING RCC 100

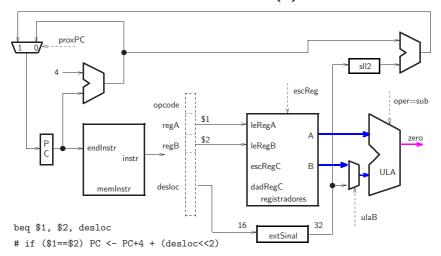
Arquitetura II — MIPS 2007-1

MIPS - desvios (i)



Arquitetura II — MIPS 2007-1

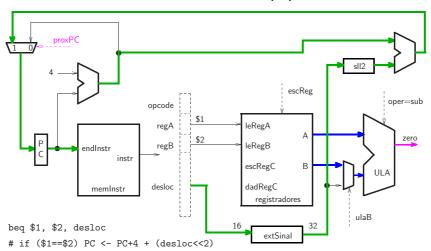
MIPS - desvios (ii)



HEPR Dist RCC 100

Arquitetura II — MIPS 2007-1

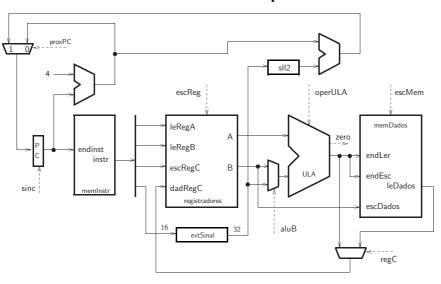
MIPS - desvios (iii)



HEPR Dinf RCC 110

Arquitetura II — MIPS 2007-1

MIPS - CPU completa



LIFPR Dinf RCC 111

MIPS - exercícios

Desenhe os diagramas completos do processador para executar as instruções:

- jump
- jump-and-link
- jump-register
- addi (add imediato)

HEPR Dist RCC 117

Arquitetura II — segmentação

2007-1

revisão: compiladores

- quem escreve um compilador deseja
- * regularidade
- * ortogonalidade
- * composabilidade
- conjuntos de instruções ortogonais quanto a
- * operações
- * tipos de dados
- * modos de endereçamento
- * completude

uma, ou condições de desvio \Longrightarrow eq lt todas as soluções, condições de desvio \Longrightarrow eq ne lt gt le ge mas não algumas

HEPR DIM RCC 113

Arquitetura II — segmentação

2007-1

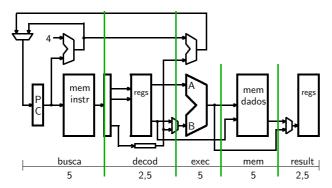
Modelo sequencial

Conjunto de instruções de processadores "comuns" define um modelo seqüencial de execução:

cada instrução é completamente executada e altera o estado do processador antes do início da próxima instrução

Este modelo facilita muito a vida do programador!

Implementação mono-ciclo

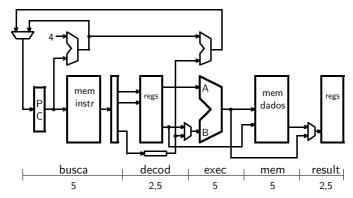


- * busca busca instrução que será executada
- * decodificação decodifica instrução e acessa conteúdo dos registradores
- * execução executa operação de ULA ou cálculo de endereço
- * memória acesso à memória de dados
- * resultado armazena resultado no bloco de registradores

HEPR Dist RCC 116

Arquitetura II — segmentação 2007-1

Implementação multi-ciclo



Instrução completa em alguns ciclos de relógio: CPI>1

$$egin{array}{ll} T_{
m ciclo} & \geq & \max(T_{
m busca}\,, T_{
m decod}\,, T_{
m ex}\,, T_{
m mem}\,, T_{
m res}) \ & \geq & {f 5} \; {
m unidades} \; {
m de} \; {
m tempo} \end{array}$$

HEPR Dist RCC 116

Arquitetura II — segmentação 2007-1

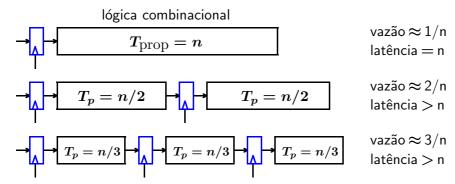
Implementação multi-ciclo

$$\mathsf{CPI} = \sum_{j=0}^{n} \mathsf{CPI}_{j} \times \mathsf{F}_{j}$$

$$\frac{\mathsf{instr}}{\mathsf{ALU}} \quad \begin{array}{c} \mathsf{freq}[\%] \quad \mathsf{ciclos} \\ \mathsf{ALU} \quad \quad 40 \quad 4 \\ \mathsf{load} \quad \quad 25 \quad 5 \\ \mathsf{store} \quad \quad 10 \quad 4 \\ \mathsf{desvios} \quad \quad 25 \quad 4 \\ \hline \quad \quad \quad \mathsf{CPI} \quad 4.25 \\ \end{array}$$

IIFPR Dinf RCC 117

Segmentação

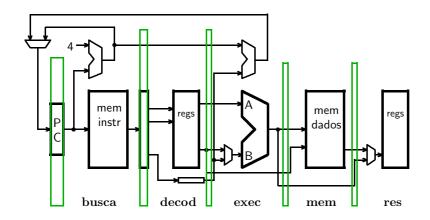


- Segmentação: adicionar registradores entre unidades funcionais
- vazão é proporcional ao número de estágios
- latência aumenta por causa dos registradores

HEPR Dist RCC 11

Arquitetura II — segmentação 2007-1

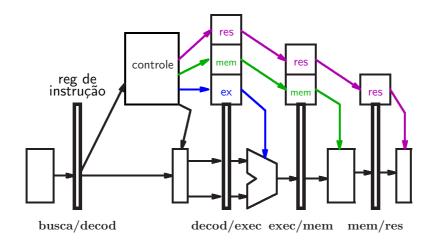
Implementação com segmentos



HEPR DIM RCC 110

Arquitetura II — segmentação 2007-1

Implementação com segmentos - controle



LIEPR Dinf RCC 120

Segmentação

- Circuito de dados
- * componentes diversos com tempos de propagação distintos
- * instruções usam mais componentes que os disponíveis
 - ⇒ riscos estruturais

structural hazards

- dependências de dados entre instruções nos diversos estágios
- * ∃ dependências entre instruções nos estágios
 - ⇒ dependências de dados

data dependencies

- o escalonamento de instruções pode ser afetado por instruções nos outros estágios
- * ∃ relacionamento entre instruções nos estágios
 - ⇒ dependências de controle

control dependencies

LIEPR Dinf BCC

Arquitetura II — segmentação

2007-1

Desempenho da segmentação

$$\mathsf{ganho} \ = \ \frac{\mathsf{CPI} \ \mathsf{sem} \ \mathsf{pipeline}}{\mathsf{CPI} \ \mathsf{com} \ \mathsf{pipeline}} \times \frac{\mathsf{ciclo} \ \mathsf{sem} \ \mathsf{pipeline}}{\mathsf{ciclo} \ \mathsf{com} \ \mathsf{pipeline}}$$

HEPR Dinf RCC

122

Arquitetura II — segmentação

2007-1

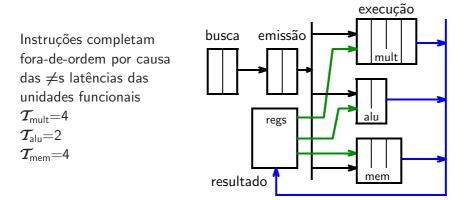
Desempenho da segmentação (cont.)

Tipos de processadores segmentados linear busca decod exec mem res busca emissão execução regs regs regs tipo Diagrací

Pipeline não-linear

2007-1

Arquitetura II - segmentação



Barramento de resultado deve ser reservado quando instrução é emitida para execução (reg de deslocamento) → se barramento estará ocupado, segura execução

HEPR Dinf RCC 12

Arquitetura II — segmentação 2007-1

Riscos com segmentação

Riscos são condições que levam a comportamento incorreto se medidas apropriadas não forem tomadas

• Riscos estruturais structural hazards

* quando duas instruções **diferentes** necessitam do **mesmo** recurso no **mesmo** ciclo

• riscos com dados data hazards

- * quando duas instruções **diferentes** usam **mesmo** local de armazenamento;
- * resolução do risco deve garantir aparência de que instruções executaram na ordem seqüencial correta
- riscos de controle control hazards
- * quando uma instrução determina **quais** instruções serão executadas a seguir (desvios, saltos, funções)

Riscos com segmentação

Riscos

Arquitetura II - segmentação

- * potenciais violações de dependências no programa
- * implementação deve garantir que dependências não são violadas
- Resolução de riscos
- * estática: compilador/programador garante corretude
- * dinâmica: hardware verifica conflitos em tempo de execução
- Inter-travamentos no pipeline

interlocks

- * mecanismo para resolução dinâmica de riscos
- * deve detectar e resolver dependências em tempo de execução

HEDR Dief RCC 10

Riscos estruturais

- Duas instruções diferentes necessitam do mesmo recurso no mesmo ciclo
- * exemplo: processador em que loads e stores usam mesma porta de memória que busca de instruções

- * nenhuma instrução completa no oitavo ciclo (bolha)
- * Arquitetura Harvard versus Arquitetura Princeton

HEPR Dinf RCC 128

Arquitetura II — segmentação 2007-1

Resolução de riscos estruturais

Riscos estruturais podem ser reduzidos se:

- cada instrução usa recurso somente uma vez,
- sempre no mesmo estágio, e
- durante um ciclo.
- RISCs são projetados para satisfazer estes requisitos:
- ightharpoonup por causa do pipeline, memória deve ser capaz de atender ≥ 1 requisição/ciclo 1 instr/ciclo + 1/4 ld-st/ciclo

Riscos com dados

DATA HAZARDS:

Quando duas instruções **diferentes** usam **mesmo** local de armazenamento

Deve parecer que instruções executam na ordem següencial correta

```
i: r1 ← r4 + r5
j: r2 ← r1 - r9  r1 foi produzido por i
k: r1 ← r6 ⊕ r3  valor de r1 em i é sobre-escrito resultado de i;j;k é o mesmo que i;k;j?
```

Convenção: **nome do risco** é a ordem do programa que **deve ser preservada pela implementação** (segment, superescalar)

HEPR Disf RCC 120

Arquitetura II — segmentação 2007-1

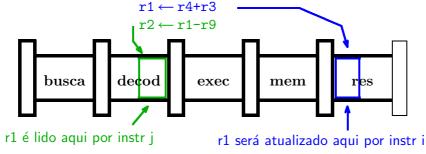
Riscos com dados - RAW

• Read-After-Write (RAW) instr **j** tenta ler operando r1 ANTES que instr **i** escreva resultado

i:
$$r1 \leftarrow r4 + r3$$

j: $r2 \leftarrow r1 - r9$

 Risco decorre de uma dependência de dados, causada pela comunicação entre as duas instruções add e sub através de r1



HEPR Dinf RCC

Riscos com dados - WAR

2007-1

Write-After-Read (WAR)
 instr j escreve operando em r2 ANTES que instr i leia-o

i:
$$r4 \leftarrow r2 + r3$$

j: $r2 \leftarrow r1 - r9$

Arquitetura II - segmentação

- Risco decorre de uma anti-dependência dependência artificial causada pelo re-uso do nome r2
- anti-dependências não ocorrem num pipeline linear de 5 estágios:
- ⊳ leituras ocorrem sempre no estágio 2 (decod)
- ⊳ escritas ocorrem sempre no estágio 5 (result)
- → MAS anti-dependências ocorrem em pipelines não-lineares p.ex CRAY e super-escalares

Riscos com dados - WAW

• Write-After-Write (WAW) instr **j** escreve operando em r3 ANTES que instr **i** escreva-o

i:
$$r3 \leftarrow r4 + r5$$

j: $r3 \leftarrow r8 - r9$

- Risco decorre de uma dependência de saída dependência artificial causada pelo re-uso do nome r3
- dependências de saída não ocorrem num pipeline linear de 5 estágios
- ⊳ escritas em regs ocorrem somente no estágio 5 (result)
- ⊳ MAS dependências de saída ocorrem em pipelines não-lineares
- Read-After-Read (RAR) não é risco

HEPR Dinf RCC 133

Arquitetura II — segmentação 2007-1

Dependências e riscos

dependências são características do programa o programador escreveu o código numa certa ordem

riscos decorrem de características da implementação o projetista do processador escolheu truques "perigosos"

HEPR Dinf RCC

Arquitetura II — segmentação 2007-1

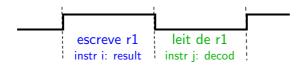
Solução simples para RAW (parcial)

- Circuito detecta risco, e então insere bolha:
 atrasa instrução j até ocorrer escrita do resultado de i
- * solução simples segura instrução dependente na busca
- * desempenho ruim por causa dos ciclos desperdiçados

i: $r1 \leftarrow r2+r3$ B D Ex M R j: $r2 \leftarrow r1-r9$ B ** ** D Ex M F

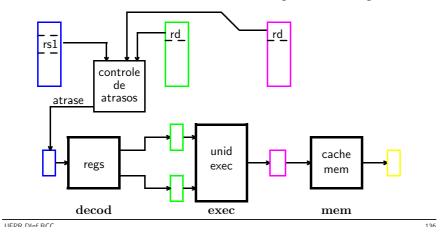
stall

 solução pressupõe que, em cada ciclo, registradores são atualizados e então são lidos



Controle de atrasos (stalls)

- Compara com estágios posteriores
- * if (rs1(decod)==rd(exec) || rs1(decod)==rd(mem)) { bolha }
- ★ mesmo para rs2
- ★ nem todos conflitos são riscos: st não escreve reg, addi não lê reg...

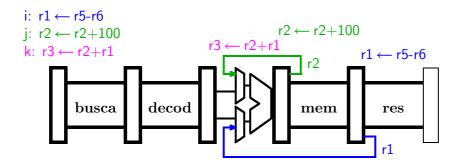


Arquitetura II — segmentação

2007-1

Adiantamento

Ao invés de atrasar, adianta resultado para entradas da ULA: usa controle de atrasos para decidir se deve adiantar, e usa multiplexadores para escolher fonte do resultado.



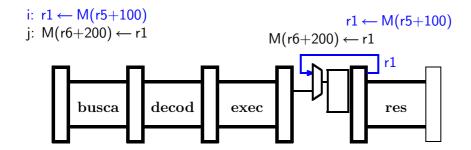
Adiantamento = forwarding, bypassing, short-circuiting

HEPR Dinf RCC 12:

Arquitetura II — segmentação 2007-1

Adiantamento

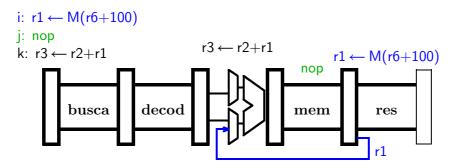
Adiantamento para o estágio de memória load seguido de store



LIEPR Dinf RCC 138

Adiantamento

Adiantamento para o estágio de memória: load seguido por add



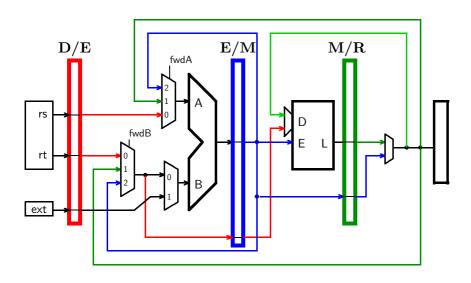
- Risco deve ser detectado por hardware e bolha inserida (nop)
- * desempenho cai por causa da bolha
- Compilador deve tentar preencher bolha com instrução "boa"
- ★ load delay slot usado no CdI MIPS-I e eliminado em MIPS-II por que?

HEPR Dist RCC 130

Arquitetura II — segmentação

Adiantamento - circuito completo

2007-1

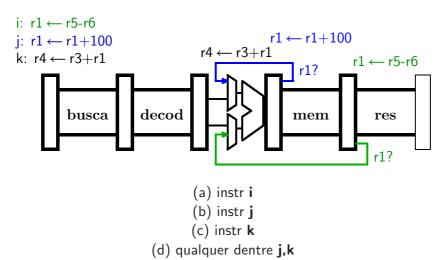


HEPR Dist RCC

Arquitetura II — segmentação 2007-1

Pergunta:

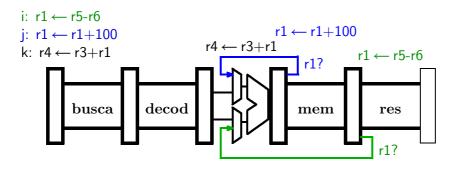
Quem fornece r1 para instrução k?



LIEPR Dinf RCC

Resposta:

Quem fornece r1 para instrução k?



(a) instr i

(b) instr j é a que mantém a ordem seqüencial

(c) instr k

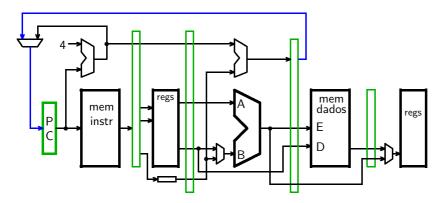
(d) qualquer dentre j,k

HEPR Dist RCC 1.

Arquitetura II — segmentação 2007-1

Riscos de Controle

Quando uma instrução determina **quais** instruções serão executadas a seguir (desvios, saltos, funções)



Desvio é efetuado em **mem**, após comparação na ULA ⇒ as instruções no caminho não-tomado devem ser anuladas

HEPR Dinf RCC 14

2007-1

Riscos de Controle

		1	2	3	4	5	6	7	8	9
i:	beq r1,r0,8	В	D	Е	М	R				
i+1:	sub r2,r8,r9		В	D	Ε	*	*			anulada
i+2:	lw rx, C(ry)			В	D	*	*	*		anulada
i+3:	sw ri,D(rj)				В	*	*	*	*	anulada
i+8:	add r3,r4,r5					В	D	Ε	Μ	R
i+9:	add r6,r7,r8						В	D	Ε	B R

Se CPI=1, 30% das instruções são desvios, penalidade de 3 ciclos \Longrightarrow CPI_{novo}=1.9

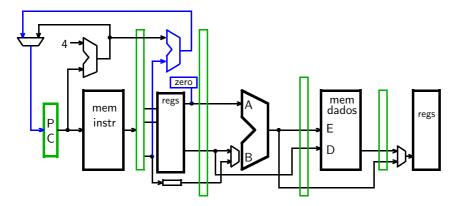
• Desempenho pode ser melhorado:

Arquitetura II — segmentação

- * descobre direção antes: compara com zero em decod
- * computa endereço de destino: adiciona somador em decod
- → penalidade de UM ciclo ao invés de três

Riscos de Controle

Desvio é efetuado em **decod**, após comparação com ZERO



Se CPI=1, 30% das instruções são desvios, penalidade de 1 ciclo \Longrightarrow CPI_{novo}=1.3

LIEPR Dinf BCC

Arquitetura II — segmentação 2007-1

Tratamento de Riscos de Controle

- Tratamento de riscos de controle é MUITO importante pelo impacto no desempenho
- * de 1/3 a 1/6 de todas instruções são desvios
- * penalidade de um, dois ou três ciclos
- * pipelines com mais estágios \iff desempenho pior
- Atrasa instruções até decidir se desvia
- Previsão de desvios → prever direção do desvio
- * reduz/elimina penalidade se previsão correta
- * pode aumentar penalidade quando erra
- * Técnicas:

previsão estática - segue sempre mesmo caminho previsão dinâmica - depende do comportamento do programa

HEPR Dinf RCC

Arquitetura II — segmentação 2007-1

Previsão de desvios

• Previsão Estática

⊳ prevê sempre não-tomado MIPS, SPEC95, $\approx 40\%$ MIPS, SPEC95, $\approx 60\%$

⊳ prevê sempre tomado

⊳ prevê que para-trás é tomado

> previsão associada a opcodes

delayed branches

• Previsão dinâmica

⊳ será visto mais tarde

Desvios Atrasados

Sempre executa a próxima instrução
 OK com pipelines simples

i: $PC \leftarrow PC+8 \text{ if } r1 == 0$

i+1: r2 ← r8 - r9 sempre é executada

. . .

i+8: $r3 \leftarrow r4 + r5$

- É necessário que branch delay slot seja preenchido
- ⊳ preenche com instrução antes do desvio (i-1)
- ▶ preenche com instrução de destino
 - * SE é seguro executar instrução de destino
 - * ajuda somente no caso do desvio tomado
- ⊳ preenche com instrução após o desvio (i+2)
 - * SE é seguro executar instrução após o desvio
 - * ajuda somente no caso do desvio não-tomado

HEPR Dist RCC 12

Arquitetura II — segmentação

2007-1

Desvios Atrasados, Loads Atrasados

Arquiteturas que expõem riscos do pipeline ao software, como nos casos de *delayed load*s e *branch delay slot*s resultam em programas com número significativo de nops inseridos pelo compilador.

Estas soluções misturam arquitetura (ISA) com implementação (pipeline linear)

Isso é uma má idéia©

HEPR Dist RCC

Arquitetura II — desvios e excessões

2007-1

resumo - Riscos de Controle

- Riscos de controle tem efeito devastador no desempenho
- Mecanismos simples de previsão de desvios:
- ★ segura fluxo até decidir direção do desvio (penalidade 1-2 ciclos)
- ★ prediz que desvio é sempre não-tomado (penalidade 1 ciclo, ≈40% das vezes)
- ★ prediz que desvio é sempre tomado (penalidade 1 ciclo, ≈60% das vezes)
- desvio atrasado "esconde" penalidade quando desvio não-tomado
- * ocupar o espaço pode ser difícil
- \star mistura arquitetura com implementação

HEPR Dinf RCC 15(

revisão - Riscos

• modelo seqüencial de execução:

cada instrução é completamente executada e altera o estado do processador antes do início da próxima instrução

- Riscos são potenciais violações de dependências no programa
- implementação deve garantir que dependências não são violadas
- Resolução de riscos:
- * estática: compilador/programador garante corretude
- * dinâmica: hardware verifica conflitos em tempo de execução
- Inter-travamento no pipeline

interlocks

2007-1

- * mecanismo para resolução dinâmica de riscos
- * deve detectar e resolver dependências em tempo de execução

HEPR Dinf RCC

Arquitetura II — desvios e excessões

2007-1

mais revisão - Riscos

• Risco estrutural

structural hazards

- Risco de dados

data hazards

- duas instruções diferentes usam mesmo local de armazenamento
 - aparência de execução na ordem següencial correta
- Risco de controle

control (flow) hazards

 □ uma instrução determina quais instruções serão executadas a seguir (desvios, saltos, funções) pior impacto no desempenho

HEPR Dinf RCC 15

Arquitetura II — desvios e excessões

2007-1

Dependências de Controle

- Desempenho dos processadores é limitado pelos desvios
- Maioria das dependências de dados podem ser eliminadas pelo compilador + adiantamento;
- Dependências de controle não podem ser eliminadas por causa do comportamento dinâmico dos programas...
- Regra dos 85/60: medições indicam que aproximadamente
- * 85% dos desvios "para trás" são tomados
- * 60% dos desvios "para frente" são tomados
- solução: adivinhar destino do desvio para antecipar a busca da próxima instrução porque fazer alguma coisa é melhor que não fazer nada

Riscos de controle e previsão de desvios

Previsão Dinâmica de Desvios

- Passado recente pode auxiliar a prever futuro próximo: usar a história dinâmica dos desvios para prever sua direção
- * se desvio será tomado ou não-tomado
- * qual o endereço de destino
- porque é dinâmico, previsor pode se ajustar a
- * fases distintas da execução do programa
- * comportamento de desvios específicos

HEPR Dist RCC 154

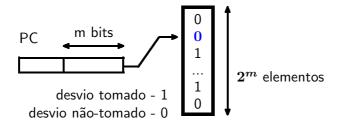
Arquitetura II — desvios e excessões

2007-1

Previsão Dinâmica de Desvios

Tabela com $\mathbf{2}^m$ bits no estágio de busca

- acessa tabela com m bits do PC
- * tabela contém 1 se desvio foi tomado na última execução
- st tabela contém $oldsymbol{0}$ se desvio não foi tomado na última execução
- prevê que nesta execução tomará mesmo caminho que na anterior
- atualiza o bit quando errar a previsão



HEPR Dist RCC 155

Arquitetura II — desvios e excessões

2007-1

Previsão Dinâmica de Desvios

Exemplo: loop executa três voltas e então prossegue:

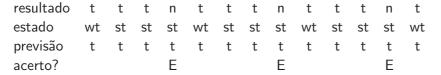
Na previsão com um bit, ocorrem duas previsões erradas a cada mudança c.r.a história recente

resultado t t t n t t t n t t t n t t previsão t t t T N t t T N t t T N acerto?

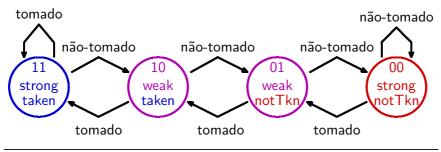
Com este método são duas as previsões erradas a cada vez que seqüência muda...

Previsores melhores: contadores de 2 bits

Contador de dois bits implementa histerese e melhora previsão:



Previsões erradas caem para a metade.



HEPR Dist RCC 15

Arquitetura II — desvios e excessões

2007-1

Previsores melhores: correlação de previsões

Desvios em instruções distintas podem ser relacionados:

if
$$(a < 0) a = 0$$
;
if $(b > 0) b = 0$;
if $(a != b) \{ ... \}$

Se as duas primeiras condições são verdadeiras, a terceira será falsa; portanto a idéia é:

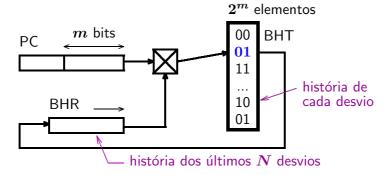
- manter história de todos os desvios recentes
 história aproxima caminho seguido pelo programa
- Branch History Register (BHR) é um registrador de deslocamento que mantém resultado dos últimos N desvios
- PC e BHR são combinados para acessar tabela de previsão
- taxas de acerto ficam entre 80% e 90%

HEPRING 153

Arquitetura II — desvios e excessões

2007-1

Previsores melhores: correlação de previsões



Apontador para tabela de previsão é uma combinação dos conteúdos do PC e do *Branch History Register* soma, sobreposição, ou concatenação

Gshare: um BHR de m-bits, ou-exclusivo do BHR com PC

Previsores ainda melhores: mais de um previsor

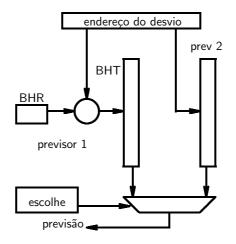
Desvios tem comportamentos distintos
→ usar previsores diferentes em desvios diferentes.

escolha:

se previsor escolhido errou e outro acertou inverte escolha

previsores híbridos, competitivos (tournament)

Gshare + contador de 2 bits



HEPR Dist RCC 160

Arquitetura II — desvios e excessões

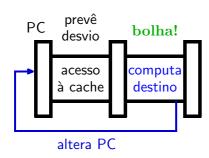
2007-1

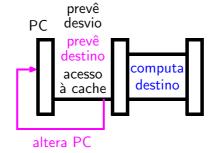
Previsão de endereço de destino

Tabela de endereços de destino (Branch Target Buffer [BTB])

permite prever endereço de destino

se acerta previsão, evita inserção de bolha no pipeline





HEPR Dist RCC

Arquitetura II — desvios e excessões

2007-1

Tabela de endereços de destino

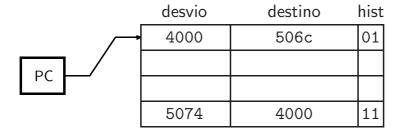
A tabela de endereços de destino é uma cache especial:

- índice e etiqueta são comparados com endereço da instrução de desvio
- conteúdo do BTB é o endereço de destino (cache de destinos)
- bits de previsão mantidos junto com endereço de destino
- geralmente associado à cache de instruções

Branch Target Cache:

- * armazena instrução de destino junto com endereço
- * permite eliminar o desvio do fluxo de instruções branch folding

Tabela de endereços de destino



HEPR Dist RCC 163

Arquitetura II — desvios e excessões

2007-1

Tabela de endereços de destino

Administração da tabela:

- armazena desvios não-tomados?
- o que fazer quando inicia execução (tabela vazia)?
- o que fazer com conflitos na tabela (tabela cheia)?
- o que fazer quando errar a previsão?
- quanto custa (CPI) errar uma previsão?

HEPR Dinf RCC

Arquitetura II — desvios e excessões

2007-1

Previsão de chamada e retorno de funções

- Destino de saltos indiretos para retorno de função (jr)
 é muito difícil de prever
- * uma função pode ter vários endereços de retorno estáticos e/ou dinâmicos
- Solução: implementar uma pilha de retorno de funções no estágio de busca do processador
- * salva endereço de retorno na pilha instrução jal
- * desempilha ao retornar e prevê destino instrução jr r31
- O quê fazer com saltos indiretos que ∉ à funções? instrução jr

Resumo de Previsão de Desvios

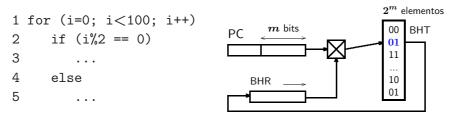
- Previsão é útil em pipelines simples
- é imprescindível em processadores super-escalares
- processadores modernos usam
- * previsores sofisticados
- * tabela de endereco de destino
- * pilha de retorno de funções

HEPR Dist RCC 166

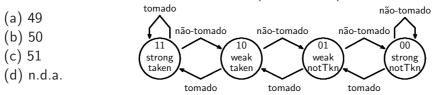
Arquitetura II — desvios e excessões

2007-1

Pergunta:



Quantos erros de previsão na linha **2** para previsor BHR+BHT? BHT iniciaizado com 00 (não-tomado)



HEPR Dist RCC 16:

Arquitetura II — desvios e excessões 2007-1

Resposta:

Estados: 11-St 10-wt 01-wn 00-Sn

BHR com 8 bits, insere resultado na posição menos-significativa BHT com 256 posições, indexada com 8 bits do PC

PC = 0x200 (para simplificar), operação do PC com BHR é XOR Estruturas inicializadas com 0 em todas as posições

1 for(i=0; i<100; i++)
2 if (i%2 == 0)
3 ...
4 else
5 ...
quantos erros de previsão na linha 2?
OU
quantas voltas do laço para treinar
o previsor?

Resposta: (cont)

```
pos BHT nova novoBHR
i dsv
      BHR
                                         acerto?
  nt
      0000.0000 0 00 -> 00 0000.0000 0
0
                                          sim
1
  t
      0000.0000 0 00 -> 01 0000.0001 1
                                          não
2
  nt 0000.0001 1 00 -> 00 0000.0010 2
                                          sim
3
      0000.0010 2 00 -> 01 0000.0101 5
  t
                                          não
4
  nt 0000.0101 5 00 -> 00 0000.1010 a
                                           sim
5
       0000.1010 a 00 -> 01 0001.0101 15
  t
                                          não
  nt 0001.0101 15 00 -> 00 0010.1010 2a
6
7
       0010.1010 2a 00 -> 01 0101.0101 55
                                          não
  nt 0101.0101 55 00 -> 00 1010.1010 aa
8
                                          sim
      1010.1010 aa 00 -> 01 0101.0101 55
9
                                          não - 5 erros
10 nt 0101.0101 55 00 -> 00 1010.1010 aa
                                          sim
      1010.1010 aa 01 -> 10 0101.0101 55
11 t
                                          sim <--acerto
```

HEPR Disf RCC 16

Arquitetura II — desvios e excessões

2007-1

Interrompemos a nossa programação...

HEPR Dist RCC 117

Arquitetura II — desvios e excessões

2007-1

Interrupções e Excessões

- Interrupções
- * geralmente com causa externa ao processador (assíncrona)
- * não são relacionadas a uma instrução específica
- * Exemplos: interrupção de dispositivo; falta de energia
- Excessões
- * relacionadas à execução de uma instrução específica (síncrona)
- * Exemplos: overflow, falta de página
- Traps
- * relacionadas à execução de uma instrução específica
- * rotina (hw+sw) de tratamento de excessões

2007-1

Eventos Síncronos vs Assíncronos

- Síncrono evento relacionado com a seqüência de instruções ocorre durante a execução de uma instrução
- * deve parar instrução que está executando correntemente
- * falta de página durante instrução load ou store
- * excessões aritméticas
- ★ chamadas de sistema (instruções trap ou syscall)
- Assíncrono evento sem relação com seqüência de instruções causada por um evento externo
- ★ desnecessário interromper instruções que estão executando
- ★ interrupções por dispositivos de E/S
- ★ interrupções por eventos catastróficos (erro na memória)
- Semi-síncrono
- ★ evento externo que pode interromper seqüência de instruções para garantir serviços (alta disponibilidade)

HEPR Dinf RCC 17

Arquitetura II — desvios e excessões 2007-1

Interrupções e Excessões – arquitetura

- Registrador de Interrupção
- * vetor de bits indicando quais interrupções/excessões ocorreram
- Registrador de máscara

Arquitetura II - desvios e excessões

- * vetor de bits indica quais inters/excessões estão desabilitadas
- * escrever no registrador de máscara é instrução privilegiada
- * alguns bits podem ser atualizados em modo usuário (overflow)
- * algumas interrupções/excessões não são mascaráveis

HEPR Dist RCC 17

Interrupções/Excessões Precisas

Uma interrupção ou excessão é considerada **precisa** se existe uma única instrução (ou ponto de interrupção) tal que todas as instruções anteriores àquela tenham alterado o estado, e nenhuma instrução após (e incluindo) aquela tenham modificado o estado.

Isso significa que a execução pode ser re-iniciada a partir do ponto de interrupção e resultados corretos serão produzidos.

Ver artigo de Smith & Pleszkun para soluções. IEEE-TC 37:5 1988

Interrupções Precisas

- Interrupções Precisas facilitam MUITO o trabalho do SO na continuação do programa, ou na depuração
- Excessão/trap
- ⊳ todas instruções antes da causadora completam
- > nenhuma das instruções após a causadora completam
- ▶ instrução causadora ou completou ou não iniciou
- ▶ PC aponta instrução causadora
- Interrupção
- ⊳ Mesmo que excessão, mas ∄ instrução causadora
- ▶ Deve parecer que ocorreu entre duas instruções
- ▶ PC aponta para instrução que seria executada

HEPR Disf RCC 177

Arquitetura II — desvios e excessões

2007-1

Tratamento de interrupções/traps

Quando ocorre interrupção/trap:

- efetua salto para rotina do SO
- * vetor de tratadores, em geral amarrado no cj de instruções
- computa endereço de retorno
- salva informação de estado essencial
- * PC
- * CCs (condition codes)
- * PSW (processor status word)
- ◆ troca modo de execução: usuário ⇒ sistema

HEPR DIM RCC 177

Arquitetura II — desvios e excessões

2007-1

Implementação de interrupções/traps

- Precisão é importante
- Interrupções e excessões simultâneas
- ▶ busca falta de página, referência desalinhada, protection fault
- ▶ decod instrução ilegal ou privilegiada
- ▶ exec excessões de aritmética (overflow, divisão por zero)
- ⊳ mem falta de página, referência desalinhada, protection fault
- ▶ res nenhuma

Falta de Página

 \mathbf{B} \mathbf{D} $\mathbf{E}\mathbf{x}$ \mathbf{M} \mathbf{R} falta de página (mem) \mathbf{B} \mathbf{D} $\mathbf{E}\mathbf{x}$ \mathbf{M} falta de página (busca) \mathbf{B} \mathbf{D} $\mathbf{E}\mathbf{x}$ \mathbf{M} anulada \mathbf{B} D Ex M R anulada \mathbf{B} Ex M R 🗕 anulada salta para tratador

início do tratador \longrightarrow B D Ex M R

- Instrução anterior completou res//mem
- anula todas instruções subseqüentes
- impede que instrução causadora altere estado (se for load)

Arquitetura II — desvios e excessões

LIEPR Dinf BCC

2007-1

Excessão de Aritmética

Ex M R \mathbf{B} \mathbf{D} \mathbf{B} D $\mathbf{E}\mathbf{x} \mathbf{M} \mathbf{R}$ excessão (exec) \mathbf{B} Ex M R \mathbf{B} \mathbf{D} $\mathbf{E}\mathbf{x} \mathbf{M} \mathbf{R}$ anulada \mathbf{B} \mathbf{D} $\mathbf{E}\mathbf{x} \ \mathbf{M}$ anulada salta para tratador

início do tratador B D Ex M R

- Instruções anteriores podem completar
- anula todas instruções subseqüentes

HEPR Dist RCC 17

Arquitetura II — desvios e excessões

2007-1

Excessões Múltiplas

falta de página instrução ilegal falta de página overflow

usca	decod	exec	mem	result		_	
	busca	decod	exec	mem	result		
•		busca	decod	exec	mem	result	
			busca	decod	exec	mem	result

- Pipeline pode resolver
- * estado c.r.a excessões passa pelos estágios junto com instrução
- * mantém valor do PC junto a cada instrução
- * não trata excessão até instrução chegar a result
- Interrupção é detectada na busca e instrução NOP é inserida
- Quando instrução chega a result
- * salva PC \longrightarrow EPC; endereço do vetor de interrupções \longrightarrow PC
- * transforma instruções nos outros estágios em NOPs

Complicações no Conjunto de Instruções

- Estado espalhado em vários recursos (CCs)
- Atualização antes de **mem/res** (auto-incremento)
- Instruções que completam fora de ordem
- * outros mecanismos serão estudados adiante
- Interrupções ocorrem a qualquer momento
- O caso frequente não é o único/mais importante
- * Desempenho não é o parâmetro mais crítico
- * Os casos raros devem funcionar corretamente

LIEPR Dinf BCC

Arquitetura II - paralelismo no nível de instrução

2007-1

Resumo - previsão de desvios e excessões

- Previsão de desvios é imprescindível
- * contador de 2 bits com histerese para prever cada desvio
- * tabela(s) com 2^m previsores + registrador de deslocamento
- * registrador de história dos desvios

branch history register

* tabela com endereços de destino

branch target buffer

- * pilha de retorno de funções
- Interrupções e excessões complicam muito projeto de pipelines
- * precisão é importante mas implementação é cara (hardware)
- * drenar o pipeline custa caro...
- * detecta evento cedo e trata-o tarde em result

HEDD DINE DCC

Arquitetura II — paralelismo no nível de instrução

Revisão – previsão de desvios e excessões

- Previsão de desvios é imprescindível pelo efeito no desempenho
- * contador de 2 bits com histerese para prever cada desvio
- * tabela(s) com 2^m previsores + registrador de deslocamento
- * registrador de história dos desvios
- * tabela com endereços de destino

branch history register

branch target buffer

- * pilha de retorno de funções
- Interrupções e excessões complicam muito projeto de pipelines
- * precisão é importante mas implementação é cara (hardware)
- * drenar os segmentos custa caro...
- * detecta evento cedo e trata-o tarde em result

LIEPR Dinf RCC

Superpipelining

tempo de CPU = núm instr x CPI x período do relógio

Se aumentar velocidade do relógio (e reduzir período),
pode haver ganho de desempenho

→ implementar estágios com menor latência e em maior número
→ reduz período mas aumenta CPI

Com um pouco de sorte (uh?) relógio compensa CPI

Exemplos:

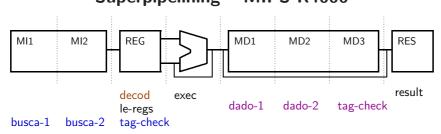
MIPS 4000 usa dois estágios na busca e três em MEM/cache Pentium III tem pipeline com 10 estágios e Pentium IV tem pipeline com 20 estágios, com dois ciclos só para transmitir bits através da CPU

HEPR Dist RCC 15

Arquitetura II — paralelismo no nível de instrução

2007-1

Superpipelining – MIPS R4000



busca dura **2** ciclos – verifica acerto na L1-l enquanto decodifica desvios tomados causam três bolhas

mem dura 3 ciclos por causa da verificação de acerto
 → não pode escrever valor até estar certo do acerto na L1-D
 só pode usar valor de load depois de dois ciclos
 usa valor do load antes de verificar etiqueta?

HEPRING 181

Arquitetura II — paralelismo no nível de instrução

2007-1

Superpipelining

Ganhos:

freqüência do relógio mais alta moda até há pouco permite ligar CPU rápida à memória lenta interf c/ mem segmentada

Perdas:

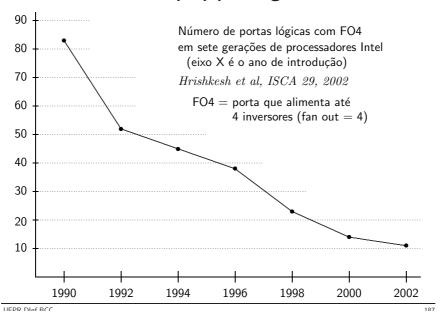
CPI mais alto
penalidade maior nos desvios
penalidade maior nas faltas nas caches (mais ciclos)

penalidade maior no uso do valor de load penalidade maior nas excessões

maior complexidade:

circuitos de adiantamento, bloqueios...

Superpipelining



Arquitetura II — paralelismo no nível de instrução

2007-1

Desempenho de Pipelines

$$\mathsf{ganho} \ = \ \frac{\mathsf{CPI} \ \mathsf{sem} \ \mathsf{pipeline}}{\mathsf{CPI} \ \mathsf{com} \ \mathsf{pipeline}} \times \frac{\mathsf{ciclo} \ \mathsf{sem} \ \mathsf{pipeline}}{\mathsf{ciclo} \ \mathsf{com} \ \mathsf{pipeline}}$$

Desempenho ótimo

ignorando latches

ganho = número de estágios

Causa das Perdas:

• dependências estruturais

escalonamento de instruções resolve

• dependências de dados

escalonamento e adiantamento resolvem

• dependências de controle

previsão resolve pprox 85% dos casos

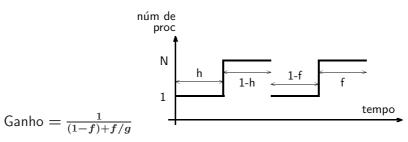
• Lei de Amdahl: ganho é limitado pelo pior componente

HEPR Dist RCC

Arquitetura II — paralelismo no nível de instrução

2007-1

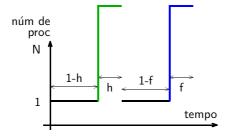
Lei de Amdahl



- h fração de tempo em código serial
- ullet 1-h fração de tempo que pode ser paralelizada (N processadores)
- f fração que pode ser vetorizada
- $ullet \ g$ ganho durante f

Lei de Amdahl

$$\lim_{g\to\infty}\frac{1}{(1-f)+f/g}=\frac{1}{(1-f)}$$



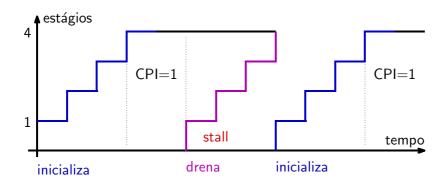
Desempenho é limitado pela parte que **não pode** ser paralelizada (h) ou vetorizada (f)

HEPR Dist RCC 100

Arquitetura II — paralelismo no nível de instrução

2007-1

Desempenho de Pipelines I

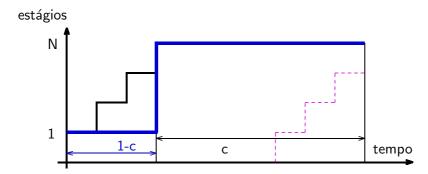


LIEDR Dinf RCC 10

Arquitetura II — paralelismo no nível de instrução

2007-1

Desempenho de Pipelines II



- ullet c fração de tempo com pipeline cheia
- ullet 1-c fração de tempo em que ocorrem bolhas
- ullet quando c é um pouquinho abaixo de 100%, a queda no desempenho é enorme!

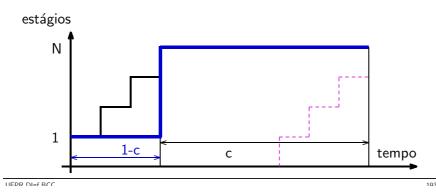
ightarrow a fração 1-c deve ser minimizada

LIFPR Dinf RCC

Desempenho de Pipelines III

Se executar mais de uma instrução por ciclo, a fração 1-c pode ser reduzida pelo fator de escalaridade ${\it S}$

$$\mathsf{Ganho} = \frac{1}{\frac{(1-c)}{S} + \frac{c}{g}}$$



Arquitetura II — paralelismo no nível de instrução

2007-1

Pipeline para obter CPI pprox 1

- Gargalo de Flynn
- * emissão de uma instrução por ciclo limita CPI=IPC=1
- * riscos + overhead \rightarrow CPI ≥ 1 IPC ≤ 1
- * ganhos cada vez menores com superpipelining
- solução: emitir mais de uma instrução por ciclo

	1	2	3	4	5	6	7	8	9
inst0	В	D	Е	М	R				
inst1	В	D	Ε	M	R				
inst0 inst1 inst2		В	D	Ε	M	R			
inst3				Ε					

- * paralelismo no nível de instrução (PNI) instruction-level paralelism (ILP) Fischer81
- * primeiro super-escalar: IBM América ightarrow RS/6000 ightarrow POWER1

HEPR Dinf RCC

104

Arquitetura II — paralelismo no nível de instrução

2007-1

Processadores Escalares e Vetoriais

- Processador escalar opera em uma palavra de memória: escalar
- * processadores escalares são todos aqueles vistos até agora
- processador vetorial opera em um grupo de palavras: vetor
- * The CRAY-1 Computer System, R M Russel, CACM 21:1, jan78
- * certas instruções usam vetores como operandos
- * registradores contém vetores (64 elementos/registrador)
- * acessos à memória para carregar ou escrever vetores inteiros
- * unidades funcionais são pipelines para elementos de vetores
- * ver Apêndice G em www.mkp.com/CA3

Paralelismo no nível de instrução

Qual a diferença entre dependência e risco?

HEPR Dinf RCC

Arquitetura II — paralelismo no nível de instrução

2007-1

Paralelismo no nível de instrução (PNI)

PNI é uma propriedade do software (e não do hardware)

quanto paralelismo existe entre as instruções de um programa? depende do software hw pode ou não explorá-lo

Inúmeras maneiras de explorar PNI:

- ⊳ pipelining (Apêndice A)
- ⊳ superescalar (Cap 3)
- ⊳ execução fora de ordem (escalonamento dinâmico Cap 3)
- ⊳ escalonamento pelo compilador (escalonamento estático Cap 4)

HEPR Dist RCC

Arquitetura II — paralelismo no nível de instrução

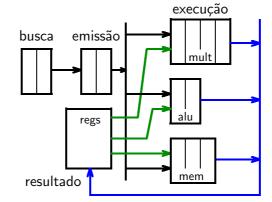
2007-1

Pipelines não-lineares

Quatro fases: busca;

decod+emissão ;
execução ;

resultado



Estágio(s) de emissão envia(m) instruções para unidade funcional latência da execução depende da unidade funcional

$$\mathcal{T}_{ ext{mult}} > \mathcal{T}_{ ext{mem}} > \mathcal{T}_{ ext{alu}}$$

Superescalar: implementação base

- superescalar com escalonamento estático, em-ordem

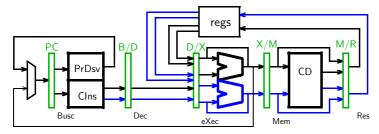
- ⊳ exemplos: MIPS 8000, Sun UltraSPARC, Alpha 21264

HEPR Dist RCC 100

Arquitetura II — paralelismo no nível de instrução

2007-1

Pipeline de 5 estágios com emissão dupla



- O que é necessário para
- * buscar duas instruções por ciclo?
- * decodificar duas instruções por ciclo?
- * executar duas instruções de ALU no mesmo ciclo?
- * acessar a cache de dados duas vezes no mesmo ciclo?
- * escrever dois registradores no mesmo ciclo?
- e se forem 4 ou 8 instruções num ciclo?

HEPR DInf RCC

Arquitetura II — paralelismo no nível de instrução

2007-1

Busca com largura N

- O que é necessário na busca de N instruções em um ciclo?
- se as instruções são seqüenciais
- * e no mesmo bloco (\mathbf{k}) da cache \rightarrow nada
- * e em blocos diferentes → cache intercalada + rede de combinação

		0	1	2	3	4	5	6	7			
	bloco k	Х	у	Z	W					\rightarrow	xyzw	
,	bloco i			Χ	У							rede combina
	bloco j					Z	W			\longrightarrow	xyzw	os 2 blocos

- se as instruções não são seqüenciais
- * dois acessos em série: $acesso1 \rightarrow prevê_destino \rightarrow acesso2$
- desvios no meio de um bloco: fácil se desvios não-tomados (NT)
- * acesso serial + previsão em paralelo
- * se previsão é tomado (T) → descarta instrs após desvio

Decodificação com largura N

- O que é necessário para decodificar N instruções em um ciclo?
- decodificar as instruções
- * fácil se instruções de tamanho fixo (múltiplos decodificadores ||s)
- * difícil, porém possível, se tamanho variável → x86
- ler operandos dos registradores
- * 2N portas de leitura no bloco de registradores
- * na verdade, menos que 2N porque muitos valores são adiantados
- como fica a lógica de controle dos atrasos (stalls)?

HEPR Dist RCC 20

Arquitetura II - paralelismo no nível de instrução

2007-1

Decodificação com largura N

- lógica de controle de atrasos num pipeline simples:
- * rs1(D)==rd(D/X) || rs1(D)==rd(X/M) || rs1(D)==rd(M/R)
- * mesmo para rs2
- * com adiant completo: rs1(D)==rd(D/X) && opc(D/X)==load
- dobrando a largura de emissão, quadruplica lógica de atrasos
- ⋆ não são só 2 instruções em Decod, mas 2 instr em todos estágios
- \star rs1(D1)==rd(D/X1) && opc(D/X1)==load
- ★ rs1(D1)==rd(D/X2) && opc(D/X2)==load
- ★ repetir para rs1(D2), rs2(D1), rs2(D2)
- ★ testar dependência da segunda instr na primeira: rs1(D2)==rd(D1)
- num pipeline de largura N, circuito de atrasos cresce com N^2 mesmo vale para adiantamento, só que pior...
- ullet |lógica de controle de atrasos $|\propto N^2$

HEPRING 201

Arquitetura II — paralelismo no nível de instrução

2007-1

Execução com largura N

- O que é necessário para executar N instruções em um ciclo?
- múltiplas unidades funcionais. N de cada tipo?
- ⋆ N ULAs? Pode ser, ULAs são pequenas
- * N divisores de ponto flutuante? Não, circuito é enorme e divPF é infreqüente
- tipicamente, usa combinação proporcional ao uso
- ★ RS/6000: 1 ULA/endereços/desvios + 1 PF
- ★ Pentium: 1 ULA complexa + 1 ULA simples
- \star PentiumII: 1 ULA/PF + 1 ULA + 1 load + 1 store + 1 desvios
- * Alpha 21164: 1 ULA/PF/desvios + 2 ULA + 1 load/store
- ullet circuito de adiantamento $\propto N^2$
- * lógica de controle é pequena porque variáveis tem 5 bits (regs)
- * circuitos de dados é gigantesco (32 ou 64 bits por caminho)
- * layout da fiação é infernal, MUXes são enormes e lentos
- * menos horrível se agrupar unidades funcionais em clusters

Interface de memória com largura N

- O que é necessário para acessar memória 2 vezes no mesmo ciclo?
- cache de dados com múltiplos bancos (detalhes mais tarde)
 necessita lógica de detecção de conflitos (2 refs ao mesmo banco)
 necessita lógica de detecção de riscos RAW
- aproximadamente 20% das instruções são loads e 15% stores
- ⊳ para largura N, são necessárias
 0,2N portas de leitura e
 0,15N portas de escrita na memória

HEPR Dist RCC 205

Arquitetura II — paralelismo no nível de instrução

2007-1

Gravação de resultados com largura N

- O que é necessário para escrever 2 registradores no mesmo ciclo?
- apenas mais uma porta de escrita no bloco de registradores
 tudo o que deveria já foi feito nos estágios anteriores...
- MAS o tratamento de excessões é ainda mais complicado deve-se usar buffer de re-ordenação à lá Smith&Plezkun

HEPR Dist RCC

rquitetura II — paralelismo no nível de instrução

2007-1

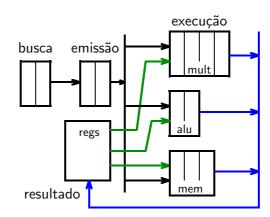
Modelo de pipeline

Quatro fases:

busca ;
decod+emissão ;
execução ;
resultado

cada fase pode ser um pipeline com mais de um estágio

nos interessa a **emissão** (*issue*)



Ordem de execução dos programas

• modelo seqüencial de execução:

cada instrução é completamente executada e altera o estado do processador antes do início da próxima instrução

- → modelo habitual de raciocinar sobre execução de programas
- modelos com paralelismo:

paralelismo explícito \longrightarrow multiprocessadores, threads, vetores paralelismo implícito \longrightarrow hardware esconde execução paralela

HEPR Dist RCC 203

Arquitetura II - paralelismo no nível de instrução

2007-

Ordem

- Emissão de instruções para execução
- * emissão em ordem o que vimos até agora
- emissão fora-de-ordem processador examina várias instruções para decidir quais podem seguir para execução
- ordem de execução
- * execução em ordem o que vimos até agora
- * execução fora-de-ordem se operandos e unidade funcional disponíveis então executa
- ordem de escrever resultado (writeback)
- * resultados gravados em ordem o que vimos até agora
- * resultados gravados fora-de-ordem grava resultado quando não viola nenhuma dependência

HEPR Dinf RCC 20

Arquitetura II — paralelismo no nível de instrução

2007-1

Ordem de emisão

- Emissão em ordem
 o que vimos até agora
 emissão pára se ∃ qualquer tipo de dependência ou risco
 risco estrutural, de dados, ou de controle
- emissão fora-de-ordem processador examina várias instruções para decidir quais podem seguir para execução lookahead
- * se instrução bloqueia por causa de conflito por recurso processador segue buscando até encontrar instrução independente
- * devem ∃r recursos para executar instruções independentes
- * ordem sequencial de execução deve ser obedecida mesmo se instruções executam e completam fora-de-ordem

Ordem de execução

• Execução em ordem

- o que vimos até agora
- * conflitos por recursos seguram execução de instruções + novas
- * idem para dependências de dados ou de controle
- execução fora-de-ordem

explora paralelismo entre instruções

- * se operandos e unidade funcional disponíveis então executa
- * dependências de dados devem ser respeitadas
- * como descobre e resolve riscos RAW, WAR, WAW?

algoritmo do placar algoritmo de Tomasulo Register Update Unit

scoreboarding, Thornton CDC6600 1961 IBM 360/91 1967 Sohi IEEE-TC 39:3 1990

LIEPR Dinf BCC

Arquitetura II - paralelismo no nível de instrução

2007-1

Ordem de gravar resultado

- Resultados gravados em ordem
- o que vimos até agora
- * instrução lenta segura outras que não dependem do seu resultado
- * facilita recuperação de excessões e implementação do SO mas desempenho é ruim
- resultados gravados fora-de-ordem
- * grava resultado quando não viola nenhuma dependência
- * o que acontece se instrução mais-nova grava resultado antes que excessão seja detectada em instrução mais-velha?

Smith & Plezkun IEEE-TC 37:5 1988

HEPR Dinf RCC

2007-1

Arquitetura II — paralelismo no nível de instrução

Emissão em-ordem + resultado em-ordem

pipeline decodifica 2/ciclo executa 2/ciclo escreve 2/ciclo

i1 executa em 2 ciclos i3 e i4 competem por UF i5 depende de result de i4 i5 e i6 competem por UF

decodif ciclo 1 i1 i2 2 i3 i43 **i4** 4 **i4** 5 **i**6 6 **i6** 7 8

execução i1 i2 i1 i3 **i**4 i5 **i**6

resultado i2 i1 i3 **i**4 **i**6 i5

Emissão em-ordem + resultado fora-de-ordem

pipeline i1 executa em 2 ciclos decodifica 2/ciclo i3 e i4 competem por UF executa 3/ciclo i5 depende de result de i4 escreve 2/ciclo i5 e i6 competem por UF

ciclo	dec	odif
1	i1	i2
2	i3	i4
3		i4
4	i 5	i6
5		i6
6		
7		

execução					
i1	i2				
i1		i3			
		i4			
	i5				
	i6				

resultado					
i2					
i1	i3				
i4					
i5					
i6					

TIEPR DINFRCC 21

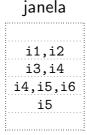
Arquitetura II - paralelismo no nível de instrução

2007

Emissão fora-de-ordem + resultado fora-de-ordem

pipeline i1 executa em 2 ciclos decodifica 2/ciclo i3 e i4 competem por UF executa 3/ciclo i5 depende de result de i4 escreve 2/ciclo i5 e i6 competem por UF

ciclo	dec	odif
1	i1	i2
2	i3	i4
3	i5	i6
4		
5		
6		



execução					
i1	i2				
i1		i3			
	i6	i4			
	i5				

result	resultado					
i2						
i1	i3					
i4	i 6					
i5						

HEPR Dinf RCC 21

Arquitetura II — paralelismo no nível de instrução

2007-1

Processador Super-escalar

Processador possui dois ou mais Circuitos de Dados paralelos

Estágio de decodificação escolhe tuplas de instruções e as despacha de acordo com disponibilidade de unidades funcionais

Inter-travamento entre estágios e unidades funcionais resolve dependências de dados e de controle

Algoritmo do Placar ou Algoritmo de Tomasulo: detalhes adiante resolvem "problema dos gladiadores" \Longrightarrow que gladiador luta contra qual, e quando

Unidade de reordenamento enfilera resultados de acordo com dependências entre resultados e operandos

Processador Superescalar

Dependências de controle são resolvidas com execução especulativa

Instruções nos dois lados do desvio são executadas especulativamente;

→ quando decide, anula efeito das instruções do caminho errado

Registradores fantasma mantém valores da execução especulativa; quando resolve desvio copia de regs fantasma para regs visíveis

TIEPP DIMERCC 21

Arquitetura II — algor de Tomasulo 2007-

Resumo

- Superpipelining
- * aumenta freqüência do relógio & reduz lógica em cada estágio
- * ganho: desempenho ∝ freqüência do relógio
- * perda: CPI aumenta por causa das penalidades (desvios, depend)
- Processadores superescalares
- * recursos replicados para usar paralelismo no nível de instrução
- * algoritmo para garantir dependências de dados e controle algoritmo do placar ou de Tomasulo
- * PNI é característica do SW HW pode ou não suportar PNI>1
- Ordem modelo sequencial deve ser respeitado pelo HW
- * emissão em-ordem VS fora-de-ordem
- * execução em-ordem VS fora de ordem
- * gravar resultados em-ordem VS fora-de-ordem

HEPR Dist RCC 21

Arquitetura II — algor de Tomasulo 2007

Revisão

- Superpipelining
- * aumenta freqüência do relógio & reduz lógica em cada estágio
- * ganho: desempenho ∝ freqüência do relógio
- * perda: CPI aumenta por causa das penalidades (desvios, depend)
- Processadores superescalares
- * recursos replicados para usar paralelismo no nível de instrução
- * algoritmo para garantir dependências de dados e controle algoritmo do placar ou de Tomasulo
- * PNI é característica do SW HW pode ou não suportar PNI>1
- Ordem modelo seqüencial deve ser respeitado pelo HW
- * emissão em-ordem VS fora-de-ordem
- * execução em-ordem VS fora de ordem
- * gravar resultados em-ordem VS fora-de-ordem

Trecho de código comum

```
for (i=0; i<TAM; i++)
                            // Multiplica elmtos de
      A[i] = B[i] * C[i];
                            // B e C, armazena em A
      addi r1,r0,TAM
                            # tamanho dos vetores
      or r2, r0, r0
                            # indice nos vetores
      la r4, endB
                            # ender vetor B
      la r5. endC
                           # ender vetor C
      la r6, endA
                           # ender vetor A
loop: lw r3, 0(r4+r2)
                           # carrega B[i]
      lw r7, 0(r5+r2)
                            # carrega C[i]
      mul r7,r7,r3
                            # r7 ←r7*r3
      addi r1,r1,-1
                           # decr contador
      sw r7, 0(r6+r2)
                            # guarda A[i]
      addi r2,r2,8
                            # acerta indice
      bne r1,r0,loop
                            # terminou?
LIEPR Dinf RCC
```

Arquitetura II - algor de Tomasulo

2007-1

Ordem de execução

```
lw r3, 0(r4+r2)
                   # carrega B[i]
lw r7, 0(r5+r2)
                  # carrega C[i]
mul r7, r7, r3
                   # r7 ←r7*r3
addi r1,r1,-1
                  # decr contador
```

• Problema:

- * dois LWs emitidos para unidade de memória
- * MUL é bloqueado por causa do risco RAW em r3 e r7
- * ADDI pára porque MUL parou ADDI pára por causa da execução em-ordem

HEPR Dinf RCC

Arquitetura II — algor de Tomasulo

2007-1

Ordem de execução – soluções

- Escalonamento estático (sw)
- * compilador re-ordena as instruções
- * compilador pode usar algoritmos poderosos
- * hardware pode ser simples
- * escalonamento não se adapta à execução: faltas nas caches
- Escalonamento dinâmico (hw)
- * hardware re-ordena instruções
- * hw trata eventos desconhecidos em tempo de compilação
- * software é mais portável
- * hardware é (muito) mais complexo
- * é usado em muitos processadores super-escalares

LIEPR Dinf RCC

Riscos em registradores

- Emissão fora-de-ordem não pode violar riscos RAW que são as únicas dependências verdadeiras
- riscos WAR e WAW podem inibir reordenação desnecessariamente dependências falsas:

```
lw r2, 8(r4)
mul r2,r2,r5
add r5,r6,r9
se r5 é modificado fora-de-ordem, MUL pode ler operando
```

HEPR Dist RCC 22

Arquitetura II - algor de Tomasulo

mas não todos eles

2007-1

Riscos em registradores

Compiladores podem evitar alguns riscos WAR e WAW

```
add r2, r8, r4
bne r1,r0, else
add r2, r9, r4
else: add r6, r2, r5
```

- * não há risco WAW se desviar há risco WAW se seguir (em r2)
- * instrução pode ter risco consigo mesma no escalonamento dinâmico de loops

HEPR Dinf RCC

224

Arquitetura II — algor de Tomasulo

2007-1

Re-nomeação de registradores em harwdare

- Considere registradores do CdI como nomes para valores não como locais de armazenagem
- hardware pode re-nomear registradores para evitar riscos WAR e WAW

```
de:
```

```
lw r2, 8(r4)
mul r2,r2,r5
add r5,r6,r9
para:
lw t1, 8(r4)
mul t2,t1,r5
add t3,r6,r9
```

Compilador deve garantir que usos subsequentes de r2 e r5 usem t1 e t3 \longrightarrow mesmo através de desvios

Algoritmo de Tomasulo

- usado no IBM 360/91, em 1967
- * modelo rápido para aplicações científicas
- * unidades de PF segmentadas somador e multiplicador (divisão efetuada no multiplicador)
- escalonamento dinâmico nas unidades de ponto flutuante

HEPR Dief RCC 22

Arquitetura II - algor de Tomasulo

Arquitetura II — algor de Tomasulo

2007-1

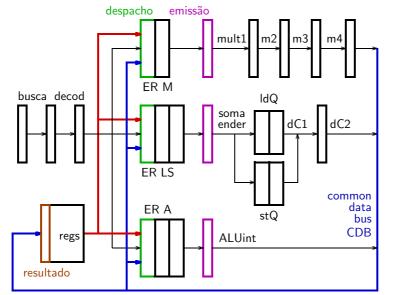
2007-1

Algoritmo de Tomasulo Generalizado

- Melhora paralelismo em processador segmentado
- aplica-se algoritmo em todas as unidades segmentadas
 → trata instruções de memória como as demais
- Usa etiquetas (tags) para identificar valores valores podem ser nomeados por tags e nomes de registradores
- estações de reserva (ER) implementam controle distribuído
- * uma ER por unidade funcional
- * tag identifica resultado da instrução na ER
- Common Data Bus (CDB) difunde todos resultados etiqueta da unidade funcional é transmitida junto com valor

HEPR Dist RCC 22

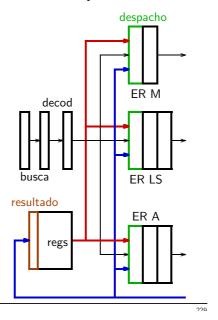
Algoritmo de Tomasulo Generalizado



Algoritmo de Tomasulo - despacho

Despacho:

recebe próx instr de busca/decod há estação de reserva livre? não - bloqueia emissão sim – despacha para emissão: copia regs prontos para ER copia tags dos regs não-prontos para ER



Arquitetura II — algor de To

2007-1

Algoritmo de Tomasulo - emissão

Emissão:

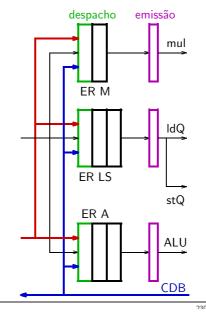
LIEPR Dinf BCC

se operandos estão prontos emite operação para Unid Func senão, espera operando ser transmitido no CDB (tag, valor)

DESVIOS:

nenhuma instrução executa até que desvios anteriores sejam resolvidos

→ excessões precisas



HEPR Dinf RCC

Arquitetura II — algor de Tomasul

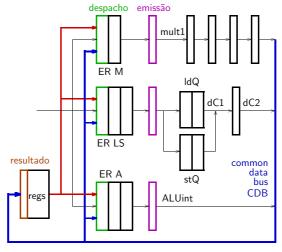
Algoritmo de Tomasulo – finalização

Finalização:

se CDB está livre. então transmite resultado no CDB senão, bloqueia enquanto espera para transmitir

ADIANTAMENTO:

CDB adianta valor para ins dependentes 1 ciclo entre res \rightarrow uso



todas estações de reserva com tags iguais à tag no CDB (resultado) recebem valor

Algoritmo de Tomasulo - EdR

- Estações de Reserva
- * controlam a execução das instruções
- * detecção de riscos distribuída
- * tag de 4 bits identifica cada elemento das estações de reserva
- todos os receptores de valor usam tag para identificar dado que passa pelo CDB
- no despacho
- * ocorre mapeamento entre valores (resultados) e tags
- * isso feito, "nomes" de registradores são desnecessários

HEPR Dinf RCC 23

Arquitetura II — algor de Tomasulo

2007-

Algoritmo de Tomasulo - Estações de Reserva

- Premissa: tag=0 → dado/valor disponível
- campos das estações de reserva
- **Op** opcode (operação)
- Qj, Qk etiquetas das fontes
 - $\mathbf{Q}\mathbf{j}\,|\,\mathbf{Q}\mathbf{k}{=}0\longrightarrow\mathbf{V}\mathbf{j}\,|\,\mathbf{V}\mathbf{k}$ contém valor pronto
- Vj, Vk valores dos operandos (fontes)
- **B** (busy) em uso
- A campo de endereço nas EdR da unidade de memória
- campos do bloco de registradores
- Qi etiqueta da unidade funcional que produzirá valor
 - **Qi**=0 → registrador contém valor pronto
- Vi valor/conteúdo do registrador

HEPR Dinf RCC 23

Arquitetura II — algor de Tomasulo

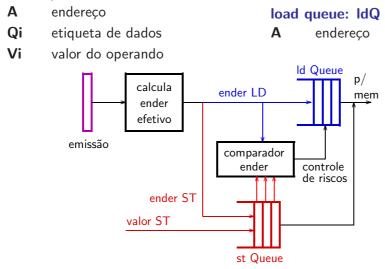
2007-1

Loads e Stores

- Loads e Stores tratados como demais instruções
- * usam unidade de memória
- * endereços passam através de uma load queue
- stores: cálculo do ender efetivo é separado do acesso à memória
- * cálculo do ender efetivo despachado para ER
- * acesso à memória despachado para fila de escrita
- * unidade de memória envia endereços para fila de escrita
- * valor é associado ao endereço pela tag na fila de escrita
- resolução de riscos de memória
- * emissão de loads e stores ocorre em-ordem
- * compara ender dos loads com ender anteriores na fila de escrita bloqueia se há endereços repetidos RAW? loads podem ultrapassar stores se endereços diferem

Loads e Stores

store queue: stQ



HEPR Dist RCC 235

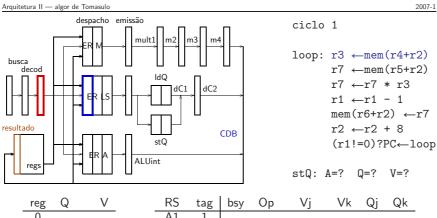
Arquitetura II — algor de Tomasulo

2007-1

Exemplo

addi r1,r0,TAM # tamanho dos vetores or r2,r0,r0 # indice nos vetores la r4, endB # ender vetor B 1000 la r5, endC # ender vetor C 2000 la r6, endA # ender vetor A 3000 loop: lw r3, 0(r4+r2)# carrega B[i] lw r7, 0(r5+r2)# carrega C[i] # r7 ←r7*r3 mul r7,r7,r3 addi r1,r1,-1 # decr contador sw r7, 0(r6+r2)# guarda A[i] addi r2,r2,8 # acerta indice bne r1,r0,loop # terminou?

HEPR Dist RCC 226

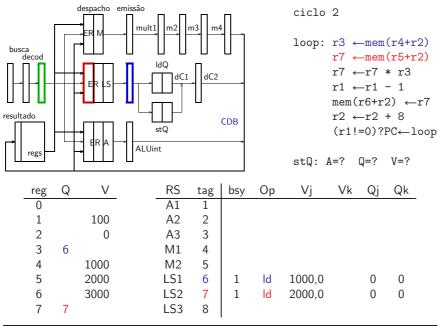


reg	5 (Q	V	R5	tag	bsy	Оp	٧J	VK	QJ	Qĸ
0				A1	1						
1			100	A2	2						
2			0	А3	3						
3		6		M1	4						
4			1000	M2	5						
5			2000	LS1	6	1	ld	1000,0		0	0
6			3000	LS2	7						
7				LS3	8						
						•					

LIEPR Dinf RCC 237

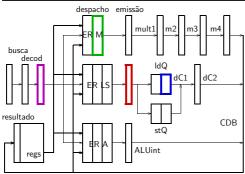
2007-1

2007-1



HERR Dist RCC 23

Arquitetura II — algor de Tomasulo



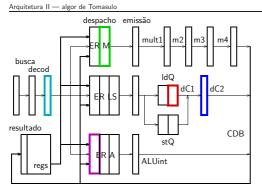
loop: r3 ←mem(r4+r2)
r7 ←mem(r5+r2)
r7 ←r7 * r3
r1 ←r1 - 1
mem(r6+r2) ←r7
r2 ←r2 + 8
(r1!=0)?PC←loop

ciclo 3

stQ: A=? Q=? V=?

reg	Q	V	RS	tag	bsy	Ор	Vj	Vk	Qj	Qk
0			A1	1						
1		100	A2	2						
2		0	А3	3						
3	6		M1	4	1	mul			6	7
4		1000	M2	5						
5		2000	LS1	6	1	ld	1000 + 0		0	0
6		3000	LS2	7	1	ld	2000,0		0	0
7	4		LS3	8						

HEPR Dief RCC 23



loop: r3 ←mem(r4+r2)

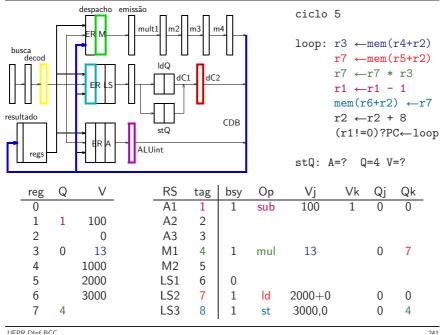
ciclo 4

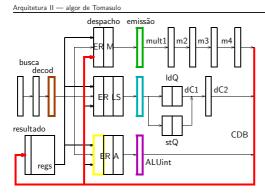
r7 ←mem(r5+r2) r7 ←r7 * r3 r1 ←r1 - 1 mem(r6+r2) ←r7 r2 ←r2 + 8 (r1!=0)?PC←loop

stQ: A=? Q=? V=?

reg	Q	V		RS	tag	bsy	Ор	Vj	Vk	Qj	Qk
0			-	A1	1	1	sub	100	1	0	0
1	1	100		A2	2						
2		0		А3	3						
3	6			M1	4	1	mul			6	7
4		1000		M2	5						
5		2000		LS1	6	1	ld	1000 + 0		0	0
6		3000		LS2	7	1	ld	2000+0		0	0
7	4			LS3	8						

LIEPR Dinf RCC 240





ciclo 6

2007-1

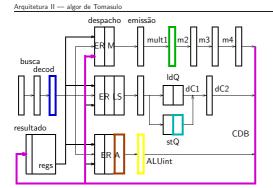
2007-1

loop: r3 ←mem(r4+r2) r7 ←mem(r5+r2) r7 ←r7 * r3 r1 ←r1 - 1 mem(r6+r2) ←r7 r2 ←r2 + 8 (r1!=0)?PC←loop

stQ: A=3000 Q=4 V=?

reg	Q	V	RS	tag	bsy	Op	Vj	Vk	Qj	Qk
0			A1	1	1	sub	100	1	0	0
1	1	100	A2	2	1	add	0	8	0	0
2	2	0	А3	3						
3	0	13	M1	4	1	mul	13	11	0	0
4		1000	M2	5						
5		2000	LS1	6	0					
6		3000	LS2	7	0					
7	4		LS3	8	1	st	3000,0		0	4
					•					

HEPR Dist RCC

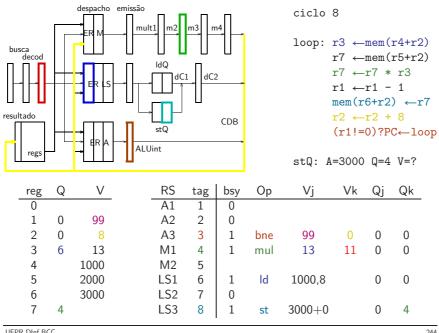


ciclo 7

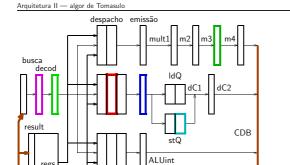
loop: r3 ←mem(r4+r2) r7 ←mem(r5+r2) r7 ←r7 * r3 r1 ←r1 - 1 mem(r6+r2) ←r7 r2 ←r2 + 8 (r1!=0)?PC←loop

stQ: A=3000 Q=4 V=?

reg	Q	V		RS	tag	bsy	Ор	Vj	Vk	Qj	Qk
0			•'	A1	1	0					<u></u>
1	0	99		A2	2	1	add	0	8	0	0
2	2	0		А3	3	1	bne	99	0	0	0
3	0	13		M1	4	1	mul	13	11	0	0
4		1000		M2	5						
5		2000		LS1	6	0					
6		3000		LS2	7	0					
7	4			LS3	8	1	st	3000+0		0	4
						,					



HEPR Dinf RCC



ciclo 9

2007-1

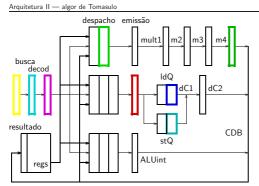
2007-1

loop: $r3 \leftarrow mem(r4+r2)$ $r7 \leftarrow mem(r5+r2)$ r7 ←r7 * r3 r1 ←r1 - 1 $mem(r6+r2) \leftarrow r7$ $r2 \leftarrow r2 + 8$ (r1!=0)?PC \leftarrow loop

stQ: A=3000 Q=4 V=?

reg	Q	V	RS	tag	bsy	Ор	Vj	Vk	Qj	Qk
0			A1	1	0					
1	0	99	A2	2	0					
2	0	8	А3	3	0					
3	6	13	M1	4	1	mul	13	11	0	0
4		1000	M2	5						
5		2000	LS1	6	1	ld	1000,8		0	0
6		3000	LS2	7	1	ld	2000,8		0	0
7	7		LS3	8	1	st	3000+0		0	4

HEPR Dinf RCC

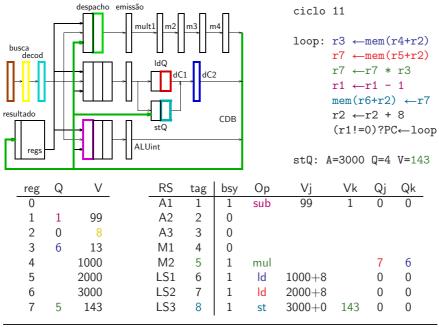


ciclo 10

loop: $r3 \leftarrow mem(r4+r2)$ $r7 \leftarrow mem(r5+r2)$ r7 ←r7 * r3 2x r1 ←r1 - 1 $mem(r6+r2) \leftarrow r7$ $r2 \leftarrow r2 + 8$ (r1!=0)?PC←loop

stQ: A=3000 Q=4 V=?

reg	Q	V	F	RS	tag	bsy	Ор	Vj	Vk	Qj	Qk
0			-	\ 1	1	0					
1	0	99	A	۹2	2	0					
2	0	8	A	43	3	0					
3	6	13	N	/ 11	4	1	mul	13	11	0	0
4		1000	N	/ 12	5	1	mul			7	6
5		2000	L	S1	6	1	ld	1000 + 8		0	0
6		3000	L	S2	7	1	ld	2000,8		0	0
7	5		L	S3	8	1	st	3000+0		0	4

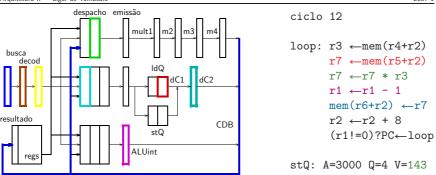


HEPR Dist RCC

Arquitetura II — algor de Tomasulo

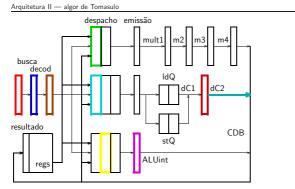
2007-1

2007-1



reg	Q	V	RS	tag	bsy	Op	Vj	Vk	Qj	Qk
0			A1	1	1	sub	99	1	0	0
1	1	99	A2	2	0					
2	0	8	А3	3	0					
3	0	12	M1	4	0					
4		1000	M2	5	1	mul	12		0	7
5		2000	LS1	6	1	st	3000,8		0	5
6		3000	LS2	7	1	ld	2000 + 8		0	0
7	5	143	LS3	8	1	st	3000+0	143	0	0

HEPR Dinf RCC 2.

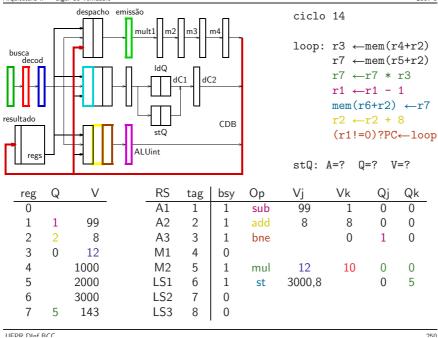


loop:	$r3 \leftarrow mem(r4+r2)$
	$r7 \leftarrow mem(r5+r2)$
	r7 ←r7 * r3
	r1 ←r1 - 1
	$mem(r6+r2) \leftarrow r7$
	$r2 \leftarrow r2 + 8$
	$(r1!=0)$?PC \leftarrow loop

ciclo 13

stQ: A=? Q=? V=?

reg	Q	V	ı	RS	tag	bsy	Ор	Vj	Vk	Qj	Qk
0				A1	1	1	sub	99	1	0	0
1	1	99		A2	2	1	add	8	8	0	0
2	2	8		А3	3	0					
3	0	12	1	М1	4	0					
4		1000	1	M2	5	1	mul	12		0	7
5		2000	L	.S1	6	1	st	3,000,8		0	5
6		3000	L	.S2	7	1	ld	2000 + 8		0	0
7	5	143	L	.S3	8	0					



Arquitetura II — algor de Tomasulo

2007-1

Algoritmo de Tomasulo – resumo

- Estações de Reserva
- * emissão fora-de-ordem quando operandos disponíveis
- Re-nomeação de registradores
- * elimina riscos WAW e WAR importante com poucos regs
- * permite desenrolar laços dinamicamente
- * lógica é relativamente complexa
- Common Data Bus
- * transmite resultados para múltiplos destinos
- pode ser duplicado, com mais hw * gargalo central
- excessões precisas não implementadas
- * poderia usar buffer de re-ordenação para atualizar registradores

HEPR Dinf RCC

Arquitetura II — algor de Tomasulo

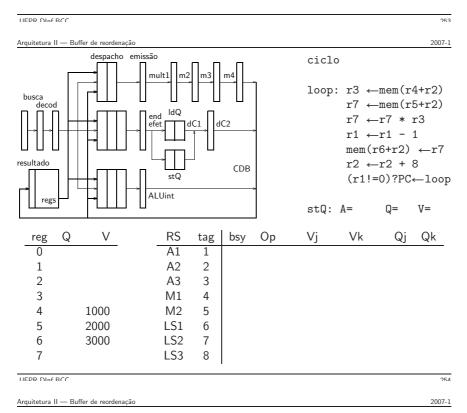
2007-1

Algoritmo de Tomasulo – tratamento de riscos

PROBLEMA	SOLUÇAO
risco estr: Estação de Reserva	bloqueia no despacho
risco estr: unidade funcional	bloqueia na ER
RAW	resolve com as etiquetas
WAR	copia oper para ER no despacho
WAW	renomeação de registradores

Exercícios

- 1. Re-escreva o código do laço para minimizar as bolhas.
- 2. Desenrole o loop para minimizar as bolhas. Qual o ganho com relação ao exercício 1?
- 3. Considerando a latência da multiplicação (4 ciclos), quantas vezes o loop deve ser desenrolado para minimizar as bolhas?
- 4. Complete a seqüência de diagramas após o ciclo 9, em que o destino do desvio é decidido (veja diagrama do ciclo 10). Simule mais uma volta do loop.



revisão - Algoritmo de Tomasulo

- Estações de Reserva
- * emissão fora-de-ordem quando operandos disponíveis
- * emissão continua adiante dos desvios/jumps
- Re-nomeação de registradores elimina riscos WAW e WAR
- * permite desenrolar laços dinamicamente
 múltiplas iterações usam lugares ≠s para valores/regs
- Common Data Bus transmite resultados para múltiplos destinos
- Lógica de detecção de riscos é distribuída pelas ER e CDB
- * se múltiplas instruções esperam por um mesmo resultado, com outro já disponível, então valor no CDB libera todas p/ execução
- escalonamento é dinâmico porque grafo de dependências é construído durante execução

Superescalar + Fora-de-Ordem + Especulação

- Superescalar + Fora-de-Ordem + Especulação conceitos que funcionam bem (melhor) quando combinados
- ◆ CPI ≥1 consegue com superescalar
- superescalar sofre mais riscos? por que? resolve com escalonamento dinâmico (exec/reslt fora-de-ordem)
- dependências RAW causam problemas?
 resolve com janela de instruções grande
- desvios impedem que janela fique cheia?
 resolve com especulação

HEPR Dist RCC 25

Arquitetura II - Buffer de reordenação

2007-1

Excessões precisas e especulação

- Especulação adivinha futuro e executa tentativamente
- importante na previsão de desvios:
 prevê destino e inicia execução naquele caminho
- se especulação foi errada,
 deve voltar e re-iniciar execução de onde previu erradamente
 mesmo comportamento que em excessões precisas
- mesma técnica para excessão precisa e especulação:
 completar em-ordem in-order commit

HEPR Dist RCC 25

Arquitetura II — Buffer de reordenação

2007-1

Especulação e excessões precisas

- Excessões tem semântica seqüencial
- von Neumann
- * todas instruções antes da excessão completam
- * todas instrs após excessão devem parecer nunca ter iniciado
- * mesmas condições que para um desvio previsto erradamente
- qual a dificuldade com excessões/interrupções precisas?
- * completar fora-de-ordem → deve desfazer escritas anteriores
- * em-ordem \sim nenhuma escrita posterior à instr de desvio, é confirmada até que desvio complete
- * fora-de-ordem pode ocorrer
- fora-de-ordem:

excessões precisas e recuperação de especulação errada são o mesmo problema → mesma solução

Algoritmo de Tomasulo e excessões precisas

- Algoritmo de Tomasulo faz:
- * emissão em-ordem
- * execução fora-de-ordem
- * resultados fora-de-ordem
- necessário resolver ordenação dos resultados para que excessões sejam precisas
- * "ponto de quebra" definido:
 instruções mais velhas completaram (escreverem resultado)
 instruções mais novas não alteram estado
- habilidade de abortar e re-iniciar cada instrução
 - → estado preciso

HEPR Dist RCC 25

Arquitetura II - Buffer de reordenação

2007-1

writeback

Qual o problema com estado preciso?

 problema no estágio de escrita/resultado mistura o que deveria estar separado:

* difunde valores para Estações de Reserva, e adianta valores para instruções → po

ra instruções → pode ser fora-de-ordem

* escreve valores nos registradores

→ melhor se é em-ordem

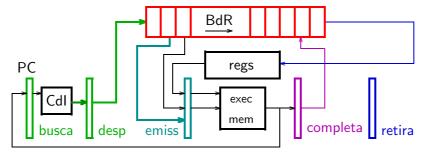
- P: qual a solução para todo problema de funcionalidade?
 R: adicionar um nível de indireção
- * como é feito para execução fora-de-ordem separa DECOD em
 - (1) DESPACHO em-ordem e (2) EMISSÃO fora-de-ordem
- * usa um buffer de instruções para separar os dois estágios placar ou estações de reserva

HEPR Dinf RCC

Arquitetura II — Buffer de reordenação

2007-1

Buffer de Re-ordenação (BdR)



buffer de instruções ≡ Buffer de Reordenação

 mantém resultados prontos a caminho dos regs e da cache de dados

pode ser combinado com as ERs (como na figura) ou separado

divide estágio de RESULTADO em COMPLETA e RETIRA

Completa e Retira BdR regs exec mem completa retira

completa (CM)

commit

- \star valores prontos são gravados no BdR fora-de-ordem
- ★ estágio fora-de-ordem → riscos resultam em espera (wait)
- retira (RT)

retire, graduate

- \star BdR escreve no bloco de registradores em-ordem
- ★ estágio em-ordem → riscos resultam em bloqueios (stalls)

HEPR DIAFRCC 29

Arquitetura II — Buffer de reordenação

2007-1

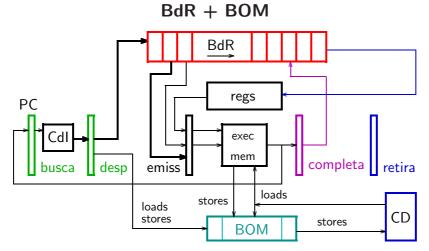
Buffer de Ordenação de Memória

BdR escreve nos regs em-ordem - como fica escrever em memória?

- mesmo que antes → escreve na cache no estágio de memória
- * ruim!!! memória imprecisa é pior que regs imprecisos
- Buffer de ordenação de memória (BOM)
- * fila de escrita, store buffer, load/store queue
- * STORES em COMP (mas não retirados) escrevem no BOM
- st em RET, STORE na cabeça do BOM é escrito em memória
- * LOADs procuram no BOM e na cache em paralelo adiantamento do BOM se referencia mesmo endereço

HEPR Dinf RCC 26

Arquitetura II — Buffer de reordenação 2007-1



Esta é uma versão simplificada, porém realista, do PentiumPro (P6)

Tomasulo + BdR

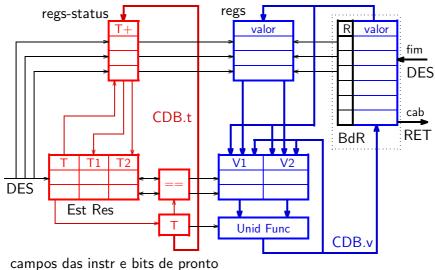
- adicionar BdR ao algoritmo de Tomasulo
- * combinação de BdR e ER é chamada de *Register Update Unit* (RUU) ou "método de Sohi"
- * BdR e ER separadas são chamados de estilo P6
- Exemplo: P6-simples
- * BdR e ER separados
- * organização: 1 ALU, 1 load, 1 store, 2 PF (latência de 3 ciclos)

HEPR Dist RCC 26

Arquitetura II — Buffer de reordenação

2007-1

Organização no Estilo P6



etiquetas (tags) valores

Arquitetura II — Buffer de reordenação

2007-1

Estruturas de dados

- Estações como em Tomasulo simples
- Buffer de Reordenação (BdR)
- * cabeça, fim mantém ordem seqüencial
- * R registrador destino da instr
- * V valor produzido pela instrução
- etiqueta

HEPR Dinf RCC

- * #BdR (elmto do BdR) (Tomasulo usava #ER)
- status dos registradores
- * T+ bit + indica "pronto no BdR"
- * tag = $0 \rightarrow$ resultado pronto em registrador (preciso)
- * tag $\neq 0 \longrightarrow$ resultado não está pronto
- * tag = $N+ \rightarrow$ resultado pronto no BdR #N

Estruturas de dados

hd		Bdl			tat						
tl	#	instr	R	V	ender	EM	EX	CM		reg	T+
	1	ldf f0,X(r1)								f0	
	2	mulf f4,f0,f2								f2	
	3	stf f4,Z(r1)								f4	
	4	add r1,r1,8								r1	
	5	ldf f0,X(r1)								CI	OB .
	6	mulf f4,f0,f2								Τ	V
	7	stf f4,Z(r1)									•

Est Resv										
#	UF	bsy	Т	Vj	Vk	Qj	Qk			
1	ALU									
2	load									
3	store									
4	FP1									
5	FP2									

LIEPR Dinf BCC

Arquitetura II — Buffer de reordenação

Segmentos do P6

Estrutura do processador: BUS, DES, EMI, EXE, COM, RET Busca, Despacho, Emissão, Exec, Completa, Retira

```
• DES (despacho)
```

```
(ER, BdR, BOM, cheios)? (bloqueia):
   { aloca elementos em ER, BdR, BOM,
      atribui #BdR à etiqueta da ER
      atribui #BdR ao status do reg destino, bit + desligado
      copia regs/valores prontos para ER
```

EXE (execução)

libera elemento da ER

(em Tomasulo fazia na finalização, pode fazer em EXE porque #ER não é mais etiqueta)

HEPR Dinf RCC

Arquitetura II - Buffer de reordenação

2007-1

Segmentos do P6

Estrutura do processador: BUS, DES, EMI, EXE, COM, RET

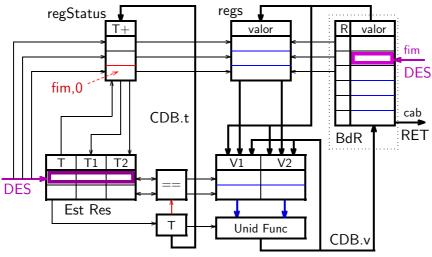
COM (completa)

```
(CDB ocupado) ? (espera) :
                                                            wait
   { escreve valor no elmto do BdR indicado pela etiqueta na ER,
      marca elmto do BdR como pronto,
      liga bit + (pronto-no-BdR) no status do reg destino
```

• RET (retira)

```
(cabeça do BdR não está pronto) ? (bloqueia) :
                                                            stall
   { escreve valor da cabeça do BdR no bloco de regs,
      se STORE escreve cabeça do BOM na cache,
      trata quaisquer excessões,
      libera elementos do BdR e BOM
```

P6 DESpacha (i)



aloca elmtos no BdR e ER (etiqueta dest da ER = #BdR) regStatus aponta para #BdR, bit pronto-no-BdR=0 (+)

HEPR Dinf RCC 27

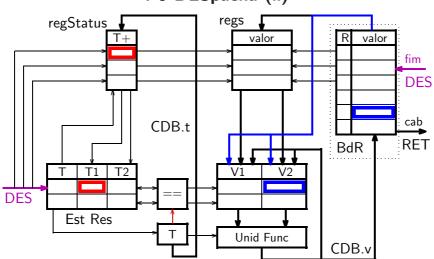
Arquitetura II — Buffer de reordenação

Arquitetura II — Buffer de reordenação

2007-1

2007-1

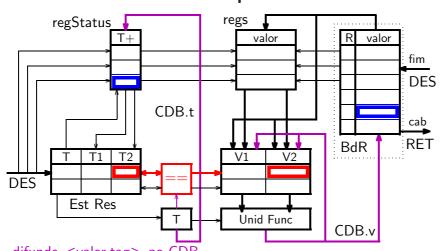
P6 DESpacha (ii)



verifica tags de operandos do regStatus (tag==0: copia dos regs) tag !=0: copia tag p/ ER tag ==n+: copia tag do BdR[n]

HEPR Dief RCC 27

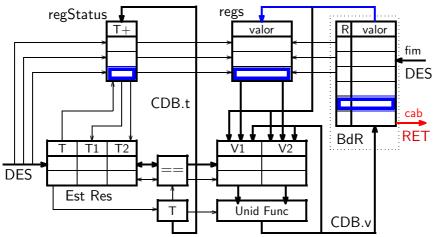
P6 COMpleta



 $\begin{array}{l} {\sf difunde} < \!\! {\sf valor,tag} \!\! > {\sf no} \ {\sf CDB} \\ {\sf escreve} \ {\sf result} \ {\sf no} \ {\sf BdR}, \ {\sf liga} \ {\sf bit} + {\sf em} \ {\sf regStatus} \end{array}$

'casa' tags, grava CDB.v nas ER das instr dependentes

P6 RETira



bloqueia até que instr na cabeça do BdR tenha completado escreve valor da cabeça do BdR nos regs (na cache se STORE), limpa elmto de regStatus libera elmto do BdR

HEPR Dinf RCC 2

Arquitetura II — Buffer de reordenação

2007-1

rStat eg T+

CDB

B-1

Exemplo P6 - ciclo 1

h	d		BdReord + BOrdMem										
t		#	instr	R	V	ender	EM	EX	CM		reg		
h,	,t	1	ldf f0,X(r1)	f0		&X[0]					f0		
		2	mulf f4,f0,f2								f2		
		3	stf $f4,Z(r1)$								f4		
		4	add r1,r1,8								r1		
		5	ldf f0,X(r1)										
		6	mulf f4,f0,f2								Т		
		7	stf $f4,Z(r1)$										

 Est Res

 oper
 T
 Vj
 Vk
 Qj
 Qk

 Idf
 B-1
 [r1]

aloca BdR[1] aloca EstRes[2] ajusta regStat[f0]

HEPR Dist RCC 27

Arquitetura II — Buffer de reordenação

UF

ALU

load

store

FP1

FP2

1

2

3

4

5

2007-1

Exemplo P6 - ciclo 2

hd		Bd	Reor	d +	BOrdMe	m			rS	tat
tl	#	instr	R	V	ender	EM	EX	CM	reg	T+
h	1	ldf f0,X(r1)	f0		&X[0]	c2			f0	B-1
t	2	mulf f4,f0,f2	f4						f2	
	3	stf f4,Z(r1)							f4	B-2
	4	add r1,r1,8							r1	
	5	ldf f0,X(r1)							CI	DB
	6	mulf f4,f0,f2							Τ	V
	7	stf f4,Z(r1)								

Est Res

#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk
1	ALU							
2	load	1	ldf	B-1		[r1]		
3	store							
4	FP1	1	mulf	B-2		R[f2]	B-1	
5	FP2							

aloca BdR[2] aloca EstRes[4] ajusta regStat[f4]

LIEPR Dinf RCC 276

Exemplo P6 - ciclo 3

hd		Bd	Reor	d +	BOrdMe	m				rS	tat
tl	#	instr	R	V	ender	EM	EX	CM	•	reg	T+
h	1	ldf f0,X(r1)	f0		&X[0]	c2	c3		•	f0	B-1
	2	mulf f4,f0,f2	f4							f2	
t	3	stf f4,Z(r1) &Z[0]							f4	B-2	
	4	add r1,r1,8								r1	
	5	Idf f0,X(r1)								CI	DВ
	6	mulf f4,f0,f2								Τ	V
	7	stf f4,Z(r1)									

Est Res

#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk
1	ALU							
2	load	0	_	_		ldf em	EXE(_
3	store	1	stf	B-3		R[r1]	B-2	
4	FP1	1	mulf	B-2		R[f2]	B-1	
5	FP2							

aloca BdR[3] aloca ER[3] libera EstRes[2]

HEPR Dinf RCC

Arquitetura II — Buffer de reordenação

2007-1

Exemplo P6 - ciclo 4

hd		В	dReo	rd + E	30rdMen	n			r	Stat
tl	#	instr	R	V	ender	EM	EX	CM	reg	T+
h	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	с3	с4	f0	B-1+
	2	mulf f4,f0,f2	f4			c4			f2	
	3	stf f4,Z(r1)			&Z[0]				f4	B-2
t	4	add r1,r1,8	r1						r1	B-4
	5	ldf f0,X(r1)							C	.DB
	6	mulf f4,f0,f2							Τ	V
	7	stf f4, $Z(r1)$							B-1	[f0]

Est Res

#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk
1	ALU	1	add	B-4	R[r1]	8		
2	load	0						
3	store	1	stf	B-3		R[r1]	B-2	
4	FP1	1	mulf	B-2	CDB.V	R[f2]	b-1	
5	FP2							

aloca BdR[4] aloca ER[1]

ldf completa, [f0] passa no CdB, marca f0 pronto no BdR: B-1+

HEPR Dinf RCC

Arquitetura II — Buffer de reordenação

2007-1

Exemplo P6 - ciclo 5

hd		В	dRed	ord +	BOrdMei	m				Stat
tl	#	instr	R	V	ender	EM	EX	CM	reg	g T+
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	с3	c4	f0	B-5
h	2	mulf f4,f0,f2	f4			c4	c5,1		f2	
	3	stf $f4,Z(r1)$			&Z[0]				f4	B-2
	4	add r1,r1,8	r1			c5			r1	B-4
t	5	ldf f0,X(r1)	f0							CDB
	6	mulf f4,f0,f2							Т	V
	7	stf f4,Z(r1)							-	

Est Res

#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk
1	ALU	1	add	B-4	R[r1]	8		
2	load	1	ldf	b-5				B-4
3	store	1	stf	B-3		R[r1]	B-2	
4	FP1	0	_	_		mulf e	m EX	EC
5	FP2							

aloca BdR[5] aloca ER[2] libera ER[4] ldf retira, grava [f0] em regs

Exemplo P6 - ciclo 6

hd		В	dRec	ord +	BOrdMei	m			rS	tat
tl	#	instr	R	V	ender	EM	EX	CM	reg	T+
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	c3	c4	f0	B-5
h	2	mulf f4,f0,f2	f4			c4	c5,2		f2	
	3	stf f4,Z(r1)			&Z[0]				f4	B-6
	4	add r1,r1,8	r1			c5	с6		r1	B-4
	5	Idf f0,X(r1)	f0						CI	OB .
t	6	mulf f4,f0,f2	f4						Т	V
	7	stf f4,Z(r1)								

Est Res

#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk
1	ALU	0	_	_		add ei	n EXE	C
2	load	1	ldf	B-5				B-4
3	store	1	stf	B-3		R[r1]	B-2	
4	FP1	0						
5	FP2	1	mulf	B-6		R[f2]	B-5	

aloca BdR[6] aloca ER[5] libera ER[1]

HEPR Dinf RCC

Arquitetura II — Buffer de reordenação

2007-1

Exemplo P6 - ciclo 7

hd		В	dRec	ord + I	BOrdMei	n			•	r	Stat
tl	#	instr	R	V	ender	EM	EX	CM	•	reg	T+
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	с3	с4		f0	B-5
h	2	mulf f4,f0,f2	f4			c4	c5,3			f2	
	3	stf f4,Z(r1)			&Z[0]					f4	B-6
	4	add r1,r1,8	r1	[r1]		c5	с6	c7		r1	B-4+
	5	Idf f0,X(r1)	f0		&X[1]	c7				C	.DB
t	6	mulf f4,f0,f2	f4							Т	V
	7	stf f4, $Z(r1)$		∄ ER livre						B-4	[r1]

Est Res

#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk
1	ALU	0						
2	load	1	ldf	B-5		CDB.V		b-4
3	store	1	stf	B-3		R[r1]	B-2	
4	FP1	0						
5	FP2	1	mulf	B-6		R[f2]	B-5	

add completa

[r1] passa no CDB marca r1 pronto no BdR: B-4+

bloqueia DESPacho porque ∄ ER livre para stf

HEPR Dinf RCC

2007-1

Arquitetura II — Buffer de reordenação

Exemplo P6 - ciclo 8

hd		BdReord + BOrdMem										
tl	#	instr	R	V	ender	EM	EX	CM	reg			
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	c3	c4	f0			
h	2	mulf f4,f0,f2	f4	[f4]		c4	c5	с8	f2			
	3	stf f4,Z(r1)			&Z[0]	с8			f4			
	4	add r1,r1,8	r1	[r1]		c5	с6	с7	r1			
	5	ldf f0,X(r1)	f0		&X[1]	с7	с8					
t	6	mulf f4,f0,f2	f4						Т			
	7	stf f4, $Z(r1)$		ÆΕ	R livre				B-2			

f2 f4 B-6 r1 B-4+ CDB T V B-2 [f4]

rStat g T+

B-5

	Π.
-cT	RAG

#	UF	bsy	oper	Τ	Vj	Vk	Qj	Qk
1	ALU	0						
2	load	0	_	—		ldf em	EXE(2
3	store	1	stf	B-3	CDB.V	R[r1]	b- 2	
4	FP1	0						
5	FP2	1	mulf	B-6		R[f2]	B-5	

mulf completa

[f4] passa no CDB bloqueia DESP: /3 ER livre para stf

bloqueia add em RETira: retira em-- ordem

Exemplo P6 - ciclo 9

hd		BdReord + BOrdMem												
tl	#	instr	R	V	ender	EM	EX	CM	reg					
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	c3	с4	f0					
	2	mulf f4,f0,f2	f4	[f4]		c4	c5	с8	f2					
h	3	stf f4,Z(r1)			&Z[0]	с8	с9		f4					
	4	add r1,r1,8	r1	[r1]		c5	с6	с7	r1					
	5	ldf f0,X(r1)	f0	[f0]	&X[1]	c7	с8	с9						
	6	mulf f4,f0,f2	f4			с9			Т					
t	7	stf $f4,Z(r1)$			&Z[1]				B-5					

rS	Stat
reg	T+
f0	B-5+
f2	
f4	B-6
r1	B-4+
С	DB
Т	V
B-5	[f0]

	Est Res											
#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk				
1	ALU	0										
2	load	0										
3	store	1	stf	B-7		B-4.v	B-6					
4	FP1	0										
5	FP2	1	mulf	B-6	CDB.V	R[f2]	b- 5					

Idf completa
[f0] passa no CDB
libera BdR[3]
. (stf em EX)
e aloca BdR[7]
bloqueia add:
RETira em-ordem

HEPR Dinf RCC

283

Arquitetura II — Buffer de reordenação

2007-1

Exemplo P6 - ciclo 10

hd		BdReord + BOrdMem									
tl	#	instr	R	V	ender	EM	EX	CM			
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	с3	с4			
	2	mulf f4,f0,f2	f4	[f4]		c4	c5	с8			
h	3	stf $f4,Z(r1)$			&Z[0]	c8	с9	c10			
	4	add r1,r1,8	r1	[r1]		c5	с6	с7			
	5	ldf f0,X(r1)	f0	[f0]	&X[1]	c7	c8	с9			
	6	mulf f4,f0,f2	f4			с9	c10,1				
t	7	stf f4,Z(r1)			&Z[1]						

rS	rStat							
reg	T+							
f0	B-5+							
f2								
f4	B-6							
r1	B-4+							
C	.DB							
Т	V							

Est Res

#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk
1	ALU	0						
2	load	0						
3	store	1	stf	B-7		B-4.v	B-6	
4	FP1	0						
5	FP2	0				mulf e	m EXE	EC

libera ER[5] bloqueia add bloqueia ldf RETira em-ordem

HEPR Dinf RCC

2007-1

Arquitetura II — Buffer de reordenação

Exemplo P6 - ciclo 11

hd		BdReord + BOrdMem										
tl	#	instr	R	V	ender	EM	EX	CM				
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	c3	c4				
	2	mulf f4,f0,f2	f4	[f4]		c4	c5	с8				
	3	stf f4,Z(r1)			&Z[0]	c8	с9	c10				
h	4	add r1,r1,8	r1	[r1]		c5	сб	с7				
	5	ldf f0,X(r1)	f0	[f0]	&X[1]	c7	c8	с9				
	6	mulf f4,f0,f2	f4			с9	c10,2					
t	7	stf f4, $Z(r1)$			&Z[1]							

rS	Stat
reg	T+
f0	B-5+
f2	
f4	B-6
r1	B-4+
C	.DB
Т	V

Est Res

#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk
1	ALU	0						
2	load	0						
3	store	1	stf	B-7		B-4.v	B-6	
4	FP1	0						
5	FP2	0						

retira stf e libera BdR[3] bloqueia add bloqueia ldf RETira em-ordem

rStat $\mathsf{T}+$

B-5+

B-6

B-4+ CDB Τ

reg f0

f2 f4

r1

Exemplo P6 - ciclo 12

hd		BdReord + BOrdMem									
tl	#	instr	R	V	ender	EM	EX	CM			
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	с3	с4			
	2	mulf f4,f0,f2	f4	[f4]		c4	c5	с8			
	3	stf $f4,Z(r1)$			&Z[0]	с8	c9	c10			
h	4	add r1,r1,8	r1	[r1]		c5	с6	c7			
	5	ldf f0,X(r1)	f0	[f0]	&X[1]	c7	c8	с9			
	6	mulf f4,f0,f2	f4			с9	c10,3				
t	7	stf f4,Z(r1)			&Z[1]						

Est Res												
#	UF	bsy	oper	Т	Vj	Vk	Qj	Qk				
1	ALU	0										
2	load	0										
3	store	1	stf	B-7		B-4.v	B-6					
4	FP1	0										
5	FP2	0										

retira add **EM-ORDEM** (bloq desde c7)

e libera BdR[4]

bloqueia ldf RETira em-ordem

LIEPR Dinf BCC

Arquitetura II - Buffer de reordenação

2007-1

Estado preciso com BdR

- Razão de ser do BdR é manter estado preciso
- como?
- 1. espera até que condição precisa chegue ao estágio onde RETira
- 2. limpa conteúdo do BdR, ER, e regStatus

enche de zeros

- 3. re-inicia
- funciona porque zeros são o estado correto para re-iniciar
- *~0 em BdR/ER \longrightarrow elementos estão vazios
- * tag=0 na tabela de status dos regs \rightarrow registrador atualizado
- e porque escritas nos registradores e na cache ocorrem em RETira
- exemplo: falta de página em STORE

HEPR Dinf RCC

Arquitetura II — Buffer de reordenação

Exemplo de Estado Preciso no P6 - ciclo 9

hd		BdReord + BOrdMem									
tl	#	instr	R	R V ender		EM	EX	CM			
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	c3	c4			
	2	mulf f4,f0,f2	f4	[f4]		c4	c5	c8			
h	3	stf f4,Z(r1)			&Z[0]	c8	c9				
	4	add r1,r1,8	r1	[r1]		c5	с6	c7			
	5	ldf f0,X(r1)	f0	[f0]	&X[1]	c7	c8	c9			
	6	mulf f4,f0,f2	f4			с9					
t	7	stf f4,Z(r1)			&Z[1]						

rStat $\mathsf{T}+$ reg f0 B-5+ f2 f4 B-6 B-4+ r1 CDB Т V [f0] B-5

Est Res

	UF	bev	onor	т	Vi	Vk	Qi	Ωk
#	UF	bsy	oper	ı	٧J	VK	۷J	Qk
1	ALU	0						
2	load	0						
3	store	1	stf	B-7		B-4.v	B-6	
4	FP1	0						
5	FP2	1	mulf	B-6	CDB.V	R[f2]	b- 5	

FALTA de PÁGINA em stf em c9

LIEPR Dinf RCC

Exemplo de Estado Preciso no P6 - ciclo 10

hd		Во	dReo	rd + E	3OrdMen	n			rStat		
tl	#	instr	R	V	ender	EM	EX	CM	re	eg	T+
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	c3	с4	f	0	
	2	mulf f4,f0,f2	f4	[f4]		c4	c5	c8	f	2	
h,t	3	stf f4,Z(r1)			&Z[0]	с8	с9	c10	f	4	
	4	add r1,r1,8							r	1	
	5	ldf f0,X(r1)								CI	OB
	6	mulf f4,f0,f2							-	Γ	V
	7	stf f4.Z(r1)									

Est Res										
UF	bsy	oper	Т	Vj	Vk	Qj	Qk			
ALU	0									
load	0									
store	0									
FP1	0									

instrução com excessão (stf) chega à cabeça do BdR LIMPA TUDO!

LIEPR Dinf BCC

1

2

3

4

5

FP2

Arquitetura II — Buffer de reordenação

Arquitetura II — Buffer de reordenação

0

2007-1

2007-1

Exemplo de Estado Preciso no P6 – ciclo 11

hd		Во			rS1	tat					
tl	#	instr	R	V	ender	EM	EX	CM	•	reg	T+
	1	ldf f0,X(r1)	f0	[f0]	&X[0]	c2	c3	c4		f0	
	2	mulf f4,f0,f2	f4	[f4]		c4	c5	с8		f2	
h,t	3	stf f4,Z(r1)			&Z[0]					f4	
	4	add r1,r1,8								r1	
	5	Idf f0,X(r1)							_	CE)B
	6	mulf f4,f0,f2								Τ	V
	7	stf $f4,Z(r1)$							_		

Est Res UF Vk Qj Qk bsy oper ALU 1 0 2 0 re-inicia de onde parou load 3 B-3 R[f4] R[r1] store 1 stf 4 FP1 0 FP2 5 0

LIEPR Dinf BCC

Desempenho do Estilo P6

Qual é o custo de excessões precisas?

- em geral, mesmo desempenho que Algoritmo de Tomasulo
- * pouco melhor porque libera ER antes (EXEC) e portanto incorre em menos riscos estruturais
- se BdR não for suficientemente grande, desempenho ruim
- * BdR causa riscos estruturais
- regras de projeto para BdR N=largura
- $ightarrow |\mathsf{BdR}| \geq (N \times \mathsf{núm} \; \mathsf{de} \; \mathsf{estágios} \; \mathsf{entre} \; \mathsf{DESP} \; \mathsf{e} \; \mathsf{RET})$
- $\triangleright |\mathsf{BdR}| \geq (N \times \mathcal{T} \mathsf{acerto} \mathsf{na} \mathsf{L2})$
- ⊳ qual o porquê destas regras?

Problemas com P6

- Organização foi popular (início 90s)
- * relativamente fácil de implementar corretamente
- * para recuperar estado, zera tudo
- * exemplos: PentiumPro, AMD K6, PowerPC
- Problemas com implementação de alto desempenho
- ⊳ excessiva movimentação de valores:

$$regs/BdR \rightarrow ER \rightarrow BdR \rightarrow regs$$

- ▷ ER misturam valores com etiquetas (controle) caminhos ficam longos e isso baixa freqüência do relógio

HEPR Dinf RCC 200

Arquitetura II — Buffer de reordenação

2007-1

Solução: implementação do MIPS R10K

- Separa controle (BdR/ER) dos dados (regs/Unid Funcionais)
- ▷ BdR e ER usados somente para controle e etiquetas → pequenos
- ⊳ bloco de registradores próximo das UFs, todo o resto fica ao lado

Detalhes em K C Yeager, *The MIPS R10000 Superscalar Microprocessor*, IEEE Micro 16:2, Abr1996.

TIEPR DIM RCC 200

Arquitetura II — escalonamento estático

2007-1

resumo – Especulação e Excessões Precisas

- Excessões tem semântica seqüencial
- von Neumann
- * todas instruções antes da excessão completam
- * todas instrs após excessão devem parecer nunca ter iniciado
- * mesmas condições que para um desvio previsto erradamente
- qual a dificuldade com excessões/interrupções precisas?
- * completar fora-de-ordem → deve desfazer escritas anteriores
- implementação: Tomasulo com Buffer de Reordenação Estilo P6
- ▶ instruções completam mas só escrevem resultado nos registradores
 - se ∄ excessão
- ⊳ se há excessão (ou previsão de desvio errada) zera ER, BdR, e tabela de status dos regs, e re-inicia de onde parou

revisão - Especulação e Excessões Precisas

- Excessões tem semântica seqüencial
- von Neumann
- * todas instruções antes da excessão completam
- * todas instrs após excessão devem parecer nunca ter iniciado
- * mesmas condições que para um desvio previsto erradamente
- qual a dificuldade com excessões/interrupções precisas?
- * completar fora-de-ordem → desfazer escritas anteriores
- implementação I: Tomasulo + Buffer de Reordenação Estilo Pé
- ▶ instruções completam mas só escrevem resultado nos registradores
 - se ∄ excessão
- ⊳ se (∃ excessão || previsão_desvio_errada)
 zera ER, BdR, e tabela statusRegs; re-inicia de onde parou
- implementação II: MIPS R10K
 separa controle (BdR/ER) dos dados (regs/Unid Funcionais)

LIEPR Dinf RCC

205

Arquitetura II - escalonamento estático

2007-1

Trilha do desempenho

- segmentação simples: emissão em-ordem, largura=1
- primeira extensão: super-escalar, emissão múltipla, largura>1
- segunda extensão: escalonamento de instruções para > PNI
- * 1ª opção: escalonamento dinâmico
- * 2ª opção: escalonamento estático

LIEPR Dinf BCC

2007-1

Arquitetura II — escalonamento estático

Very Long Instruction Word – VLIW

Problemas com implementação super-escalar

- ▶ largura de busca + previsão de desvios
- riangleright adiantamento $\propto N^2$
- riangleright avaliação de dependências $\propto N^2$

Alternativa: Very Long instruction Word - VLIW

- \rightarrow pipeline com emissão simples de grupo com N instruções (VLIW)
 - ⊳ compilador garante independência das instrns numa VLIW
 - ▶ processador não testa dependências "dentro" de uma VLIW

 - ight
 angle alocação \pm fixa: $1^{\underline{a}}$ é ULA, $2^{\underline{a}}$ é LOAD, etc

ULA Id beq ULA

História

- iniciou com microcódigo "horizontal"
- projetos na academia
- ▷ ELI-512 [Fischer 85]
- ⊳ IMPACT [Hwu 91]
- produtos comerciais
- ⊳ Multiflow [Colwell, Fischer 85] → falhou
- ▷ Cydrome [Rau 85] → falhou
- ⊳ EPIC, IA-64, Itanium [Colwell, Fischer, Rau 97] ~>??
- \triangleright Transmeta [Ditzel 99] \rightarrow ??

traduz x86→ VLIW

HEPR Dinf RCC 20

Arquitetura II — escalonamento estático

2007-1

VLIW per se

- VLIW puro: nenhum teste de dependência em hw nem mesmo entre grupos de VLIWs
- compilador responsável por escalonar todos os segmentos
- ⊳ é possível se conhece exatamente a estrutura da CPU e latências
- 1º problema: estrutura e latências diferem entre implementações
- ▶ TransMeta faz re-compilação em tempo de execução
- 2º problema: latências não são fixas numa implementação
- ⊳ não usar caches? uh?
- ▷ escalonar prevendo uma falta na cache? para que a cache então?

HEPR Dist RCC 20

Arquitetura II — escalonamento estático

2007-1

VLIW viável

Explicitly Parallel Instruction Computing (EPIC) ISA-64, Itanium

menos rígido que VLIW

- talvez nem é VLIW...
- palavras de instruções com tamanho variável
- ⊳ implementadas como "amontoados" com bits de dependência
- b torna código compatível com máquinas de larguras ≠
- pressupõe que hw implementa lógica de bloqueio inter-amontoados
- b torna código compatível com ≠s núm_segmentos, latências
- ▷ permite bloqueios por faltas na cache, e execução fora-de-ordem
- explora toda informação sobre paralelismo que compilador descobre
- compatível com ≠s implementações da mesma arquitetura

Escalonamento e Paralelismo no Nível de Instrução

Não faz sentido implementar um pipeline com largura $oldsymbol{N}$ se, na média,

muito menos que ${m N}$ instruções independentes executam num ciclo

desempenho é importante

mas utilização também é (sustentada/pico)

HEPR Dist RCC

Arquitetura II - escalonamento estático

2007-1

Exemplo: SAXPY

Single Precision A*X+Y

rotina de álgebra linear para resolver sistemas de equações parte dos "Livermore Loops"

HEPR Dist RCC 30

Arquitetura II — escalonamento estático

2007-1

Desempenho e Utilização - SAXPY escalar

Processador segmentado, escalar MULF 5 ciclos, ADDF 2 ciclos, ambos segmentados adiantamento completo, desvios previstos como tomados

```
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16
                B D X M R
ldf f0, X(r1)
                  B D d* E E E E
mulf f4, f0, f2
                     B p* D X M R
ldf f6, Y(r1)
                          B D d* d* d*
addf f8,f6,f4
                                        Ε
                                           Ε
                                              R
                                  p*
                             В р*
                                     р*
stf f8,Z(r1)
                                        D
                                           X M R
                                           D
                                              X M R
add r1, r1, 8
                                                       R
                                           В
                                              D
                                                 X M
ble r1,r2,loop
               d*=dep dados; p*=dep segm, bolha
```

uma iteração com 7 instruções — latência de 15 ciclos desempenho: 7 instr/15 ciclos → IPC = 0,47 utilização: 0,47 IPCreal / 1,0 IPCideal → 47%

Desempenho e Utilização - SAXPY superescalar

Processador segmentado, superescalar, largura $N{=}2$ MULF 5 ciclos, ADDF 2 ciclos, ambos segmentados

```
3 4 5 6 7 8 9 10 11 12 13 14 15
              B D X M
                         R
ldf f0, X(r1)
              B D d* d* E E E E R
mulf f4, f0, f2
              B p^* p^* D X M R
ldf f6, Y(r1)
                 B p* p* D d* d* d* d* E E
addf f8,f6,f4
                   B p* D p* p* p* p* X d* M R
stf f8,Z(r1)
                         B p* p* p* p* D p* X M R
add r1,r1,4
                         B p* p* p* p* D p* d* X M
                                                    R
ble r1,r2,loop
              d*=dep dados; p*=dep segm, bolha
```

latência ainda em 15 ciclos — por que???
utilização: 0,47 IPCreal / 2,0 IPCideal → 27%
mais riscos → mais bloqueios — por que?
cada bloqueio é mais caro — por que?

HEPR Dist RCC 30.

Arquitetura II - escalonamento estático

2007-1

Escalonamento e Emissão

Escalonamento de instruções

→ escolhe ordem de execução das instr

- importante para melhorar utilização e desempenho
- relacionado à emissão de instr→ quando é que instr executam

escalonamento estático com emissão dinâmica (ou quase)

 escalar em-ordem depende de bom escalonamento pelo compilador

HEPR Dinf RCC

Arquitetura II — escalonamento estático

2007-1

Escalonamento de Instruções

- idéia: instruções independentes entre operações lentas e uso RAW
- * evita bloqueio enquanto espera resolução do risco RAW
- * ≈ mesmo que escalonamento dinâmico (Tomasulo)
- para isso, necessita de instruções independentes
- escopo de escalonamento: região de código sob exame
- * quanto maior melhor maior # instr independentes para escolher
- * com escopo definido, escalonar é simples
- * o problema é expandir o escopo → cruzar desvios?
- escalonamento pelo compilador

▷ desenrolar loops
 ▷ software pipelining
 ▷ escalonamento de traçados trace scheduling
 resto

Escalonamento: compilador ou hardware

compilador

- > possibilita hw simples com relógio rápido
- ▶ ∄ informação sobre latências (faltas na cache)
- ⊳ recuperação de especulação errada é cara

hardware

- previsão de desvios com melhor acurácia
- ▷ informação dinâmica sobre latências (cache) e dependências
- ⊳ fácil de especular e recuperar em caso de erro
- ▷ |buffer_de_instruções| finito limita escopo para escalonamento
- ⊳ hw mais complicado (mais energia), relógio lento

HEPR Dist RCC 30

Arquitetura II - escalonamento estático

2007-1

Escalonamento por SW – desenrolar o SAXPY

Deseja-se separar operações dependentes

- * escopo de uma iteração é muito pequeno
- * cadeia mais longa de operações tem 9 ciclos
- ⊳ resultado do ldf 1 ciclo

- ⊳ não dá para esconder latência de 9 ciclos em 7 instruções, mas dá se forem 2×9 ciclos em 14 instruções
- * desenrolar laço = escalonar usando duas iterações

HEPR Dist RCC 300

Arquitetura II — escalonamento estático 2007-

Desenrolar fase 1: agrupar iterações

ldf f0, X(r1) ldf f0, X(r1)mulf f4,f0,f2 mulf f4,f0,f2 Combinar 2 iterações ldf f6, Y(r1) ldf f6, Y(r1) (no geral N) addf f8,f6,f4 addf f8,f6,f4 stf f8,Z(r1) stf f8,Z(r1)fundir controle do laco: add r1,r1,8 --- --,--,-incremento do índice --- --,--,--ble r1,r2,loop e desvio ldf f0, X(r1) ldf f0, X+8(r1)mulf f4,f0,f2 mulf f4,f0,f2 ajustar usos implícitos ldf f6, Y(r1) ldf f6, Y+8(r1) das addf f8,f6,f4 addf f8,f6,f4 variáveis de indução stf f8,Z+8(r1)stf f8,Z(r1)(r1 neste caso) add r1,r1,8 add r1,r1,8+8 ble r1,r2,loop ble r1,r2,loop

Desenrolar fase 2: escalonar segmentos

```
ldf f0, X(r1)
                                           ldf f0, X(r1)
                        mulf f4,f0,f2
                                           ldf f0, X+8(r1)
Escalonar para
                        ldf f6, Y(r1)
                                           mulf f4,f0,f2
reduzir bloqueios
                        addf f8, f6, f4
                                           mulf f4,f0,f2
por riscos RAW
                        stf f8,Z(r1)
                                           ldf f6, Y(r1)
                        ldf f0, X+8(r1)
                                           ldf f6, Y+8(r1)
                        mulf f4,f0,f2
                                           addf f8,f6,f4
                        ldf f6, Y+8(r1)
                                           addf f8,f6,f4
                        addf f8, f6, f4
                                           stf f8,Z(r1)
                        stf f8,Z+8(r1)
                                           stf f8, Z+8(r1)
                                           add r1,r1,16
                        add r1,r1,16
                        ble r1,r2,loop
                                           ble r1,r2,loop
```

LIEPR Dinf BCC

Arquitetura II — escalonamento estático

Desenrolar fase 3: re-nomear registradores

```
ldf f0, X(r1)
                                          ldf f0, X(r1)
                       mulf f4,f0,f2
                                          ldf f10, X+8(r1)
Escalonamento causa
                       ldf f6, Y(r1)
                                          mulf f4,f0,f2
riscos WAR:
                       addf f8, f6, f4
                                          mulf f14,f10,f2
é necessário
                       stf f8,Z(r1)
                                          ldf f6, Y(r1)
re-nomear
                       ldf f0, X+8(r1)
                                          ldf f16, Y+8(r1)
registradores
                       mulf f4,f0,f2
                                          addf f8, f6, f4
                       ldf f6, Y+8(r1)
                                          addf f18, f16, f14
                       addf f8, f6, f4
                                          stf f8, Z(r1)
                       stf f8,Z+8(r1)
                                          stf f18,Z+8(r1)
                       add r1,r1,16
                                          add r1, r1, 16
                       ble r1,r2,loop
                                          ble r1,r2,loop
```

LIEPR Dinf BCC

Arquitetura II — escalonamento estático

2007-1

Desempenho do SAXPY desenrolado

```
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17
               B D X M R
ldf f0, X(r1)
ldf f10, X+8(r1) B D X M R
                  BDEEEEER
mulf f4,f0,f2
                    B D E E E E R
mulf f14,f10,f2
                        B D X M R
ldf f6, Y(r1)
                           B D X M s* s* R
ldf f16, Y+8(r1)
                             B D d^* E E s^* R
addf f8,f6,f4
                                                 dep em f4
                               B p^* D E p^* E R
addf f18, f16, f14
stf f8,Z(r1)
                                    B D X M R
stf f18,Z+8(r1)
                                       B D X M R
                                          B D X M R
add r1, r1, 16
ble r1,r2,loop
                                            B D X M R
               d*=dep dados; p*=dep segm; s*=risco estrut regs
```

duas iterações com 12 instruções — latência de 17 ciclos desempenho: 12 instr/17 ciclos \rightarrow IPC = 0,71 $(\gg 47\%)$

utilização: 0,71 IPCreal / 1,0 IPCideal \rightarrow 71%

Nem tudo são flores...

Desenrolar loops é bom mas

- código é maior ($\approx N \times$)
- escalonamento pode ser ruim nas fronteiras dos laços
- não é possível tratar dependências entre iterações

```
for (i=0; i<N; i++)
```

X[i] = A*X[i-1]; /* iteração depende da anterior */

laço desenrolado mantém uma cadeia de dependências que não pode ser escalonada

o que fazer se (número_de_iterações MOD desenroladas ≠ 0) ??

HEPR Dist RCC 311

Arquitetura II - escalonamento estático

2007-1

Software pipelining

Software pipelining = Desenrolamento simbólico de laços

Se as iterações do laço são independentes, então pode extrair mais paralelismo se usar instruções de iterações distintas

• idéia:

- ⊳ inicia com laço original
- ▷ escreve novo laço com instruções de iterações distintas

ganhos:

- ⊳ maximiza distância entre resultado e uso
- ⊳ enche e esvazia segmentos uma vez por volta, ao invés de uma vez em cada iteração desenrolada

HEPR Dinf RCC 31

Arquitetura II — escalonamento estático

2007-1

Analogia com pipeline

segmentação em hw

cada ciclo contém:

estágio 3 da instr 1, est 2 da instr i+1, est 1 da instr i+2

segmentação em sw

ciclo \sim iteração "física" do laço instrução \sim iteração "lógica" original do laço estágio \sim instrução

• uma iteração "física" contém

instruções de múltiplas iterações originais

instr 3 da iteração 1, instr2 da iter i+1, instr 1 da iter i+2

Exemplo 1

```
iteração física: stf da iter original i ; ldf,mulf da iter original i+1 prólogo: inicializa pipeline → ldf,mulf da iteração 0 epílogo: completa trabalho → stf da última iteração
```

```
ldf f0, X(r1)
                           ldf f0,X-8(r1)
                                            ← prólogo
mulf f4,f0,f2
                           mulf f4,f0,f2
stf f4,X(r1)
                     loop: stf f4,X(r1)
                                            ← iteração
                           ldf f0,X(r1)
                                             ← física
add r1,r1,8
ble r1,r2,loop
                           mulf f4,f0,f2
ldf f0,X(r1)
                           add r1, r1, 8
mulf f4,f0,f2
                           ble r1,r2,loop
stf f4,X(r1)
                           stf f4,X+8(r1) \leftarrow epilogo
add r1, r1, 8
ble r1,r2,loop
```

HEPR Dist RCC 316

Arquitetura II — escalonamento estático

2007-1

Diagrama de tempo

Mesmo diagrama com nova terminologia: ciclos → iteração física (horizontal) instruções → iteração lógica (vertical) estágios → instruções (LM=ldf,mulf S=stf)

	fís1	fís2	fís3	fís4	fís5	
lóg1	LM	S				← prólogo
lóg1 lóg2		LM	S			
lóg3			LM	S		
lóg4				LM	S	← epílogo

iteração física 2 tem **stf** da iter lógica 1, e ldf,mulf da iter lógica 2

Note que numa iteração física os grupos de instruções estão na ordem invertida → OK porque grupos não são relacionados (exec //)

HEPR Dist RCC 311

Arquitetura II — escalonamento estático

2007-1

Exemplo 2

Altera estrutura para tolerar mais latência (3 iterações)

```
ldf f0, X(r1)
                          ldf f0, X-8(r1)
mulf f4,f0,f2
                          mulf f4,f0,f2
                                          \leftarrow prólogo
stf f4, X(r1)
                         ldf f0, X-4(r1)
add r1,r1,8
                    loop: stf f4,X-8(r1) ← iteração
ble r1,r2,loop
                         mulf f4,f0,f2
                                          ← física
ldf f0, X(r1)
                          ldf f0, X(r1)
mulf f4,f0,f2
                          add r1,r1,4
stf f4, X(r1)
                          ble r1,r2,loop ←
add r1,r1,8
                          stf f4, X-4(r1)
ble r1,r2,loop
                         mulf f4,f0,f2
                                          ← epílogo
ldf f0, X(r1)
                          stf f4, X(r1)
mulf f4,f0,f2
stf f4, X(r1)
                          L M S
add r1, r1, 8
                            L M S
ble r1,r2,loop
                               L M S
                                   L M S
```

Software pipelining

```
⊳ não aumenta (muito) o tamanho do código
```

⊳ pode variar grau de paralelismo para tolerar maiores latências software super-pipelining uma iteração física com iterações lógicas i, i+2, i+4

HEPR Dinf RCC 31

Arquitetura II — escalonamento estático

2007-1

Suporte de hardware para escalonamento

Instruções condicionais ("predicadas")

Dois níveis de predicação

- predicação completa: cada instrução carrega um predicado IA-64
- move condicional: registrador é copiado se condição é válida conditional move = CMOVE

```
cmoveqz r1,r2,r3 # if (r3 == 0) r1 \leftarrow r2;
```

if-conversion converte fluxo de controle em fluxo de dados → elimina desvios

pode causar montes de replicação de código

Alpha, IA-32

HEPR Dinf RCC 320

Arquitetura II — escalonamento estático

2007-1

Suporte de hardware para escalonamento

Instruções condicionais ("predicadas")

Exemplo: emissão de 2 instruções/ciclo

PRIMEIRA INSTR SEGUNDA INSTR lw r1,16(r2) add r3,r4,r5 add r6,r3,r7

beqz r10,label
lw r8,0(r10)
lw r9,0(r8)

desperdiça oportunidade de emissão porque 3° LW depende do resultado do 2° LW.

Suporte de hardware para escalonamento

Instruções condicionais ("predicadas")

Usar versão condicional da instrução load: LWC = load condicional load é efetivado se terceiro operando é zero

```
PRIMEIRA INSTR
                      SEGUNDA INSTR
lw r1,16(r2)
                      add r3,r4,r5
lwc r8,0(r10),r10
                      add r6, r3, r7
beqz r10, label
lw r9,0(r8)
```

Se a seqüência após o desvio fosse curta, todo o bloco básico poderia ser convertido para execução condicional e o desvio eliminado

LIEPR Dinf BCC

Arquitetura II - escalonamento estático

2007-1

Suporte de hardware para especulação em memória

Compilador só pode mover LOADs adiante de STORES se há certeza absoluta de que o movimento é correto

Arquitetura pode incluir instruções CHK e LDA que verificam conflitos de endereçamento

```
CHK = check address
                        LDA = load advanced
```

- ⊳ instrução CHK é inserida no local original do LOAD e LDA é movido para antes do(s) STORE(s)
- ⊳ quando load avançado LDA (especulativo) é executado, endereço efetivo é armazenado (no BdR?)
- ▶ se um STORE subsequente referencia aquele endereço antes da instrução CHK então a especulação falhou
- ⊳ se somente LOAD foi especulado, então basta repetir referência à memória quando instrução CHK é executada

HEDD DINE DCC

Arquitetura II — escalonamento estático

Suporte de hardware para especulação em memória

Compilador não consegue garantir que $X+r1 \neq Y+r2$

```
ldf f2,X(r1)
                   ldf f2,X(r1)
mulf f4,f2,f0
                   1dfa f6, Y(r2)
                                    \leftarrow ld adiantado
stf f4,X(r1)
                  mulf f4,f2,f0
ldf f6,Y(r2)
                   stf f4,X(r1)
mulf f8,f6,f0
                   chk Y(r2)
                                    ← verifica
stf f8,Y(r2)
                   mulf f8,f6,f0
                   stf f8,Y(r2)
```

Se quando executa primeiro STORE X+r1 = Y+r2, então repete ldf f6,Y(r2) antes do segundo mulf

Mecanismo é chamado memory conflict buffer e usado no IA-64

Comportamento com excessões

Compilador não pode violar comportamento c.r.a excessões

Exemplo 1: instruções condicionais não podem causar excessões se a condição é falsa, e o efeito da instr é anulado anula inclusive/especialmente efeitos excepcionais

Exemplo 2: busca antecipada não pode causar excessões

TIEPR DINFRCC 32

Arquitetura II - escalonamento estático

2007-1

Suporte de hardware para especulação em memória

O que acontece se o LOAD adiantado causa uma falta de página? Nada, se usar LDE = load especulativo

```
\begin{array}{llll} \mbox{ldf } f2, X(r1) & \mbox{ldf } f2, X(r1) \\ \mbox{mulf } f4, f2, f0 & \mbox{ldfe } f6, Y(r2) & \leftarrow \mbox{ld especulativo} \\ \mbox{stf } f4, X(r1) & \mbox{mulf } f4, f2, f0 \\ \mbox{bfnez } f4, \mbox{label} & \mbox{stf } f4, X(r1) \\ \mbox{ldf } f6, Y(r2) & \mbox{mulf } f8, f6, f0 & \leftarrow \mbox{mul especulativo} \\ \mbox{mulf } f8, f6, f0 & \mbox{bfnez } f4, \mbox{label} \\ \mbox{stf } f8, Y(r2) & \mbox{stf } f8, Y(r2) \\ \end{array}
```

- ▶ Bit de interrupção/excessão é associado ao f6
- bit é propagado para f8 → mulfe é especulativo

HEPR Dist RCC

Arquitetura II — hierarquia de memória

2007-1

resumo - Escalonamento Estático

- desenrolar laços
- ⊳ reduz a freqüência de desvios
- > aumenta tamanho do código
- software pipelining
- ⊳ não há dependências no corpo do laço
- ⊳ não reduz a freqüência de desvios
- ⊳ necessita blocos com prólogo e epílogo
- escalonamento de traçados
- ▶ pode usar para código sem laços
- ⊳ não é simples
- suporte pela arquitetura (conjunto de instruções)
- ▶ loads especulativos, loads avançados
- ▶ predicação, move-condicional

Sistema de Memória

• Sistemas de memória

- * hierarquia de memória
- * vazão e latência
- **★** DRAM e SRAM
- Memória cache
- ⋆ organização
- ★ leitura
- ★ escrita
- * otimizações
- Memória Virtual
- * endereçamento
- * paginação
- * TLB
- ⋆ excessões

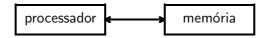
LIEPR Dinf RCC

378

Arquitetura II — hierarquia de memória

2007-1

Sistemas de Memória



Atualmente, o desempenho de computadores é limitado pela latência da memória e pela vazão até a memória

latência é o tempo de um único acesso [s] tempo de acesso à memória ≫ ciclo do processador

vazão é o número de acessos por unidade de tempo [refs/s] se uma fração m das instruções acessam a memória \longrightarrow ocorrem 1+m referências/instrução \longrightarrow CPI=1 se e só se ocorrerem 1+m referências/ciclo

HEPR Dinf RCC 320

Arquitetura II — hierarquia de memória

2007-1

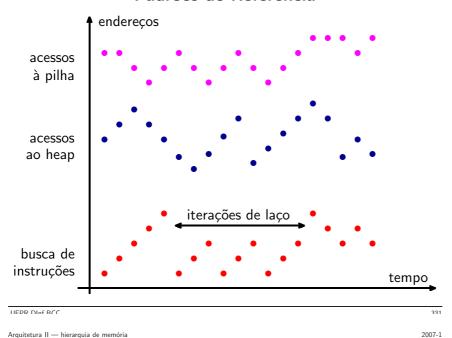
Vazão e Latência

- Latência da memória é o intervalo desde a requisição pelo processador até a disponibilidade para o processador [s]
- Vazão é a taxa de transferência de/para a memória [x/s]
 bandwidth = largura de banda
- Vazão e latência são intimamente relacionadas:
 Se R é o número de requisições que a memória pode atender simultaneamente, então

$$V = R/L$$
 $[x/s = x/s]$

Porque raramente V = 1/L ?

Padrões de Referência



Princípios da Localidade I

Referências à memória apresentam duas propriedades:

Localidade Temporal: se um local é referenciado, possivelmente, este será referenciado de novo em futuro próximo

Localidade Espacial: se um local é referenciado, possivelmente, locais próximos serão referenciados em futuro próximo

HEPR Dinf RCC 32'

Arquitetura II — hierarquia de memória

2007-1

Princípios da Localidade II

```
for (i=0; i<100; i++)
  for (j=0; j<200; j++)
    res[i] += vetor[i] * matriz[i][j];</pre>
```

localidade temporal: código, variáveis locais i e j

Princípios da Localidade III

Caches podem tirar proveito dos dois tipos de localidade:

da localidade temporal ao lembrar do conteúdo de posições acessadas recentemente

→ localidade temporal ajuda escolher o que expurgar da cache

da localidade espacial ao buscar blocos em áreas acessadas recentemente

→ localidade espacial ajuda escolher o que carregar na cache

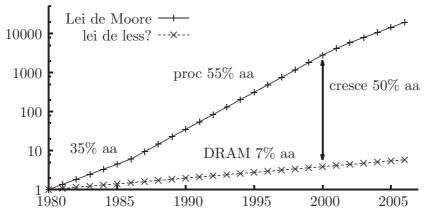
HEPR Dist RCC 33

Arquitetura II — hierarquia de memória

2007-1

Hierarquias de Memória

desempenho: processadores vs DRAM



Processador que emite quatro instruções por ciclo poderia executar 800-1000 instruções durante uma falta na cache!!!

HEPR Dinf RCC 33

Arquitetura II — hierarquia de memória

2007-1

Hierarquias de Memória

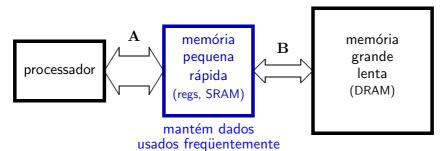
Solução para aumento na latência relativa:

esconder a latência

interpondo memória pequena mas com alta velocidade entre processador e memória DRAM

Memória cache permite esconder latência da memória porque padrões de referência são muito previsíveis → localidade

Hierarquia de Memória



tamanho: registradores \ll SRAM \ll DRAM por que? latência: registradores \ll SRAM \ll DRAM por que? vazão: interna ao CI \gg externa ao CI por que?

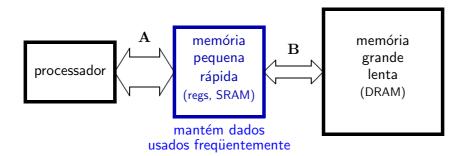
Memória rápida será efetiva se e só se tráfego em $B \ll A$

HEPR Dinf RCC 327

Arquitetura II — hierarquia de memória

2007-1

Hierarquia de Memória



No acesso ao conteúdo:

acerto se encontra $(dado \in cache) \longrightarrow latência baixa$ *hit* **falta** se não encontra $(dado \notin cache) \longrightarrow latência elevada$ *miss*

HEPR Dinf RCC 33

Arquitetura II — hierarquia de memória

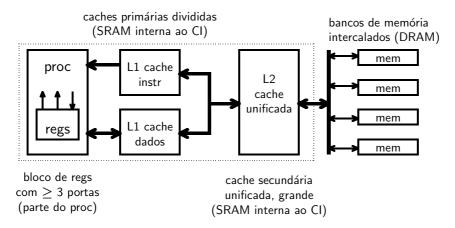
2007-1

Gerência da Hierarquia de Memória

- Gerência por software (registradores)
- ⊳ é parte do estado da CPU visível compilador/programador
- ⊳ SW tem controle completo sobre alocação de armazenadores mas HW pode fazer coisas às escondidas, como re-nomeação de registradores
- gerência por hardware (caches)
- ⊳ não é parte do estado visível
- ► HW decide automaticamente o que manter em memória rápida mas SW pode dar pistas,

como fazer busca antecipada ou marcar don't cache

Sistemas de Memória Típicos



"Lei de less": Processador super-escalar pode executar $pprox 10^3$ instruções durante falta na cache

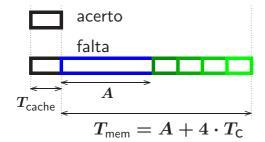
HEPR Dinf RCC 344

Arquitetura II — hierarquia de memória

2007-1

Tempo médio de acesso à memória

$$T_{\mathsf{med}} = T_{\mathsf{cache}} * (\mathsf{acertos} + \mathsf{faltas} * T_{\mathsf{mem}})$$



A = tempo de acesso à memória

 $T_{\mathsf{C}} = \mathsf{tempo}$ para copiar da memória para a cache (4 palavras por bloco)

Custo de acerto: \approx ciclo do processador penalidade por falta: \gg ciclo do processador

HEPR Dinf RCC

Arquitetura II — hierarquia de memória

2007-1

Memória - vazão

- Vazão: número de acessos por unidade de tempo [acessos/s]
- ⊳ afeta tempo de carga dos blocos da cache e operações de E/S
- vazão de pico ≤ largura/unidade_de_tempo
 vazão de pico é aquela impossível de ser sustentada
- vazão efetiva ≈ largura/(unidade_de_tempo + overhead)
- \triangleright se fração M das instruções referencia memória
 - ightarrow 1 + M referências à memória/instrução
 - ightarrow CPI=1 só se executa $\geq 1+M$ referências/ciclo

Memória – latência

• Latência: tempo para um único acesso

[s]

- ⊳ afeta todas as operações na memória
 - → afeta penalidade nas faltas em cache
- ⊳ não pode ser diminuída só com \$\$\$

mas pode ser escondida

⊳ no limite: velocidade da luz

 $\simeq 0.7c$ no silício/alumínio ≈ 1 palmo/ns

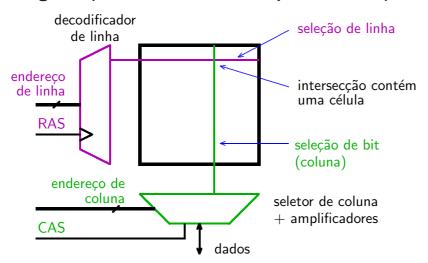
c = velocidade da luz (30cm/ns no vácuo)

HEPR Dist RCC 343

Arquitetura II — hierarquia de memória

2007-1

Organização da Memória Principal - introdução



HEPR Dinf RCC 344

Arquitetura II — hierarquia de memória

2007-1

Organização da Memória Principal - introdução

Memória principal usa DRAM
 Dynamic Random Access Memory

▶ latência: afeta penalidade nas faltas em cache

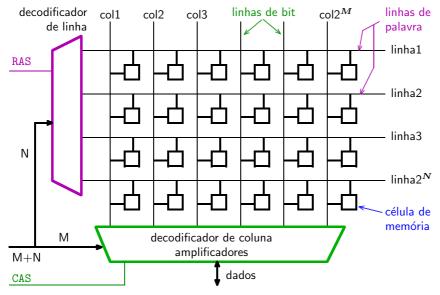
tempo de acesso: tempo entre requisição e dados prontos tempo de ciclo: tempo entre duas requisições

⊳ organizada como matriz (quadrado/retângulo) endereços multiplexados e tranferidos em 2 partes linhas, controlado por → RAS = row address strobe colunas, controlado por → CAS = column address strobe porque memória é organizada numa matriz 2D

• Cache usa SRAM Static Random Access Memory

ightharpoonup sem refresh (6 transistores/célula ao invés de 1) tamanho de células SRAM/DRAM pprox 4-8 custo e tempo de ciclo DRAM/SRAM pprox 8-16





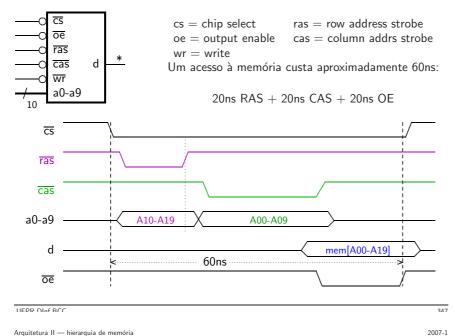
HEPR Dist RCC 34

Arquitetura II — hierarquia de memória

2007-1

2007-1

Sistemas de Memória



Memória Dinâmica

- DRAM é otimizada para densidade e não para velocidade
- ▷ célula com um transistor que se comporta como capacitor
- ⊳ bits armazenados como carga elétrica no capacitor
- ightharpoonup capacitor se descarrega na leitura \longrightarrow leitura destrutiva

- ciclo dura aprox dobro do tempo de acesso
- ⊳ é necessário pré-carregar as linhas de dados antes de acessar

Célula de Memória Dinâmica

Escrita:

- 1. ativa seleção de bit (bit line)
- 2. seleciona linha

Leitura:

- 1. carrega linha de bit até Vdd/2
- 2. seleciona linha
- 3. linha de bit e capacitor dividem carga
- 4. amplifica diferença de voltagem amplificador sente diferença de ${f 10}^6$ elétrons

5. escreve e reforça valor

transistor tipo P: $g (g=1) \longrightarrow (s \leftrightarrow d)$

seleção de linha

transistor

de passagem

capacitor

=bit $\{1,0\}$

Refresh:

1. lê conteúdo de cada célula

TIEDR Dinf RCC 340

seleção de bit

Arquitetura II — hierarquia de memória

Célula de Memória Dinâmica - escrita/leitura

Escrita:

- 1. força valor na linha de bit
- 2. seleciona linha
- 3. capacitor mantém valor por 60ms, então refresca

Leitura:

- 1. carrega linha de bit até Vdd/2
- 2. seleciona linha
- 3. linha de bit e capacitor dividem carga
- 4. amplificador detecta valor (1/0)
- 5. re-escreve valor

Refresh:

igual a leitura

Arquitetura II — hierarquia de memória

bit I linha 1

2007-1

linha

HEDR DIAFROC

DRAM - histórico

	capacidade	acesso [ns]	ciclo [ns]
ano	[Mbits]	mesma col	nova col
1980	0,064	150	250
1983	0,256	120	220
1986	1	100	190
1989	4	80	160
1992	16	60	120
1996	64	50	100
2000	256	40	90
2002	512	5	50
2004	1024	3	45
	60% aa	7% aa	

capacidade cresce 60% aa (melhor desde 96) tempo de acesso reduz 7% aa

Desempenho da Memória Principal

- Tempo de Ciclo \gg tempo de acesso (\approx 2:1) (leit/escr)
- Tempo de Ciclo (leit/escr)
- → com que freqüência pode iniciar um acesso?
- Tempo de Acesso (leit/escr)
- → quanto tempo demora para acessar o conteúdo, após iniciar acesso?
- Tempo de ciclo limita vazão da memória



LIEPR Dinf RCC

352

Arquitetura II - hierarquia de memória

2007-1

Memória Estática - SRAM

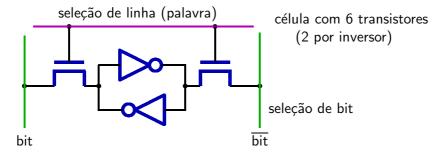
- Organização retangular (linha/coluna)
- ▷ linhas de endereços e dados não são multiplexadas
- ullet optimizada para 1° velocidade, e 2° densidade
- ▶ 4-6 transistores por célula
- bits armazenados em flip-flops/latches
- ⊳ estática → sem refresh
- tempo de acesso = duração do ciclo
- densidade/desempenho
- ▶ 1/4 da densidade de DRAM
- em 1998, capacidade de 1Mbit, tempo de acesso de 7 a 20ns

HEPR Dinf RCC

2007-1

Arquitetura II — hierarquia de memória

Célula de Memória Estática (Static RAM)



Escrita:

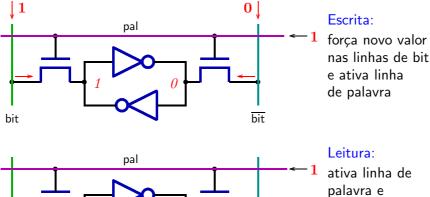
1. ativa linhas de bit bit=1, $\overline{\text{bit}}$ =0 (FF \leftarrow 1)

2. seleciona linha/palavra

Leitura:

- 1. carrega linhas de bit até Vdd/2
- 2. seleciona linha/palavra
- 3. célula puxa uma das linhas para 0
- 4. amplificador detecta diferença entre bit e bit

Célula de Memória Estática – escrita/leitura



bit bit o

bit armazenado

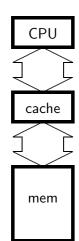
célula força

HEPR Dinf RCC

Arquitetura II — hierarquia de memória

2007-1

Memória Principal Simples (simplória?)



Parâmetros:

4 ciclos para emitir endereço

56 ciclos para acessar cada palavra

4 ciclos para transmitir palavra para cache/CPU

Penalidade para carregar bloco de 4 palavras: $(4 + 56 + 4) \times 4$ palavras = 256 ciclos

Simples:

CPU, cache, barramento e memória mesma largura (32 bits)

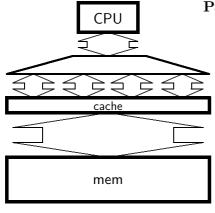
Barato mas desempenho ruim!

HEPR Dinf RCC

Arquitetura II — hierarquia de memória

2007-1

Memória Principal Larga



Parâmetros:

4 ciclos para emitir endereço 56 ciclos para acessar 4 palavras

4 ciclos para transf p cache/CPU

Penalidade para carregar bloco de 4 palavras: (4 + 56 + 4) = 64 ciclos

Largo:

CPU/mux 1 palavra mux/cache, barramento, memória N palavras (Alpha com 64 e 256 bits)

Desempenho bom mas caro!

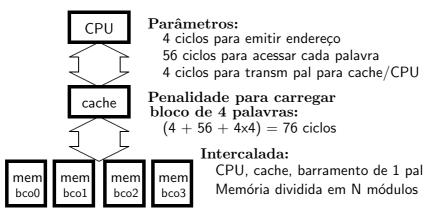
L1 com largura 1 e L2 com largura 4, remove o MUX do caminho entre L1 e CPU

HEPR Dinf RCC

357

2007-1

Memória Principal Intercalada



Boa relação de custo/desempenho

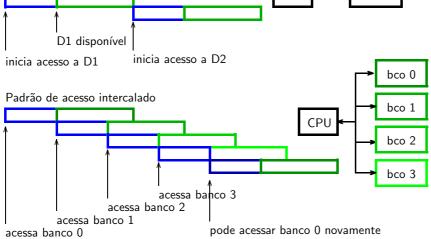
Arquitetura II — hierarquia de memória

2007-1

Intercalar para aumentar a vazão

Padrão de acesso SEM intercalar

CPU mem



HEPR Dinf RCC 35

Arquitetura II — hierarquia de memória

Intercalar para aumentar a vazão



- Palavras vizinhas alocadas em bancos distintos
- ▷ alocação definida pelos bits menos-significativos do endereço
- Bancos podem ser acessados concorrentemente
- ▷ se número de bancos ≫ tamanho do bloco da cache
- \triangleright acessos a endereços ≠s não-contíguos podem ser emitidos para bancos distintos: 1000,1004,1008,100c || 1014,1018,101c,1020 pressupõe que ∃ ≥ 8 bancos
- ullet número de bancos $\geq \#$ ciclos para acessar palavra num banco

Memória Intercalada

- Intercalação simples
- ⊳ necessita circuito de controle em cada banco
- ⊳ efetua acesso em paralelo
- ⊳ bom para sistemas com cache porque unidade de transfer é bloco
- bom com escrita preguiçosa (escreve todo o bloco sujo)
- Intercalação complexa (super-bancos)
- - → super-bancos endereçados individualmente
 - → interface de memória com pipeline

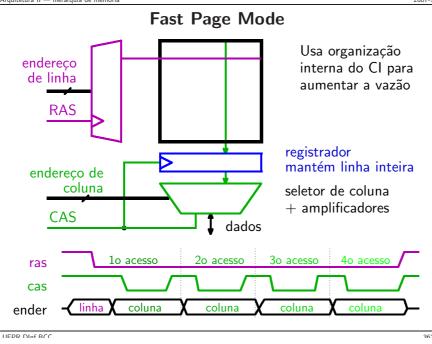
número do super-banco	deslocamento no super-banco	
	núm do banco	deslocamento no banco

HEPR DIst RCC 361

Arquitetura II - hierarquia de memóri

2007-1

2007-1



Fast Page Mode

• Fast Page Mode

Arquitetura II - hierarquia de memóri

- inventado para adaptadores de vídeo
- > registrador mantém toda uma linha da matriz
- Extended Data Out acrescenta latch entre mux e saída
 → permite mudar CAS mais cedo
- Synchronous DRAM acrescenta sinal de relógio para acelerar sinalização entre processador e controlador de memória
- Double Data Rate RAM como SDRAM, mas transmite dados nas duas bordas do relógio

Pergunta:

Considere um sistema de memória com tempo de acesso (mesma linha) de 14 ciclos. Quantos bancos são necessários para se obter o melhor desempenho ao menor custo?

- a) 4 bancos
- b) 8 bancos
- c) 16 bancos
- d) 32 bancos

HEPR DIst RCC 36.

Arquitetura II — hierarquia de memória

2007-1

Resposta:

Considere um sistema de memória com tempo de acesso (mesma linha) de 14 ciclos. Quantos bancos são necessários para se obter o melhor desempenho ao menor custo?

- a) 4 bancos
- b) 8 bancos
- c) 16 bancos # bancos $\geq \#$ ciclos, # bancos $= 2^B$ d) 32 bancos

HEPR Dinf RCC 36

Arquitetura II — hierarquia de memória

2007-1

Pergunta:

Considere um sistema com memória intercalada em 16 bancos de 32 bits. O bloco da cache é de quatro palavras.

Quantas referências concorrentes podem ser sustentadas?

LIEPR Dinf RCC 366

Resposta:

Considere um sistema com memória intercalada em 16 bancos de 32 bits. O bloco da cache é de quatro palavras.

Quantas referências concorrentes podem ser sustentadas?

Depende das combinações de tipos de acesso:

→ a quatro blocos de cache

(faltas de leitura, cache não-bloqueante, sem conflitos nos bancos)

→ escritas a 16 palavras consecutivas

(se barramento suportar)

ightarrow leitura de dois blocos, e escritas $n \leq 8$ palavras (se cache não-bloqueante e barramento suportarem) etc...

HEPR Dist RCC 367

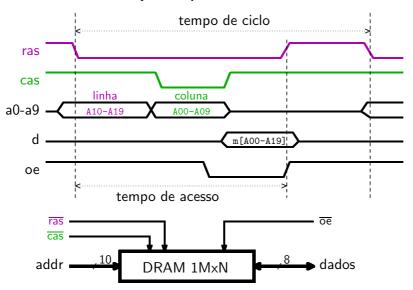
Arquitetura II — hierarquia de memória

Arquitetura II — hierarquia de memória

2007-1

2007-1

Temporização de Leitura



HEDR Dinf RCC 36

Parâmetros de Temporização de Leitura

 $\mathcal{T}_{\mathsf{RAC}}$ intervalo entre descida de RAS e dados válidos

- ★ quotado como velocidade da DRAM (tempo de acesso)
- * valor típico para DRAM de 4Mb é \mathcal{T}_{RAC} =60ns

 $\mathcal{T}_{\mathsf{RC}}$ intervalo mínimo entre o início de um acesso a uma linha e o próximo (tempo de ciclo)

- \star au_{RC} =110ns para DRAM de 4Mb com au_{RAC} =60ns
- \mathcal{T}_{CAC} intervalo mínimo entre descida do CAS e dados válidos
 - ★ T_{CAC} =15ns para DRAM de 4Mb com T_{RAC} =60ns

 $\mathcal{T}_{\mathsf{PC}}$ intervalo mínimo entre o início de um acesso a uma coluna e o próximo acesso

 \star T_{PC} =35ns para DRAM de 4Mb com T_{RAC} =60ns

Memória Intercalada – vazão entre L1 e L2

- Processador com relógio de 1GHz, 2 instruções por ciclo
- ★ 0.02 = taxa de faltas de instruções na L1i

1.25 refs/instrução

 \star 0.10 = taxa de faltas de dados na L1d

qual a demanda na L2?

 $1*0.02+0.25*0.10=0.045\;\text{faltas/instr}$

2 instr/ns * 0.045 faltas/instr = 0.09 faltas/ns

- a cada falta carrega 4 pals da L2 na L1d+L1i: $0.09 \text{ bloco/ns} * 4 \text{ pals/bloco} = 0.36 \text{ pals/ns} \longrightarrow 2.78 \text{ ns/pal}$
- 0.09 bloco/ns * 4 pals/bloco * 4 byte/pal → 694.4 Mbyte/s

bancos	SRAM 32bits, $\mathcal{T}_{ac}{=}12$ ns	OK?
4	$16b / (12+4) \; ns = 1.0 \; Gbyte/s$	> 0.694 Gb/s
8	32b / (12+8) ns = 1.6 Gbyte/s	\gg 0.694 Gb/s
16	64b / (12+16) ns = 2.28 Gbyte/s	$\gg 0.694 \text{ Gb/s}$

HEPR Dist RCC

Arquitetura II — hierarquia de memória

2007-1

Memória Intercalada - latência na L2

- Processador com relógio de 1GHz, 2 instruções por ciclo
- ★ 0.02 = taxa de faltas de instruções na L1i

1.25 refs/instrução

 \star 0.10 = taxa de faltas de dados na L1d

qual a demanda na L2?

1 * 0.02 + 0.25 * 0.10 = 0.045 faltas/instr * 2 = 0.09 faltas/ns

- a cada falta carrega 4 pals da L2 na L1d+L1i: $0.09 \text{ bloco/ns} * 4 \text{ pals/bloco} = 0.36 \text{ pals/ns} \longrightarrow 2.78 \text{ ns/pal}$
- quantos bancos para suprir demanda/vazão?

bancos	SRAM 32bits,	$oldsymbol{\mathcal{T}}_{ac}{=}12$ ns	OK?
4	3.0 ns/pal	$\gg 1$ ns	×
8	1.5 ns/pal	$>1{\sf ns}$	×
16	0.75 ns/pal	$< 1\mathrm{ns}$	\checkmark
	× vazão OK m	nas taxa xfer >	> ciclo

HEPR Dist RCC

Arquitetura II — hierarquia de memória

2007-1

Memória Intercalada - vazão entre L2 e RAM

- Processador com relógio de 1GHz, 2 instruções por ciclo
- ★ 0.02 = taxa faltas instruções L1i

 $0.10={\sf taxa}$ faltas dados L1d

 \star 0.075 = taxa de faltas de dados na L2

1.25 refs/instrução

qual a demanda na RAM?

(1*0.02+0.25*0.10)*0.075 = 0.003375 faltasL2/instr 2 instr/ns * 0.003375 faltas/instr = 0.00675 faltas/ns

• a cada falta carrega 8 pals da RAM para a L2:

6.75 bloco/ μ s * 32 bytes/bloco = 216 bytes/ μ s \rightarrow 216 Mb/s

bancos	$T_{RAC}=30$ ns, $T_{RC}=70$ ns	OK?
4	16b / (70+16) ns = 186 Mb/s	< 216 Mb/s
8	32b / (70+32) ns = 313 Mb/s	> 216 Mb/s
16	64b / (70+64) ns = 477 Mb/s	≫ 216 Mb/s
4 ns	nor ciclo do harramento (250MHz). DR	AM 32hits

HEPR Dinf RCC 37%

Balanceamento de Vazão entre CPU e Memória

- Requisitos de vazão pelo processador
- * processador super-escalar com ciclo de 1ns
- * sem cache de dados
- * 1 acesso à memória (64 bits) por ciclo do processador
- Vazão suportada pela memória
- ★ DRAM com ciclo de 90ns
- ★ Cls com 256Mbits

```
\longrightarrow 128 bancos para prover 1 palavra por ciclo
```

```
256Mb x 1 \longrightarrow 64x1 x 128 = 8192 chips \longrightarrow 256 Gbytes 64Mb x 4 \longrightarrow 16x4 x 128 = 2048 chips \longrightarrow 64 Gbytes
```

• Estes números não são razoáveis...

HEPR Disf RCC 375

Arquitetura II — otimização de caches

2007-1

Balanceamento ainda

- Adicionar uma cache de dados
- ★ taxa de faltas de 5%
- ★ blocos com 4 palavras de 64 bits
- * escrita preguiçosa
- ★ 25% das faltas são em blocos sujos
 - → demanda de aprox 1 pal a cada 4 ciclos demanda diminui por fator de 4
- Capacidade mínima ainda é impraticável:

```
256Mb x 1 \longrightarrow 64x1 x 128 / 4 = 2048 chips \longrightarrow 64 Gbytes
64Mb x 4 \longrightarrow 16x4 x 128 / 4 = 512 chips \longrightarrow 16 Gbytes
```

- Solução: adicionar mais um nível de cache mais caro que 16G?
- Outra Solução: usar Cls mais largos (x8, x16, x32)
- Outra Solução: usar CIs mais avançados fast page mode, EDO, Sdram, DRdram...

HEPRING 274

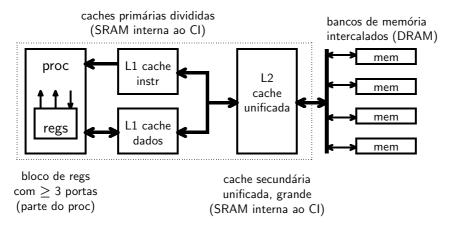
Arquitetura II — otimização de caches

2007-1

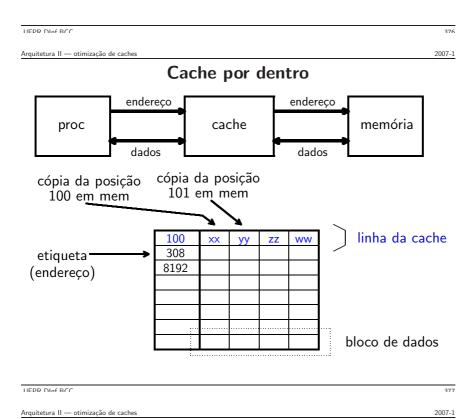
Cache

- Sistemas de memória
- Memória cache
- * organização
- * leitura
- * escrita
- * otimizações
- * busca antecipada
- Memória Virtual
- \star endereçamento
- ⋆ paginação
- ⋆ TLB
- * excessões

Sistemas de Memória Típicos



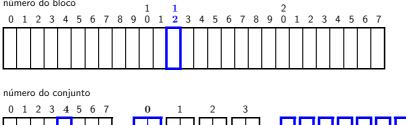
"Lei de Less": Processador super-escalar pode executar $pprox 10^3$ instruções durante falta na cache



Algoritmo para leitura

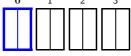
```
cache recebe endereço;
se objeto existe na cache (acerto hit)
entrega para processador;
senão (falta miss)
lê bloco da memória
espera o tempo de acesso à memória
entrega ao processador
e atualiza cache
```

Alocação de blocos



mapeamento direto

bloco 12 mapeia somente no bloco 4 12 % 8 = 4



associativa por conjuntos (binária)

bloco 12 mapeia no conjunto 0 12 % 4 = 0

associativa (total)

bloco 12 mapeia em qualquer conjunto

HEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

Mapeamento Direto

dado para processador

vál etiqueta bloco de dados

cache

dados para processador

endereço da referência

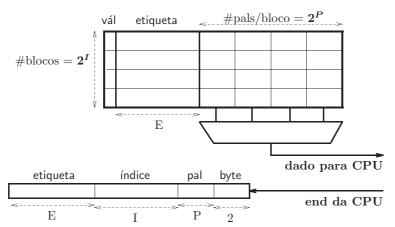
referência

etiqueta índice pal byte

endereço

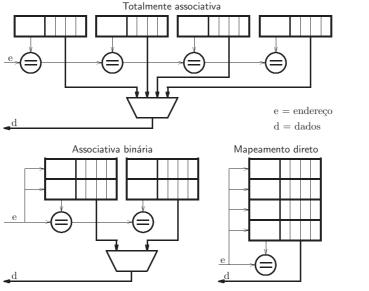
Arquitetura II — otimização de caches 2007-1

Organização da cache



tamanho da cache = num_blocos * tam_bloco * tam_palavra Quanto mais blocos \Longrightarrow menor a taxa de faltas Blocos com 4-16 palavras tiram vantagem de localidade espacial

Associatividade I



LIEPR Dinf BCC

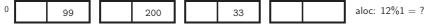
Arquitetura II — otimização de caches

2007-1

Associatividade II

M[12] = 99; M[13] = 200; M[14] = 33

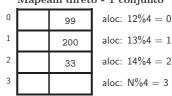
Associativa - 4 conjuntos



Associativa binária - 2 conjuntos



aloc: 12%2 = 0aloc: 13%2 = 1



 ${\bf Mapeam~direto~-~1~conjunto}$ aloc: 12%4 = 0

> aloc: 14%4 = 2aloc: N%4 = 3

Escolha de vítima para reposição: com mapeamento direto ∄ escolha...

1) escolhe o bloco usado no passado mais distante (LRU)

2) escolhe a esmo

HEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

Taxas de Faltas e de Acertos

Se encontra na cache — acerto Se não encontra na cache — falta

taxa de acertos
$$\stackrel{\triangle}{=} \frac{\text{número de acertos}}{\text{número de referências}}$$

taxa de faltas $\stackrel{\triangle}{=} 1$ — taxa de acertos

Pergunta:

Qual a taxa de faltas se |conj_dados| ≫ capacidade_da_cache?

- (a) baixa
- (b) média
 - (c) alta
- (d) faltas independem dos tamanhos relativos

HEPR Dief RCC 38

Arquitetura II — otimização de caches

2007-1

Resposta:

Qual a taxa de faltas se |conj_dados| ≫ capacidade_da_cache?

- (a) baixa
- (b) média
- (c) alta
- (d) faltas independem dos tamanhos relativos

Qual a taxa de faltas se |conj_dados| ≪ capacidade_da_cache?

HEPR Dinf RCC 386

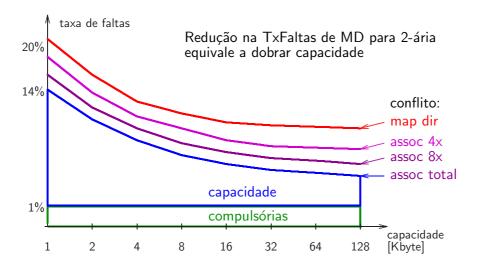
Arquitetura II — otimização de caches

2007-1

Causas de faltas

- faltas compulsórias linhas são tocadas pela primeira vez.
- ⊳ Solução: blocos com muitas palavras para fazer busca antecipada implícita;
- faltas por capacidade conj. de trabalho ≫ tamanho da cache.
- ⊳ Solução: aumentar tamanho da cache;
- faltas por conflito mapeamento de endereços ≠s em mesmo bloco da cache.
- ⊳ Solução: aumentar associatividade.
- > ocorrem mesmo quando a cache não está cheia;

Causas de faltas



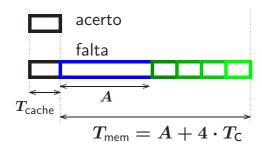
HEPR Dist RCC 38

Arquitetura II — otimização de caches

2007-1

Tempo médio de acesso à memória I

- ullet 1 ciclo por acerto na cache $(T_{\sf cache})$
- ullet penalidade: 2 a 100 ciclos para acessar memória (T_{mem})
- tMed = T_{cache} + (faltas · T_{mem})



A = tempo de acesso à memória

 $T_{\mathsf{C}} = \mathsf{tempo}$ para copiar da memória para a cache (4 palavras por bloco)

Custo de acerto: ≈ ciclo do processador penalidade por falta: ≫ ciclo do processador

HEPR Dinf RCC 3:

Arquitetura II — otimização de cache

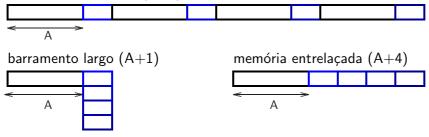
2007-1

Tempo médio de acesso à memória II

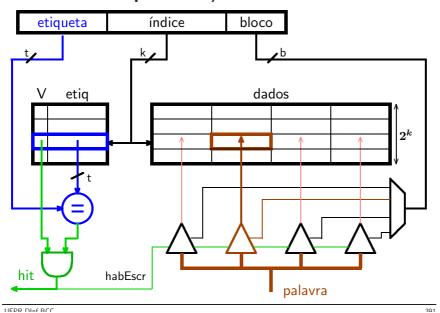
Tempo de acesso não pode ser (muito) reduzido tempo de transferência depende de projeto do barramento:

- * barramento estreito: uma palavra por ciclo (\$↓);
- * barramento/memória largos: mais de uma palavra por ciclo; (\$\epsilon);
- * memória entrelaçada: acesso paralelo e transferência serializada;

barramento estreito 4(A+1)



Implementação da Escrita

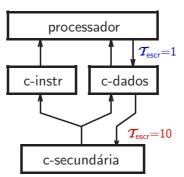


Arquitetura II — otimização de caches

2007-1

Políticas de escrita

- Escrita forçada (write-through)
- > propaga escrita até memória
- ▷ escritas completam na velocidade da memória
- Escrita preguiçosa (write-back)
- → acumula escritas na cache
 → bit sujo
- ⊳ falta em leitura custa propagação até memória



HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

Políticas de escrita

• Escrita forçada

write-through

- ★ cada escrita é propagada até a memória
- * escritas ocorrem na velocidade da memória

• Escrita preguiçosa

write-back

- * acumula escritas na cache
- ★ a cada escrita, o bit sujo é ligado
- * bloco permanece sujo na cache até que seja substituído
- ★ na substituição, se bloco vítima está sujo então o bloco inteiro é enviado para atualizar memória
- ★ falta na leitura pode causar escrita do bloco sujo

• Fila de Escrita

write-buffer

- \star imprescindível com escrita forçada
- \star com escrita forçada, fila tem uma palavra de largura
- ★ com escrita preguiçosa, fila tem um bloco de largura

O que fazer quando ocorre uma falta na escrita?

Aloca espaço na escrita

- write-allocate
- * espaço é alocado na cache para bloco faltante, e então é atualizado
- * se uma palavra no bloco foi atualizada, outras também o serão...
- ★ se cache com escrita preguiçosa, falta provoca até duas transações:
 - 1) se bloco está sujo, expurga-o
 - 2) carrega bloco faltante
- Não aloca espaço na escrita

no-write-allocate

- ★ não é alocado espaço na cache para bloco faltante
- * se não ocorreu falta de leitura, bloco pode não ser necessário...
- * bloco faltante é atualizado diretamente na memória
- * fila de escrita é imprescindível
- Combinações comuns
- ★ escrita forçada & não-alocação de espaço
- * escrita preguiçosa & alocação de espaço na escrita

fetch-on-write

LIEPR Dinf BCC

Arquitetura II — otimização de caches

2007-1

Fila de escrita I

Referências de escrita bloqueiam processador até que acesso à nível mais baixo da hierarquia complete

10% das referências são escritas: escritas na cache secundária custam ≈ 10 ciclos $o \mathcal{T}_{\text{mem}} = 0.9 * 1 + 0.1 * 10 = 1.9$ ciclos

solução: fila de escritas para desacoplar velocidade do processador da velocidade da memória

HEPR Dinf RCC

Arquitetura II — otimização de cache

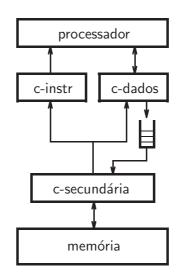
2007-1

Fila de escrita II

Fila reduz tempo médio da referência de escrita

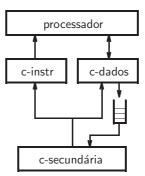
cada elemento da fila contém um par <endereço, valor>

controlador de cache efetua atualização da memória



Fila de escrita III

- Processador executa sw \$5,24(\$8)
- ⊳ se há espaço na fila, escrita completa em um ciclo
- ⊳ senão, processador bloqueia até abrir espaço na fila, enquanto escrita anterior é propagada até L2
- Fila tem capacidade para 2-16 registros
- fila fica cheia na entrada de funções com muitos parâmetros...



HEPR Dist RCC 30:

Arquitetura II — otimização de caches

2007-1

revisão - Revisão de Caches

- Localidade Temporal
- * objeto será referenciado novamente no futuro próximo: pilha, código
- LocalidadeEspacial
- * objetos vizinhos serão referenciados no futuro próximo: vetores, instruções
- Taxa de Acertos = núm_acertos / núm_referências
- \star Taxa de Faltas = 1 Taxa_de_acertos
- 3 Cs: faltas compulsórias, por conflitos, por capacidade
- ★ compulsórias: são compulsórias...
- ★ conflitos: endereços mapeiam no mesmo bloco ↑associatividade
- ★ capacidade: |conj de dados| > |cache|; ↑tamanho
- Políticas de escrita
- * escrita forçada, escrita preguiçosa, fila de escrita
- * nas faltas: aloca espaço, não-aloca espaço

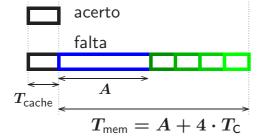
HEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

Tempo médio de acesso à memória

$$T_{ extsf{med}} = T_{ extsf{cache}} + (extsf{faltas} \cdot T_{ extsf{mem}})$$



A = tempo de acesso à memória

 $T_{\mathsf{C}} = \mathsf{tempo}$ para copiar da memória para a cache (4 palavras por bloco)

Custo de acerto: ≈ ciclo do processador penalidade por falta: ≫ ciclo do processador

Desempenho de caches

tempo de CPU =
$$\#Instr \times (CPI_{exec} + CPI_{mem}) \times ciclo$$

$$CPI_{mem} = \frac{numRefs}{instr} \times (T_{acerto} + txFaltas \cdot penalidade)$$

Tempo médio de acesso à memória

 $T_{\mathsf{mem}} = T_{\mathsf{acerto}} + \mathsf{txFaltas} \cdot \mathsf{penalidade}$

HEPR Dist RCC 40

Arquitetura II — otimização de caches

2007-1

Melhoria no desempenho de caches

Tempo médio de acesso a memória = tempo de acerto + taxa de faltas x penalidade pela falta

Para melhorar o desempenho é necessário:

- reduzir a taxa de faltas
- reduzir a penalidade
- reduzir o tempo de acerto

Estratégia de projeto mais simples é projetar a cache primária maior possível, sem elongar o ciclo de relógio, e sem estágios adicionais no pipeline fica mais complicado com emissão fora-de-ordem

HEPR Dist RCC 40:

Arquitetura II — otimização de caches

2007-1

Reduzir a taxa de faltas

- Faltas compulsórias: primeira referência a um bloco (a frio)
- ⊳ faltas que ocorreriam numa cache infinita
- Faltas por capacidade: cache não comporta conjunto de trabalho
- ▶ faltas que ocorreriam mesmo com políticas perfeitas de alocação e reposição
- Faltas por conflito: colisões no mapeamento de blocos na cache
- ▶ faltas que não ocorreriam se a cache fosse totalmente associativa e reposição fosse com LRU
- ⊳ podem ocorrer com cache não-cheia

Reduzir a taxa de faltas – parâmetros de projeto

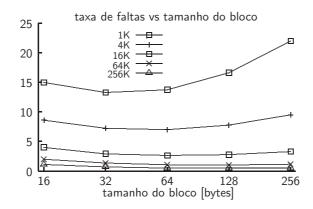
- Maior capacidade
- ▷ reduz faltas por capacidade e por conflito
- Tamanho do bloco
- ▷ localidade espacial reduz faltas compulsórias
- ▶ menor número de blocos pode aumentar conflitos
- bloco maior pode aumentar penalidade (tempo de preenchimento)
- Associatividade
- ⊳ reduz faltas por conflito (até associatividade 4-8)
- ▶ pode aumentar tempo de acesso

HEPR DINFRCC ADD

Arquitetura II — otimização de caches

2007-1

1. Reduzir faltas - Tamanho do Bloco



- localidade espacial reduz faltas compulsórias
- para mesmo tamanho, menor número de blocos aumenta conflitos
- bloco maior aumenta tempo de preenchimento >penalidade

HEPR DIAF RCC

Arquitetura II — otimização de caches 2007-1

2. Reduzir faltas - Associatividade

Regra 2:1

taxa de faltas de uma cache com mapeamento direto de tamanho $oldsymbol{N}$

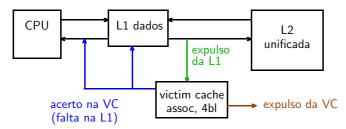
taxa de faltas de cache com associatividade binária de tamanho N/2

Achtung: desempenho é medido pelo tempo de execução!! maior associatividade pode aumentar tempo de acesso

Em [Hill88] sugere-se que relação entre tempo de acerto de assoc-binária / mapeamentoDireto é
1,02 para caches internas e
1,10 para caches externas

dados estarão disponíveis somente após comparação da etiqueta

3. Reduzir faltas - Cache de Vítimas



Cache de Vítimas é uma cache associativa pequena (4-8 blocos), ligada à cache com mapeamento direto para conter linhas recentemente expurgadas

[Jouppi90]: CdV com 4 elementos remove de 20% a 95% dos conflitos numa cache de 4Kbytes com mapeamento direto;

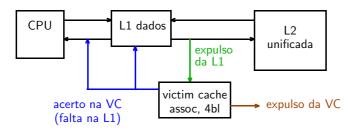
Usado no Alpha e máquinas da HP

HEPR Dist RCC 40

Arquitetura II — otimização de caches

2007-1

3.1 Reduzir faltas – Cache de Vítimas



- Procura na L1;
- se falta, procura da cache de vítimas;
- se encontrou, troca linha por aquela que vai ser expulsa da L1;
- se falta na CdV, vítima da L1 → CdV;
- Efeito combinado: tempo de acesso do mapeamento direto com redução nas faltas por conflito.

HEPR Dist RCC

Arquitetura II — otimização de cache

2007-1

4. Reduzir faltas – Busca antecipada por hardware

- Busca antecipada de instruções
- ⊳ Alpha 21064 busca dois blocos numa falta
- ⊳ bloco extra é colocado num stream buffer
- ▷ numa falta, procura no stream buffer
- Busca antecipada de dados
- ▶ 1 stream buffer para dados resolve 25% das faltas
 numa cache de 4KB; 4 stream buffers resolvem 43% [Jouppi90]
- ⊳ para aplicações científicas, 8 streams resolveram 50-70% das faltas com duas caches de 64KB, assoc-quaternária [Palacharla94]
- Busca antecipada necessita de sistema de memória com vazão sobrando e que possa ser utilisada sem penalidade

5. Reduzir faltas – Busca antecipada por software

- Busca antecipada de dados
- * carregar dados em registrador (HP PA-RISC load)
- * carregar dados para a cache (MIPS-4, PowerPC)
- ★ instruções de busca antecipada não podem causar faltas de página
 → é uma forma de execução especulativa
- Executar instruções de busca antecipada custa tempo
- ★ custo da busca antecipada é menor que os ganhos pelos acertos?
- * processadores super-escalares gulosos facilitam despacho destas instruções
- Busca antecipada causa poluição na cache
- * carregar blocos em stream buffers e não na cache
- \star numa falta, carrega stream buffers com blocos após o faltante
- * numa falta, procura na cache e nos SBs
- \star 1 SB é bom para instruções \geq 1 SBs para dados (pprox 4..8)

HEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

6. Reduzir faltas – Otimização pelo compilador

- Instruções
- > reordenar funções na memória para reduzir faltas por conflito
- Dados
- ▶ agrupar vetores: acessos a vetor de elementos compostos tem melhor localidade espacial que 2 vetores
- ▶ troca de loops: mudar aninhamento dos loops para acessar dados na ordem de armazenamento em memória
- ▶ fusão de loops: combinar 2 loops independentes com comportamento similar e sobreposição de variáveis
- blocagem: melhorar a localidade temporal ao acessar "blocos" de dados repetidamente ao invés de percorrer linhas/colunas inteiras

HEPRING RCC 411

Arquitetura II — otimização de caches

2007-1

Exemplo: agrupamento de vetores

```
/* Antes: 2 vetores */
int val[TAM];
int key[TAM];

/* Depois: 1 vetor de estruturas */
struct agrup {
  int val;
  int key;
}
struct agrup doisVetores[TAM];
```

Redução de conflitos entre val e key: melhora localidade espacial

Exemplo: troca de loops

```
for (k=0; k < 100; k++)
                                     /* Antes */
2
      for (j=0; j < 100; j++)
3
         for (i=0; i < 5000; i++)
            x[i][j] = 2 * x[i][j];
                                    /* |passo|=100 */
   for (k=0; k < 100; k++)
                                     /* Depois */
5
6
      for (i=0; i < 5000; i++)
7
         for (j=0; j < 100; j++)
8
            x[i][j] = 2 * x[i][j]; /* |passo| = 1 */
```

Acessos seqüenciais ao invés de passadas de 100 palavras

→ melhora localidade espacial

passada ≡ *stride*

2007-1

TIEDR Dief RCC A1'

Exemplo: fusão de loops

Arquitetura II — otimização de caches

```
/* Antes: duas faltas por acesso em a e em c */
1 for (i=0; i < N; i++)
2
      for (j=0; j < N; j++)
         a[i][j] = 1/b[i][j] * c[i][j];
3
4
  for (i=0; i < N; i++)
5
      for (j=0; j < N; j++)
         d[i][j] = a[i][j] + c[i][j];
/* Depois: uma falta por acesso em a e em c */
7
   for (i=0; i < N; i++)
8
      for (j=0; j < N; j++) {
9
         a[i][j] = 1/b[i][j] * c[i][j];
         d[i][j] = a[i][j] + c[i][j];
            /* re-uso e melhor localidade espacial */
```

HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

Exemplo: blocagem

- Dois loops internos:
- ★ lê todos os NxN elementos de z[]
- \star lê N elementos de uma linha de y[] repetidamente
- \star escreve N elementos de 1 coluna de x []
- Faltas por capacidade são função de N e tamanho da cache
- \star núm faltas $=2N^3+N^2$, se não ocorrerem conflitos...
- Idéia: computar sub-matriz que cabe na cache

IIFPR Dinf RCC

Exemplo: blocagem

```
for (i=0; i < N; i++)
                                            /* Antes */
2
       for (j=0; j < N; j++) {
3
          for (r=0.0, k=0; k < N; k++)
4
              r = r + y[i][k] * z[k][j];
5
          x[i][j] = r;
6
           j
                                                       j
   Z
                                              Х
                        У
 k^{\ 3}
        não tocado
                       tocado há muito
                                           tocado há pouco
```

LIEPR Dinf RCC

2007-1

Arquitetura II — otimização de cache

Exemplo: blocagem

```
for (jj=0; jj < N; jj = jj+B)
                                        /* Depois */
   for (kk=0; kk < N; kk = kk+B)
   for (i=0; i < N; i++)
4
      for (j = jj; j < min(jj+B,N); j++) {
         for (r = 0.0, k = kk; k < min(kk+B,N); k++)
5
6
            r = r + y[i][k] * z[k][j];
7
         x[i][j] = x[i][j] + r;
                                        /* x[i][j] parcial */
      }
8
```

- B é chamado de fator de blocagem
- ullet Faltas por capacidade caem de $2N^3+N^2$ para $2N^3/B+N^2$
- Faltas por conflito talvez diminuem

LIEPR Dinf RCC

Arquitetura II — otimização de cache

2007-1

Exemplo: blocagem

```
/* Depois */
   for (jj=0; jj < N; jj = jj+B)
   for (kk=0; kk < N; kk = kk+B)
   for (i=0; i < N; i++)
4
      for (j = jj; j < min(jj+B,N); j++) {
5
         for (r = 0.0, k = kk; k < min(kk+B,N); k++)
6
            r = r + y[i][k] * z[k][j];
7
         x[i][j] = x[i][j] + r; 
                                       /* x[i][j] parcial */
  Z
                                          Х
                      У
 k
```

Melhoria no desempenho de caches

Tempo médio de acesso a memória = tempo de acerto + taxa de faltas x penalidade pela falta

Para melhorar o desempenho é necessário:

- reduzir a taxa de faltas
- reduzir a penalidade nas faltas
- reduzir o tempo de acerto

HEPR Dist RCC 41

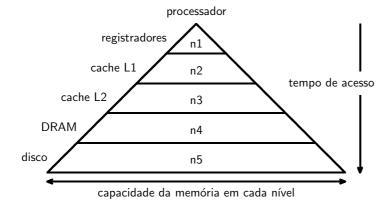
Arquitetura II — otimização de caches

2007-1

A. Reduzir a penalidade – hierarquia de caches

Memória não pode ser simultaneamente grande e rápida;

→ hierarquia com caches maiores mais longe do processador

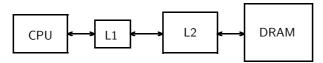


HEPR Dist RCC 410

Arquitetura II — otimização de caches

2007-1

A.a Reduzir a penalidade – hierarquia de caches



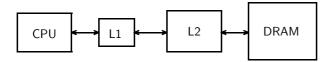
 $\mbox{Taxa de faltas local} = \mbox{faltas na cache} \; / \; \mbox{acessos à cache} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas na cache} \; / \; \mbox{referências pelo processador} \\ \mbox{Taxa de faltas pelo processador} \\ \mbox{Taxa de faltas global} = \mbox{faltas pelo processador} \\ \mbox{Taxa de faltas pelo processador} \\ \mbox{Taxa de$

$$T_{mem} = tA_{L1} + F_{L1} \times P_{L1}$$

= $tA_{L1} + F_{L1} \times [tA_{L2} + F_{L2} \times P_{DRAM}]$

onde tA é o tempo de acerto, $m{F}$ é taxa de faltas, e $m{P}$ é a penalidade por falta

A.b Reduzir a penalidade - hierarquia de caches



Propriedade de Inclusão:

Conteúdo das caches maiores *inclui* o conteúdo das menores —> cache menor contém cópias dos dados na cache maior

- ▶ qual a relação entre associatividade, tamanho e inclusão?
- > qual a relação entre bits de tag/índice e assoc, tamanho e inclusão?
- ▶ [Baer, Wang88] On the Inclusion Properties for Multi-Level Cache Hierarchies

Caches *exclusivas* trocam linhas numa falta \rightarrow L1 \leftarrow novo: L2 \leftarrow velho

TIEPR Dinf RCC 42

Arquitetura II — otimização de caches

B. Red penalidade - prioridade para faltas de leitura

- Escrita forçada com fila de escrita causa risco RAW
- ▷ processador tenta ler dado que está na fila de escrita
- ⊳ solução ruim: espera até que fila esvazie e então continua LOAD

 → aumenta penalidade nas faltas de leitura
- ⊳ solução melhor: verifica se ∃m dependências; se não, LOAD prossegue sem bloquear
- Escrita preguiçosa: falta na leitura substitui bloco sujo
- ⊳ solução ruim: escreve bloco sujo e então busca bloco faltante
- ⊳ solução melhor:
 - 1. escreve bloco sujo num buffer/fila,
 - 2. lê bloco faltante e entrega ao processador,
 - 3. então escreve bloco sujo
- → rápido porque CPU prossegue assim que ler bloco faltante

HEPR DInf RCC

Arquitetura II — otimização de caches 2007-

C. Reduzir a penalidade – minimizar tempo de carga

- Tempo de preenchimento do bloco é longo
- ...mas não precisa esperar até que bloco seja preenchido
- Early restart assim que palavra requisitada chegar da memória, entrega ao processador, que continua a executar
- Critical word first busca palavra requisitada primeiro e a entrega ao processador assim que chegar da memória
- > processador continua enquanto bloco está sendo preenchido
- ⊳ processador requisita p2:
 p0 p1 p2 p3
- ho **p2** é entregue, e bloco é preenchido com **p2** ightarrow p3 ightarrow p0 ightarrow p1
- ▷ localidade espacial indica que próximo acesso será na palavra seguinte, que pode não ter chegado ainda

D. Reduzir a penalidade – cache não-bloqueante

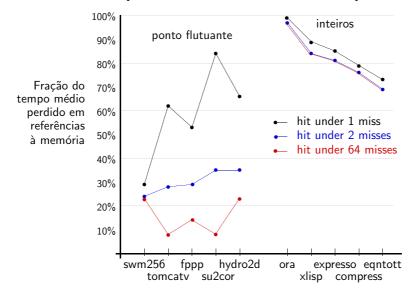
- Cache não-bloqueante entrega dados enquanto trata de uma falta
 processador com execução fora-de-ordem
- ▶ [Kroft81] Lock-up Free Instruction Fetch/Prefetch Cache Organization
- acerto sob falta (hit under miss) reduz penalidade efetiva porque cache fornece dados, ao invés de bloquear processador na falta
- acerto sob múltiplas faltas (miss under miss) reduz penalidade ainda mais ao sobrepor tratamento de múltiplas faltas
- > necessita de memória com múltiplos bancos
- ▶ PentiumPro permite até 4 faltas pendentes

HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

D.1 Reduzir a penalidade – cache não-bloqueante



HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

Melhoria no desempenho de caches

Tempo médio de acesso a memória = tempo de acerto + taxa de faltas \times penalidade pela falta

Para melhorar o desempenho é necessário:

- reduzir a taxa de faltas
- reduzir a penalidade nas faltas
- reduzir o tempo de acerto

I. Reduzir o tempo de acerto

caches simples e pequenas

Tempo de acerto curto se:

Alpha 21164 tem 8KB para dados e 8KB para instruções, 96KB cache secundária, unificada, associatividade ternária, e relógio muito rápido!

LIEPR Dinf BCC

Arquitetura II — otimização de cache

2007-1

II. Reduzir o tempo de acerto - evitar tradução de endereços

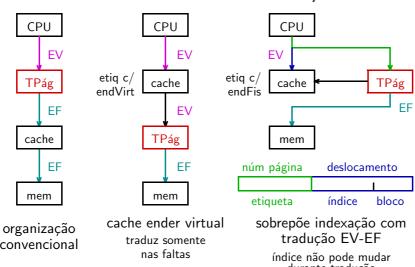
- Cache endereçada com: endereço virtual vs endereço físico
- a cada troca de contexto deve expurgar a cache
- * senão ocorrem falsos-acertos
- * processo que entra no processador tem montes de faltas compulsórias
- \star conseqüência da troca de contexto pode durar 10^5 ciclos
- * solução: usar identificador de processo na etiqueta na cache
- como trata de sinônimos?
- ★ dois endereços virtuais que mapeiam no mesmo endereço físico
- * solução: garantir que indexação da cache ocorre com bits que coincidem no EV e EF

EVirtual	núm página virtual	deslocamento
ender na cache	—índice—	
EFísico	núm página física	deslocamento

LIEPR Dinf BCC

2007-1

II.i Reduzir T de acerto – evitar tradução de ender



convencional

durante tradução |cache| = |página|

LIEPR Dinf RCC

III.i Reduzir o tempo de acerto – escritas em pipeline

Escrita consome dois ciclos

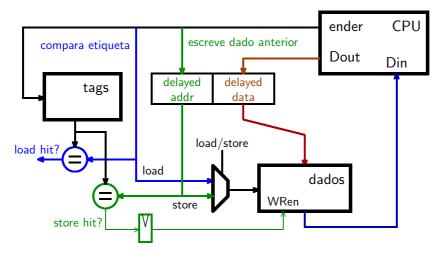
- um para comparar etiqueta e outro para escrever na cache
- armazena dados e endereço em buffer
- escreve dados na cache durante comparação de etiqueta na próxima escrita
- risco RAW entre leitura e escrita atrasada → adiantamento

HEPR Dist RCC 43

Arquitetura II — otimização de caches

2007-1

III.ii Reduzir tempo de acerto – escritas em pipeline



HEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

Resumo – melhorar desempenho de caches

 $T_méd_acesso_mem = T_acerto + tx_faltas \times penalidade/falta$

Para melhorar o desempenho é necessário:

★ reduzir a taxa de faltas

- > parâmetros de projeto: núm,tam blocos, associatividade
- > cache de vítimas, stream buffers, busca antecipada
- Dotimização do código: agrupamento, troca de índices, fusão, blocagem

* reduzir a penalidade nas faltas

- ▶ hierarquia de caches
- ▷ priorizar faltas na leitura (critical word first, early restart)

* reduzir o tempo de acerto

- ▷ não traduzir endereço virtual para físico
- ▷ escritas em pipeline

Cache

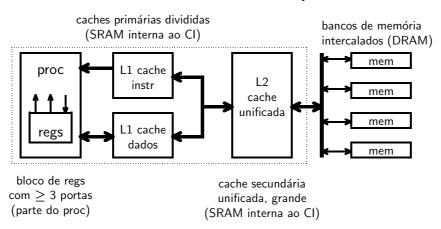
- Sistemas de memória
- Memória cache
- * organização
- * leitura
- * escrita
- * otimizações
- * busca antecipada
- Busca antecipada
- Memória Virtual
- * endereçamento
- ⋆ paginação
- ⋆ TLB
- ⋆ excessões

HEPR Dinf RCC 433

Arquitetura II — otimização de caches

2007-1

Sistemas de Memória Típicos



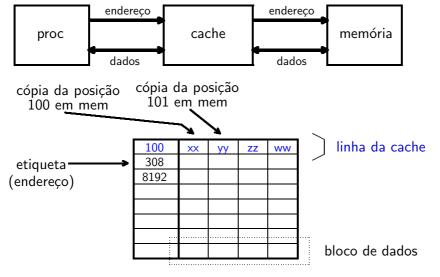
"Lei de Less": Processador super-escalar pode executar $pprox \mathbf{10^3}$ instruções durante falta na cache

HEPR DInf RCC 424

Arquitetura II — otimização de caches

2007-1

Cache por dentro



Algoritmo para leitura

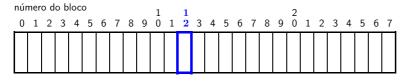
```
cache recebe endereço;
se objeto na cache (acerto hit)
entrega para processador;
senão (falta miss)
lê bloco da memória
espera
entrega ao processador
e atualiza cache
```

HEPR Dist RCC 433

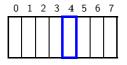
Arquitetura II — otimização de caches

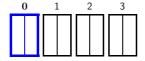
2007-1

Alocação de blocos



número do conjunto







mapeamento direto

bloco 12 mapeia somente no bloco 4 12 % 8 = 4

associativa por conjuntos (binária)

bloco 12 mapeia no conjunto 0 12 % 4 = 0 associativa (total)

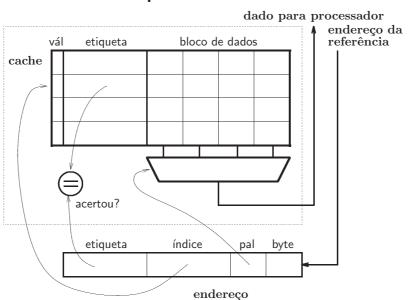
bloco 12 mapeia em qualquer conjunto

IIEPR Dinf RCC

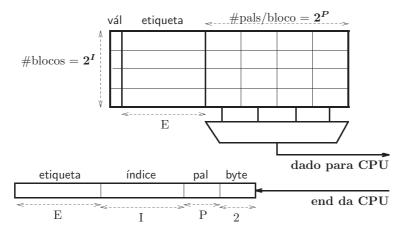
Arquitetura II — otimização de cache

2007-1

Mapeamento Direto



Organização da cache III



tamanho da cache = num_blocos * tam_bloco * tam_palavra Quanto mais blocos ⇒ menor a taxa de faltas Blocos com 4-16 palavras tiram vantagem de localidade espacial

HEPR Dinf RCC Arquitetura II — otimização de caches 2007-1 Associatividade II M[12] = 99; M[13] = 200; M[14] = 33Associativa - 4 conjuntos aloc: 12%1 = ?99 33 ${\bf Mapeam\ direto\ -\ 1\ conjunto}$ Associativa binária - 2 conjuntos aloc: 12%4 = 0aloc: 12%2 = 099 aloc: 13%4 = 1aloc: 13%2 = 1aloc: 14%4 = 233 aloc: N%4 = 3

Escolha de vítima para reposição: com mapeamento direto ∄ escolha...

1) escolhe o bloco usado no passado mais distante (LRU)

2) escolhe a esmo

LIEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

Taxas de Faltas e de Acertos

Se encontra na cache — **acerto** Se não encontra na cache — **falta**

taxa de acertos
$$\stackrel{\triangle}{=} \frac{\text{número de acertos}}{\text{número de referências}}$$

taxa de faltas $\stackrel{\triangle}{=} 1$ — taxa de acertos

Pergunta:

Qual a taxa de faltas se |conj_dados| ≫ capacidade_da_cache?

- (a) baixa
- (b) média
 - (c) alta
- (d) faltas independem dos tamanhos relativos

HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

Resposta:

Qual a taxa de faltas se |conj_dados| ≫ capacidade_da_cache?

- (a) baixa
- (b) média
 - (c) alta
- (d) faltas independem dos tamanhos relativos

Qual a taxa de faltas se |conj_dados| ≪ capacidade_da_cache?

HEPR Dinf RCC 441

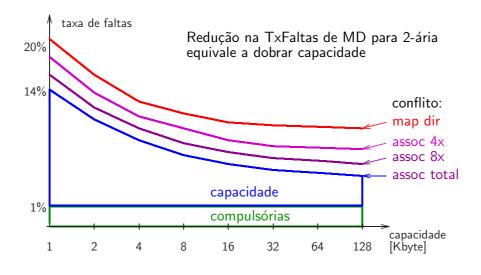
Arquitetura II — otimização de caches

2007-1

Causas de faltas

- faltas compulsórias linhas são tocadas pela primeira vez.
- \star Solução: blocos com muitas palavras para fazer busca antecipada implícita;
- ullet faltas por capacidade conj. de trabalho \gg tamanho da cache.
- * Solução: aumentar tamanho da cache;
- faltas por conflito mapeamento de endereços no mesmo bloco da cache.
- * Solução: aumentar associatividade.
- ★ Faltas por conflito ocorrem mesmo quando a cache não está cheia; estas se devem aos conflitos de mapeamento nos blocos da cache

Causas de faltas



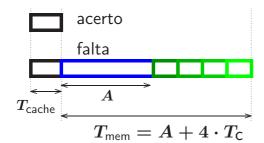
HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

Tempo médio de acesso à memória I

- 1 ciclo por acerto na cache (tCache)
- penalidade: 2..100 ciclos para acessar memória (tMem)
- tMed = tCache * (acertos + faltas * tMem)



A = tempo de acesso à memória

 $T_{\mathsf{C}} = \mathsf{tempo}$ para copiar da memória para a cache (4 palavras por bloco)

Custo de acerto: ≈ ciclo do processador penalidade por falta: ≫ ciclo do processador

HEPR Dist RCC

Arquitetura II — otimização de cache

2007-1

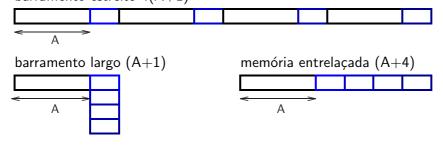
Tempo médio de acesso à memória II

Tempo de acesso não pode ser (muito) reduzido

tempo de transferência depende de projeto do barramento:

- ★ barramento estreito: uma palavra por ciclo (\$↓);
- ★ barramento/memória largos: mais de uma palavra por ciclo; (\$↑);
- ★ memória entrelaçada: acesso paralelo e transferência serializada;

barramento estreito 4(A+1)



2007-1

Implementação da Escrita

Arquitetura II — otimização de caches

2007-1

etiqueta índice bloco

t dados

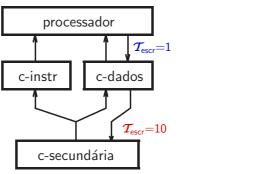
hit habEscr palavra

Políticas de escrita

- Escrita forçada (write-through)
- ⋆ propaga escrita até memória

Arquitetura II — otimização de cache

- * escritas com velocidade da memória
- Escrita preguiçosa (write-back)
- \star acumula escritas na cache \longrightarrow bit sujo
- \star falta em leitura custa propagação até memória



LIFPR Dinf RCC 450

Políticas de escrita

Escrita forçada

write-through

- * cada escrita é propagada até a memória
- * escritas ocorrem na velocidade da memória

• Escrita preguiçosa

write-back

- * acumula escritas na cache
- * a cada escrita, o bit sujo é ligado
- * bloco permanece sujo na cache até que seja substituído
- na substituição, se bloco vítima está sujo então bloco inteiro enviado para atualizar memória
- * falta na leitura pode causar escrita do bloco sujo

Fila de Escrita

write-buffer

- * imprescindível com escrita forçada
- * com escrita forçada, fila tem uma palavra de largura
- * com escrita preguiçosa, fila tem um bloco de largura

TIEPR DIst RCC

Arquitetura II — otimização de caches

2007-1

O que fazer quando ocorre uma falta na escrita?

Aloca espaço na escrita

write-allocate

- * espaço é alocado na cache para bloco faltante, e então é atualizado
- ★ se uma palavra no bloco foi atualizada, outras também o serão...
- ★ se cache com escrita preguiçosa, falta provoca até duas transações:
 - 1) se bloco está sujo, expurga-o
 - 2) carrega bloco faltante

Não aloca espaço na escrita

no-write-allocate

- * não é alocado espaço na cache para bloco faltante
- * se não ocorreu falta de leitura, bloco pode não ser necessário...
- * bloco faltante é atualizado diretamente na memória
- * fila de escrita é imprescindível
- Combinações comuns
- ★ escrita forçada & não-alocação de espaço
- ★ escrita preguiçosa & alocação de espaço na escrita

fetch-on-write

HEPR DInf RCC

Arquitetura II — otimização de cache

2007-1

Fila de escrita I

Referências de escrita bloqueiam processador até que acesso à nível mais baixo da hierarquia complete

10% das referências são escritas; escritas na cache secundária custam ≈ 10 ciclos $\Longrightarrow T_{\rm mem} = 0.9*1 + 0.1*10 = 1.9$ ciclos

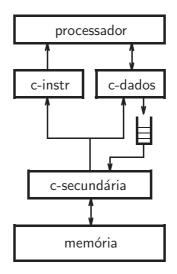
solução: fila de escritas para desacoplar velocidade do processador da velocidade da memória

Fila de escrita II

Fila reduz tempo médio da referência de escrita

cada elemento da fila contém um par <endereço, valor>

controlador de cache efetua atualização da memória



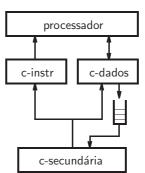
HEPR Dist RCC 454

Arquitetura II — otimização de caches

2007-1

Fila de escrita III

- Processador executa sw \$5,24(\$8)
- * na fase MEM, processador insere na fila par <(24+\$8), \$5>
- ★ se há espaço na fila, escrita completa em um ciclo
- ★ senão, processador bloqueia até abrir espaço na fila, enquanto escr anterior é propagada até L2
- Fila tem capacidade para 2-16 registros
- ★ fila enche na entrada de funções com muitos parâmetros...



HEPR Dinf RCC 455

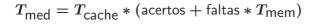
Arquitetura II — otimização de caches

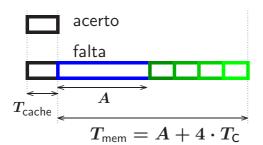
2007-1

revisão - Revisão de Caches

- Localidade Temporal
- \star objeto será referenciado novamente no futuro próximo: pilha, código
- LocalidadeEspacial
- * objetos vizinhos serão referenciados no futuro próximo: vetores, instruções
- Taxa de Acertos = núm_acertos / núm_referências
- ★ Taxa de Faltas = 1 Taxa_de_acertos
- 3 Cs: faltas compulsórias, por conflitos, por capacidade
- ★ compulsórias: são compulsórias...
- ★ conflitos: endereços mapeiam no mesmo bloco ↑associatividade
- * capacidade: |conj de dados| > |cache|; ↑tamanho
- Políticas de escrita
- * escrita forçada, escrita preguiçosa, fila de escrita
- ★ nas faltas: aloca espaço, não-aloca espaço

Tempo médio de acesso à memória





A = tempo de acesso à memória

 $T_{\mathsf{C}} = \mathsf{tempo}$ para copiar da memória para a cache (4 palavras por bloco)

Custo de acerto: \approx ciclo do processador penalidade por falta: \gg ciclo do processador

LIEPR Dinf BCC

457

Arquitetura II — otimização de caches

2007-1

Desempenho de caches

tempo de CPU
$$= \#Instr \times (CPI_{exec} + CPI_{mem}) \times ciclo$$

$$\mathsf{CPI}_{\mathsf{mem}} \ = \ \frac{\mathsf{numRefs}}{\mathsf{instr}} \times (T_{\mathsf{acerto}} + \mathsf{txFaltas} \cdot \mathsf{penalidade})$$

Tempo médio de acesso à memória

$$T_{\mathsf{mem}} = T_{\mathsf{acerto}} + \mathsf{txFaltas} \cdot \mathsf{penalidade}$$

HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

Melhoria no desempenho de caches

Tempo médio de acesso a memória = tempo de acerto + taxa de faltas \times penalidade pela falta

Para melhorar o desempenho é necessário:

- reduzir a taxa de faltas
- reduzir a penalidade
- reduzir o tempo de acerto

Estratégia de projeto mais simples é projetar a cache primária maior possível, sem elongar o ciclo de relógio, e sem estágios adicionais no pipeline fica mais complicado com emissão fora-de-ordem

Reduzir a taxa de faltas

- Faltas compulsórias: primeira referência a um bloco (partida a frio)
- ▶ faltas que ocorreriam numa cache infinita
- Faltas por capacidade: cache não comporta conjunto de trabalho
- ▶ faltas que ocorreriam mesmo com políticas perfeitas de alocação e reposição
- Faltas por conflito: colisões no mapeamento de blocos na cache
- ▶ faltas que não ocorreriam se a cache fosse totalmente associativa e reposição fosse com LRU
- ▶ podem ocorrer com cache não-cheia

HEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

Reduzir a taxa de faltas – parâmetros de projeto

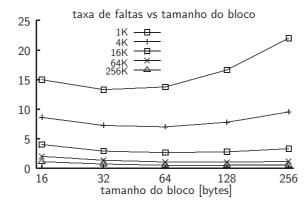
- Maior capacidade
- ⊳ reduz faltas por capacidade e por conflito
- Tamanho do bloco
- ▶ localidade espacial reduz faltas compulsórias
- ⊳ bloco maior pode aumentar penalidade (tempo de preenchimento)
- Associatividade
- ⊳ reduz faltas por conflito (até associatividade 4-8)
- ▶ pode aumentar tempo de acesso

HEPR DIAF RCC

Arquitetura II — otimização de caches

2007-1

1. Reduzir faltas - Tamanho do Bloco



- localidade espacial reduz faltas compulsórias
- para mesmo tamanho, menor número de blocos aumenta conflitos
- bloco maior aumenta tempo de preenchimento >penalidade

2. Reduzir faltas - Associatividade

Regra 2:1

taxa de faltas de uma cache com mapeamento direto de tamanho N =

taxa de faltas de cache com associatividade binária de tamanho N/2

Achtung: desempenho é medido pelo tempo de execução!! maior associatividade pode aumentar tempo de acesso

Em [Hill88] sugere-se que relação entre tempo de acerto de assoc-binária / mapeamentoDireto é
1,02 para caches internas e
1,10 para caches externas

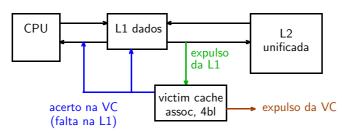
lembre que dados estão disponíveis só após comparação da etiqueta

HEPR Disf RCC 465

Arquitetura II — otimização de caches

2007-1

3. Reduzir faltas – Cache de Vítimas



Cache de Vítimas é uma cache pequena (4-8 blocos), associativa, ligada à cache com mapeamento direto para conter linhas recentemente expurgadas

[Jouppi90]: CdV com 4 elementos remove de 20% a 95% dos conflitos numa cache de 4Kbytes com mapeamento direto;

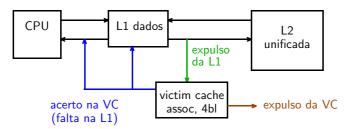
Usado no Alpha e máquinas da HP

HEPR Dinf RCC 464

Arquitetura II — otimização de caches

2007-1

3.1 Reduzir faltas - Cache de Vítimas



- Procura na L1;
- se falta, procura da cache de vítimas;
- se encontrou, troca linha por aquela que vai ser expulsa da L1;
- se falta na CdV, vítima da L1 → CdV;
- Efeito combinado: tempo de acesso do mapeamento direto com redução nas faltas por conflito.

4. Reduzir faltas – Busca antecipada por hardware

- Busca antecipada de instruções
- * Alpha 21064 busca dois blocos numa falta
- * bloco extra é colocado num stream buffer
- * numa falta, procura no stream buffer
- Busca antecipada de dados
- ★ 1 stream buffer para dados resolve 25% das faltas numa cache de 4KB; 4 stream buffers resolvem 43% [Jouppi90]
- ★ para aplicações científicas, 8 streams resolveram 50-70% das faltas com duas caches de 64KB, assoc-quaternária [Palacharla94]
- Busca antecipada necessita de sistema de memória com vazão sobrando e que possa ser usada sem penalidade

HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

5. Reduzir faltas – Busca antecipada por software

- Busca antecipada de dados
- * carregar dados em registrador (HP PA-RISC load)
- ★ carregar dados para a cache (MIPS-4, PowerPC)
- ★ instruções de busca antecipada não podem causar faltas de página → forma de execução especulativa
- Executar instruções de busca antecipada custa tempo
- * custo da busca antecipada é menor que os ganhos pelos acertos?
- * processadores super-escalares gulosos facilitam despacho destas instruções
- Busca antecipada causa poluição na cache
- * carregar blocos em stream buffers e não na cache
- \star numa falta, carrega stream buffers com blocos após o faltante
- * numa falta, procura na cache e nos SBs
- \star 1 SB é bom para instruções > 1 SBs para dados (\approx 4..8)

HEPR DIAF RCC

Arquitetura II — otimização de caches

2007-1

6. Reduzir faltas – Otimização pelo compilador

- Instruções
- * reordenar funções na memória para reduzir faltas por conflito
- ★ verificar perfil de execução para descobrir conflitos, e então re-ligar aplicativos
- Dados
- ★ **agrupar vetores:** vetor de elmtos compostos tem melhor localidade espacial que 2 vetores
- ★ troca de loops: mudar aninhamento dos loops para acessar dados na ordem de armazenamento em memória
- * fusão de loops: combinar 2 loops independentes com comportamento similar e sobreposição de variáveis
- * blocagem: melhorar a localidade temporal ao acessar "blocos" de dados repetidamente ao invés de percorrer linhas/colunas inteiras

Exemplo: agrupamento de vetores

```
/* Antes: 2 vetores */
int val[TAM];
int key[TAM];

/* Depois: 1 vetor de estruturas */
struct agrup {
   int val;
   int key;
}
struct agrup doisVetores[TAM];
```

Redução de conflitos entre val e key: melhora localidade espacial

LIEPR Dinf BCC

Arquitetura II — otimização de caches

2007-1

Exemplo: troca de loops

```
for (k=0; k < 100; k++)
                                    /* Antes */
2
      for (j=0; j < 100; j++)
3
         for (i=0; i < 5000; i++)
4
            x[i][j] = 2 * x[i][j]; /* passo=100 */
5
   for (k=0; k < 100; k++)
                                     /* Depois */
6
      for (i=0; i < 5000; i++)
7
         for (j=0; j < 100; j++)
            x[i][j] = 2 * x[i][j]; /* passo=1 */
8
```

Acessos seqüenciais ao invés de passadas de 100 palavras

melhora localidade espacial
passada == stride

HEPR Dinf RCC

470

Arquitetura II — otimização de caches

2007-1

Exemplo: fusão de loops

```
/* Antes: duas faltas por acesso em a e em c */
1 for (i=0; i < N; i++)
2
      for (j=0; j < N; j++)
3
         a[i][j] = 1/b[i][j] * c[i][j];
4 for (i=0; i < N; i++)
      for (j=0; j < N; j++)
         d[i][j] = a[i][j] + c[i][j];
/* Depois: uma falta por acesso em a e em c */
7
   for (i=0; i < N; i++)
8
      for (j=0; j < N; j++) {
         a[i][j] = 1/b[i][j] * c[i][j];
9
         d[i][j] = a[i][j] + c[i][j];
Α
            /* re-uso e melhor localidade espacial */
```

Exemplo: blocagem

- Dois loops internos:
- ★ lê todos os NxN elementos de z[]
- ★ lê N elementos de uma linha de y[] repetidamente
- ★ escreve N elementos de 1 coluna de x[]
- Faltas por capacidade são função de N e tamanho da cache
- \star núm faltas $=2N^3+N^2$, se não ocorrerem conflitos...
- Idéia: computar sub-matriz que cabe na cache

HEPR Dist RCC 47"

Arquitetura II — otimização de caches

Arquitetura II — otimização de caches

2007-1

2007-1

Exemplo: blocagem

IIFPR Dinf RCC 474

H&P QA Fig-5.21

HEPR Dinf RCC 475

Arquitetura II — otimização de caches

2007-1

Exemplo: blocagem

```
1 for (jj=0; jj < N; jj = jj+B)  /* Depois */
2 for (kk=0; kk < N; kk = kk+B)
3 for (i=0; i < N; i++)
4  for (j = jj; j < min(jj+B-1 , N); j++) {
5  for (r = 0, k = kk; k < min(kk+B-1, N); k++)
6  r = r + y[i][k] * z[k][j];
7  x[i][j] = r;
8 }</pre>
```

- B é chamado de fator de blocagem
- ullet Faltas por capacidade caem de $2N^3+N^2$ para $2N^3/B+N^2$
- Faltas por conflito talvez diminuem

HEPR Dinf RCC 471

Arquitetura II — otimização de caches

2007-1

Exemplo: blocagem

2007-1

HEPR Dinf RCC 478

Arquitetura II — otimização de caches 2007-1

H&P QA Fig-5.22

HEPR Dist RCC A76

Melhoria no desempenho de caches

Tempo médio de acesso a memória = tempo de acerto + taxa de faltas \times penalidade pela falta

Para melhorar o desempenho é necessário:

• reduzir a taxa de faltas

Arquitetura II — otimização de caches

- reduzir a penalidade nas faltas
- reduzir o tempo de acerto

LIFPR Dinf RCC 480

A. Reduzir a penalidade – hierarquia de caches

Memória não pode ser simultaneamente grande e rápida; ⇒ hierarquia com caches maiores mais longe do processador

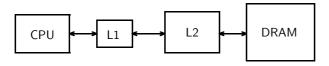


HEPR Dist RCC

Arquitetura II — otimização de caches

2007-1

A.a Reduzir a penalidade – hierarquia de caches



Taxa de faltas local = faltas na cache / acessos à cache
Taxa de faltas global = faltas na cache / referências pelo processador

$$T_{mem} = tA_{L1} + F_{L1} \times P_{L1}$$

= $tA_{L1} + F_{L1} \times [tA_{L2} + F_{L2} \times P_{DRAM}]$

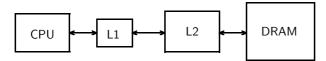
onde tA é o tempo de acerto, F é taxa de faltas, e P é a penalidade por falta

HEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

A.b Reduzir a penalidade – hierarquia de caches



Propriedade de Inclusão:

Conteúdo das caches maiores *inclui* o conteúdo das menores ⇒ cache menor contém cópias dos dados na cache maior

- ★ qual a relação entre associatividade, tamanho e inclusão?
- ★ qual a relação entre bits de tag/índice e assoc, tamanho e inclusão?
- ★ [Baer,Wang88] On the Inclusion Properties for Multi-Level Cache Hierarchies

Caches *exclusivas* trocam linhas numa falta $\Longrightarrow L1 \leftarrow novo; L2 \leftarrow velho$

A.c Reduzir a penalidade – hierarquia de caches

Hierarquia no Itanium-2

- L1: 2×16 KB, 4 conjuntos, linhas com 64 bytes, quatro portas (2-load + 2-store), latência de 1 ciclo
- L2: 256KB, 4 conjuntos, linhas com 128 bytes, quatro portas (2-load + 2-store), latência de 5 ciclos
- L3: 3MB, 12 conjuntos, linhas com 128 bytes, uma porta de 32 bytes, latência de 12 ciclos

HEPR Dinf RCC 48

Arquitetura II — otimização de caches

2007-1

B. Reduzir a penalidadeprioridade para faltas de leitura

- Escrita forçada com fila de escrita causa risco RAW
- * processador tenta ler dado que está na fila de escrita
- ★ solução ruim: espera até que fila esvazie e então prossegue com LOAD
 ⇒ aumenta penalidade nas faltas de leitura
- ★ solução melhor: verifica se ∃ dependências; se não, LOAD prossegue sem bloquear
- Escrita preguiçosa: falta na leitura substitui bloco sujo
- * solução ruim: escreve bloco sujo e então busca bloco faltante
- ★ solução melhor:
 - 1. escreve bloco sujo num buffer/fila,
 - 2. lê bloco faltante e entrega ao processador,
 - 3. então escreve bloco sujo
- \star menos perda de tempo porque CPU prossegue assim que ler bloco faltante

HEPR DIAF RCC

Arquitetura II — otimização de caches

2007-1

C. Reduzir a penalidade – minimizar tempo de carga

- Tempo de preenchimento do bloco é longo
- ...mas não precisa esperar até que bloco seja preenchido
- Early restart assim que palavra requisitada chegar da memória, entrega ao processador, que continua a executar
- Critical word first busca palavra requisitada primeiro e a entrega ao processador assim que chegar da memória
- * processador continua enquanto bloco está sendo preenchido
- * técnica é útil com blocos grandes
- ★ processador requisita p2:
 p0 p1 p2 p3
- \star **p2** é entregue, e bloco é preenchido com p2 \to p3 \to p0 \to p1
- \star localidade espacial indica que próximo acesso será na palavra seguinte, que pode não ter chegado ainda

D. Reduzir a penalidade – cache não-bloqueante

- Cache não-bloqueante entrega dados enquanto trata de uma falta
- * processador com execução fora-de-ordem
- ★ memória com múltiplos bancos ⇒ cache & DRAM
- * [Kroft81] Lock-up Free Instruction Fetch/Prefetch Cache Organization
- acerto sob falta (hit under miss) reduz penalidade efetiva fornecendo dados, ao invés de bloquear processador
- acerto sob múltiplas faltas (miss under miss) reduz penalidade ainda mais ao sobrepor tratamento de múltiplas faltas
- ★ complexidade do controlador da cache é MUITO maior por causa dos acessos concomitantes à memória
- * necessita de memória com múltiplos bancos
- ★ PentiumPro permite até 4 faltas pendentes

HEPR Dinf RCC

Arquitetura II — otimização de caches

2007-1

Melhoria no desempenho de caches

Tempo médio de acesso a memória = tempo de acerto + taxa de faltas x penalidade pela falta

Para melhorar o desempenho é necessário:

- reduzir a taxa de faltas
- reduzir a penalidade nas faltas
- reduzir o tempo de acerto

HEPR Dinf RCC 488

Arquitetura II — otimização de caches

2007-1

I. Reduzir o tempo de acerto

- caches simples e pequenas

Tempo de acerto curto se:

- cache pequena ⇒ reduz tempo de acesso
- cache simples \implies mapeamento direto

Alpha 21164 tem 8KB para dados e 8KB para instruções, 96KB cache secundária, unificada, associatividade ternária,

e ⇒ relógio muito rápido!

II. Reduzir o tempo de acertoevitar tradução de endereços

Cache endereçada com: endereço virtual vs endereço físico

- a cada troca de contexto deve expurgar a cache
- * senão ocorrem falsos-acertos
- * processo que entra no processador tem montes de faltas compulsórias
- \star conseqüência da troca de contexto pode durar 10^5 ciclos
- * solução: usar identificador de processo na etiqueta na cache
- como trata de sinônimos?
- * dois endereços virtuais que mapeiam no mesmo endereço físico
- \star solução: garantir que indexação da cache ocorre com bits que coincidem no EV e EF

EVirt	núm página virtual	deslocamento
ender. cache	—índice—	
EFís	núm página física	deslocamento

TIEPR Dist RCC

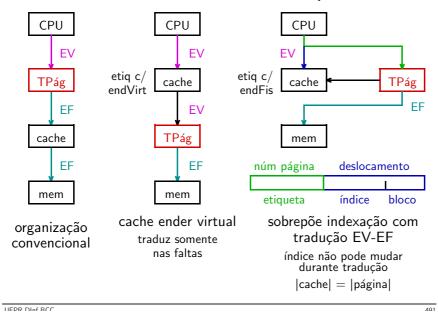
Arquitetura II — otimização de caches

Arquitetura II — otimização de cache

2007-1

2007-1

II.i Reduzir T de acerto – evitar tradução de ender's

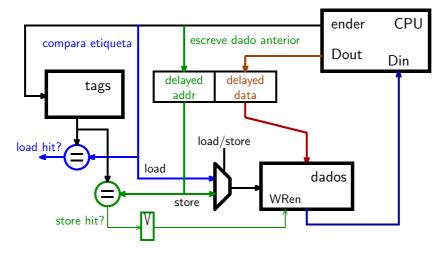


III.i Reduzir o tempo de acerto – escritas em pipeline

Escrita consome dois ciclos

- um para comparar etiqueta e outro para escrever na cache
- armazena dados e endereço em buffer
- escreve dados na cache durante comparação de etiqueta na próxima escrita
- risco RAW entre leitura e escrita atrasada ⇒ adiantamento

III.ii Reduzir tempo de acerto – escritas em pipeline



HEPR Dist RCC 403

Arquitetura II - memória principal e virtual

2007-1

Resumo - melhorar desempenho de caches

 $T_méd_acesso_mem = T_acerto + tx_faltas \times penalidade/falta$

Para melhorar o desempenho é necessário:

* reduzir a taxa de faltas

- De parâmetros de projeto: núm, tam blocos, associatividade
- > cache de vítimas, stream buffers, busca antecipada
- Do otimização do código: agrupamento, troca de índices, fusão, blocagem

* reduzir a penalidade nas faltas

- ▷ hierarquia de caches
- ▷ priorizar faltas na leitura (critical word first, early restart)

* reduzir o tempo de acerto

- > caches simples e pequenas
- ▷ não traduzir endereço virtual para físico
- ▷ escritas em pipeline

HEPR Dist RCC

Arquitetura II — memória principal e virtual

2007-1

Revisão - melhorar desempenho de caches

 $T_méd_acesso_mem = T_acerto + tx_faltas \times penalidade/falta$

Para melhorar o desempenho é necessário:

* reduzir a taxa de faltas

- ▷ parâmetros de projeto: núm,tam blocos, associatividade
- > cache de vítimas, stream buffers, busca antecipada
- $\, \rhd \,$ otimização do código: agrupamento, troca de índices, fusão, blocagem

* reduzir a penalidade nas faltas

- ▶ hierarquia de caches
- ▷ priorizar faltas na leitura (critical word first, early restart)

* reduzir o tempo de acerto

- ▷ não traduzir endereço virtual para físico
- ▷ escritas em pipeline

Sistema de Memória

- Sistemas de memória
- Memória cache
- * organização
- ★ leitura
- ★ escrita
- * otimizações
- * busca antecipada
- Memória Virtual
- ⊳ endereçamento e proteção
- ▶ paginação
- ⊳ excessões

HEPR Dinf RCC A0

Arquitetura II - memória principal e virtual

2007-1

Tipos de loci em memória



- Endereço de linguagem de máquina como especificado pelo montador
- ◆ Endereço virtual → endereço efetivo conj de instruções especifica tradução dos endereços gerados pelo montador para endereços usados pelo ligador
- Endereço físico

sistema operacional especifica mapeamento de endereço virtual para o número de uma posição na memória física

HEPR DInf RCC

Arquitetura II — memória principal e virtual

2007-1

Espaço de Endereçamento – três modelos

- Primeiro Modelo (até ≈1960)
- * Programador supõe que espaço de endereçamento é contínuo espaço se extende de 0x0000 0000 a 0xffff ffff
- * programa deve ser carregado sempre no endereço físico inicial
- * memória pode conter somente um programa em execução
- Segundo Modelo (até 1972 IBM 370, idéia de 1962)
- * EdE contínuo, > que memória física → overlays
- * Registrador de tradução para fazer relocação do programa
 → programa pode ser carregado em qualquer endereço físico
- * mais de um programa carregado em memória para execução
- Terceiro Modelo (após 1972 IBM 370)
- * EdE > que memória física, tradução de ender transparente
- * muitos programas carregados na memória

Endereços Absolutos

No início da década de 50, computador da moda era o EDSAC, e endereço virtual = endereço na memória física

Só um programa executa na máquina, com acesso irrestrito a todos seus recursos (RAM+E/S)

Endereços no programa dependem de **onde** o programa está carregado na memória

É fácil escrever sub-rotinas com código independente de posição?

HEPR Dist RCC 400

Arquitetura II - memória principal e virtual

2007-1

Tradução Dinâmica de Endereços

• Motivação:

operações de E/S nas máquinas antigas eram lentas e processador era envolvido em todas as transferências (polling)

Produtividade maior se dois ou mais programas executassem concorrentemente → multiprogramação

- Programas independentes de posição:
 facilitam a programação e gerenciamento da memória
 → é necessário um registrador_base
- Proteção:

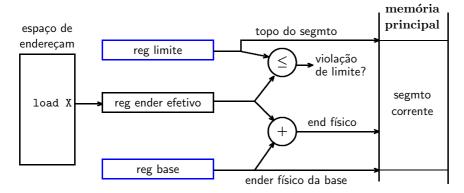
programas independentes não deveriam afetar-se acidentalmente
→ é necessário um registrador_limite

HEPR Dist RCC 50

Arquitetura II — memória principal e virtual

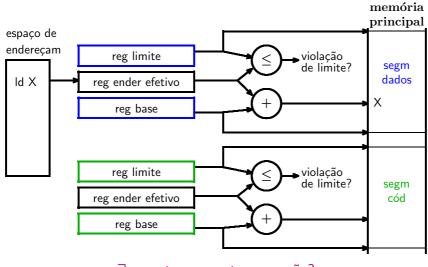
2007-1

Tradução com registradores base e limite



Registradores base e limite são visíveis somente quando processador executa em modo supervisor

Áreas Separadas para Programa e Dados



∃m vantagens nesta separação?

HEPR NINF RCC

Arquitetura II — memória principal e virtual

2007-1

Fragmentação da Memória

- Programas com segmentos de tamanho fixo devem ser acomodados na memória que está disponível
- à medida que programas/usuários entram e saem do sistema, memória fica com buracos onde estavam os programas
 - → fragmentação
- De quando em quando, programas devem ser re-locados/movidos para abrir espaço
 - → compactação

SO da Microsoft $^{\rm TM}$ obriga usuários a fazer compactação de disco, um problema que foi resolvido (não pela primeira vez) no Unics, em 1979. *That IS innovation* $^{\rm TM}$.

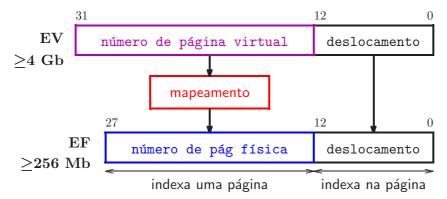
HEPR DIst RCC 50

Arquitetura II — memória principal e virtual

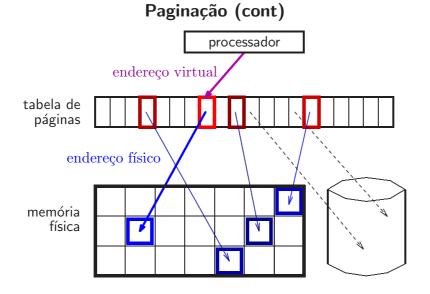
2007-1

Paginação

Espaço de endereçamento dividido em páginas de 4-8 Kbytes Tabela de Páginas mapeia endereços virtuais em endereços físicos



Quais as vantagens das páginas de tamanho fixo e da indireção através da tabela de páginas?



HEPR Dinf RCC 505

Arquitetura II - memória principal e virtual

2007-1

Tabela de Páginas

- Qual o tamanho da tabela de páginas?
- Quantas tabelas de páginas são necessárias?
- O espaço necessário para as tabelas de páginas é proporcional ao espaço de endereçamento, ao núm de usuários ...
- ▷ espaço necessário é grande!
- ⊳ grande demais para manter em registradores (muitas páginas)
- Alocar TP em registradores especiais, só para o usuário corrente
 pode não ser factível para TPs grandes
- Alocar TP na memória principal
- ⊳ necessita uma referência para buscar endereço base da página e...
- - → dobra o número de referências à memória

HEPR DIAF RCC 50

Arquitetura II — memória principal e virtual

2007-1

Paginação sob demanda

- Paginação reduz fragmentação externa (buracos entre programas)
- Problema:

acomodar programa com muitos dados em memória pequena?

- ▶ usar overlays programador deve controlar quais trechos de código/dados estão residentes em RAM
 eeeeca!
- ▶ programação com overlays é complicada!
- Paginação sob demanda:

"uma página da memória secundária é trazida para a memória primária sempre que for (implicitamente) requisitada pelo processador"

Atlas 1962

> memória primária atua como cache para memória secundária

DRAM como Cache de Disco

Se página virtual não está em memória física,

→ deve ser copiada do disco

→ alguma página deve ser ejetada para abrir espaço localidade recomenda:

vítima é a que foi usada no passado mais distante (LRU)

Exemplo: 11, 10, 12, 9, 11, 7, 11, 13 LRU c.r.a 13? 11?

Escrita é preguiçosa:

- páginas semente de leitura (código) são substituídas
- páginas com atualizações são marcadas sujas e, antes de substituídas, disco deve ser atualizado

HEPR Dinf RCC

Arquitetura II — memória principal e virtual

2007-1

Cache vs Paginação sob Demanda



cache	paginação	
linha na cache	quadro (moldura)	
bloco na cache 32-64 bytes	página 4-8K bytes	
faltas na cache 1-20%	faltas de página $<$ 0.001 $\%$	
acerto na cache 1 ciclo	acerto na TP 100 ciclos	
falta na cache 100 ciclos	falta de página ${f 10^6}$ ciclos	
falta tratada em HW	falta tratada espec em SW	

HEPR Dist RCC 500

Arquitetura II — memória principal e virtual

2007-1

Sistemas modernos de Memória Virtual

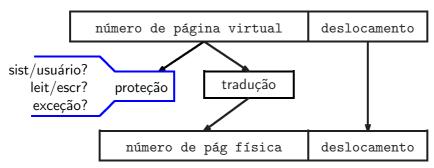
• Proteção e privacidade

- ⊳ vários usuários,
- ⊳ cada um com seu espaço de endereçamento privativo,
- ▶ e um ou mais espaços de endereçamento compartilhados tabela de páginas ≈ espaço de endereçamento
 mapeamento
- paginação sob demanda
- De capacidade de executar programas maiores que a memória física

• melhor utilização dos recursos

- ⊳ menor fragmentação externa
- custo: tradução de endereço a cada referência à memória

Tradução de Endereços e Proteção



 cada referência a dado ou instrução necessita de tradução e verificação de proteção

proteção obtida através da indireção

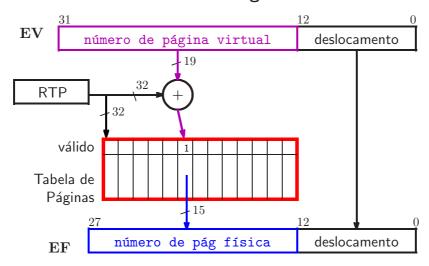
• um bom sistema de memória virtual deve ser rápido (≤ 1 ciclo) e usar espaço eficientemente

HEPR DINF RCC 51

Arquitetura II — memória principal e virtual

2007-1

Tabela de Páginas



HEPR Dinf RCC 51

Arquitetura II — memória principal e virtual

2007-1

Tabela de Páginas Linear

Registrador de Tabela de Páginas (RTP) aponta para início da Tabela de Páginas do processo

A cada troca de contexto o SO atualiza o RTP para que este aponte a tabela de páginas do processo de usuário.

Cada elemento da tabela de páginas contém:

nPF número da página física residente em memória

nPD número da página em disco (se página foi movida para disco)

nil página não-existente

stat bits de status e proteção/uso used, dirty, {RO,RW,EX}

2007-1

Arquitetura II — memória principal e virtu

Tamanho da TP linear

Com endereços de 32 bits, páginas de 4 Kbytes, 4 bytes/elemento:

ightarrow 2²⁰ elementos = 4 Mbytes por processo

ightarrow 4 Gbytes de swap para conter espaço de endereço completo

Páginas maiores!

mais fragmentação interna penalidade por falta maior última página meio vazia mais tempo para ler do disco

Processadores de 64 bits

ightarrow mesmo páginas de 1 Mbyte implicam em tabelas com ${f 2}^{44}$ elementos de 8 bytes (pprox 35 TBytes)

Há salvação? Como?

LIEPR Dinf BCC Arquitetura II - memória principal e virtual 2007-1 Tabela de Páginas Hierárquica 11 0 21 31 desloc p2 desl índice no índice no nível 1 nível 2 Raiz da TP corrente p2 RTP р1 TP nível1 TP nível 2 pág em mem fis pág inexistente (não foi mapeada) pág em mem secund páginas HEPR Dinf RCC

Cache de Mapeamentos I

A cada referência, processador consulta TP
para descobrir endereço físico do objeto

→ faz uma referência à Tabela de Página para obter endereço
e então faz referência ao objeto...

→ para cada referência, DOIS acessos à memória...

Solução:

manter mapeamentos numa Cache de Mapeamentos ou Translation Lookaside Buffer (TLB) ou Translation Buffer (TB, no Vax)

Translation Lookaside Buffer

Processador procura mapeamento na TLB; se encontra, completa referência; senão, busca mapeamento da Tabela de Páginas (em mem física) e guarda na TLB para uso futuro

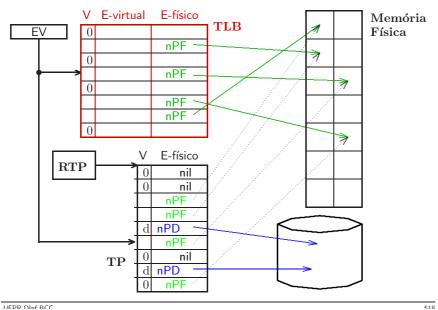
Se encontra na TLB \rightarrow tradução em \leq 1 ciclo TLB hit se falta na TLB \rightarrow procura na TP (30-100 ciclos) TLB miss

HEPR Dist RCC 51

Arquitetura II - memória principal e virtual

2007-1

Cache de Mapeamentos II



Arquitetura II — memória principal e virtual

2007-1

Parâmetros de projeto de TLBs

- Tamanho da TLB: 32 a 128 blocos
- tamanho de bloco: 1..2 mapeamentos por bloco
- associatividade alta
- - → maior número de conflitos
- ▶ TLBs grandes (256-512 elmtos) tem associatividade 4-8

Pergunta:

Qual a cobertura da TLB?

a)
$$|TLB| \times |pag|$$

b)
$$|TLB| \times |Tab_pag|$$

c)
$$|pag| \times |Tab_pag|$$

d) nenhuma das acima

HEPR Dief RCC

Arquitetura II — memória principal e virtual

2007-1

Resposta:

Qual a cobertura da TLB? ou, qual o maior espaço de endereçamento mapeado pela TLB?

a)
$$|TLB| \times |pag|$$

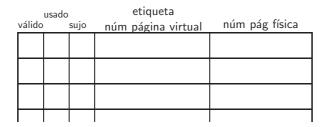
- b) $|TLB| \times |Tab_pag|$
- c) $|pag| \times |Tab_pag|$
- d) nenhuma das acima

HEPR Dist RCC

 $\underline{\mathsf{Arquitetura}\;\mathsf{II}-\mathsf{mem\'oria}\;\mathsf{principal}\;\mathsf{e}\;\mathsf{virtual}}$

2007-1

Campos da Cache de Mapeamentos (TLB)



válido == 1 se mapeamento é válido usado == 1 se página foi referenciada recentemente (LRU) sujo == 1 se ocorreu uma ou mais escritas na página

Se página virtual não está em memória, NÃO pode haver um mapeamento da página na TLB

LIEPR Dinf RCC 522

Tratamento de Faltas na TLB

• Em software MIPS, Alpha

- ▶ falta na TLB causa exceção, SO caminha pela TP hierárquica e re-carrega elemento da TLB
- ⊳ instrução que causou a falta na TLB é re-iniciada (busca)
- ⊳ caminhada ocorre em modo privilegiado (end fís), sem tradução

• Em hardware

SparcV8, x86, PowerPC

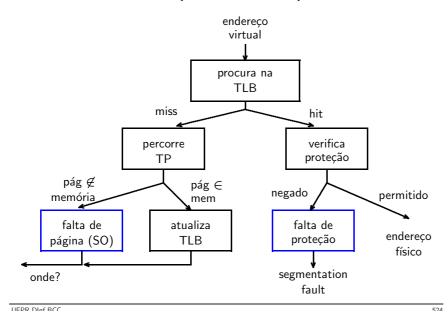
- □ unidade de gerenciamento de memória (MMU) caminha pela TP
 e re-carrega elemento da TLB
- ⊳ instrução que causou a falta na TLB é re-iniciada (busca)
- ⊳ se página (dados/código, ou TPs 2º nível) não está em memória, MMU sinaliza uma falta de página na instrução que causou falta na TLB

Por que a instrução causadora deve ser buscada novamente?

HEPR Dist RCC 593

Arquitetura II — memória principal e virtual 2007-1

Tradução de Endereços



Arquitetura II — memória principal e virtual 2007-1

TLB & Cache - excessões



Faltas & Excessões

Causas: endereço ilegal fora do espaço de endereçamento ou desalinhado end $\%4 \neq 0$

ou falta na TLB

ou falta na tabela de páginas page fault

TLB & Cache

Falta na TLB pode ocorrer

- na busca de instrução
- (i) processador esvazia pipeline; (ii) re-carrega TLB; e
- (iii) busca instrução novamente custo: drenar pipeline + carga da TLB
- referência a dados (ld e st)
- (i) instrs à frente completam; instr não pode alterar estado: ld r1,0(r1)
- (ii) processador re-carrega TLB; e (iii) busca instrução novamente

custo: drenar pipeline + carga da TLB + estágios até MEM



HEPR Dinf RCC 526

Arquitetura II — memória principal e virtual

2007-1

Reposição de Páginas

Se página virtual não está em memória SO assume controle e requisita cópia ao controlador de disco SO escolhe vítima para ser substituída pela nova página;

→ se vítima estiver suja, SO deve atualizar disco

Localidade implica em vítima ser página usada no passado mais distante (LRU)

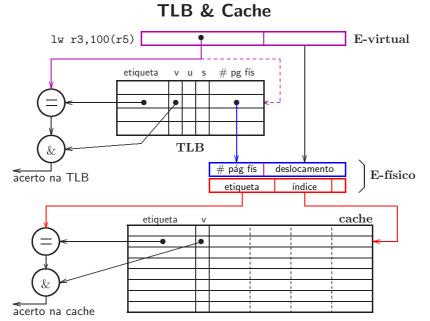
Como implementar LRU para 4 conjuntos com 3 bits?

Pode ser pseudo-LRU...

HEPR Dief RCC 597

Arquitetura II — memória principal e virtual

2007-1



Resumo - memória física e virtual

- Aumento da vazão

- Redução da Latência
- b fast page mode
 b fast page mode
 c fast page mode
- Memória Virtual
- ▶ Proteção vs utilização dos recursos proteção obtida através da indireção
- ⊳ |TP| grande → tabela de página hierárquica (multi-nível)
- ⊳ relação entre Mem Virtual e tratamento de excessões falta de página é evento freqüente → tratamento rápido

LIEPR Dinf BCC

520

Arquitetura II - multiprocessadores

2007-1

Processamento Paralelo & Multiprocessadores

- Motivação
- Tipos de máquinas paralelas
- Coerência entre caches
- Sincronismo
- Ordenação de operações em memória

HEPR Dinf RCC

530

Arquitetura II — multiprocessadores

2007-1

Motivação

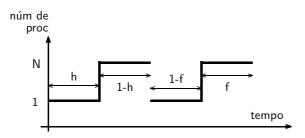
Problema:

algumas computações são muito grandes para uniprocessadores Por **grande** entende-se muito tempo ou muita memória

Solução:

- a) projetar processadores mais rápidos
- \rightarrow de 1985-2005 o crescimento de desempenho foi enorme
 - → processadores modernos já são MUITO complexos
- b) usar vários processadores trabalhando em paralelo
- → progresso no desenvolvimento de software tem sido lento
- → progresso bom sw de servidores e algumas aplics embutidas
 - → base de sw para desktops é enorme...

Complicador: Lei de Amdahl



Ganho de 80 com 100 processadores só se 0, 25% do trabalho é serial fazer Exercício 6.1

TIEDR Dief RCC 53

Arquitetura II — multiprocessadores 2007-1

Complicador: Custos da cooperação

• Cooperação entre processadores envolve:

ightharpoonup comunicação \longrightarrow processadores esperam ruim

- Custo da comunicação depende de:
- ⊳ sw/hw para empacotar mensagens
- > rede que transporta mensagens
- Evento de comunicação custa vários micro-segundos
- ightarrow processador @ 1GHz ightarrow 1 μ s=1000 ciclos

HEPR Dinf RCC 53

Arquitetura II — multiprocessadores 2007-1

Aplicações

- Processamento paralelo
- ⋆ paralelismo real num único programa
- \star compartilhamento de dados pode ser freqüente/intenso
- Processamento de Transações
- ★ paralelismo entre transações independentes
- ★ ênfase em produtividade (throughput)
- Sistemas Operacionais
- ⋆ programa paralelo grande que executa por muito tempo
- * SO paralelo é geralmente otimizado "a mão"
- ★ compartilhamento de dados in-frequente
 - ightarrow estruturas de dados protegidas com granularidades variadas
- estas formas são distintas da exploração de paralelismo entre as instruções de uma única linha de execução single-thread vs multiple-threads

HEPR Dinf RCC 53/

Taxonomia de Flynn

Taxonomia de Flynn, publicada 1966, e baseada no número de fluxos de instruções e número de fluxos de dados

- SISD Single Instruction, Single Data
- ▷ uniprocessador
- SIMD Single Instruction, Multiple Data
- processadores vetoriais/sistólicos
- MISD Multiple Instruction, Single Data
- poucos exemplos práticos
- MIMD Multiple Instruction, Multiple Data

HEPR Dist RCC 538

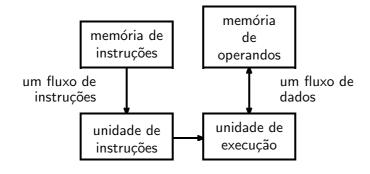
Arquitetura II — multiprocessadores

2007-1

(processadores sistólicos ?)

Taxonomia de Flynn - SISD

- SISD Single Instruction stream, Single Data stream
- ▷ ...o bom e velho uniprocessador



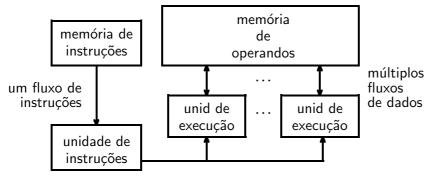
HEPR DIst RCC 53

Arquitetura II — multiprocessadores

2007-1

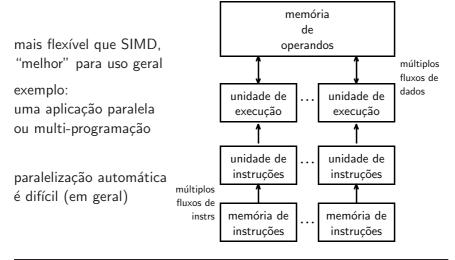
Taxonomia de Flynn - SIMD

- SIMD Single Instruction stream, Multiple Data streams
- > armazenagem de instruções e operandos usualmente separada
- ⊳ modelo de programação é data parallel
- ⊳ paralelização automática pode ser possível —geralmente é



Taxonomia de Flynn - MIMD

• MIMD - Multiple Instruction streams, Multiple Data streams



TIEDR Dinf RCC 538

Arquitetura II - multiprocessadores

2007-1

Eixos de Projeto

- Modelos de programação
- ⊳ sequencial (SISD)

shared memory

message passing

• Interface binária da aplicação

camada de sw que mapeia modelo de programação no hardware

- b troca de mensagens → comunicação com send() recv()
- Hardware
- ⊳ memória compartilhada ← comunicação "grudada" na memória

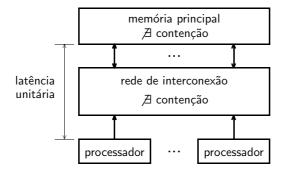
HEPR Dist RCC 53

Arquitetura II — multiprocessadores

2007-1

Modelo PRAM

- Parallel RAM = perfeição (≈ irreal)
- > memória completamente compartilhada
- ▶ latência unitária
- ⊳ largura de banda infinita → não há contenção
- ▶ localização dos dados irrelevante



Perfeição não é realizável

- Latência cresce se tamanho do sistema cresce
 ▷ veloc da luz é fixa → latência ∝ distância
- Vazão é limitada pela organização da memória e da rede de interconexão
- estas levam a uma divisão entre arquiteturas MIMD com Uniform Memory Access – UMA

 ϵ

Non-Uniform Memory Access -NUMA

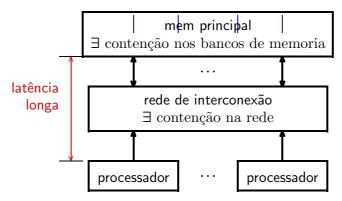
HEPR Dist RCC 54:

Arquitetura II - multiprocessadores

2007-1

UMA - Uniform Memory Access

- latência no acesso à memória é a mesma para todos processadores mas pode ser longa
- latências crescem com tamanho do sistema aumentar escala é difícil



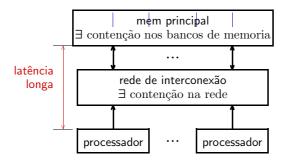
HEPR Dinf RCC 54/

Arquitetura II — multiprocessadores

2007-1

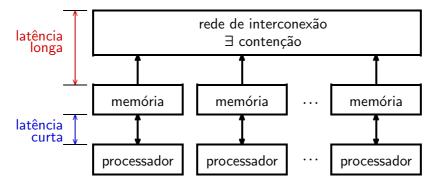
UMA - Uniform Memory Access

- Localização dos dados é irrelevante (não importa)
- contenção limita vazão → na rede e na memória
- geralmente usam caches
- usado em pequenos multiprocessadores = MPs simétricos Symmetric MultiProcessors ou SMPs



NUMA - Non-Uniform Memory Access

- Latência pequena no acesso à memória local
- latência grande no acesso à memória remota $P \rightarrow M_{loc} \rightarrow R \rightarrow M_{rem}$
- vazão para memória local pode ser mais alta que para remota
- contenção na rede e no acesso à memória



HEPR Dist RCC

Arquitetura II - multiprocessadores

2007-1

Multiprocessadores NUMA Non-Uniform Memory Access

- Memória logicamente compartilhada mas fisicamente distribuída
- > pode ser tratado como memória compartilhada
- desempenho depende fortemente da localização dos dados
- Multicomputadores
- ⊳ cada processador tem espaço de endereçamento privativo
- ⊳ comunicação através de troca de mensagens (explícita)

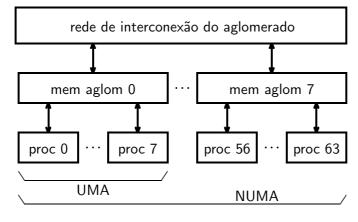
HEPR Dist RCC

Arquitetura II — multiprocessadores

2007-1

Aglomerados (clusters)

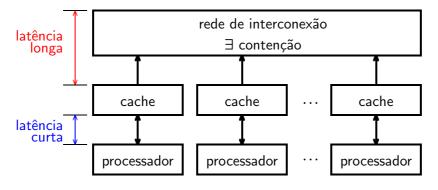
- Nós UMA pequenos num sistema NUMA grande
- híbrido? aglomerado de aglomerados?



2007-1

COMA - Cache Only Memory Architecture

- Dados nas caches migram para onde são necessários localização dos dados importa muito pouco
- COMA é uma forma de NUMA
- área de pesquisa quente (90-94), só uma implementação (KSR-1)



HEPR Dist RCC 547

Arquitetura II — multiprocessadores

Pergunta

Com 100 processadores interligados por uma rede perfeita e ideal, qual fração do código pode ser executada serialmente, se desejo ganho de {50,75,100}, com relação a 1 processador?

HEPR Dinf RCC 54

Arquitetura II — multiprocessadores 2007-1

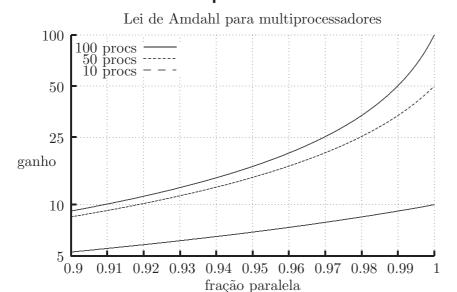
Resposta I

Com 100 processadores interligados por uma rede perfeita e ideal, qual fração do código pode ser executada serialmente, se desejo ganho de {50,75,100}, com relação a 1 processador?

$$\mathsf{Ganho}_{\mathsf{total}} = \frac{1}{(1 - \mathsf{Frac}_{\mathsf{melhor}}) + (\mathsf{Frac}_{\mathsf{melhor}} \, / \, \, \mathsf{Ganho}_{\mathsf{melhor}})}$$

$$50 = 1/(1 - 0.99) + (0.99/100)$$
 1%
 $75 = 1/(1 - 0.9965) + (0.9965/100)$ 0,35%
 $100 = 1/(1 - 0.999999...) + (0.999999.../100)$ $\approx 0\%$

Resposta II



HEPR Dist RCC 550

Arquitetura II — multiprocessadores

2007-1

O que é a fração serial?

- Inicialização
- ⊳ de estruturas de dados
- ⊳ criação de processos nos processadores
- ⊳ distribuição de trabalho entre os processadores
- Sincronização
- ⊳ contenção no acesso às regiões críticas

semáforos

loop mais externo

- Término
- ⊳ E/S para disco

HEPR Dinf RCC

551

Arquitetura II — multiprocessadores

2007-1

Comunicação Inter-processos

Na maioria das aplicações paralelas interessantes, processos paralelos (tarefas) devem comunicar-se

Método de comunicação leva à divisão entre: memória logicamente distribuída→ troca de mensagens memória logicamente compartilhada→ memória compartilhada

Tipos de Máquinas Paralelas

Duas religiões em computação paralela:

- memória logicamente distribuída → multicomputadores
- memória logicamente compartilhada → multiprocessadores
- comunicação explícita vs comunicação implícita

Nos dois casos, há uma rede no caminho crítico

na verdade, multiprocessadores são multicomputadores travestidos: **sempre** há uma rede de interconexão;

→ eventos de comunicação baseados em mensagens...

HEPR Dinf RCC 55

Arquitetura II - multiprocessadores

2007-1

Métricas de Desempenho na Comunicação largura de banda/vazão

- A vazão pode ser limitada por
- ⊳ processador, memória, rede de interconexão
- a largura de banda pode ser medida em dois pontos:
- ▷ na rede → bisection bandwidth:
 - no pior diâmetro (c.r.a vazão global), corta os fios e mede vazão
- ⊳ em cada nó da rede → na interface com a rede
- mecanismo de comunicação afeta a largura de banda
- → ocupação de recursos

utilização

(# interrupções → intervalo entre msgns)

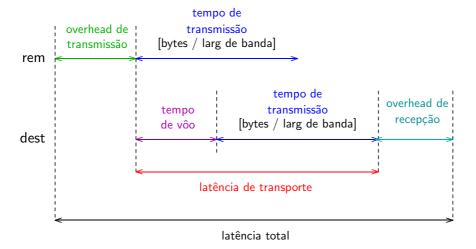
Arquitetura II — multiproces

HEPR Dinf RCC

LIEPR Dinf RCC

2007-1

Métricas de Desempenho na Comunicação latência



Métricas de Desempenho na Comunicação latência

- latência = Tx_ovhd
 + tempo_de_vôo + tempo_de_transmissão
 + Rx_ovhd
- \triangleright tempo_de_vôo \propto velocidade da luz (\approx 0,65c), distância

- latência afeta desempenho e facilidade de programação
- ⊳ programação: alocação de dados para aumentar localidade
- há forte relação entre overhead e ocupação...

HEPR Dinf RCC 55

Arquitetura II — multiprocessadores

2007-1

Métricas de Desempenho na Comunicação granularidade

- esconder latência pela superposição de comunicação e computação
- ⊳ superposição geralmente dificulta a programação
- ▶ granularidade ≡ computação / comunicação

As três métricas (vazão, latência, granularidade) são influenciadas pela aplicação: tamanho dos ítens (msgns) afeta vazão e latência padrões de comunicação podem causar congestionamento

HEPRING RCC 55:

Arquitetura II — multiprocessadores

2007-1

Mem compartilhada vs Mem distribuída

Memória compartilhada

▶ hw mais complicado

- (caches são coerentes)
- ⊳ 'facilidade' de programar com modelo de memória compartilhada
- ▶ fácil de programar quando padrões de comunicação são complexos
 - ou mudam dinamicamente durante execução
- baixa latência quando transfere ítens pequenos loads e stores
- ▶ MP com memória compartilhada emula troca de mensagens com facilidade

Mem compartilhada vs Mem distribuída

• Memória distribuída

- ▶ hardware mais simples
- ▶ padrões de comunicação devem ser explicitados com precisão
- ▶ programador deve estruturar programa para minimizar custos da comunicação
- ⊳ sincronização associada aos eventos de comunicação
- ▷ emular memória compartilhada é difícil e/ou ineficiente
 SO envolvido em proteção,
 nomeação e
 envio/recepção de msgns

HEPR Dist RCC

Arquitetura II - multiprocessadores

2007-1

Multicomputadores

- Arquitetura memória distribuída
- mais detalhes em breve

- ⊳ cada nó tem identificador único, processos em cada nó tbém
- Efeitos na latência
- ▶ rede é pipeline para mensagens

SE implementação permite

▶ latência = tempo_por_nó * núm_de_nós

+

comprimento_da_msgm / bytes_por_un_de_tempo

- - → comunicação infrequente e com msgns longas
 - \rightarrow granularidade relativamente grande (comput/comun $\gg 1$)

HEPR DInf RCC 54

Arquitetura II — multiprocessadores

2007-1

NUMA

Multiprocessadores

- Arquitetura memória logicamente compartilhada
- > todos os nós enxergam toda a memória
- ▶ memória local ao nó + memória remota ao nó
- Efeitos na latência
- ▶ Referências a objetos compartilhados não distinguem localização, se local ou remota
 diferença no desempenho
- ▶ mesmo modelo de programação de uniprocessador
- Exemplos:
- ▷ Servidores da Sun, Silicon Graphics, IBM
- ▶ Kendall Square Research KSR-1 (COMA) faliu em 95

Resumo

- Vazão através da rede (bisecção), e na interface processador-rede
- latência = Tx_ovhd
 + tempo_de_vôo + tempo_de_transmissão
 + Rx_ovhd
- granularidade ≡ computação / comunicação
- multiprocessadores: vários processadores interligados espaço de endereçamento compartilhado comunicação implícita, granularidade pequena/média
- multicomputadores: máquinas independentes cooperam espaços de endereçamento disjuntos comunicação explícita, granularidade grande

HEPR Dist RCC 56

Arquitetura II - coerência entre caches

2007-1

Processamento Paralelo & Multiprocessadores

- Motivação
- Tipos de máquinas paralelas
- Coerência entre caches
- Sincronismo
- Ordenação de operações em memória

HEPR Dinf RCC

Arquitetura II — coerência entre caches

2007-1

Revisão - Eixos de Projeto

- Modelos de programação
- ⊳ sequencial (SISD)
- → memória compartilhada (MIMD)
- b troca de mensagens (MIMD)

shared memory message passing

• Interface binária da aplicação

camada de sw que mapeia modelo de Programação no hardware

- b troca de mensagens → comunicação com send(), recv()
- Hardware
- → memória compartilhada ← comunicação via sist de memória

HEPR Dinf RCC 56.

Revisão

- Vazão da rede (bisecção), e na interface processador-rede
- latência = Tx_ovhd
 + tempo_de_vôo + tempo_de_transmissão
 + Rx_ovhd
- granularidade ≡ computação / comunicação
- multiprocessadores: vários processadores interligados espaço de endereçamento compartilhado comunicação implícita, granularidade pequena/média
- multicomputadores: máquinas independentes cooperam espaços de endereçamento disjuntos comunicação explícita, granularidade grande

HEPR Dist RCC 56

Arquitetura II — coerência entre caches

2007-1

Execução atômica

/* Quais os valores de "print c" e "print d"? */

Comandos 'atômicos' em C **não são** executados atomicamente pelo processador

```
a = a + 1; \equiv lw r1,0(r2)
addi r1,r1,1
sw r1,0(r2)
```

HEPR Dinf RCC

Arquitetura II — coerência entre caches

2007-1

Sincronização e Atomicidade

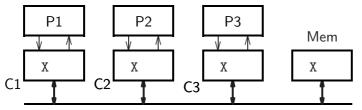
- Operações atômicas devem ser serializadas pelos mecanismos de escrita em memória
- Problemas decorrentes de disputas:
- ▶ latência sem contenção
- ⊳ serialização se há contenção
 - → dificulta escalar para sistemas maiores

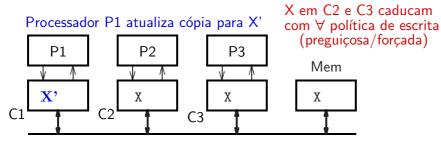
Barramento é um meio de comunicação compartilhado comunicação por **difusão** broadcast

→ comunicação é serializada através do barramento

Coerência de caches

Inicialmente, três cópias idênticas de X em C1, C2 e C3





HEPR Dist RCC 566

Arquitetura II — coerência entre caches

2007-1

Coerência de caches

P1, P2 e P3 carregam X em suas caches; P1 atualiza sua versão; C2 e C3 ficam com sua cópia desatualizada. → isso ocorre **mesmo** que as caches usem escrita forçada (ou escrita preguiçosa)

Informalmente, **coerência entre caches** é um método para garantir que os acessos à memória sejam coerentes, apesar das caches.

HEPR Dinf RCC 566

Arquitetura II — coerência entre caches

2007-1

Causas de In-coerência

- Compartilhamento de dados para escrita
 é a causa considerada mais comumente
- migração de processos processo que sai deixa dados para trás
- entrada/saída geralmente corrigida com invalidação de blocos nas caches

Não é necessária uma noção de tempo_absoluto:
o que é necessário é a coerência 'aparente',
ao invés de coerência absoluta

→ a incoerência temporária entre memória
e cache com escrita preguiçosa pode ser aceitável

Soluções para Problemas de Coerência

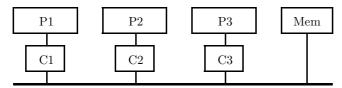
- Evitar caches privativas mover caches para sistema de memória
- Dados compartilhados não podem ficar na cache
- ▶ possivelmente, a solução mais simples por software
- Invalidar blocos em momentos estratégicos
- ⊳ por exemplo, após seções críticas
- Usar uma tabela global
- ▶ tabela mantém localização dos dados compartilhados
- ⊳ já foi tentado, mas não cresce em escala
- Usar caches que espionam o barramento em uso comum
- Usar diretórios distribuídos junto com a memória

HEPR Dist RCC 57

Arquitetura II — coerência entre caches

2007-1

Métodos de Coerência para Barramentos



As caches usam escrita preguiçosa

(tipicamente)

- ▶ porque estas minimizam tráfego no barramento
- Os blocos nas caches podem estar num destes estados (típicos)
- ▷ INVÁLIDO: conteúdo do bloco não pode ser usado
- \triangleright VÁLIDO: não está sujo, compartilhado (há ≥ 1 cópia)
- ⊳ SUJO: única cópia suja
- ▶ RESERVADO: não-sujo, única cópia (aumenta eficiência)

HEPR Dinf RCC 57

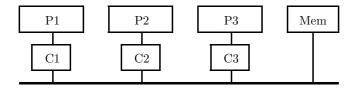
Arquitetura II — coerência entre caches

2007-1

Métodos de Coerência para Barramentos

- Espionagem (snooping)
- b todas as caches observam todo o tráfego no barramento
- ⊳ etiquetas com duas portas: CPU e barramento
- O que ocorre nas escritas?

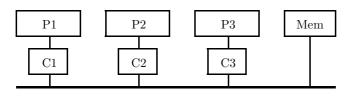
invalidar cópias nas outras caches: protocolos de invalidação atualizar cópias nas outras caches: protocolos de atualização



Protocolo de Invalidação

atividade	atividade	cache	cache	memória
no processador	no barramento	C1	C2	ender X
				0
P1 lê X	falta em X	0		0
P2 lê X	falta em X	0	0	0
P1 faz X=1	invalidação para X	1	-	0
P2 lê X	falta em X *	1	1	1

^{*} P1 responde à falta em C2 e atualiza memória



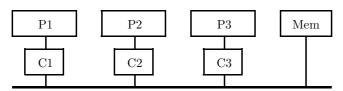
HEPR Dist RCC 574

Arquitetura II — coerência entre caches

2007-1

Protocolo de Atualização

atividade no processador	atividade no barramento	cache C1	cache C2	memória ender X
				0
P1 lê X	falta em X	0		0
P2 lê X	falta em X	0	0	0
P1 faz X=1	broadcast de X=1	1	1	1
P2 lê X	acerto em X	1	1	1



Há mais tráfego no barramento, se P2 não usar X novamente.

HEPR Dinf RCC 57

Arquitetura II — coerência entre caches

2007-1

Protocolo de Invalidação Simplificado

- Caches usam escrita preguiçosa
- Estados
- ▷ INVÁLIDO: conteúdo do bloco não pode ser usado
- ightarrow COMPART-ilhado: não está sujo, compartilhado (há ≥ 1 cópia)
- ▶ EXCLUSIVO: única cópia suja
- a cada transação no barramento, o controlador de cache:
- > verifica se bloco está na cache
- ▷ se estiver na cache, efetua mudança de estado cfe ME (adiante)

LIEPR Dinf RCC 576

Protocolo de Invalidação Simplificado (cont.)

H-P QA fig 6.11

HEPR Dist RCC 57

Arquitetura II — coerência entre caches

2007-1

Protocolo de Invalidação Simplificado (cont.)

H-P QA

fig 6.11

HEPR Dist RCC

Desempenho de Protocolos para Barramentos

- Tipos de faltas nas caches:

Arquitetura II — coerência entre caches

- Faltas por coerência são faltas adicionais causadas pelo protocolo de coerência
- Falso compartilhamento
- ▷ a unidade de coerência é um bloco de cache
- ⊳ pode ocorrer que um bloco inteiro seja compartilhado, mas palavras individuais não são compartilhadas

Escalabilidade

Multiprocessadores construídos com barramentos escalam bem até 10-20 processadores

barramento é o gargalo !!!

por que não usar mais barramentos? — dividir tráfego entre barramento de endereços & coerência, e barramento de dados

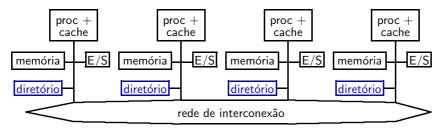
mesmo assim, limite é lá pelos 20 processadores...

HEPR Dist RCC 580

Arquitetura II — coerência entre caches

2007-1

Coerência de Caches com Diretórios



Para construir MPs grandes e escaláveis:

Informação sobre compartilhamento é mantida em diretórios distribuídos fisicamente pela memória

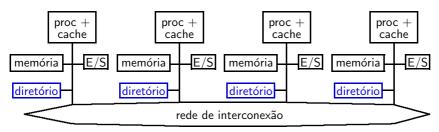
A informação de estado é mantida nas caches e nos diretórios em memória

HEPR DIst RCC 58:

Arquitetura II — coerência entre caches

2007-1

Coerência de Caches com Diretórios



Não há um barramento para serializar acessos de escrita

Transações do protocolo de coerência usam mensagens e recibos (ACKs)

Coerência de Caches com Diretórios

O espaço de endereçamento é alocado estaticamente entre os processadores

Endereço especifica número do nó e endereço no nó nó endereço nó de residência nó em cuja memória o endereço foi alocado, nó contém diretório que controla estado do bloco nó local nó que mantém cópia de bloco em sua cache nó é local c.r.a processador que está referenciando o bloco nó remoto nó que também mantém cópia do bloco cuja cópia deve ser invalidada

HEPR Dist RCC 588

Arquitetura II - coerência entre caches

2007-1

Coerência de Caches com Diretórios

Informação sobre compartilhamento mantida em diretório na memória:

- estados de blocos em memória:
- \triangleright shared: ≥ 1 processador/es tem cópia/s; memória atualizada
- ▶ uncached: não há cópias em caches
- ▶ exclusive: 1 processador tem cópia; memória está desatualizada
- elementos do diretório apontam para procs. com cópias
- ▷ p.ex. com vetor de bits um bit por processador
- tipos de apontadores:
- → mapa completo N processadores, N apontadores (bits)
- ⊳ limitado número fixo de apontadores (4-8)
- ⊳ lista lista encadeada de blocos de caches

HEPR DIAF RCC 55

Arquitetura II — coerência entre caches

2007-1

Coerência de Caches com Diretórios

- bloco no estado uncached → memória atualizada
- ⊳ falta na leitura: envia dado a quem pede;

marca bloco como shared; apontador=dono

▶ falta na escrita: envia dado a quem pede;

marca bloco como exclusive; apontador=dono

- bloco no estado shared → memória atualizada
- ⊳ falta na leitura: envia dado da memória a quem pede;

apontador=dono,novo

⊳ falta na escrita: envia dado da memória a quem pede;

invalida todas as outras cópias;

marca bloco como exclusive; apontador=dono

Coerência de Caches com Diretórios

- bloco no estado **exclusive** → memória des-atualizada
- ⊳ falta na leitura: pede dado a quem o possui;

atualiza memória;

marca bloco como **shared**; apontador=dono,novo

▶ Write-Back: marca bloco como uncached;

apontador=NIL

⊳ falta na escrita: pede dado a quem o possui;

marca bloco como **exclusive**; apontador=dono

HEPR Dinf RCC

Arquitetura II — coerência entre caches

Arquitetura II — coerência entre caches

2007-1

Coerência de Caches com Diretórios - ME cache

H-P QA

fig 6.29

HEPR DInf RCC 58

Coerência de Caches com Diretórios - ME diretório

H-P QA

fig 6.30

LIEPR Dinf RCC 588

Coerência de Caches com Diretórios - exemplo

Considere três processadores, P1, P2, P3, e a següência de eventos

- 1. bloco X não está em nenhuma cache
- 2. os três processadores efetuam leitura de X
- 3. P3 atualiza X (fica bloqueado)
- 4. acerto em C3, mas P3 não é dono do bloco
- 5. C3 requisita exclusividade (para poder escrever)
- 6. memória envia invalidações para C1 e C2
- 7. memória recebe respostas (acks), marca linha como exclusiva informa a C3 que agora é a dona do bloco
- 8. C3 atualiza bloco, marca-o como sujo P3 prossegue

LIEPR Dinf RCC

- KRU

Arquitetura II - sincronização e consistência

2007-1

Resumo

- Coerência de caches: aparência de tempo absoluto mantém ordem aparente das operações na memória
- Protocolos com espionagem em meio compartilhado (difusão)
- ⊳ invalidação invalida cópias
- Protocolos com comunicação ponto-a-ponto
- ⊳ invalidação
- ⊳ diretório distribuído mantém localização das cópias
 - * mapa completo
 - * mapa parcial
 - * lista encadeada

HEPR Dinf RCC

Arquitetura II — sincronização e consistência

Processamento Paralelo & Multiprocessadores

- Motivação
- Tipos de máquinas paralelas
- Coerência entre caches
- Sincronismo
- Ordenação de operações em memória

HEPR Dinf RCC 50:

Revisão

- Coerência de caches: aparência de tempo absoluto mantém ordem aparente das operações na memória
- Protocolos com espionagem em meio compartilhado (difusão)
- ⊳ invalidação invalida cópias
- Protocolos com comunicação ponto-a-ponto (redes)
- ▶ invalidação
- - * mapa completo
 - * mapa parcial
 - * lista encadeada

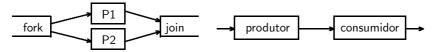
HEPR Dinf RCC 50

Arquitetura II — sincronização e consistência

2007-1

Sincronismo

- Sincronismo é necessário sempre que há processos concorrentes num sistema
 mesmo em uni-processador
- Forks e Joins: os processos de um programa paralelo podem ter que esperar até que certos eventos tenham ocorrido
- Produtor-Consumidor: o processo consumidor deve esperar até que o processo produtor tenha produzido dados
- Uso exclusivo de recurso: o sistema operacional deve garantir que somente um processo de cada vez use um recurso



HEPRING SO

rquitetura II — sincronização e consistência

2007-1

Execução atômica

/* Quais os valores de "print c" e "print d"? */

Comandos 'atômicos' em C não são executados atomicamente pelo processador

$$a = a + 1; \equiv lw r1,0(r2)$$

addi r1,r1,1
 $sw r1,0(r2)$

HEPR Dinf RCC 50.

Comunicação e Sincronização

Sincronização **pode** ser programada com loads e stores normais:

```
/* proc 1 */
                                /* proc 2 */
       flag1=0
                                       flag2=0
       . . .
                                       . . .
       flag1=1
                                       flag2=1
LAB1: if(flag2 == 1)
                              LAB2: if(flag1 == 1)
       goto LAB1
                                       goto LAB2
       /* seção crítica */
                                       /* seção crítica */
       flag1=0
                                       flag2=0
```

MAS é difícil de programar e depurar \longrightarrow este exemplo pode bloquear

HEPR Dinf RCC 50

Arquitetura II — sincronização e consistência

2007-1

Semáforos (locks)

Um **semáforo** é um T.A.D. composto de um inteiro não-negativo, e das operações: *E W Dijkstra, 1965*

P(s): if s > 0 decrementa s de 1; senão espera V(s): incrementa s de 1 e acorda um dos processos a esperar

P e V **devem** ser executados atomicamente, sem interrupções, ou acessos a **s** por outros processos

```
Processo( ) P(s) \hspace{1cm} \text{Valor inicial de } \textbf{s} \hspace{0.1cm} \text{determina número} \\ \hspace{0.1cm} \textit{região crítica} \hspace{1cm} \text{máximo de processos na região crítica.} \\ V(s) \hspace{1cm} \text{Valor inicial de } \textbf{s} \hspace{0.1cm} \text{determina número} \\ \hspace{0.1cm} \text{máximo de processos na região crítica.} \\ \text{V(s)} \hspace{0.1cm} \text{Valor inicial de } \textbf{s} \hspace{0.1cm} \text{determina número} \\ \hspace{0.1cm} \text{máximo de processos na região crítica.} \\ \text{V(s)} \hspace{0.1cm} \text{Valor inicial de } \textbf{s} \hspace{0.1cm} \text{determina número} \\ \hspace{0.1cm} \text{máximo de processos na região crítica.} \\ \text{V(s)} \hspace{0.1cm} \text{V(s)} \hspace{0.1cm} \text{Valor inicial de } \textbf{s} \hspace{0.1cm} \text{determina número} \\ \text{Maximo de processos na região crítica.} \\ \text{V(s)} \hspace{0.1cm} \text{V(s)} \hspace{0.1c
```

HEPR Dinf RCC 50

Arquitetura II — sincronização e consistência

2007-1

Exemplos de Primitivas de Sincronização (hw)

Test_and_set e Reset

Exemplos de Primitivas de Sincronização (hw)

Fetch_and_Add

HEPR Dinf RCC

Arquitetura II — sincronização e consistência

2007-1

Sincronização em modo usuário

Sincronização com spin locks:

→ baixa latência para obter trava
 → região crítica pode ficar travada por pouco tempo

```
li r2, 1
lockit: swap r2, 0(r1) # troca atômica
bnez r2, lockit # se falhou, repete
# região crítica
...
unlock: st r0, 0(r1) # libera
```

Em sistema com caches:

variável com a trava mantida na cache

swap ≡ escrita → montes de tráfego de coerência

HEPR Dinf RCC

500

Arquitetura II — sincronização e consistência

2007-1

Sincronização em modo usuário

Em sistema com caches:

variável com a trava mantida na cache

swap ≡ escrita → montes de tráfego de coerência

Alternativa: espera em loop até liberar, então tenta obter trava

```
tryit: ld r4, 0(r1) # em cache como SHARED # qdo outro proc escreve: # SHARED \rightarrow INVAL bnez r4, tryit # espera até ficar livre li r2, 1 # livre! talvez ainda está? lockit: swap r2, 0(r1) # troca atômica bnez r2, tryit # se falhou, repete
```

Exemplos de Primitivas de Sincronização - LL+SC

Load_linked + Store_conditional

não-bloqueante

- par de instruções usa registrador especial para manter flag, endereço e resultado do store_conditional
- load_linked lê valor e escreve endereço em reg de endereçamento global
- store_conditional escreve valor **se** reg de endereço global não foi alterado; retorna indicação de sucesso (1) ou falha (0)
- opera sob a premissa de que máquina tem barramento

HEPR Dist RCC 60:

Arquitetura II - sincronização e consistência

2007-1

Exemplos de Primitivas de Sincronização - LL+SC

- \star load_linked lê valor e escreve endereço em reg de endereçamento global
- * store_conditional escreve valor se reg de endereço global não foi alterado; retorna indicação de sucesso (1) ou falha (0)

TIEDR DIN RCC 60'

Arquitetura II — sincronização e consistência

2007-1

Exemplos de Primitivas de Sincronização - LL+SC

- * load_linked lê valor e escreve endereço em reg de endereçamento global
- * store_conditional escreve valor **se** reg de endereço global não foi alterado; retorna indicação de sucesso (1) ou falha (0)

Exemplo: troca atômica de r4 com mem[r1] atomic swap

```
try: mov r3, r4  # carrega valor para troca

l1 r2, 0(r1)  # load linked: r2\leftarrow mem[r1]

sc r3, 0(r1)  # store conditional:

# OK: mem[r1]\leftarrow r3; r3\leftarrow 1

# NO: r3\leftarrow 0

beqz r3, try  # se SC falhou, repete

mov r4, r2  # escreve valor trocado
```

Desempenho de locks

Instruções com read-modify-write atômico bloqueantes

test_and_set, fetch_and_ Φ , swap

VS

Instruções com read-modify-write atômico não-bloqueantes

load_linked+store_conditional

1/5

Protocolos baseados em load's e store's normais

Desempenho depende de vários fatores inter-relacionados: latência sem contenção

grau de contenção

serialização se há contenção

impede aumentar escala

caches

execução de load's e store's fora de ordem

HEPR Dinf RCC

PI IA

Arquitetura II — sincronização e consistência

2007-1

Atomicidade, Consistência e Ordenação de Eventos

- Atomicidade
- ★ informalmente: acessos a uma variável acontecem um de cada vez, e afetam todas as cópias imediatamente
- Consistência Seqüencial
- ★ informalmente: extensão intuitiva do modelo de memória dos uniprocessadores
- Consistência Relaxada
- \star informalmente: somente operações de sincronismo são consistentes \longrightarrow melhor desempenho
- Ordenação de Eventos

= característica da implementação

▶ ordenação forte – consistência seqüencial

⊳ ordenação fraca – consistência relaxada

coerência de caches → mesmo endereço consistência de caches → endereços distintos

HEPR Dinf RCC

....

Arquitetura II — sincronização e consistência

2007-1

Ordenação de Eventos na Memória

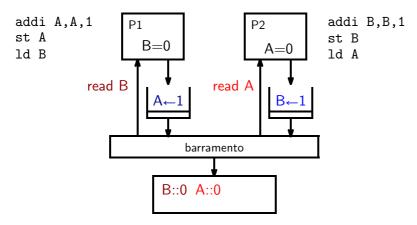
 Atualizações podem ser re-ordenadas pelo sistema de memória Exemplo:

```
P1: A=0; P2: B=0; ... A=1; B=1; L1: if (B==0) ... L2: if (A==0) ...
```

- intuitivamente, é impossível que ambos A e B sejam zero MAS isso pode ocorrer se as atualizações forem re-ordenadas pelo sistema de memória
- num Multi-Proc, as regras que ordenam a atualização em memória devem ser cuidadosamente definidas e atendidas pela

devem ser cuidadosamente definidas e atendidas pela implementação

Causas de Inconsistência – Filas de Escrita



leituras podem ultrapassar valores que estão esperando vez na fila porque endereços são distintos

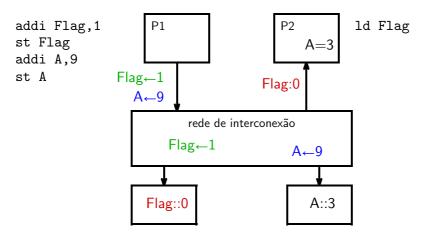
HEPR Dist RCC 607

Arquitetura II — sincronização e consistência

2007-1

2007-1

Causas de Inconsistência - Atrasos na Rede



tempo de viagem é ∞ à distância entre fonte e destino

TIEDR DIN RCC 603

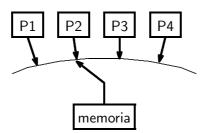
Arquitetura II — sincronização e consistência

2007-1

Consistência Sequencial

Um sistema é seqüencialmente consistente se o resultado de uma execução é o mesmo caso as operações de todos os processadores fossem executadas em alguma ordem seqüencial, e as operações de cada processador individual aparecem na ordem especificada pelo seu programa.

Leslie Lamport, 1979



2007-1

Consistência Sequencial

Um sistema é seqüencialmente consistente se o resultado de uma execução é o mesmo caso as operações de todos os processadores fossem executadas em alguma ordem seqüencial, e as operações de cada processador individual aparecem na ordem especificada pelo seu programa.

Leslie Lamport, 1979

P1:
$$\alpha, \beta, \gamma$$

P2: $\mathcal{P}, \mathcal{Q}, \mathcal{R}$
P3: $\clubsuit, \heartsuit, \spadesuit, \diamondsuit$

ordem global: $\clubsuit, \alpha, \beta, \mathcal{P}, \heartsuit, \mathcal{Q}, \gamma, \mathcal{R}, \spadesuit, \diamondsuit$

Consistência Seqüencial =

interleaving arbitrário que preserva a ordem de referências à memória de programas seqüenciais

HEPR Dist RCC 611

Arquitetura II — sincronização e consistência

Atomicidade

• Definição:

um acesso pelo processador \mathcal{P} à variável X é atômico se nenhum outro processador pode acessar qualquer cópia de X enquanto o acesso pelo processador \mathcal{P} estiver ocorrendo.

- atomicidade no hardware pode ser relaxada
- ⋆ o hardware pode parecer que tem comportamento atômico sem que todas as cópias sejam atualizadas simultaneamente
- Atomicidade é importante para sincronismo baseado em operações read-modify-write
- Problemas com ordenação podem ocorrer mesmo com atomicidade
- * atomicidade é uma propriedade de acessos à variáveis individuais, e não garante ordenação de atualizações de variáveis distintas.
- * operações atômicas são serializadas pelos mecanismos de escrita em memória

HEPR Dinf RCC

Arquitetura II — sincronização e consistência 2007-1

Ordenação de Execução (de Programas)

 Considere diferentes sequências de execução de instruções inicialmente A=B=C=0

• Só podemos observar indiretamente a seqüência real de eventos, com base no que é impresso

BC, AC, AB representa os valores impressos

- ⊳ se processo é executado em ordem, um conjunto de seqüências é possível: acbdef [10, 10, 11], ecdfab [11, 01, 01], etc
- ⊳ se instruções podem executar fora de ordem, outras seqüências são possíveis: adbcfe [00, 10, 11] e indesejáveis?

Ordenação de Execução (de Programas)

- Cada interleaving sugere que um conjunto específico de valores será impresso se os acessos à memória são atômicos
- ⊳ se acessos à memória não são atômicos, diferentes processadores podem perceber seqüências distintas:

[01, 10, 01] P1: $e \rightarrow c$ P2: $a \rightarrow e$ P3: $c \rightarrow a$

- atomicidade + operações de memória ordenadas
 → ordenação forte
- num sistema que é seqüencialmente consistente os eventos são fortemente ordenados

HEPR Dinf RCC 61

Arquitetura II — sincronização e consistência

2007-1

Coerência entre caches (revisitada)

- Informalmente, coerência entre caches é um método para fazer com que escritas em qualquer local na memória sejam percebidos na mesma ordem pelos outros processadores

HEPR Dinf RCC 61

Arquitetura II — sincronização e consistência

2007-1

Consistência Relaxada

- Consistência seqüencial pode diminuir desempenho porque
 a consistência é forte
 e as restrições de ordenamento tbém
- Consistência relaxada: pode relaxar

ordem de execução ou atomicidade da escrita

- ⊳ relaxa ordem de escritas para leituras
- ⊳ relaxa ordem de escritas para escritas
- ⊳ relaxa ordens de leituras para leituras, e leituras para escritas
- ▷ lê escritas pelos outros adiantadamente
- ▶ lê escritas próprias antecipadamente
- relaxamento é geralmente determinado pela implementação
 em hardware
 vários modelos já foram propostos
- ∃m mecanismos que permitem ignorar relaxamento

Consistência Relaxada

- Consistência seqüencial pode diminuir desempenho
 porque a consistência é forte
 e as restrições de ordenamento tbém
- Consistência relaxada: pode relaxar
 ordem de execução ou atomicidade da escrita
- * relaxa ordem de escrita para leitura
- * relaxa ordem de escrita para escrita
- * relaxa ordens de leitura para leitura, e leitura para escrita
- ★ lê escrita pelos outros adiantadamente
- * lê escritas próprias antecipadamente
- relaxamento é geralmente determinado pela implementação em hw
- ∃ mecanismos que permitem ignorar relaxamento
- ▷ cercas (fence operations)
- ⊳ semântica das instruções de sincronização explícitas

HEPR DINF RCC 61

Arquitetura II — sincronização e consistência

2007-1

Cercas em memória

Cercas são instruções para

fences

forçar serialização dos acessos à memória

Processadores com modelos de ordenação fracos nos quais load's e store's para endereços distintos podem ser re-ordenados devem prover instruções para forçar a serialização dos acessos à memória

Cercas são operações custosas mas paga-se o custo da serialização somente quando necessário

HEPR Dinf RCC 61

Arquitetura II — sincronização e consistência

2007-1

Ordenação Fraca

Ordenação fraca:

não se presume nada sobre ordenação de operações **entre** pontos de sincronização explícita

- regras informais:
- ▶ instruções de sincronização notificam hardware (controladores de caches)
- ▶ ações para manter consistência ocorrem somente nas instruções de sincronização

2007-1

Programas propriamente sincronizados

Variáveis de sincronização (mutexes) são disjuntas das variáveis com dados

acessos à variáveis compartilhadas para escrita são protegidos por regiões críticas → não há corridas exceto pelos semáforos/locks

Em geral, não é possível provar que um programa é livre de corridas no acesso aos dados.

TIEPR DINF RCC
610

Arquitetura II — sincronização e consistência 2007-1

Programas propriamente sincronizados

Modelos de ordenação relaxados permitem a re-ordenação de instruções pelo compilador, ou processador, desde que a re-ordenação não atravesse uma cerca

Além disso, o processador não pode executar especulativamente, ou fazer busca antecipada, atravessando cercas.

HEDRING RCC 620

Resumo

- \bullet Primitivas de sincronização em barramento (que serializa acessos) test_and_set, fetch_and_ Φ
- Primitivas de sincronização distribuídas (sem serializar acessos)
 load-linked e store-conditional
- oerência de caches → mesmo endereço

Arquitetura II - entrada e saída

- consistência de memória → endereços distintos
- Consistência seqüencial: a ordem interna dos acessos de cada processo é respeitada
- Consistência fraca: ordem de acessos não é respeitada
- Cercas garantem pontos de sincronização ordenados

Entrada e Saída

- Desempenho
- Tipos e Características de Dispositivos
- Arquitetura do Sistema de E/S
- **Barramentos**
- ▶ Processamento de E/S
- Discos
- ⊳ Sistemas de Discos de Alto Desempenho

HEPR Dinf RCC 69

Arquitetura II — entrada e saída 2007-1

Desempenho

I/O certainly has been lagging in the last decade Seymour Cray, 1976

Also, I/O needs a lot of work
David Kuck, 1988

If, ..., performance of CPUs improves as 55% per year and I/O did not improve, every task would become I/O bound

H&P, 2001

HEPR Dist RCC 623

Arquitetura II — entrada e saída 2007-1

Quem se importa com desempenho de E/S?

- Desempenho de CPU cresce 55% ao ano
- Desempenho de sist de E/S é limitado por partes mecânicas (discos)

melhora <10% aa em operações/segundo ou MBytes/segundo

• Lei de Amdahl:

ganho de desempenho é limitado pela parte mais lenta

ightharpoonup 10x CPU e 10% E/S \longrightarrow desempenho cresce 5x perde 50%

 \triangleright 100x CPU e 10% E/S \longrightarrow desempenho cresce 10x perde 90%

- Gargalo em E/S
- ⊳ fração cada vez menor do tempo na CPU

Quem se importa com desempenho de E/S?

- Se desempenho de E/S é limitado por partes mecânicas,
 - e desempenho da eletrônica segue a Lei de Moore,
- **então** é melhor aumentar complexidade da eletrônica para melhorar desempenho global dos sistemas:
- reduzir número de fios + aumentar veloc da comunicação USB
- aumentar complexidade do controlador de disco caches+DMA
- melhorar a saída do adaptador gráfico mesmo monitor

HEPR Dist RCC 695

Arquitetura II — entrada e saída 2007-1

Desempenho de E/S

There is an old network saying:

Bandwidth problems can be cured with money.

Latency problemas are harder because the speed of light is fixed, you can't bribe God

David Clark

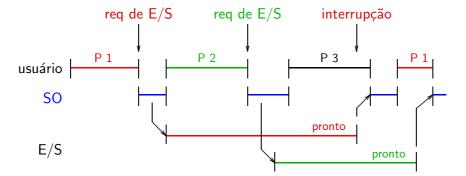
- Produção (throughput)
- ▶ largura de banda ou operações por segundo
- Latência
- ▶ tempo decorrido depende de concorrência
 latência de operações escondida se concorre com computação

HEPR Dist RCC 621

Arquitetura II — entrada e saída 2007-1

Concorrência entre E/S e Computação

E/S concorre com computação de maneiras complexas



Desempenho

$$\mathcal{T}_{\mathsf{tarefa}} = \mathcal{T}_{\mathsf{cpu}} + \mathcal{T}_{\mathsf{ES}} - \mathcal{T}_{\mathsf{concorr}}$$

Exemplo 1: E/S completamente escondido: 10 = 10 + 4 - 4

Exemplo 2: acelera CPU 2x, qual é o novo desempenho?

a)
$$\mathcal{T}_{\mathsf{tarefa}} = 5 + 4 - 4 = 5$$
 melhor \leftarrow esconde

b)
$$\mathcal{T}_{\mathsf{tarefa}} = 5 + 4 - 0 = 9$$
 pior \leftarrow expõe

c)
$$\mathcal{T}_{\mathsf{tarefa}} = 5 + 4 - 2 = 7$$
 médio \leftarrow esconde um pouco

HEPR Dinf RCC 69

Arquitetura II — entrada e saída 2007-1

Exemplos de Aplicações

- Supercomputadores
- * taxa de transferência é importante
 - → muitos MBytes/seg para arquivos grandes
- Processamento de Transações
- * taxa de serviço é importante
- * acessos aleatórios aos dados
 - → taxa de acessos a disco por segundo
- Sistemas de Arquivos para Multi-programação (ex. Unix)
- * tratamento de arquivos pequenos é importante
- * acessos seqüenciais
- * muitas criações e remoções

HEPR Dist RCC 620

Arquitetura II — entrada e saída 2007-1

Características dos Dispositivos

- Comportamento
- ⊳ entrada (lê uma vez só)
- ⊳ saída (escreve uma vez só)
- ▷ armazenagem (lê muitas vezes, também escreve)
- contra-parte
- ▶ humano
- taxa de transferência
- ⊳ taxa de pico

Tipos de Dispositivos

dispositivo	comportamento	parceiro	taxa [bytes/seg]
teclado	entrada	humano	10
mouse	entrada	humano	20
scanner	entrada	humano	400 K
monitor gráf	saída	humano	60 M
rede local	E/S	computador	0.5-10 M
fita	armazenagem	computador	2 M
disco	armazenagem	computador	2-10 M

HEPR Dinf RCC

Arquitetura II — entrada e saída 2007-1

Arquitetura do Sistema de E/S

Hierarquia de vias:

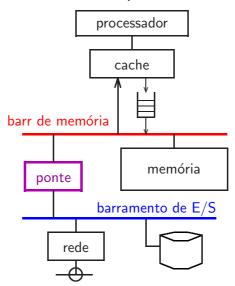
largura de banda é menor a medida em que desce na hierarquia barramentos distintos em cada nível

Processamento de E/S:

controlado por programa **ADM**

processadores de E/S

Arquitetura II — entrada e saída



2007-1

HEPR Dinf RCC

Características dos Barramentos

Opção	alto desempenho	baixo custo
linhas de dados, endereços ≠s	sim	não
largura do barram de dados	largo	estreito
tamanho das transferências	rajadas	palavra
número de mestres	mais de um	um
barramento em pipeline	sim	não
operação	síncrona	assíncrona

LIEPR Dinf RCC

Características dos Barramentos

- Barramento ≡ meio de comunicação compartilhado → broadcast
- Barramento de E/S: dispositivos com grande faixa de valores de vazão e latência
- ▶ Mestre é capaz de iniciar transferência
- ▶ Escravo responde aos comandos do mestre
- Num barramento em pipeline, há dois tipos de transações: comandos/requisições e respostas

```
transação = requisição → processamento → resposta
```

HEPR DINF RCC 63

Arquitetura II — entrada e saída

2007-1

Tipos de Barramentos – CPU-mem

- Barramento que interliga CPU e memória
- * desempenho muito elevado:

800MByte/s a 64Gbyte/s

largura de 64 a 256 bits,

relógio de 100 a 400 MHz

pode usar as duas bordas do relógio (double data rate)

- * projeto altamente otimizado
- * número de dispositivos é limitado pelo projeto # pequeno!
 menor comprimento → maior velocidade
- * tendência: maior velocidade e menor largura
- Barramento entre CPU e cache
- * desempenho mais alto ainda
- Barramento de E/S

HEPR Dinf RCC

Arquitetura II — entrada e saída

2007-1

Tipos de Barramentos - E/S

- Barramento que liga CPU e memória
- Barramento entre CPU e cache
- Barramento de E/S
- * compatibilidade com algum padrão é importante (PCI, SCSI)
- * projetado para tolerar número variado de dispositivos
- * projetado para tolerar diversidade em vazão e latência
- * maior comprimento → menor velocidade
- * recentemente:

usa mais eletrônica para simplificar (menos fios) e/ou

maior desempenho (protocolos mais complexos)

Método de Comutação do Barramento

- Barramento com comutação de circuito circuit switched
- * similar a um telefonema
- * posse do barramento é mantido até que transação completa
- * protocolo mais simples
- * latência dos dispositivos afeta utilização do barramento
- Barramento com comutação de pacotes

packet switched

- * similar ao envio de uma carta
- * barramento é liberado após envio do comando outro mestre pode utilizar barramento até que resposta seja emitida
- * protocolo mais complexo mas permite segmentação
- * circuito de controle mais complexo, mas utilização do barramento é maior

HEPR DInf RCC 63

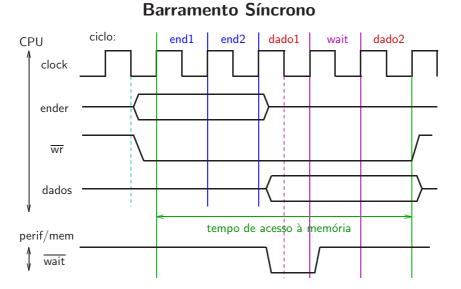
Arquitetura II — entrada e saída 2007-1

Barramento Síncrono vs Assíncrono

- Barramento Síncrono
- * maior velocidade de operação → maior vazão
- * projeto elétrico complexo e restrições fortes na temporização
- * maior desempenho porque não tem realimentação sinais percorrem barramento em somente um sentido: mestre → escravo
- Barramento Assíncrono

LIEDR Dinf RCC 63

Arquitetura II — entrada e saída 2007-1



Barramento Síncrono vs Assíncrono

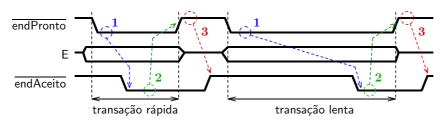
- Barramento Síncrono
- Barramento Assíncrono
- * menor velocidade de operação → menor vazão
- * projeto elétrico/temporal mais simples
- * pior desempenho porque tem realimentação sinais percorrem todo o barramento em dois sentidos: mestre(comando)→ escravo; escravo(resposta)→ mestre
- * a cada evento, sinais devem ser sincronizados aos relógios internos dos dispositivos
 - → relógios ≠s em dispositivos distintos (freqüência e fase)

HEPR Dist RCC 640

Arquitetura II — entrada e saída 2007-1

Barramento Assíncrono

Transação de Endereçamento

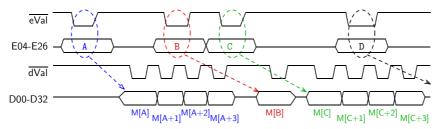


SEQÜÊNCIA DE EVENTOS

- 1: mestre inicia transação enderPronto escravo responde com enderAceito 2: mestre responde com enderPronto 3: escravo completa transação com enderAceito
- HEDR Dinf RCC 641

Arquitetura II — entrada e saída 2007-1

Split Transaction Bus



- Barramento é liberado após envio do comando outro mestre pode usar barram até que resposta seja emitida
- protocolo: pended packet switched
- maior utilização do barramento
- pacotes com identificador de transação

Pergunta

Num sistema de E/S, a CPU é o mestre ou o escravo no barramento?

- a) mestre
- b) escravo
- c) mestre e escravo
- d) mestre ou escravo
- e) nenhuma das acima

HEPR Dist RCC 641

Arquitetura II — entrada e saída 2007-1

Resposta

Num sistema de E/S, a CPU é o mestre ou o escravo no barramento?

- a) mestre
- b) escravo
- c) mestre e escravo
- d) mestre ou escravo, dependendo se transação é leitura ou escrita e) nenhuma das acima

HEPR Dinf RCC 64

Arquitetura II — entrada e saída 2007-1

Pergunta

Qual o tipo de barramento que permite maior nível de sobreposição de comunicação (através do barramento) com computação (pela CPU e/ou pelos dispositivos)?

- a) barramento síncrono
- b) barramento assíncrono
- c) split transaction síncrono
- d) split transaction assíncrono
 - e) nenhum dos acima

LIEPR Dinf RCC 645

Resposta

Qual o tipo de barramento que permite maior nível de sobreposição de comunicação (através do barramento) com computação (pela CPU e/ou pelos dispositivos)?

- a) barramento síncrono
- b) barramento assíncrono
- c) split transaction síncrono
- d) split transaction assíncrono
 - e) nenhum dos acima

As transações divididas permitem maior superposição entre as transações dos diversos pares mestre-escravo...

HEPR Dinf RCC 64

Arquitetura II — entrada e saída

2007-1

Processamento de E/S (i)

- E/S controlada por programa
- * CPU gerencia explicitamente todas as transferências loop: load→ store
- * pode causar enorme desperdício de tempo de CPU
- * não é muito usado há casos especiais interessantes
 ▷ aplicações embutidas de baixo custo
 - ▶ SO verifica todos dispositivos lentos de uma vez
- Acesso Direto à Memória
- Processadores de E/S dedicados

HEPR Dinf RCC

Arquitetura II — entrada e saída

2007-1

Processamento de E/S (ii)

- E/S controlada por programa
- Acesso Direto à Memória Direct Memory Access = DMA
- Processador inicializa controlador de DMA com endereço inicial, tamanho,
 - tipo da transferência
- * controlador de DMA dispara transferência de um bloco e gera interrupção após completar transferência
- * bom para volumes grandes e periféricos gulosos
- Processadores de E/S dedicados

Arquitetura II — entrada e saída 2007

Processamento de E/S (iii)

- E/S controlada por programa
- Acesso Direto à Memória
- Processadores de E/S dedicados
- * controladores capazes de executar 'programas' de E/S podem ser quase tão sofisticados quanto a CPU
- * podem ser de uso geral ou dedicado
- * e.g. canais de E/S da IBM

HEPR Dinf RCC 64

Arquitetura II — entrada e saída

2007-1

Comunicação com Processadores de E/S

Início/controle das operações

- Periféricos mapeados como memória
- * acesso em faixa de endereços de E/S para disparar operação
- * dispositivos no mesmo espaço de endereçamento que a memória
- * faixa de endereços de E/S são protegidos pelo SO
- * endereços não podem ser carregados nas caches (em geral)
- * dispositivo modelado por estrutura de dados
 → programação simples
- Periféricos acessados com instruções especiais

HEPR Dinf RCC

Arquitetura II — entrada e saída

2007-1

Comunicação com Processadores de E/S

Início/controle das operações

- Periféricos mapeados como memória
- Periféricos acessados com instruções especiais
- * instruções especiais iniciam operações de E/S in, out
- * instruções são previlegiadas executam somente em modo supervisor
- * difícil de modelar dispositivos em linguagem de alto nível

Arquitetura II — entrada e saída 2007-

Comunicação com Processadores de E/S

Notificação de término das operações

• E/S por programa

polling

- > processador fica em loop até que bit de status mude
- ⊳ latência elevada
- Interrupção
- > quando operação completa, interrompe processador
- ⊳ baixa latência

HEPR Dist RCC 65'

Arquitetura II — entrada e saída

2007-1

2007-1

Gerenciamento de E/S

- SO gerencia recursos de E/S
- ★ serializa requisições para recursos compartilhados
- ⋆ aloca espaço em disco, etc
- aplicativos invocam serviços do SO chamadas de sistema
- SO converte requisições em invocações de drivers
- ★ driver contém programa específico para dispositivo → comandos, registradores, etc
- driver comanda dispositivo

Arquitetura II — entrada e saída

⋆ operações específicas para dispositivo ou barramento



HEPR Dinf RCC

Proteção de E/S

- Instruções de E/S são previlegiadas
- ⊳ somente o SO pode executar instruções de E/S
- Dispositivos mapeados em memória são protegidos contra acessos em modo usuário,
 - e seus endereços/valores não são armazenados em cache
- ⊳ somente o SO pode acessar áreas de E/S mapeadas em memória
- Código do SO verifica cada chamada de sistema por usuários
- ⊳ se operação é permitida, então é executada

HEPR Dinf RCC 65.

resumo - Sistemas de E/S

- ullet $\mathcal{T}_{\mathsf{tarefa}} = \mathcal{T}_{\mathsf{cpu}} + \mathcal{T}_{\mathsf{ES}} \mathcal{T}_{\mathsf{concorr}}$
- Hierarquia de barramentos desempenho $des(CPU-cache) \gg des(cache-mem) \gg des(E/S)$
- Tipos de barramento:
- ★ comutação de circuitos × comutação de pacotes
- ★ síncrono × assíncrono × split transaction
- Processamento de E/S
- * por programa
- * acesso direto a memória
- * canais de E/S
- Mapeamento do espaço de endereçamento
- * como memória × como E/S
- Notificação de término: interrupção × polling

HEPR Dist RCC

Arquitetura II — discos 2007-1

Entrada e Saída

- Desempenho
- Tipos e Características de Dispositivos
- Arquitetura do Sistema de E/S
- Discos
- ▶ Mecanismo, componentes de gravação e de posicionamento
- ▶ Controlador
- ▶ RAID Sistemas de Discos de Alto Desempenho

fonte:

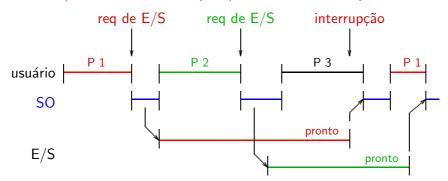
[RW94] An introduction to Disk Drive Modeling, Ruemmler & Wilkes, IEEE Computer 27(3):17-28, Mar 1994

HEPR DIst RCC 654

Arquitetura II — discos 2007-1

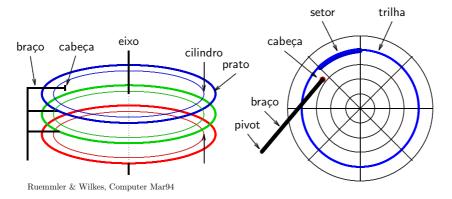
revisão - Concorrência entre E/S e Computação

E/S concorre com computação de maneiras complexas



$$\mathcal{T}_{\mathsf{tarefa}} = \mathcal{T}_{\mathsf{cpu}} + \mathcal{T}_{\mathsf{E/S}} - \mathcal{T}_{\mathsf{concorr}}$$

Discos Magnéticos



HEPR DIst RCC 658

Arquitetura II — discos 2007-1

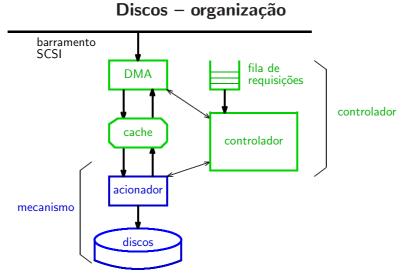
Discos Magnéticos – Parâmetros Típicos

Característica	mín	máx
diâmetro [polegadas]	1,0	* 3,5
capacidade formatada [GB]	1	>100
discos/pratos	1	20
trilhas por superfície	6.000	25.000
setores por trilha	100	600
bytes por setor	512	4.096
velocidade [rpm]	3.600	15.000
cache [MB]	0,125	≥4
taxa transferência [MB/s]	2,5-5	27-40

^{*} tamanho típico (+popular) em 2000

HEPR Dinf RCC 650

Arquitetura II — discos 2007-1



Ruemmler & Wilkes, Computer Mar
94 $\,$

Mecanismo - componentes de gravação

- Diâmetro: 1.0, 1.3, 2.5, 3.5 ,8 polegadas
- densidade linear de gravação: 100 Kbpi (bits/inch)
- densidade de trilhas: 20 Ktpi
- efeito combinado: densidade por área cresce 60% aa
- velocidade 3.600, 7.200, 10.000, 15.000 cresce 12% aa
- uma cabeça ativa por vez, taxa de leitura ≥ 100 Mbps
- conteúdo de um setor
- ▷ número do setor:
- ⊳ espaço;
- ▶ informação do setor com Cód Det Corr Erros;
- ▷ espaço;

> · · ·

LIEPR Dinf RCC

661

Arquitetura II — discos

2007-1

Mecanismo - componentes de posicionamento (i)

- Densidade é tão alta que noção de cilindro é quase irrelevante
- busca consiste de
- riangle aceleração até atingir $1/2~V_{
 m m\acute{a}x}$
- $riangleright V_{ ext{máx}}$ em distâncias longas
- ⊳ desaceleração até trilha desejada
- ⊳ estabilização da cabeça sobre a trilha (1-3 ms)
- \triangleright seeks longos: $\mathcal{T} \propto \operatorname{distância} (V_{\max})$
- \triangleright seeks médios: $\mathcal{T} \propto \sqrt{\mathrm{dist\hat{a}ncia}}$
- recalibrar posição a cada 15-30 min, durante 500-800 ms
- acompanhamento de trilhas

HEPR Dinf RCC

Arquitetura II — discos

Mecanismo – componentes de posicionamento (ii)

- Densidade é tão alta que noção de cilindro é quase irrelevante
- busca: aceleração; curso; desaceleração; estabilização
- recalibrar posição a cada 15-30 min, durante 500-800 ms
- acompanhamento de trilhas
- \triangleright troca de cabeças \rightarrow reposicionar braço (≈ 0.5 -1.5 ms)
- \triangleright troca de trilha \rightarrow estabilizar braço (\approx 1-3 ms)

 $(\approx 0.75 \text{ ms até estabilizar para escrever})$

Mecanismo – leiaute dos dados

- Disco visto pelo SO como vetor linear de blocos (256-1024 bytes)
- controlador mapeia vetor nos setores físicos $1D \sim 3D$ $V[i] \sim d[$ setor, trilha, superfície]
- #bits cresce ≈ linearmente com comprimento da trilha zoneamento: número de setores depende do raio
 ∃ 3-50 zonas com mesmo número de setores / zona
- deslocamento de setores nas trilhas:

setor0 de cada trilha deslocado para esconder tempo de reposicionamento track skewing

- trilhas/setores sobressalentes: referências a setores danificados são re-mapeadas para setores/trilhas de reserva
- ★ na formatação pula endereço da trilha com defeito slip sparing
- \star em uso re-mapeia endereço do setor/trilha para sobressalentes

TIEDR Dief RCC 6

Arquitetura II — discos 2007-1

Controlador (i)

- Funções do controlador SCSI

- > transferir dados entre disco e cliente
- operação do controlador custa 0.3-1 ms (caindo lentamente) eletrônica segue Lei de Moore mas funcionalidade mais complexa com o tempo
- interface com barramento
- buffer usado como cache

HEPR Dinf RCC 66

Arquitetura II — discos 2007-1

Controlador (ii)

- Funções do controlador SCSI
- operação do controlador custa 0.3-1 ms (caindo lentamente) eletrônica segue Lei de Moore mas funcionalidade cada vez mais complexa
- interface com barramento

- buffer usado como cache

Controlador - cache (leitura)

Políticas da cache: read-ahead ≈ busca antecipada

- usa × ignora dados em acerto parcial?
- ignora cache em requisições grandes? segmentação dos acessos
- mantém × descarta bloco após transferência?
- on-arrival read-ahead: assim que chegar na trilha, lê trilha toda
- read-ahead agressivo: atravessa trilhas e/ou cilindros
- read-ahead conservador: pára no final de trilha/cilindro
- cache associativa: particionar cache para ≠s seqüências entrelaçadas
- como cache do disco interfere com buffer cache do Unix?

HERR DIM RCC 663

Arquitetura II — discos 2007-1

Controlador - cache (escrita)

Cache pode corromper sistema de arquivos se faltar energia

- controlador avisa que completou operação após escrever na cache
- isso deve ser evitado na escrita de metadados
- se cache tem bateria, problema desaparece o que fazer no boot?
- ganhos adicionais
- > controlador pode otimizar escrita de vários blocos
- ▶ mesmos problemas que fila de escrita RAW, WAW
- como cache do disco interfere com buffer cache do Unix?

Cache com fila de comandos: controlador pode otimizar operações porque conhece geometria do disco

HEPR DInf RCC 668

Arquitetura II — discos 2007-1

Operação de Discos

- Comandos
- ▶ Latência/atraso no controlador + tempo na fila (OS)
- ⊳ 0,5ms se não encontra na cache, 0,1ms se encontra na cache
- Seek/busca movimentação do braço entre trilhas/cilindros
- ⊳ move a cabeça até a trilha desejada
- ⊳ tempo depende da posição inicial da cabeça
- Latência rotacional
- \triangleright espera até que setor desejado passe sob a cabeça ($\approx 1/2$ volta)
- \triangleright na média, 0.5/rpm \rightarrow 0,5/(7200rpm/60spm)=4,2ms
- Transferência de dados
- b taxa de transferência entre 2 e 40 MByte/s

Desempenho

- Tempo médio de acesso
 - = tempo médio de movimentação do braço (seek)
 - + latência rotacional média
 - + tempo de transferência
 - + tempo do controlador
- Exemplo: 7200 rpm, 10MByte/s

tempo médio de movimentação do braço: 10ms

tempo do controlador: 0,5ms

tempo para ler bloco de 4Kbytes (uma página)

$$10ms + 0.5/(7200rpm/60spm) + 4KB/10MB/s + 0.5ms$$

$$10ms + 4.2ms + 0.4ms + 0.5ms = 14.65ms$$

HEPR Dief RCC 670

Arquitetura II — discos 2007-1

Desempenho (cont)

	distância	
	# trilhas	fração
	0	24%
Localidade:	15	23%
discos exibem localidade	30	8%
→ em "acessos locais"	45	4%
tempo de busca (seek) cai em $1/3$	60	3%
	75	3%
Cache:	90	1%
buffer em memória (<i>Unix buffer cache</i>)	105	3%
e na unidade de disco	120	3%
→ latência cai para hit+transferência	135	2%
	150-165	(3+2)%
	180-195	(3+3)%
	unix time	e-sharing

HEPR Dist RCC

2007-1

Tendências

Arquitetura II — discos

Capacidade
$$\nearrow 100\%$$
 aa $(2x / 1.0 \text{ anos})$
Taxa de Transferência $\nearrow 40\%$ aa $(2x / 2.0 \text{ anos})$
Rotação+posicionamento $\searrow 8\%$ aa $(0.5x / 10 \text{ anos})$
MB/\$ $\nearrow 100\%$ aa $(2x / 1.5 \text{ anos})$

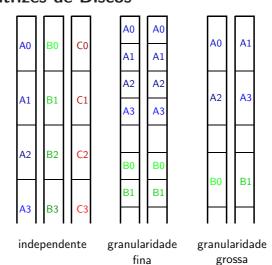
Matrizes de Discos

Conjunto de discos individuais: cada disco com seu braço/cabeça

Distribuição dos dados: endereçamento independente

listras de blocos pequenos listras de

blocos grandes



HEPR Dinf RCC 673

Arquitetura II — discos 2007-1

Matrizes de Discos (i)

- Endereçamento independente
- ⊳ software/usuário distribui os dados
- ▷ balanceamento de carga entre discos pode ser problemático
 - → é difícil de sincronizar acessos em paralelo
- Listras de blocos pequenos

fine-grain striping

• Listras de blocos grandes

coarse-grain striping

• Granularidade escolhida em função da aplicação e tipo de carga

HEPR Dinf RCC 65

Arquitetura II — discos 2007-

Matrizes de Discos (ii)

- Endereçamento independente
- Listras de blocos pequenos

fine-grain striping

- ▷ |bloco| = um bit, um byte, ou um setor
- ▶ #discos*|bloco| define menor quantidade de dados acessível
- ▶ balanceamento de carga perfeito: 1 requisição atendida por vez
- \triangleright taxa efetiva de transferência \approx N vezes melhor que um disco só
- b tempo de acesso pode aumentar,
 - a não ser que discos sejam sincronizados
- Listras de blocos grandes

coarse-grain striping

- * paralelismo na transferência de grandes volumes de dados
- ★ concorrência para transferências pequenas
- * balanceamento de carga pela aleatoriedade
- Granularidade escolhida em função da aplicação e tipo de carga

HEPR Dinf RCC 67^t

Matrizes de Discos (iii)

- Endereçamento independente
- Listras de blocos pequenos

fine-grain striping

- \star |bloco| = um bit, um byte, ou um setor
- * #discos*|bloco| define menor quantidade de dados acessível
- ★ balanceamento de carga perfeito: 1 requisição atendida por vez
- \star taxa efetiva de transferência pprox N vezes melhor que um disco só
- * tempo de acesso pode aumentar, a não ser que discos sejam sincronizados
- Listras de blocos grandes

coarse-grain striping

- ⊳ paralelismo na transferência de grandes volumes de dados
- ⊳ balanceamento de carga pela aleatoriedade
- Granularidade escolhida em função da aplicação e tipo de carga

 HEPP Disf RCC
 676

 Arquitetura II — discos
 2007-1

Mecanismos de Redundância

- Falhas em discos são parcela grande de falhas de hardware
- → striping aumenta o número de arquivos perdidos por falha
- Replicação dos dados
 - espelhamento dos discos
 - → permite leituras em paralelo
 - → escritas devem ser sincronizadas
- Proteção com paridade
 - → usar disco para manter a paridade

Arquitetura II — discos 2007-

Redundant Arrays of Inexpensive Disks - RAIDs

- Conjuntos de discos pequenos e baratos resultam em alto desempenho e alta confiabilidade
 - D = número de discos de dados no conjunto
 - V = número de discos de verificação no conjunto
- Nível 1: discos espelhados (D=1, V=1)
- ▷ desperdício é elevado
- Nível 2: código de detecção de erros (D=10, V=4)
- ▶ mesmo tipo de código de detecção de erros usado com DRAMs
- ▶ agrega bits atualizados com bits que permanecem; recomputa a paridade
- > re-escreve todo o conjunto, incluindo a verificação

Redundant Arrays of Inexpensive Disks - RAIDs

- Nível 3: paridade com bits entrelaçados (D=4, V=1)
- ⊳ disco com falha é identificado facilmente pelo controlador
- ⊳ não é necessário código especial para descobrir disco em falha
- Nível 4: paridade com blocos entrelaçados

- ⊳ escrita deve atualizar disco com dados e paridade

HEPR Dinf RCC 67

Arquitetura II — discos 2007-1

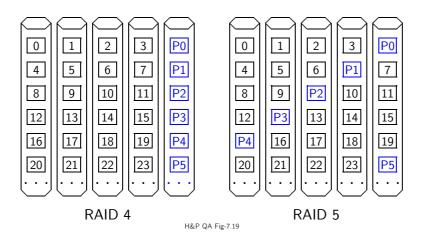
Redundant Arrays of Inexpensive Disks - RAIDs

- Nível 5: paridade distribuída com blocos entrelaçados
- ⊳ paridade é distribuída pelos discos no conjunto
- > atualizações distintas de paridade vão para discos distintos
- Nível 6: matriz bi-dimensional
- → matriz de dados é bi-dimensional, com paridade nas linhas e nas colunas
- ▷ permite recuperação de duas falhas

HEPR DIst RCC 68

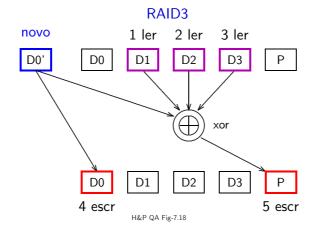
Arquitetura II — discos 2007-1

RAID 4/5 – Blocos de Paridade



RAID - Atualização "Pequena"

Qual é o número de operações de leitura/escrita nos discos individuais para efetuar escrita de poucos dados (= atualização pequena)?

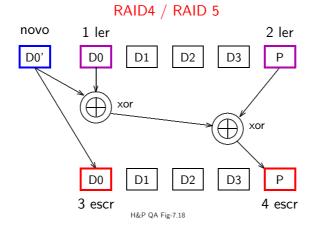


HEPR Dief RCC 68°

Arquitetura II — discos 2007-1

RAID - Atualização "Pequena"

Qual é o número de operações de leitura/escrita nos discos individuais para efetuar escrita de poucos dados (= atualização pequena)?



HEPR Dinf RCC 68

Arquitetura II — redes de interconexão

resumo - Discos

2007-1

- Tempo médio de acesso
 - = tempo médio de movimentação do braço (seek)
 - + latência rotacional média
 - + tempo de transferência
 - + tempo do controlador
- Cache no controlador para tirar proveito de localidade
- * falta de energia durante escrita de metadados corrompe sist de arquivos
- * mesmos problemas que fila de escrita (riscos RAW a WAW)
- RAID usar discos baratos para aumentar desempenho – striping e melhorar confiabilidade – paridade

Redes e Clusters

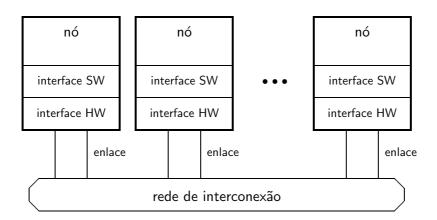
- Protocolos
- Escopo, Tipos de Redes, Meios Físicos
- Espaço de projeto
- ▶ topologias
- ⊳ comutação
- ⊳ sincronicidade
- Roteamento
- Controle de Tráfego

HEPR Dist RCC 688

Arquitetura II — redes de interconexão

2007-1

Modelo de rede

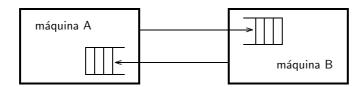


HEPR DIst RCC 68

Arquitetura II — redes de interconexão

2007-1

Modelo de comunicação



- Quem transmite?
- quando transmite?
- como interpreta mensagem?
- como garante que mensagem chegou ao destino?
- como garante que mensagem está correta?
- protocolo determina as respostas

Protocolos

Protocolo deve contemplar:

- endereçamento → máquinas + processos
- proteção entre processos na mesma máquina
- tipos de mensagem: request, reply, acknowledgement
- entrega confiável
- ▶ mensagem perdida? retransmissão + temporizadores + buffers
- minúcias:
- > duplicação de mensagens
- > ordenação de mensagens

HEPR Dinf RCC

Arquitetura II — redes de interconexão

2007-1

2007-1

Anatomia de uma mensagem

ender fonte nome de máquina

ender destino

processo fonte nome de processo ou porta

processo destino

tipo de msgm {req, rsp, ack, crtlFluxo}

núm seqüência seqüenciamento, duplicações

controle de fluxo buffer overflow/underflow

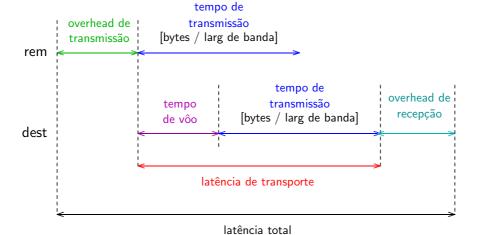
checksum detecção de erros

carga O QUE REALMENTE INTERESSA

HEPR Dinf RCC 680

Arquitetura II — redes de interconexão

Parâmetros de Desempenho



Equação do desempenho

latência total

= overhead_transmissão	proc 100% ocupado com envio
+ tempo_de_vôo	veloc da luz + atrasos hw
+ mensagem / largura_de_banda	[bits / bits/s]=[s]
+ overhead_recepção	proc 100% ocupado recebendo

 largura de banda medida no fio (inclui cabeçalhos) 	[bit/s]
ullet tempo de vôo é aprox 20cm/ns $+$ atrasos nos repetidores	[s]
• ovhead TX inclui tempo nos componentes de SW e HW	[s]
• ovhead RX geralmente maior que de TX (+interrupção)	[s]
• tamanho da mensagem inclui cabeçalhos	[bit]

HEPR DInf RCC 60

Arquitetura II — redes de interconexão

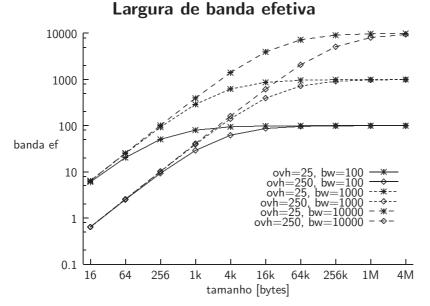
2007-1

Largura de banda efetiva

- em redes de pequeno diâmetro, tempo de vôo pode ser ignorado
- em rede de longa distância, overheads de TX e de RX podem ser ignorados
- latência total pode ser simplificada
 latência total ≈ overheads + |mensagem| / largura_de_banda
- largura de banda efetiva entregue pela rede é
 larg de banda efetiva = |mensagem| / latência_total

HEPR Dinf RCC 600

des de interconexão 2007-



moral da estória: overhead DEVE ser pequeno

2007-1

Escopo e Tipos de Redes

escopo	tipo	
interna ao Cl	crossbar, completamente conectada	
motherboard	barramento	
System Area Network	barramento, anel, malha, cubo	SAN
Local Area Network	barramento, anel	LAN
Wide Area Network	malha esparsa, ponto-a-ponto	WAN

TIEDR Dief RCC 60.

Meios Físicos

- Flat cable
- par trançado

Arquitetura II — redes de interconexão

- cabo coaxial
- fibra ótica
- ⊳ mono-modo

- ondas de rádio

HEPR DIst RCC 601

Arquitetura II — redes de interconexão 2007-1

Espaço de Projeto

• Critérios de desempenho

- ⊳ latência tempo de viagem de uma mensagem
- ⊳ largura de banda quanto tráfego a rede pode tratar
- ⊳ custo rede representa qual fração do custo de HW
- ⊳ funcionalidade combinação de msgns, tolerância a falhas

Espaço de Projeto (cont)

• Parâmetros de projeto

⊳ topologia – estática × dinâmica

com sub-classes

- ⊳ modo de operação síncrono × assíncrono
- ▷ estratégia de comutação circuitos, re-despacho, cut-through
- ▷ estratégia de controle centralizado × distribuído
- produto cartesiano destes 4 parâmetros define espaço de projeto

HEPR Dist RCC 60

Arquitetura II - redes de interconexão

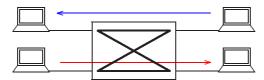
2007-1

Meio Compartilhado vs Comutado

Meio compartilhado (Ethernet em cabo coaxial)



Meio comutado (Ethernet comutada)



HEPR DIst RCC 603

Arquitetura II — redes de interconexão

2007-1

Topologias – Redes com Conexão Estática

- Redes construídas com conexões diretas entre os nós da rede
- padrão de conexões fixado na construção da rede
- grau de um nó d: número de enlaces ligados ao nó
- diâmetro da rede D: caminho mais curto maximal entre \forall 2 nós
- largura de um canal w: número de bits no canal/enlace
- diâmetro da bisecção: menor diâmetro da rede com 2 metades =s channel bisection width b: mínimo número de canais na bisecção

wire bisection width B = bw: reflete densidade da fiação

Topologias (i) - Redes com Conexão Estática

- ullet Linear N nós conectados por N-1 enlaces não é um barramento, enlaces podem ser bi-direcionais D=N-1 (elevado para N grande)
- ullet anel N nós conectados por N enlaces enlaces uni-direcionais

$$d=2$$
 $D=N$

enlaces bi-direcionais

$$d=2$$
 $D=|N/2|$

- barrel shifter anel com enlaces adicionais para vizinhos a uma distância que é potência de 2: i
 ightharpoonup j se $|j-i|=2^r$ d = 2n - 1D=n/2
- completamente conectado cada nó é ligado a todos outros nós d=N-1D=1

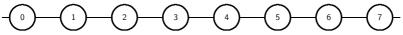
LIEPR Dinf BCC

Arquitetura II - redes de interconexão

2007-1

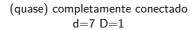
Topologias (ii) – Redes com Conexão Estática

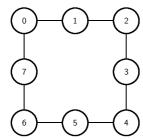
arranjo linear d=2 D=7

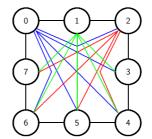


anel unidir d=2 D=8

anel bidir d=2 D=4







LIEPR Dinf BCC

2007-1

Arquitetura II — redes de interconexão

Topologias (iii) - Redes com Conexão Estática

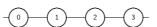
- ullet Árvore binária $N=2^k-1$ nós em k níveis $d\leq 3$ D=2(k-1)
- $d = \{1, N-1\}$ D=2• estrela – árvore de dois níveis
- árvore gorda diâmetro dos canais aumenta próximo da raiz

árvore binária estrela árvore gorda $d \le 3 D=2(k-1)$ d=1,N-1 D=2 $d \le 3 D = 2(k-1)$

Topologias (iv) - Redes com Conexão Estática

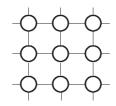
ullet Arranjo Linear é uma malha k=1 (1D) com $N{=}n$ nós

$$d=2$$
 $D=N-1$



ullet Malha – rede k-dimensional com $N{=}n^k$ nós

$$d{=}2k$$
 (nós internos) $D{=}k(n{-}1)$



ullet toróide – malha n imes n com conexões nas extremidades D = 2|n/2|d=4

HEPR Dinf RCC

Arquitetura II — redes de interconexão

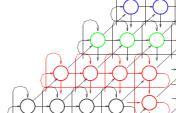
2007-1

Topologias (v) - Redes com Conexão Estática

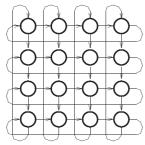
Anel (1D, k=4 n=1)



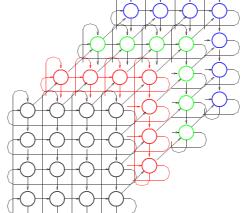
Cubo (toróide 3D, k=4 n=3)



Malha (2D, k=4 n=2)



Arquitetura II — redes de interconexão

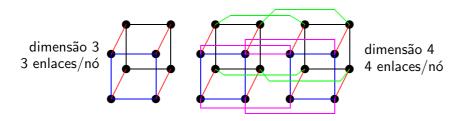


LIEPR Dinf BCC

Topologias (vi) - Redes com Conexão Estática

ullet Cubos: n-cubo k-ário – interliga N nós, via enlaces com k nós/enlace, em cubo com dimensão n

$$N = k^n \;, \quad k = \sqrt[n]{N} \;, \quad n = \log_k N$$



quando k=2 rede é hipercubo

Conexão Estática - Custo vs Desempenho

rede	$oldsymbol{d}$ grau	$oldsymbol{D}$ diâm	#enl	bissec	#nós
arranj linear	2	N-1	N-1	1	N
anel (1D)	2	$\lfloor N/2 floor$	Ν	2	Ν
totmte conect	N-1	1	N(N-1)/2	$(N/2)^{2}$	N
árvore binária	3	2(h-1)	N-1	1	$h{=}\lceil lg_2 N \rceil$
malha 2D	4	2(r-1)	2N-2r	r	$r=\sqrt{N}$
toro 2D	4	$2\lfloor r/2 \rfloor$	2N	2r	$r=\sqrt{N}$
hipercubo	n	n	nN/2	N/2	$n=lg_2N$
n-cubo k-ário	2n	$n\lfloor k/2 \rfloor$	nΝ	$2k^{n-1}$	$N=k^n$

HEPR Disf RCC 700

Arquitetura II — redes de interconexão

2007-1

Conexão Estática - Custo vs Desempenho

Rede com 64 nós

topologia	bissecção	portas	enlaces
barramento	1	1	1
anel (1D)	2	2	128
toróide (2D)	16	4	192
hipercubo	32	6	256
totalmte conect	1024	255	2080

TIEPR Disf RCC 70

Arquitetura II — redes de interconexão

2007-1

Topologias - Redes com Conexão Dinâmica

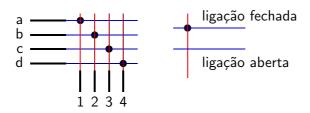
barramento

⊳ conexão entre mestre e escravo definida a cada transação

crossbar

 \triangleright é a rede mais cara \longrightarrow O(n²) comutadores

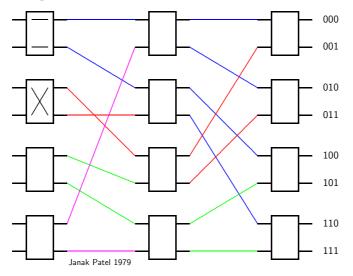
ho é a rede mais eficaz $i,j \in \{0..n\}, i
eq j \Rightarrow P_i
ightleftharpoons P_j$



• redes multi-estágio (ômega, delta, ...)

 \triangleright custo menor que crossbar \longrightarrow O(n log n) comutadores

Topologias - R Conexão Dinâmica - rede Delta



HEPR Dist RCC 700

Arquitetura II — redes de interconexão

2007-1

Topologias (i) – Redes com Conexão Dinâmica

característica	barramento	multiestágio	crossbar
latência mínima	constante	$O(log_{m{k}}n)$	constante
banda / proc	O(w/n)O(w)	O(w)O(nw)	O(w)O(nw)
complex cabeam	O(w)	$O(nw \log_k n)$	$O(n^2w)$
complex comut	O(n)	$O(n \log_k n)$	$O(n^2)$
conectividade	$1 \longrightarrow 1$	†	1

barramento: n proc, largura w

multi-est: nxn MIN, kxk comutadores de largura w

crossbar: nxn largura w

† algumas permutações e broadcast se rede livre

‡ todas permutações, uma de cada vez

HEPR Dist RCC 71

Arquitetura II — redes de interconexão

2007-1

Comutação

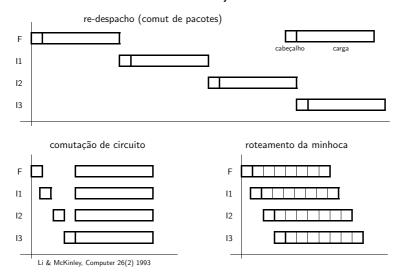
• Comutação

ho re-despacho (store-and-forward) $L_{
m rd} \propto$ núm comutadores \cdot $|{\sf msgm}|$

- imes cut-through se cabeçalho bloqueado, comutador engole msgm $L_{
 m ct} \propto f({
 m n\'um \ comutadores}) + |{
 m msgm}|/{
 m banda} \cdot {
 m tempoTX}$
- imes wormhole se cabeçalho bloqueado, msgm bloqueia caminho $L_{
 m wh}\cong L_{
 m ct}$ se rede sem tráfego

Li & McKinley, A survey of wormhole routing techniques in direct networks, Computer 26(2),1993

Comutação



HEPRING RCC 71

Arquitetura II — redes de interconexão

2007-1

Roteamento

Roteamento

- ⊳ source-based emissor define rota até destino
 - rota é parte da msgm
- ▷ circuito virtual negocia-se caminho antes de enviar nome do circuito é parte da msgm
- - * determinístico
 - * adaptativo
 - * aleatório

HEPR Dist RCC 71

Arquitetura II — redes de interconexão

2007-1

Estratégia de Controle

• Controle Centralizado

topologia	P(bloqueios)	# comutadores	distância
crossbar	$\propto 1/(N-2)$	N^2	$\propto 1$
ômega	>crossbar	$N/2\log N$	$\propto \log N$
árvore bin	?	N-1	2(h-1)
barramento	$\propto N$	1	1?

árvore – altura $h = \lceil \log N \rceil$

• Controle Distribuído

- ⊳ anel, malha, toróide
- ⊳ hipercubo

resumo - Redes

- Protocolo endereços e rotas, confiabilidade, proteção
- latência total
 - = overhead_transmissão proc 100% ocupado com envio
 - + tempo_de_vôo veloc da luz + atrasos hw
 - + |mensagem| / largura_de_banda [bits / bits/s]=[s]
 - + overhead_recepção proc 100% ocupado recebendo
- Parâmetros de projeto
- * topologia ligação estática × dinâmica
- * modo de operação síncrono × assíncrono
- * estratégia de comutação circuitos, re-despacho, cut-through
- * estratégia de controle centralizado × distribuído
- topol: ligações estáticas cubo k-ário n-dim, completam conectado
- topol: ligações dinâmicas barramento, multi-estágio, crossbar

HEPR Dinf RCC 71

Arquitetura II — clusters 2007-1

revisão - Redes

- Protocolo endereços e rotas, confiabilidade, proteção
- latência total
 - = overhead_transmissão proc 100% ocupado com envio
 - + tempo_de_vôo veloc da luz + atrasos hw

 - + overhead_recepção proc 100% ocupado recebendo
- Parâmetros de projeto
- * topologia ligação estática × dinâmica
- * modo de operação síncrono × assíncrono
- * estratégia de comutação circuitos, re-despacho, cut-through
- * estratégia de controle centralizado × distribuído
- topol: ligações estáticas cubo k-ário n-dim, completam conectado
- topol: ligações dinâmicas barramento, multi-estágio, crossbar

HEPRING RCC 716

Arquitetura II — clusters 2007-1

Ficou faltando...

- Onde é a interface com hierarquia de memória?
- $CPU \rightarrow multiprocessador \times rede \rightarrow multicomputador$
- há um processador adicional dedicado à tratar comunicação?
- * como particiona as tarefas?
- * qual o custo da comunicação entre CPU_proc e CPU_com?
- o que fazer em HW e o que fazer em SW?
- * roteamento
- * checksum
- * ...
- proteção provida por SO x usuário

Aglomerados de PCs

- Cluster construído com PCs + rede local + GNU/Linux
- em (algumas) aplicações com granularidade grande:
- * preço/desempenho entre $3..50 \times$ o de sistemas dedicados (MPPs)
- * desempenho entre 1/3..1/2 do de sistemas dedicados (MPPs)
- componentes de um nó = componentes de um PC
 disco pode ser problemático xterms no DInf sem disco!!!
- já usamos labs do DInf como clusters
- referência: Cluster Computing White Paper

TIEPR Dief RCC 718

Arquitetura II — clusters 2007-1

HW de rede

- problemas com padrão de comunicação complexos só atingem bom desempenho em MPPs
- problemas trivialmente paralelos n\u00e3o chegam a usar toda a banda dispon\u00edvel (Ethernet)
- maioria dos problemas entre estes dois extremos

HEPR Dinf RCC 71

Arquitetura II — clusters 2007-

Exemplos

Beowulf 2⁺ Ethernet's, sw roteia tráfego entre as 2⁺ interfaces físicas transparentemente à aplicação banda 'suficiente', latência $\approx 100 \mu s$

Myrinet rede comutada (\longrightarrow multi-estágio) banda 0,5–2Gbps, latência \approx 2,6–3,2 μ s (MPI)

Infiniband barramento serial unidirecional, rede comutada, RDMA banda 2×2 Gbps (2,4,8), pode agregar 4,12 canais (2–96Gbps) latência 160ns/comutador, ponta-a-ponta $\approx 6\mu$ s

VIA (extinto?) comunicação entre espaços de endereçamento de usuário, com pouca intervenção do SO banda ≥ 1 Gbps, latência $\approx 10 \mu$ s, Zero-copy

Software

- ferramentas de programação
- * linguagens C, C++, fortran, ada ...
- * bibliotecas PVM (1989), MPI transformam PCs em clusters
- * depuradores depuração de código e de desempenho
- gerência de recursos
- * instalação, administração
- * escalonamento de hw e sw
- * armazenamento
- * diagnóstico, monitoramento, disponibilidade
- funções do SO
- * multiplexar processos nos recursos físicos

resource mgmt

* prover abstrações úteis

beautification

TIEPR DIAF RCC 7

Arquitetura II — clusters 2007

Ambientes e modelos de programação

Ambientes/modelos implementados em bibliotecas e middleware

processos seqüenciais comunicantes (troca de mensagens)
 Calculus of Communicating Systems, Robin Milner
 Comunicating Sequential Processes, C A R Hoare

Message Passing Interface ou MPI

padrão de fato

• Parallel Virtual Machine ou PVM

já foi o padrão

data-parallel (SIMD)

difícil de implementar c/ eficiência

message-driven

fast_messages + splitC

memória compartilhada distribuída

Shrimp

• e muitos outros...

HEPR Dinf RCC 77

Arquitetura II — clusters 2007-

Características de SOs para aglomerados

Facilidade de gerenciamento

single system image - interface comum esconde diversidade

- estabilidade
- * robusteza contra processos errados/maliciosos
- * recuperação de falhas por reconfiguração dinâmica
- * usabilidade sob carga
- desempenho, especialmente das partes críticas:
 escalonamento, gerenciamento de memória, arquivos e E/S,
 rede
- extensibilidade dos nós, SO aberto
- escalabilidade intra-nós: depende dos nós individuais
 inter-nós: é limitada pela rede
- heterogeneidade mesmo SO em ≠s processadores

HEPR DIst RCC 723

mesmas APIs em ≠s computadores

Funcionalidades: SO x middleware

- SOs são complexos → alterações afetam partes indesejadas
- funções no *middleware* podem ser ex-portadas para outros SOs exemplo: memória compartilhada distribuída (DSM) em mw

TIEPR Dist RCC 79.

Arquitetura II — clusters

2007-1

2007-1

Single System Image

- Extensões do SO que permitem compartilhamento transparente de recursos
- 2 variantes:
- * administração do sistema (mw?)
- * escalonamento de tarefas + uso transparente de recursos remotos
- SSI implementado no SO implica em compartilhamento de estado entre TODOS os nós de um cluster
- → manter informação atualizada é caríssimo!!

HEPR Dist RCC 75

Arquitetura II — clusters

Distributed Shared Memory

- Permite compartilhamento de memória (RAM) entre todos os nós
- paradigma de programação paralela para usar mem compartilhada entre processos de uma aplicação
- geralmente implementado em bibliotecas porque o uso normal é programar com consistência seqüencial forte e isso resulta em desempenho ruim
- programação com consistência relaxada é aceitável em aplicações paralelas mas complexo demais para uso normal

Comunicação no nível de usuário

- \bullet Aglomerados de alto desempenho ${\rm N\tilde{A}O}$ usam Ethernet mas Myrinet, Infiniband, ou ...
- para uso ser vantajoso, adaptador de rede deve ser acessado com pouco *overhead*
- todas (?) as ações de comunicação devem ocorrer em espaço/modo
 - de usuário, tipicamente com DMA em modo usuário
- problema: multiplexar tráfego de vários processos na mesma interface física de forma segura
- gerenciamento de filas de TX e RX pelo usuário
 SO estabelece conexão e pega/paga interrupções circuito virtual?
- * menos trocas de contexto e menos cópias
- * Remote DMA (RDMA): acesso via DMA à memória de outros nós do aglomerado → baixa latência+overhead

HEPR Dist RCC 79

Arquitetura II — clusters 2007-1

E/S paralela

- Problemática na carga e salvamento de arquivos de dados
- usualmente, alguns nós concentram $E/S \rightarrow$ tornam-se gargalos
- é saudável ter algum espaço local a cada nó

HEPR Dist RCC 75