

# 高等计算机体系结构 第十五讲: 习题课

栾钟治  
北京航空航天大学 计算机学院 中德联合软件研究所  
2020-06-22

1

## 期末考试

- 6月26日上午9:50-12:20
  - 试题提前10分钟通过课程中心作业区发布
  - 自备白色答题纸答题, 写清题号不用抄题
  - 答题完毕后拍照上传课程中心作业区, 最迟提交时间为12:30
  - 有任何提交问题请及时沟通联系
  - 同时开通qq群课堂, 所有同学进入课堂, 不用开视频, 交卷后方可退出群课堂

2

2

## 作业1——指令集体系结构 (ISA) 和折衷 参考答案

3

## 1 指令集体系结构 (ISA)

比较五种不同风格的指令集代码序列的内存效率。不同的体系结构类型有:

1. 零地址的机器是一种基于栈的机器, 它的所有操作都通过存储在操作数栈上的值进行。允许以下操作: •PUSH M •POP M •OP

2. 单地址的机器使用一个累加器来执行计算。允许以下操作: •LOAD M •STORE M •OP M

3. 双地址的机器有两个操作数来源, 对这两个来源的操作数执行操作并将结果存回其中一个源。允许以下操作: •OP M1, M2

4. 三地址的机器, 通常有两个操作数来源, 执行操作后的结果存回不同于两个操作数来源的第三个目的地址。

(a) 对于一台操作数和结果目的地址都是内存地址的三地址机器, 允许如下操作: • OP M3, M1, M2

(b) 对于一台源和目的都是寄存器的三地址机器, 使用内存操作将值载入寄存器 (MIPS就是这种机器的例子)。允许如下操作: •OP R3, R1, R2 •LD R1, M •ST R2, M

对以上5种不同类型的指令集, 假设  
每条指令都有一个操作码和一组操作数,

•所有的操作码均为 1 字节 (8 bits)

•所有的寄存器操作数均为 1 字节 (8 bits)

•所有的内存地址均为 2 字节 (16 bits)

•所有的数据操作数均为 4 字节 (32 bits)

•所有指令的长度均为字节的整数倍

内存带宽消耗 = 传输的代码量 + 传输的数据量

传输的数据量 = 涉及的数据数量 × 4 Bytes

内存访问没有其他优化, 变量初始值都在内存中

(a) 将下边的高级语言片段翻译成前述5种结构适用的代码序列

A = B + C;

B = A + C;

D = A - B;

(b) 请计算这5种结构对应的指令序列在执行时的取指令字节数和内存数据访问 (读和写) 字节数。

(c) 从代码尺寸的角度哪一种结构最高效?

(d) 从内存总带宽的需求 (代码+数据) 角度哪一种结构最高效?

4

### 参考答案：

代码的大小：每条指令都有一个操作码和一组操作数，

- 所有的操作码均为 1 字节 (8 bits)
- 所有的寄存器操作数均为 1 字节(8 bits)
- 所有的内存地址均为 2 字节(16 bits)
- 所有的数据操作数均为 4 字节(32 bits)
- 所有指令的长度均为字节的整数倍

内存带宽：

内存带宽消耗=传输的代码量（代码大小）+ 传输的数据量

传输的数据量=涉及的数据数量x 4 Bytes

代码片段：

A = B + C;

B = A + C;

D = A - B;

以下我们用I-Bytes表示传输的代码量，用D-Bytes表示传输的数据量

(a) 和	指令集体系结构	操作码	操作数	I-Bytes	D-Bytes	总字节数
	零地址	PUSH	B	3	4	
		PUSH	C	3	4	
(b)	基于栈的机器，它的所有操作都通过存储在操作数栈上的值进行。	ADD	1	1		
	• PUSH M – 将位于内存地址 M 处的值压入操作数栈	POP	A	3	4	
	• POP M – 弹出操作数栈并	PUSH	A	3	4	
	将值存入内存地址 M 处	PUSH	C	3	4	
	• OP – 从操作数栈中弹出两个值，对这两个值执行二进制操作 OP，结果压回到操作数栈	ADD	1	1		
		POP	B	3	4	
		PUSH	A	3	4	
		PUSH	B	3	4	
		SUB	1	1		
		POP	D	3	4	
				30	36	66

5

### 参考答案（续）：

代码的大小：每条指令都有一个操作码和一组操作数，

- 所有的操作码均为 1 字节 (8 bits)
- 所有的寄存器操作数均为 1 字节(8 bits)
- 所有的内存地址均为 2 字节(16 bits)
- 所有的数据操作数均为 4 字节(32 bits)
- 所有指令的长度均为字节的整数倍

内存带宽：

内存带宽消耗=传输的代码量（代码大小）+ 传输的数据量

传输的数据量=涉及的数据数量x 4 Bytes

代码片段：

A = B + C;

B = A + C;

D = A - B;

以下我们用I-Bytes表示传输的代码量，用D-Bytes表示传输的数据量

(a) 和	指令集体系结构	操作码	操作数	I-Bytes	D-Bytes	总字节数
	单地址	LOAD	B	3	4	
		ADD	C	3	4	
(b)	使用一个累加器来执行计算。	STORE	A	3	4	
	• LOAD M – 将存储在内存地址为 M 处的值载入累加器	ADD	C	3	4	
	• STORE M – 将累加器中的值存入内存地址为 M 处	STORE	B	3	4	
	• OP M – 对内存地址为 M 处存储的值和当前在累加器中的值执行二进制操作 OP，结果存入累加器	LOAD	A	3	4	
		SUB	B	3	4	
		STORE	D	3	4	
				24	32	56

6

### 参考答案（续）：

代码的大小：每条指令都有一个操作码和一组操作数，

- 所有的操作码均为 1 字节 (8 bits)
- 所有的寄存器操作数均为 1 字节(8 bits)
- 所有的内存地址均为 2 字节(16 bits)
- 所有的数据操作数均为 4 字节(32 bits)
- 所有指令的长度均为字节的整数倍

内存带宽：

内存带宽消耗=传输的代码量（代码大小）+ 传输的数据量

传输的数据量=涉及的数据数量x 4 Bytes

代码片段：

A = B + C;

B = A + C;

D = A - B;

以下我们用I-Bytes表示传输的代码量，用D-Bytes表示传输的数据量

(a) 和	指令集体系结构	操作码	操作数	I-Bytes	D-Bytes	总字节数
	双地址	SUB	A, A	5	12	
		ADD	A, B	5	12	
(b)	有两个操作数来源，对这两个来源的操作数执行操作并将结果存回其中一个源。	ADD	A, C	5	12	
	• OP M1, M2 – 对存储在内存地址为 M1 和 M2 的值进行二进制操作 OP，将结果存回内存地址 M1 处	SUB	B, B	5	12	
		ADD	B, A	5	12	
		ADD	B, C	5	12	
		SUB	D, D	5	12	
		ADD	D, A	5	12	
		SUB	D, B	5	12	
				45	108	153

7

### 参考答案（续）：

代码的大小：每条指令都有一个操作码和一组操作数，

- 所有的操作码均为 1 字节 (8 bits)
- 所有的寄存器操作数均为 1 字节(8 bits)
- 所有的内存地址均为 2 字节(16 bits)
- 所有的数据操作数均为 4 字节(32 bits)
- 所有指令的长度均为字节的整数倍

内存带宽：

内存带宽消耗=传输的代码量（代码大小）+ 传输的数据量

传输的数据量=涉及的数据数量x 4 Bytes

代码片段：

A = B + C;

B = A + C;

D = A - B;

以下我们用I-Bytes表示传输的代码量，用D-Bytes表示传输的数据量

(a) 和	指令集体系结构	操作码	操作数	I-Bytes	D-Bytes	总字节数
	三地址 Memory-Memory	ADD	A, B, C	7	12	
		ADD	B, A, C	7	12	
(b)	一台操作数和结果目的地址都是内存地址的三地址机器	SUB	D, A, B	7	12	
	• OP M3, M1, M2 – 对存储在内存地址为 M1 和 M2 的值执行二进制操作 OP，结果存回内存地址为 M3 处					
				21	36	57

8

### 参考答案（续）：

代码的大小：每条指令都有一个操作码和一组操作数，

- 所有的操作码均为 1 字节 (8 bits)
- 所有的寄存器操作数均为 1 字节(8 bits)
- 所有的内存地址均为 2 字节(16 bits)
- 所有的数据操作数均为 4 字节(32 bits)
- 所有指令的长度均为字节的整数倍

内存带宽：

内存带宽消耗=传输的代码量（代码大小）+ 传输的数据量

传输的数据量=涉及的数据数量x 4 Bytes

代码片段：

A = B + C;

B = A + C;

D = A - B;

以下我们用I-Bytes表示传输的代码量，用D-Bytes表示传输的数据量

指令集体系结构	操作码	操作数	I-Bytes	D-Bytes	总字节数
(a) 三地址 Load-Store	LD	R1, B	4	4	
(b) 二来源和目的都是寄存器的三地址机器，使用内存操作将值载入寄存器 (MIPS 就是这种机器的例子)。	LD	R2, C	4	4	
	ADD	R1, R1, R2	4		
	ST	R1, A	4	4	
	ADD	R3, R1, R2	4		
	ST	R3, B	4	4	
	SUB	R3, R1, R3	4		
	ST	R3, D	4	4	
• OP R3, R1, R2 - 对寄存器 R1 和 R2 中的值执行 OP 操作，将结果存回寄存器 R3					
• LD R1, M - 将内存地址为 M 处的值取出存入寄存器 R1					
• ST R2, M - 将寄存器 R2 中的值存入内存地址 M 处					
			32	20	52

9

### 参考答案（续）：

代码的大小：每条指令都有一个操作码和一组操作数，

- 所有的操作码均为 1 字节 (8 bits)
- 所有的寄存器操作数均为 1 字节(8 bits)
- 所有的内存地址均为 2 字节(16 bits)
- 所有的数据操作数均为 4 字节(32 bits)
- 所有指令的长度均为字节的整数倍

内存带宽：

内存带宽消耗=传输的代码量（代码大小）+ 传输的数据量

传输的数据量=涉及的数据数量x 4 Bytes

代码片段：

A = B + C;

B = A + C;

D = A - B;

(c) 从代码尺寸角度看，三地址Memory-Memory的机器最高效，21Bytes的代码

(d) 从内存总带宽的消耗角度，三地址Load-Store的机器最高效，52Bytes的总带宽消耗

10

## 2 性能指标

请简要回答以下问题。

- 如果在具有更高主频的处理器上运行给定程序，是否意味着相比主频较低的处理器而言它总是能够在单位时间（比如1秒钟）内执行更多的指令？

参考答案：否。主频较低的处理器可能具有更高的 IPC (每周指令数)

注：一个主频较低的处理器可能能够在周期中执行多条指令，而一个主频较高的处理器可能在一个周期中只能执行一条指令。

- 如果一个处理器执行给定程序时每秒钟能够执行更多的指令，是否意味着相比每秒执行指令数较少的处理器而言它总是能够更快地执行完这个程序。

参考答案：否。因为每秒执行指令数较多的处理器可能需要执行更多的指令

注：运行完一个程序需要执行的指令总数，不同的处理器可能会不同。

11

## 3 性能评价

评价两个实现不同指令集体系结构的处理器的潜在性能。

评价是基于运行特定基准程序（benchmark）时的性能而做出的。在实现ISA A的处理器上，最优的编译代码执行benchmark的性能是10 IPC，这一款处理器主频是500MHz。在实现ISA B的处理器上，同样最优的编译代码的执行性能是2 IPC，处理器主频是 600MHz。

- 请问实现ISA A的处理器每秒可以执行多少百万条指令(MIPS)?

参考答案：ISA A:  $10 \frac{\text{instructions}}{\text{cycle}} \times 500,000,000 \frac{\text{cycle}}{\text{second}} = 5000 \text{ MIPS}$

- 请问实现ISA B的处理器每秒可以执行多少百万条指令(MIPS)?

参考答案：ISA B:  $2 \frac{\text{instructions}}{\text{cycle}} \times 600,000,000 \frac{\text{cycle}}{\text{second}} = 1200 \text{ MIPS}$

- 哪一个更高性能的处理器： A? B? 不知道?

参考答案：不知道。对每个处理器来说，最佳的编译代码所具有的指令数量可能是不一样的。

12

## 4 固定长度和可变长度ISA

### 考虑以下两种Load/Store结构的ISA

1. 第一种是固定长度的ISA，它使用如下的指令编码。

操作码	操作数 1 (目的)	操作数 2 (源 1) Reg/Imm	操作数 3 (源 2) Reg/Imm
-----	------------	---------------------	---------------------

其中：操作码 1 byte，每个操作数均为1 byte；所有寄存器/寄存器以及寄存器/立即数的操作需要1个时钟周期，所有Load和Store操作需要4个时钟周期

2. 第二种是可变长度的ISA，它使用如下的指令编码。

操作码	操作数 1 (目的)	操作数 2 (源 1) Reg/Imm	操作数 3 (源 2) Reg/Imm (可选)
-----	------------	---------------------	--------------------------

其中：操作码 1 byte，每个操作数均为1 byte，**操作数3是可选的，由操作码隐式说明**。如果指令不需要第3个操作数，这个字段就不会使用。

可变长使得第二种ISA的译码变得复杂，所以，它的所有指令执行时间都比第一种固定长度ISA的指令多2个时钟周期。也就是说，所有寄存器/寄存器以及寄存器/立即数的操作需要3个时钟周期，所有Load和Store操作需要6个时钟周期。

操作码	操作数 1 (目的)	操作数 2 (源 1) Reg/Imm	操作数 3 (源 2) Reg/Imm
操作码	操作数 1 (目的)	操作数 2 (源 1) Reg/Imm	操作数 3 (源 2) Reg/Imm (可选)

(a) 对以上两种ISA，这段汇编代码的 ADD r3, r1, r2 // r3 = r1+r2  
尺寸（字节数）分别是多少？ SLL r3, 0x2 // r3 = r3 << 2  
MOV r5, 0xa // r5 = 0x0a  
STW r3, (r5) // MEMORY[r5] = r3

参考答案：

固定长度ISA：4x4 = 16 字节

可变长度ISA：1x4 (ADD指令) + 3x3 (其它指令) = 13 字节

(b) 对以上两种ISA，执行这一段代码序列分别需要多少时钟周期？

参考答案：

固定长度ISA：3x1 (其它指令) + 1x4 (STW 指令) = 7 时钟周期

可变长度ISA：3x3 (其它指令) + 1x6 (STW 指令) = 15 时钟周期

操作码	操作数 1 (目的)	操作数 2 (源 1) Reg/Imm	操作数 3 (源 2) Reg/Imm
操作码	操作数 1 (目的)	操作数 2 (源 1) Reg/Imm	操作数 3 (源 2) Reg/Imm (可选)

(c) 哪种ISA的代码尺寸更小？为什么？

参考答案：

可变长度ISA的代码尺寸更小，因为它使用更少的字节对给定代码段的指令进行编码。

(d) 哪种ISA的执行时间更短？为什么？

参考答案：

固定长度ISA执行时间更短，因为译码复杂性低，所以译码指令所需的时钟周期更少。

## 5 可寻址性

假如我们有64MB的内存，请计算要获得以下寻址能力所需地址的长度：

(i) 位寻址 ISA

参考答案：29 bits

(ii) 字节寻址 ISA

参考答案：26 bits

(iii) 8字节寻址 ISA

参考答案：23 bits

(iv) 32字节寻址 ISA

参考答案：21 bits

## 6 微体系结构与ISA

(a) 简要叙述微体系结构和ISA之间的区别。编译器需要知道微体系结构的什么信息才能正确的编译程序？

参考答案：

ISA层是一台计算机向软件暴露的接口，而微体系结构则是这台计算机实际的底层实现。因此，微体系结构本身及其各种改变对编译器和程序员是透明的（除了性能方面的影响），而ISA及其改变将影响编译器和程序员。  
编译器不需要知道微体系结构的任何信息就能够正确的编译程序。

(b) 判别一台机器的以下性质是微体系结构的属性还是ISA的属性：

- (i) 这台机器没有减法指令.
- (ii) 这台机器的ALU没有减法单元.
- (iii) 这台机器没有状态码.
- (iv) 在加法指令中可以指明一个5位的立即数.
- (v) 执行一条加法指令需要n个时钟周期.
- (vi) 有8个通用寄存器.
- (vii) ALU的一个输入需要一个2选1多路选择器.
- (viii) 寄存器堆有1个输入端口和2个输出端口.

参考答案：

- (i) ISA
- (ii) 微体系结构
- (iii) ISA
- (iv) ISA
- (v) 微体系结构
- (vi) ISA
- (vii) 微体系结构
- (viii) 微体系结构

## 作业2——单周期vs. 多周期微体系结构 参考答案

## 1 MIPS单周期微体系结构分析

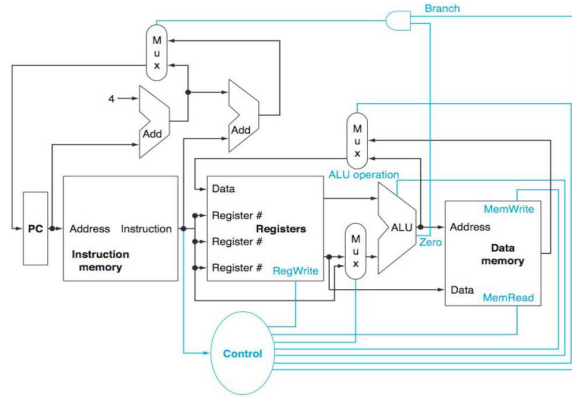
### 1.1

- ❖ 不同单元有不同的延迟时间。在图1中有七种主要单元。
- ❖ 对一条指令而言，关键路径(产生最长延迟的那条路径)上各个单元的延迟时间决定了该指令的最小延迟。
- ❖ 假设个单元的延迟时间如下表所示，回答下列3个问题。

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	控制
400ps	100ps	30ps	120ps	200ps	350ps	100ps

- (a) 对一条MIPS的与指令(AND)而言，关键路径是什么？
- (b) 对一条MIPS的装载指令(LW)而言，关键路径是什么？
- (c) 对一条MIPS的相等则分支指令(BEQ)而言，关键路径是什么？

❖ 图1



## 参考答案

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	控制
400ps	100ps	30ps	120ps	200ps	350ps	100ps

(a) 对一条MIPS的与指令(AND)而言, 关键路径是什么?

- 关键路径为: I-Mem、Regs(Read)、Mux、ALU、Mux、Regs(Write)
- 解析: 对于AND指令(and rd, rs, rt), 存在这样一条长路径: 读指令、读寄存器堆、通过ALUMux多路器、进行ALU运算、通过RegMux多路器、写寄存器堆(即I-Mem、Regs(Read)、Mux、ALU、Mux、Regs(Write))。另外一条长路径与之类似, 但是这条路径是在寄存器堆进行读操作时通过控制器的, 即: I-Mem、Control、Mux、ALU、Mux、Regs(Write), 但由于控制器的速度快于寄存器堆, 因而前者为关键路径, 其它的路径都短于这两条路径。

## 参考答案 (续)

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	控制
400ps	100ps	30ps	120ps	200ps	350ps	100ps

(b) 对一条MIPS的装载指令(LW)而言, 关键路径是什么?

- 关键路径为: I-Mem、Regs(Read)、Mux、ALU、D-Mem(Read)、Mux、Regs(Write)
- 解析: 对于LW指令, 存在这样一条长路径: 读指令、读寄存器堆获得基址、使用多路器选择立即数作为ALU的输入、使用ALU计算地址、访问数据存储器、使用多路器选择存储器输出作为寄存器的数据输入、写寄存器堆, 故有路径I-Mem、Regs(Read)、Mux、ALU、D-Mem(Read)、Mux、Regs(Write)。还有一条与之类似的长路径, 但是这条路径是通过控制器而不是寄存器堆的(用于生成ALUMux的控制信号), 由于控制器的速度快于寄存器堆, 于是前者为关键路径, 除这两条之外的路径都是比较短的路径。

## 参考答案 (续)

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	控制
400ps	100ps	30ps	120ps	200ps	350ps	100ps

(c) 对一条MIPS的相等则分支指令(BEQ)而言, 关键路径是什么?

- 由于控制器的速度快于寄存器堆, 故关键路径为: I-Mem、Regs(Read)、Mux、ALU、Mux
- 解析: 这条指令有两种长路径——决定分支条件以及计算新PC值。对于决定分支条件的路径, 需要读指令、读寄存器堆或使用控制单元、使用ALUMux、使用ALU比较两个值、使用ALU的零输出端来控制选择新PC值的多路器。对于计算新PC值的路径, 其中一条是PC值加4(Add)、加偏移量offset(Add)、选择这个值作为新的PC值(Mux); 另一条是读指令(为了取得偏移量)、使用分支加法单元和相应的多路器。但是这两条计算PC值中的路径都比决定分支条件的路径要短, 这是因为从表中可以看到指令存储器的速度要慢于执行PC+4的加法器、ALU的速度要慢于分支加法器。

## 1.2

❖ 图1中基本的单周期MIPS实现仅能实现某些指令。

- ❖ 可以在这个指令集中加入新的指令，但决定是否加入取决于给处理器的数据通路和数据通路增加的复杂度。

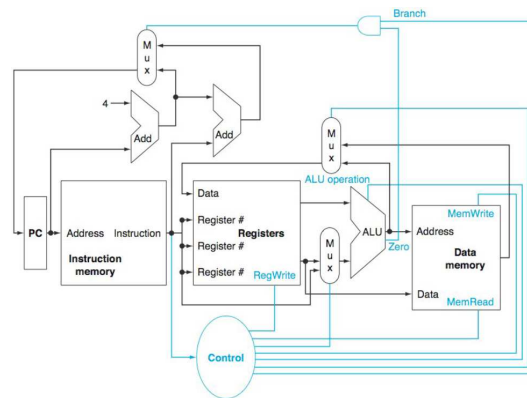
❖ 对于下表中的新指令而言，试回答下列3个问题。

指令	解释
add3 Rd,Rs,Rt,Rx	Reg[Rd]=Reg[Rs]+Reg[Rt]+Reg[Rx]

(a) 对上述指令而言, 哪些已有的单元还可以被使用?

(b) 对上述指令而言, 还需要增加哪些功能单元?

(c) 为了支持这些指令，需要在控制单元增加哪些信号？



❖ 当设计者考虑改进处理器数据通路时，往往要考虑性能与成本的折中。假设我们从图1的数据通路出发，其中指令存储器(Instruction Memory)、加法器(Add)、多路器(Mux)、ALU、寄存器堆(Registers)、数据寄存器(Data Memory)和控制单元(Control)的延迟分别为400ps、100ps、30ps、120ps、200ps、350ps和100ps，相应的成本分别为1000、30、10、100、200、2000和500。试根据表中的改进分别回答下列问题。

	改进	延迟	成本	优势
a.	更快的加法器	加法单元-20ps	每个加法单元+20	把已有的加法器用更快的加法器替代
b.	更大的寄存器堆	寄存器堆+100ps	寄存器堆+200	需要更少的load和store指令。 这将导致指令数减少5%

(d) 改进前后的时钟周期分别是多少？

(e) 改进后将获得多大的加速比？

(f) 比较改进前后的性能/价格比，进行这样的改进是否有意义？

## 参考答案

指令	解释
add3 Rd,Rs,Rt,Rx	$\text{Reg}[\text{Rd}] = \text{Reg}[\text{Rs}] + \text{Reg}[\text{Rt}] + \text{Reg}[\text{Rx}]$

(a) 对上述指令而言, 哪些已有的单元还可以被使用?

- 除分支加法器、数据存储器之外的所有单元

### 参考答案（续）

指令	解释
add3 Rd,Rs,Rt,Rx	$Reg[Rd] = Reg[Rs] + Reg[Rt] + Reg[Rx]$

(b) 对上述指令而言，还需要增加哪些功能单元？

- 寄存器堆增加一个输出端 和一个相应的读地址输入端，以读出  $Rx$ ；增加一个加法器(或为现有的ALU增加一个输入端)，加法器的一个输入端连接至寄存器堆新增的输出端，另一个输入端连接至ALU的输出端。如果采用增加一个加法器的方案，则还需增加寄存器堆数据输入选通器的一个输入端，并连接至加法器的输出端

### 参考答案（续）

指令	解释
add3 Rd,Rs,Rt,Rx	$Reg[Rd] = Reg[Rs] + Reg[Rt] + Reg[Rx]$

(c) 为了支持这些指令，需要在控制单元增加哪些信号？

- 如果增加一个加法器，则需增加寄存器堆数据输入选通器的控制信号，以实现数据通道3选1；如果为现有的ALU增加一个输入端，则需增加针对三输入端的ALU的功能控制信号定义，使其可控制新增的ADD3操作

### 参考答案（续）

	改进	延迟	成本	优势
a.	更快的加法器	加法单元-20ps	每个加法单元+20	把已有的加法器用更快的加法器替代
b.	更大的寄存器堆	寄存器堆+100ps	寄存器堆+200	需要更少的load和store指令。这将导致指令数减少5%

(d) 改进前后的时钟周期分别是多少？

- a. 由于加法单元不在关键路径上，因此对加法器的改进不影响时钟周期。
- b. 寄存器堆位于关键路径上，因而，使用更大的寄存器堆后，时钟周期变为  $1330ps + 2 \times 100ps = 1530ps$ 。
- 解析：时钟周期是由关键路径决定的，这里的关键路径为：I-Mem(读指令)、Regs(Read)(由于寄存器堆的延迟大于控制器，因而寄存器堆位于关键路径上)、Mux(选择ALU的输入)、ALU、Data Memory(Read)、Mux(选择存储器写入到寄存器堆中的数据)、Regs(Write)(数据写入寄存器堆)，该路径的延迟为  $400ps + 200ps + 30ps + 120ps + 350ps + 30ps + 200ps = 1330ps$ 。

### 参考答案（续）

	改进	延迟	成本	优势
a.	更快的加法器	加法单元-20ps	每个加法单元+20	把已有的加法器用更快的加法器替代
b.	更大的寄存器堆	寄存器堆+100ps	寄存器堆+200	需要更少的load和store指令。这将导致指令数减少5%

(e) 改进后将获得多大的加速比？

- a. 加速比由时钟周期本身的变化以及需要执行的时钟周期数目共同决定，对加法器的改进不影响时钟周期，并且，需要执行的时钟周期数目也不变，因此加速比为1.000
- b. 需要的指令数减少5%，需要的时钟周期数目也相应减少5%，同时，时钟周期由1330ps增加为1530ps，因而加速比为  $(1/0.95) \times (1330/1530) = 0.915$



### 参考答案（续）

	改进	延迟	成本	优势
a.	更快的加法器	加法单元-20ps	每个加法单元+20	把已有的加法器用更快的加法器替代
b.	更大的寄存器堆	寄存器堆+100ps	寄存器堆+200	需要更少的load和store指令。这将导致指令数减少5%

(f) 比较改进前后的性能/价格比，进行这样的改进是否有意义？

- a. 原来的处理器的总成本为  $1000(I-Mem) + 200(Regs) + 500(Control) + 100(ALU) + 2000(D-Mem) + 2 \times 30(2\text{个加法单元}) + 3 \times 10(3\text{个多选器}) = 3890$ ，更换加法器之后的总成本为  $3890 + 2 \times 20 = 3930$ ，相对成本为  $3930 / 3890 = 1.010$ ，性能/价格比为  $1.000 / 1.010 = 0.990$ ，成本增加但性能没有提升。
- b. 使用更大的寄存器堆的成本为  $3890 + 200 = 4090$ ，相对成本为  $4090 / 3890 = 1.051$ ，性能/价格比为  $0.915 / 1.051 = 0.871$ ，说明用更大的投入反而换来了性能的下降。

### 1.3

❖ 下表给出了实现处理器数据通路的逻辑单元延迟。试根据下表的两种情况分别回答下列问题。

指令存储器	加法器	多选器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
400ps	100ps	30ps	120ps	200ps	350ps	20ps	2ps

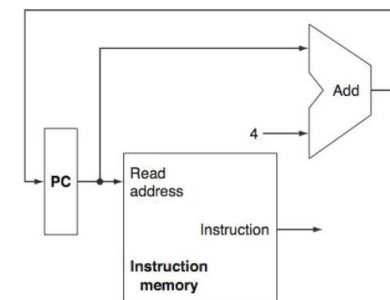
- (a) 如果处理器只需做连续取指这一件事(见图2)，那么时钟周期是多少？
  - (b) 考虑一个与图3类似的数据通路，但是假设处理器只需处理无条件相对跳转指令，那么时钟周期是多少？
  - (c) 同样考虑一个与图3类似的数据通路，但这次假设只需处理有条件相对跳转指令，那么时钟周期是多少？(请注意图3中ALU的零输出端不是与数据存储器连接，该输出与选择PC值来源的多选器的控制有关)
- 提示：图3中靠右侧的加法器延迟应当按照ALU来计算

❖ 根据下表的两种数据通路的逻辑单元，分别回答下列问题。

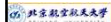
	单元
a.	执行加4的加法器(对PC)
b.	数据存储器

- (d) 哪些类型的指令需要该单元？
- (e) 对哪些类型的指令而言，该单元位于关键路径上？
- (f) 假设仅需支持beq指令和add指令，讨论该单元的延迟变化对处理器时钟周期的影响。假设其他单元的延迟不变。

❖ 图2



❖ 图3



37

### 参考答案

指令存储器	加法器	多选器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
400ps	100ps	30ps	120ps	200ps	350ps	20ps	2ps

(a) 如果处理器只需做连续取指这一件事(见图2), 那么时钟周期是多少?

- 由于指令存储器慢于加法器，因此，时钟周期决定于指令存储器的延迟，时钟周期为400ps。



38

### 参考答案 (续)

指令存储器	加法器	多选器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
400ps	100ps	30ps	120ps	200ps	350ps	20ps	2ps

(b) 考虑一个与图3类似的数据通路，但是假设处理器只需处理无条件相对跳转指令，那么时钟周期是多少？

- 关键路径为I-Mem、Sign-extend、Shift-left-2、Add(ALU)、Mux, 因此, 时钟周期为 $400\text{ps} + 20\text{ps} + 2\text{ps} + 120\text{ps} + 30\text{ps} = 572\text{ps}$ 。

39

### 参考答案 (续)

指令存储器	加法器	多选器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
400ps	100ps	30ps	120ps	200ps	350ps	20ps	2ps

(c) 同样考虑一个与图3类似的数据通路, 但这次假设只需处理有条件相对跳转指令, 那么时钟周期是多少? (请注意图3中ALU的零输出端不是与数据存储器连接, 该输出与选择PC值来源的多选器的控制有关)

- 根据题目中给出的延迟，后者为关键路径，因此时钟周期为  $400ps + 200ps + 30ps + 120ps + 30ps = 780ps$ 。
- 解析：对于有条件相对跳转指令，除存在长路径I-Mem、Sign-extend、Shift-left-2、Add(ALU)、Mux外，还存在长路径I-Mem、Registers(Read)、Mux、ALU、Mux，关键路径为这两条路径中较长的一个。

40

### 参考答案（续）

	单元
a.	执行加4的加法器(对PC)
b.	数据存储器

(d) 哪些类型的指令需要该单元？

- a. 所有指令。
- b. 与存取有关的指令，如：Lw, Sw等。

### 参考答案（续）

	单元
a.	执行加4的加法器(对PC)
b.	数据存储器

(e) 对哪些类型的指令而言，该单元位于关键路径上？

- a. 所有指令的关键路径都不会包含这个加法器，因为指令存储器的速度慢于加法器，而所有的指令都必须执行读指令这一操作。
- b. 与存取有关的指令，因为只有与存取有关的指令会用到该单元。

### 参考答案（续）

	单元
a.	执行加4的加法器(对PC)
b.	数据存储器

(f) 假设仅需支持beq指令和add指令，讨论该单元的延迟变化对处理器时钟周期的影响。假设其他单元的延迟不变。

- a. beq指令的关键路径为I-Mem、Regs(Read)、Mux、ALU、Mux，延迟为780ps，add指令的关键路径为I-Mem、Regs(Read)、Mux、ALU、Mux、Regs(Write)，延迟为980ps。这里只需讨论该单元延迟变化相比于较长的关键路径的影响，较长的关键路径为add指令，其延迟为980ps，而执行加4的加法器所在的路径为Add、Add(ALU)、Mux，其延迟为100ps + 120ps + 30ps = 250ps，若要该单元延迟变化而导致该路径成为关键路径，从而影响时钟周期，则执行加4的加法器延迟要大于980ps - 150ps = 830ps，才会影响时钟周期。
- b. 数据存储器未被beq或add指令使用，因此，它的延迟不影响时钟周期。

指令存储器	加法器	多选器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
400ps	100ps	30ps	120ps	200ps	350ps	20ps	2ps

### 1.4

❖ 本题讨论数据通路中不同的单元延迟对整个数据通路时钟周期的影响，以及指令如何利用不同的数据通路单元。根据下面的延迟情况，分别回答下列问题。

指令存储器	加法器	多选器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
500ps	150ps	100ps	180ps	220ps	1000ps	90ps	20ps

- (a) 如果仅需支持ALU类指令(如add、and等)，处理器的时钟周期是多少？
- (b) 如果仅需支持lw类指令，时钟周期是多少？
- (c) 如果必须支持add、beq、lw和sw指令，时钟周期是多少？

❖ 假设各类型指令所占比例如下表所示，试根据下表的两种情况分别回答下列问题。

add	addi	not	beq	lw	sw
30%	15%	5%	20%	20%	10%

- (d) 数据存储器平均用了多少时钟周期？
- (e) 符号扩展电路的输入平均用了多少时钟周期？在未用到该输入的其他时间，符号扩展电路在做什么？
- (f) 如果可以将数据通路上某个单元的延迟减少10%，应该减少哪个单元的延迟？改进后整个处理器的加速比是多少？

## 参考答案

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
500ps	150ps	100ps	180ps	220ps	1000ps	90ps	20ps

(a) 如果仅需支持ALU类指令(如add、and等)，处理器的时钟周期是多少？

- 时钟周期为 $500ps + 220ps + 100ps + 180ps + 100ps + 220ps = 1320ps$
- 解析：关键路径为I-Mem、Registers(Read)、Mux(选择ALU输入)、ALU、Mux(选择寄存器写入端)、Registers(Write)

## 参考答案（续）

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
500ps	150ps	100ps	180ps	220ps	1000ps	90ps	20ps

- (b) 如果仅需支持lw类指令，时钟周期是多少？
- 时钟周期为 $500ps + 220ps + 100ps + 180ps + 1000ps + 100ps + 220ps = 2320ps$
  - 解析：关键路径为I-Mem、Registers(Read)、Mux(选择ALU输入)、ALU、D-Mem(Read)、Mux(选择写入寄存器堆的数据)、Registers(Write)

## 参考答案（续）

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	符号扩展	左移两位
500ps	150ps	100ps	180ps	220ps	1000ps	90ps	20ps

(c) 如果必须支持add、beq、lw和sw指令，时钟周期是多少？

- 时钟周期为2320ps
- 解析：lw指令的关键路径最长，相比较而言，sw指令少使用了一个多路器并且不用向寄存器堆写数据，add和beq少使用了数据存储器

### 参考答案 (续)

add	addi	not	beq	lw	sw
30%	15%	5%	20%	20%	10%

(d) 数据存储器平均用了多少时钟周期?

- 平均有  $20\% + 10\% = 30\%$  的时钟周期里, 会用到数据存储器
- 解析: 只有lw和sw指令会用到数据存储器

### 参考答案 (续)

add	addi	not	beq	lw	sw
30%	15%	5%	20%	20%	10%

(e) 符号扩展电路的输入平均用了多少时钟周期? 在未用到该输入的其他时间, 符号扩展电路在做什么?

- 符号扩展电路实际上在每个周期都有计算结果, 但是它的输出在add和not指令中被忽略了, 符号扩展电路的输入只在addi指令(提供ALU需要的立即数)、beq指令(提供计算PC需要的偏移量)、lw指令和sw指令(提供寻址过程中需要的偏移量)中是需要的
- 结果为  $15\% + 20\% + 20\% + 10\% = 65\%$

### 参考答案 (续)

add	addi	not	beq	lw	sw
30%	15%	5%	20%	20%	10%

(f) 如果可以将数据通路上某个单元的延迟减少10%, 应该减少哪个单元的延迟? 改进后整个处理器的加速比是多少?

- lw指令有最长的关键路径为: I-Mem、Registers(Read)、Mux、D-Mem(Read)、ALU、Mux、Registers(Write), 决定着时钟周期的长度。
- 在路径中, 由于指令存储器的延迟值最大, 因此, 如将它的延迟从400ps减小到360ps(即减少10%), 则时钟周期由1330ps减小到1290ps, 加速比为  $1330/1290=1.031$ 。

### 1.5

❖ 本题讨论处理器时钟周期与控制单元设计之间的相互影响。根据下表的两种数据通路单元延迟情况分别回答下列问题。

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	符号扩展	左移两位	ALU控制
500ps	150ps	100ps	180ps	220ps	1000ps	90ps	20ps	55ps

- (a) 为了避免增加图4中数据通路的关键路径长度, 留给控制单元产生MemWrite信号的时间有多少?
- (b) 图4中哪个控制信号最**不**关键, 控制单元需要在多长时间内产生该信号以避免其成为关键路径?
- (c) 图4中哪个控制信号最**关**键, 控制单元需要在多长时间内产生该信号以避免其成为关键路径?

❖ 假设控制单元产生控制信号的时间如下表所示，试根据表中的两种情况回答下列问题(各部件的延迟与前面相同)。

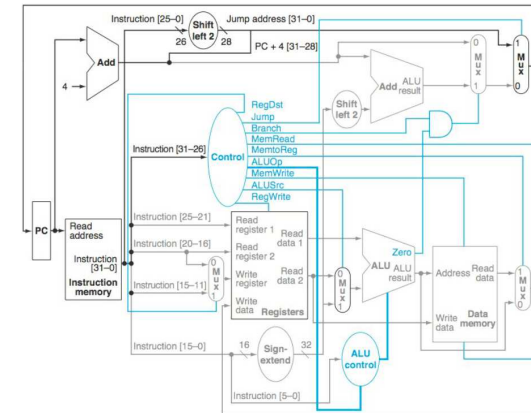
RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
1600ps	1600ps	1400ps	500ps	1400ps	400ps	1500ps	400ps	1700ps

(d) 处理器的时钟周期为多少？

(e) 如果你可以加速控制信号的产生，但加快一个控制信号5ps的代价是处理器成本增加1元。那么为了最大化性能你会加速哪些控制信号？这种性能改进的最小代价是多少？

(f) 如果一个处理器的成本已经很高，那么我们需要在维持处理器性能的同时降低其成本，而不是像第⑤问中所作的那样为提高它的性能而买单。如果你可以使用更慢的逻辑来实现对信号的控制，并且单个控制信号每减慢5ps，处理其成本就可以节省1元，那么在保持处理器性能的同时，你会减慢哪些控制信号，并且减慢多少来降低成本？

❖ 图4



### 参考答案

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	符号扩展	左移两位	ALU控制
500ps	150ps	100ps	180ps	220ps	1000ps	90ps	20ps	55ps

(a) 为了避免增加图4中数据通路的关键路径长度，留给控制单元产生MemWrite信号的时间有多少？

- 关键路径延迟为500ps + 220ps + 100ps + 180ps + 1000ps + 100ps + 220ps = 2320ps，留给控制单元产生MemWrite信号的最长时间为2320ps - 500ps - 1000ps = 820ps
- 解析：假设控制单元的延迟为0，则lw具有最长的关键路径，于是得到关键路径I-Mem、Regs(Read)、Mux、ALU、D-Mem(Read)、Mux、Regs(Write)。控制单元在读取完指令存储器之后才可以产生MemWrite控制信号，并且必须在时钟周期结束之前产生MemWrite信号，但是，由于MemWrite信号是数据存储器的写使能信号，因此在产生该信号之后还应当至少留出写存储器的时间D-Mem(Write)。

### 参考答案(续)

指令存储器	加法器	多路器	ALU	寄存器堆	数据存储器	符号扩展	左移两位	ALU控制
500ps	150ps	100ps	180ps	220ps	1000ps	90ps	20ps	55ps

(b) 图4中哪个控制信号最不关键，控制单元需要在多长时间内产生该信号以避免其成为关键路径？

- Jump信号具有最长的松弛时间，为2320ps - 500ps - 100ps = 1720ps
- 解析：所有的控制信号都必须在指令读取之后生成，同时一个信号最晚必须在时钟周期结束之前到来，对于MemWrite、RegWrite和Jump信号，由于它们都只在时钟周期的最后才会用到，因而相比于其它控制信号会拥有更长的松弛时间，由于两种情况下均是数据存储器的延迟>寄存器堆>多路器，因而Jump具有最长的松弛时间。
- 这个题目里面没有考虑PC的延迟。

### 参考答案（续）

指令存储器	加法器	多选器	ALU	寄存器堆	数据存储器	符号扩展	左移两位	ALU控制
500ps	150ps	100ps	180ps	220ps	1000ps	90ps	20ps	55ps

(c) 图4中哪个控制信号最关键，控制单元需要在多长时间内产生该信号以避免其成为关键路径？

最关键信号	产生该信号可用的时间
ALUSrc (100ps > 55ps)	220ps

- 解析：为了不影响关键路径，控制信号必须在数据到达之前产生以便不影响数据的通过，最早出现在关键路径上的控制信号是ALUOp和ALUSrc，ALUSrc的松弛时间为Regs(Read)，ALUOp的松弛时间为Regs(Read) + Mux - ALU Ctrl，拥有较小松弛时间的信号更为关键，可见二者对于ALU计算的影响取决于ALU Ctrl与Mux的延迟大小。

### 参考答案（续）

RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
720ps	730ps	600ps	400ps	700ps	200ps	710ps	200ps	800ps

RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
1600ps	1600ps	1400ps	500ps	1400ps	400ps	1500ps	400ps	1700ps

(d) 处理器的时钟周期为多少？

对时钟周期影响最大的控制信号	理想的时钟周期(由第(1)问得来)	实际的时钟周期
MemWrite (+680ps)	2320ps	3000ps

- 解析：为了便于后面的计算，我们首先计算各个控制信号的松弛时间，如下表所示：

RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
1500ps	1720ps	1620ps	500ps	1500ps	265ps	820ps	220ps	1600ps

- 若实际产生信号的时间小于松弛时间，则该信号的产生不会影响时钟周期，反之，该信号会影响时钟周期，并且显然是会使时钟周期变大，由此可以计算出新的时钟周期(下表中的数据为实际时间减去松弛时间)。

RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
100ps	-120ps	-220ps	0ps	-100ps	135ps	680ps	180ps	100ps

### 参考答案（续）

RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
1600ps	1600ps	1400ps	500ps	1400ps	400ps	1500ps	400ps	1700ps

(e) 如果你可以加速控制信号的产生，但加快一个控制信号5ps的代价是处理器成本增加1元。那么为了最大化性能你会加速哪些控制信号？这种性能改进的最小代价是多少？

对时钟周期有影响的信号	代价
RegDst (+100ps)	1195/5=239
ALUOp (+135ps)	
MemWrite (+680ps)	
ALUSrc (+180ps)	
RegWrite (+100ps)	

- 解析：应当只加速影响了时钟周期的控制信号，目标是将这些控制信号对时钟周期的影响至少变为0。

### 参考答案（续）

RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
1600ps	1600ps	1400ps	500ps	1400ps	400ps	1500ps	400ps	1700ps

(f) 如果一个处理器的成本已经很高，那么我们需要在维持处理器性能的同时降低其成本，而不是像第⑤问中所作的那样为提高它的性能而买单。如果你可以使用更慢的逻辑来实现对信号的控制，并且单个控制信号每减慢5ps，处理其成本就可以节省1元，那么在保持处理器性能的同时，你会减慢哪些控制信号，并且减慢多少来降低成本？

RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
100ps	-120ps	-220ps	0ps	-100ps	135ps	680ps	180ps	100ps

- 上表中具有最大正值的信号决定了周期，则在保证性能的前提下，可以将所有的信号都减慢到具有跟该信号相同的值，于是最小化成本的方法是按照下表减慢响应信号的速度：

RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
580ps	800ps	900ps	680ps	780ps	545ps	0ps	500ps	580ps

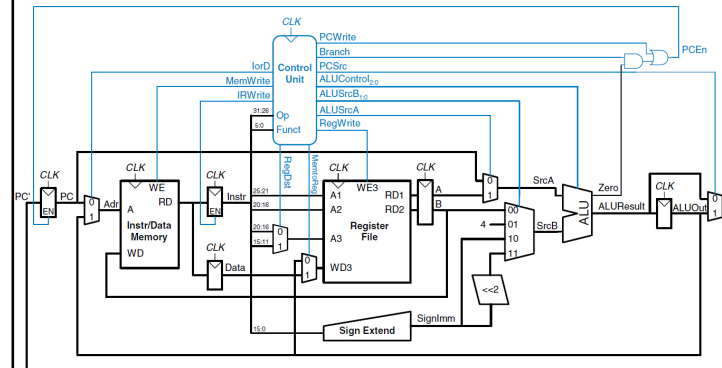
## 2 MIPS多周期微体系结构分析

### 2.1

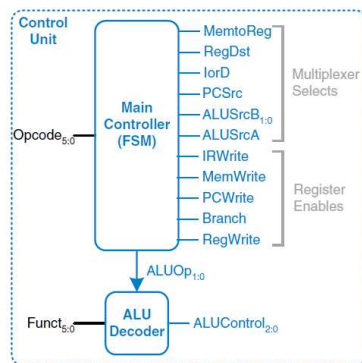
❖ 假设在多周期MIPS处理器的下列控制信号中存在固定为0缺陷，那么哪些指令将会失效？为什么？（数据通路参考图1，其中不包含j指令）

- (a) MemtoReg
- (b) ALUOp0
- (c) PCSrc

❖ 图5



❖ 图6. 图5中控制器的内部结构 (ALUOp<sub>1:0</sub>信号)



## 参考答案

### (a) MemtoReg

- lw指令将会失效，MemtoReg固定为0将导致无法选择DR作为寄存器堆写数据的来源，而lw指令需要将DR中的数据写入寄存器堆。

### (b) ALUOp0

- beq指令将会失效，ALUOp0固定为0将导致ALU无法进行beq需要的减法运算(相应的ALUOp为01)。

### (c) PCSrc

- beq指令将会失效，PCSrc固定为0将导致无法选择ALUOut的输出作为PC值的来源，而beq指令需要保存在ALUOut中的ALU的运算结果来修改PC的值。



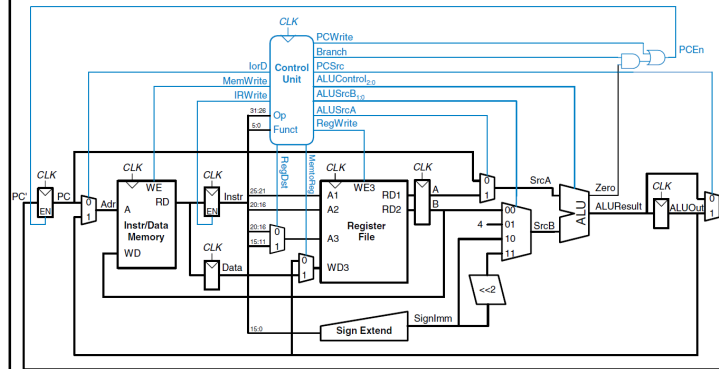
## 2.2

❖ 假设多周期MIPS处理器各个部件的延迟如下表所示(假设存储器和寄存器堆的写速度与读速度相等，数据通路参考图1)：

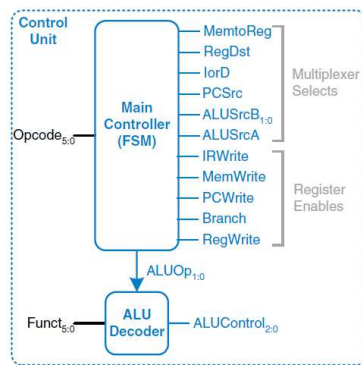
Element	Parameter	Delay(ps)
Register clk-to-Q	$t_{pcq}$	30
Register setup	$t_{setup}$	20
Multiplexer	$t_{mux}$	25
ALU	$t_{ALU}$	200
Memory read	$t_{mem}$	250
Register file read	$t_{Rfread}$	150
Register file setup	$t_{Rfsetup}$	20

- 通过提高哪个部件的速度(即减小该部件的延迟)可以对整个处理器的速度有最大的提升？
- 在避免不必要浪费的前提之下，该部件的延迟应减小到多少？
- 提升之后的处理器周期是多少？
- 有一种寄存器堆，它比现有的寄存器堆功耗低40%，但是延迟是现有寄存器堆的两倍，请分析一下使用这种寄存器堆是否有意义。

❖ 图5



❖ 图6 图5中控制器的内部结构



参考答案 (续)

Element	Parameter	Delay(ps)
Register clk-to-Q	$t_{pcq}$	30
Register setup	$t_{setup}$	20
Multiplexer	$t_{mux}$	25
ALU	$t_{ALU}$	200
Memory read	$t_{mem}$	250
Register file read	$t_{Rfread}$	150
Register file setup	$t_{Rfsetup}$	20

- 通过提高哪个部件的速度(即减小该部件的延迟)可以对整个处理器的速度有最大的提升？

▪ 由于寄存器堆的速度快于存储器( $t_{Rfread} < t_{mem}$ )且 $t_{setup} = t_{Rfsetup}$ ，于是有 $T_c = t_{pcq} + t_{mux} + \max(t_{ALU} + t_{mux}, t_{mem}) + t_{setup}$ ，由表中数据可得 $t_{ALU} + t_{mux} = 225ps < t_{mem}$ ，于是为了减小时钟周期，应当提高存储器的速度。

### 参考答案（续）

Element	Parameter	Delay(ps)
Register clk-to-Q	$t_{pcq}$	30
Register setup	$t_{setup}$	20
Multiplexer	$t_{mux}$	25
ALU	$t_{ALU}$	200
Memory read	$t_{mem}$	250
Register file read	$t_{RFread}$	150
Register file setup	$t_{RFsetup}$	20

(b) 在避免不必要浪费的前提之下，该部件的延迟应减小到多少？

- 应存储器的延迟减小到225ps。

### 参考答案（续）

Element	Parameter	Delay(ps)
Register clk-to-Q	$t_{pcq}$	30
Register setup	$t_{setup}$	20
Multiplexer	$t_{mux}$	25
ALU	$t_{ALU}$	200
Memory read	$t_{mem}$	250
Register file read	$t_{RFread}$	150
Register file setup	$t_{RFsetup}$	20

(c) 提升之后的处理器周期是多少？

- 原来处理器的时钟周期为 $30ps + 25ps + 250ps + 20ps = 325ps$ ，提升存储器速度之后处理器的时钟周期为 $30ps + 25ps + 225ps + 20ps = 300ps$ 。

### 参考答案（续）

Element	Parameter	Delay(ps)
Register clk-to-Q	$t_{pcq}$	30
Register setup	$t_{setup}$	20
Multiplexer	$t_{mux}$	25
ALU	$t_{ALU}$	200
Memory read	$t_{mem}$	250
Register file read	$t_{RFread}$	150
Register file setup	$t_{RFsetup}$	20

(d) 有一种寄存器堆，它比现有的寄存器堆功耗低40%，但是延迟是现有寄存器堆的两倍，请分析一下使用这种寄存器堆是否有意义。

- 寄存器堆的延迟变为 $t_{RFread}=300ps$
- 考虑与寄存器堆有关的两个时间，一个是读寄存器堆操作，时间为 $30ps + 300ps + 20ps = 350ps$ ；另一个是写寄存器堆操作，时间为 $30ps + 25ps + 300ps + 20ps = 375ps$ 。而原来的时钟周期为325ps，于是处理器的时钟周期将变为375ps，计算性能/功耗比为 $(325/375)/0.6 = 1.44$ ，说明功耗的下降幅度大于性能的下降幅度，使用这种寄存器堆是有一定意义的。

### 2.3

❖ 在多周期MIPS处理器上运行下面的程序需要多少个周期？这个程序的CPI是多少？

```

addi    $s0, $0, 5          # sum = 5

while:
    beq   $s0, $0, done      # if result > 0, execute the while block
    addi  $s0, $s0, -1        # while block: result = result - 1
    j     while
done:

```

### 参考答案

```
addi    $s0, $0, 5      # sum = 5

while:
    beq    $s0, $0, done  # if result > 0, execute the while block
    addi    $s0, $s0, -1   # while block: result = result - 1
    j      while

done:
```

- addi和beq都是I型指令，j是跳转指令，addi需要4个时钟周期，beq需要3个时钟周期，j需要3个时钟周期。对程序进行分析可知，第一个addi指令执行了1遍，beq指令执行了6遍(跳出循环时判断条件不成立也执行了一遍)，第二个addi指令执行了5遍，j指令执行了5遍，于是共需要  $4 + 6 \times 3 + 5 \times 4 + 5 \times 3 = 57$  个周期。
- CPI为  $57 / (1+6+5+5) = 3.35$ 。

### 作业3、4——流水线 参考答案

## 1 流水线

分别计算该程序在以下的机器上执行时花费的时钟周期数：观察以下程序：

(a) 非流水线机器  
观察以下程序：  
MULT R3, R1, R2  
ADD R5, R4, R3  
ADD R6, R4, R1  
MULT R7, R8, R9  
ADD R4, R3, R7  
MULT R10, R5, R6

(b) 采用记分板 (scoreboarding) 的流水线机器，有5个加法器、5个乘法器，没有数据转发逻辑

(c) 采用记分板 (scoreboarding) 的流水线机器，有5个加法器、5个乘法器，带数据转发逻辑

(d) 采用记分板 (scoreboarding) 的流水线机器，有1个加法器、1个乘法器，没有数据转发逻辑

(e) 采用记分板 (scoreboarding) 的流水线机器，有1个加法器、1个乘法器，带数据转发逻辑

对于所有上述的机器模型，采用以下4阶段的基本指令周期：

- 1) 取指 (1 个时钟周期)
- 2) 译码 (1 个时钟周期)
- 3) 执行
- 4) 写回 (2 个时钟周期)

MUL(5个时钟周期)、ADD(2 个时钟周期)，乘法器和加法器内部不是流水线的

4) 写回 (2 个时钟周期)

请列出为了计算对流水线的所有假设 (例如如何在流水段之间做数据转发)

### 参考答案：

(a) 非流水线机器  
 $9 + 6 + 6 + 9 + 6 + 9 = 45$  时钟周期

(b) 采用记分板 (scoreboarding) 的流水线机器，有5个加法器、5个乘法器，没有数据转发逻辑

Cycles	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28
MUL R3, R1, R2	F	D	E	E	E	E	E	E	W	W																		
ADD R5, R4, R3		F	D	I	-	-	-	-	-	I	E	E	W	W														
ADD R6, R4, R1			F	I	-	-	-	-	-	I	D	E	E	W	W													
MUL R7, R8, R9					F	D	E	E	E	E	E	E	W	W														
ADD R4, R3, R7						F	D	I	-	-	-	-	-	I	E	E	W	W										
MUL R10, R5, R6							F	I	-	-	-	-	-	I	D	E	E	E	E	W	W							

28 时钟周期(或 26时钟周期，利用寄存器堆内部的数据旁路)

### 参考答案（续）：

(c) 采用记分板（scoreboarding）的流水线机器，有5个加法器、5个乘法器，带数据转发逻辑

```
Cycles      1|2|3|4|5|6|7|8|9|10|11|12|13|14|15|16|17|18|19|20|21|22
MUL R3, R1, R2  F|D|E|E|E|E|E|W|W
ADD R5, R4, R3  F|D|-|-|-|-|E|E|W|W
ADD R6, R4, R1  F|-|-|-|-|D|E|E|W|W
MUL R7, R8, R9  F|D|E|E|E|E|E|W|W
ADD R4, R3, R7  F|D|-|-|-|-|E|E|W|W
MUL R10, R5, R6 F|-|-|-|-|D|E|E|E|E|W|W
```

22时钟周期

(d) 采用记分板（scoreboarding）的流水线机器，有1个加法器、1个乘法器，没有数据转发逻辑

```
Cycles      1|2|3|4|5|6|7|8|9|10|11|12|13|14|15|16|17|18|19|20|21|22|23|24|25|26|27|28|29
MUL R3, R1, R2  F|D|E|E|E|E|E|W|W
ADD R5, R4, R3  F|D|-|-|-|-|-|E|E|W|W
ADD R6, R4, R1  F|-|-|-|-|-|-|D|-|E|E|W|W
MUL R7, R8, R9  F|-|D|E|E|E|E|E|W|W
ADD R4, R3, R7  F|D|-|-|-|-|-|-|E|E|W|W
MUL R10, R5, R6 F|-|-|-|-|-|-|D|E|E|E|E|W|W
```

29时钟周期(或 27时钟周期，利用寄存器堆内部的数据旁路)

### 参考答案（续）：

(e) 采用记分板（scoreboarding）的流水线机器，有1个加法器、1个乘法器，带数据转发逻辑

```
Cycles      1|2|3|4|5|6|7|8|9|10|11|12|13|14|15|16|17|18|19|20|21|22|23
MUL R3, R1, R2  F|D|E|E|E|E|E|W|W
ADD R5, R4, R3  F|D|-|-|-|-|E|E|W|W
ADD R6, R4, R1  F|-|-|-|-|D|E|E|W|W
MUL R7, R8, R9  F|-|D|E|E|E|E|E|W|W
ADD R4, R3, R7  F|D|-|-|-|-|E|E|W|W
MUL R10, R5, R6 F|-|-|-|-|D|E|E|E|E|W|W
```

23时钟周期

## 2 延迟槽

一台五阶段流水线的机器，五个流水段分别是：取指、译码、执行、访存和写回。该机器采用延迟槽技术处理控制相关。无条件分支和有条件分支都在执行阶段获得分支的结果。

(a) 需要多少个延迟槽才能够确保正确的操作？

参考答案：2个

(b) 按照你在(a)中设计的延迟槽数量，下列汇编指令序列中，哪（些）条指令可以放入延迟槽？请使用合适的延迟槽填充方案重写下边的汇编指令代码。

(i) ADD R5 R4, R3  
OR R3 R1, R2  
SUB R7 R5, R6  
J X  
延迟槽  
LW R10 (R7)  
ADD R6 R1, R2  
X:

参考答案：  
ADD R5 R4, R3  
J X  
OR R3 R1, R2  
SUB R7 R5, R6  
LW R10 (R7)  
ADD R6 R1, R2  
X:

## 2 延迟槽

一台五阶段流水线的机器，五个流水段分别是：取指、译码、执行、访存和写回。该机器采用延迟槽技术处理控制相关。无条件分支和有条件分支都在执行阶段获得分支的结果。

(a) 需要多少个延迟槽才能够确保正确的操作？

参考答案：2个

(b) 按照你在(a)中设计的延迟槽数量，下列汇编指令序列中，哪（些）条指令可以放入延迟槽？请使用合适的延迟槽填充方案重写下边的汇编指令代码。

(ii) ADD R5 R4, R3  
OR R3 R1, R2  
SUB R7 R5, R6  
BEQ R5 R7, X  
延迟槽  
LW R10 (R7)  
ADD R6 R1, R2  
X:

参考答案：  
ADD R5 R4, R3  
SUB R7 R5, R6  
BEQ R5 R7, X  
OR R3 R1, R2  
NOP  
LW R10 (R7)  
ADD R6 R1, R2  
X:

## 2 延迟槽

一台五阶段流水线的机器，五个流水段分别是：取指、译码、执行、访存和写回。该机器采用延迟槽技术处理控制相关。无条件分支和有条件分支都在执行阶段获得分支的结果。

(a) 需要多少个延迟槽才能确保正确的操作？

参考答案：2个

(b) 按照你在(a)中设计的延迟槽数量，下列汇编指令序列中，哪（些）条指令可以放入延迟槽？请使用合适的延迟槽填充方案重写下边的汇编指令代码。

(iii) ADD R2 R4, R3  
OR R5 R1, R2  
SUB R7 R5, R6  
BEQ R5 R7, X  
延迟槽  
LW R10 (R7)  
ADD R6 R1, R2  
X:

参考答案：  
ADD R2 R4, R3  
OR R5 R1, R2  
SUB R7 R5, R6  
BEQ R5 R7, X  
NOP  
NOP  
LW R10 (R7)  
ADD R6 R1, R2  
X:

## 2 延迟槽

一台五阶段流水线的机器，五个流水段分别是：取指、译码、执行、访存和写回。该机器采用延迟槽技术处理控制相关。无条件分支和有条件分支都在执行阶段获得分支的结果。

(c) 你能修改流水线减少延迟槽的数量吗(不使用分支预测的技术)? 请清楚地说明你的方法并解释为什么这样能减少延迟槽。

参考答案：

将对jump和branch的目标的解析放到译码阶段。

Jump和 branch可以提前一个时钟周期得到解析，因此一个延迟槽足够保证正确的操作。

## 3 分支预测

考察以下高级语言代码段：

```
int array[1000] = { /* random values */ };
int sum1 = 0, sum2 = 0, sum3 = 0, sum4 = 0;
for (i = 0; i < 1000; i++) // 分支 1: 循环分支
{
    if (i % 4 == 0) // 分支 2: If 条件分支 1
        sum1 += array[i]; // 发生分支的路径
    else
        sum2 += array[i]; // 不发生分支的路径
    if (i % 2 == 0) // 分支 3: If 条件分支 2
        sum3 += array[i]; // 发生分支的路径
    else
        sum4 += array[i]; // 不发生分支的路径
}
```

参考答案：

分支 1:

998/1000 = 99.8%.

第一次和最后一次执行时预测错误

分支 2:

500/1000 = 50%:

分支 3:

0%. 每次执行时分支都改变方向.

(a) 当使用last-time预测器时三个分支的预测准确率分别是多少？(假设每一个分支的last-time计数器起始状态是‘不发生’) 请写出你的计算过程和依据。

## 3 分支预测

考察以下高级语言代码段：

```
int array[1000] = { /* random values */ };
int sum1 = 0, sum2 = 0, sum3 = 0, sum4 = 0;
for (i = 0; i < 1000; i++) // 分支 1: 循环分支
{
    if (i % 4 == 0) // 分支 2: If 条件分支 1
        sum1 += array[i]; // 发生分支的路径
    else
        sum2 += array[i]; // 不发生分支的路径
    if (i % 2 == 0) // 分支 3: If 条件分支 2
        sum3 += array[i]; // 发生分支的路径
    else
        sum4 += array[i]; // 不发生分支的路径
}
```

参考答案：

分支 1:

997/1000 = 99.7%. 分支头两次执行和最后一次执行(退出循环时)时预测错误

分支 2:

750/1000 = 75%. 分支重复 T N N N T N N N 模式，饱和计数器在“强不发生”和“弱不发生”之间转换(每四次预测发生一次，在分支真正发生之后)，预测结果永远是不发生。

分支 3:

500/1000 = 50%. 分支重复 T N T N 模式，每次预测饱和计数器在“强不发生”和“弱不发生”之间转换，每次预测结果是不发生，只有一半是正确的。

(b) 当使用基于2-bit饱和计数器的预测器时三个分支的预测准确率分别是多少？(假设每一个分支的2-bit计数器起始状态是‘强不发生’) 请写出你的计算过程和依据。

### 3 分支预测

考察以下高级语言代码段:

```
int array[1000] = { /* random values */ };
int sum1 = 0, sum2 = 0, sum3 = 0, sum4 = 0;
for (i = 0; i < 1000; i++) // 分支 1: 循环分支
{
    if (i % 4 == 0) // 分支 2: If 条件分支 1
        sum1 += array[i]; // 发生分支的路径
    else
        sum2 += array[i]; // 不发生分支的路径
    if (i % 2 == 0) // 分支 3: If 条件分支 2
        sum3 += array[i]; // 发生分支的路径
    else
        sum4 += array[i]; // 不发生分支的路径
}
```

参考答案:

(i) 弱不发生

分支 2:

$749/1000 = 74.9\%$

分支 3:

0%, 预测器在“弱不发生”和“弱发生”之间振荡

(ii) 弱发生

分支 2:

$749/1000 = 74.9\%$ . 分支的模式是 **T N N N T N N N**, 头四次迭代, 预测器的计数器开始于“弱发生”, 转换到“强发生”, 再转换到“弱发生”, 最后是“弱不发生”, 因此它的预测是 **T T T N**; 接下来是每组四个分支的循环, 计数器状态变化为“强不发生”“弱不发生”“强不发生”“强不发生”, 响应的预测为 **N N N N**. 因此一共有  $249 \times 3 + 2 = 749$  次正确的预测。

分支 3:

$500/1000 = 50\%$ . 分支的模式是 **T N T N**, 第一个分支计数器处于“弱发生”状态, 转换为“强发生”, 接下来是“弱发生”..., 因此预测结果是 **T T T T** ... 所以预测器准确率为50%。

注意: 对于分支 3, 预测准确率强烈地依赖于分支预测器的初始状态。

(c) 当分支2和分支3的2-bit计数器起始状态分别是 (i)‘弱不发生’; (ii)‘弱发生’ 时, 预测准确率分别是多少? 请写出你的计算过程和依据。

### 3 分支预测

考察以下高级语言代码段:

```
int array[1000] = { /* random values */ };
int sum1 = 0, sum2 = 0, sum3 = 0, sum4 = 0;
for (i = 0; i < 1000; i++) // 分支 1: 循环分支
{
    if (i % 4 == 0) // 分支 2: If 条件分支 1
        sum1 += array[i]; // 发生分支的路径
    else
        sum2 += array[i]; // 不发生分支的路径
    if (i % 2 == 0) // 分支 3: If 条件分支 2
        sum3 += array[i]; // 发生分支的路径
    else
        sum4 += array[i]; // 不发生分支的路径
}
```

(d) 当使用两层全局历史分支预测器 (2bit全局历史寄存器+每分支一张独立的模式历史表, 模式历史表的每个表项是一个2-bit饱和计数器) 时, 三个分支的预测准确率分别是多少? 假设全局历史寄存器每一位的初始状态都是‘不发生’, 模式历史表中的2-bit饱和计数器初始状态都是‘强不发生’, 计算预测准确率时, 忽略前500次循环迭代。

### 3 分支预测

参考答案:

分支 1:

$499/500 = 99.8\%$  在500次循环迭代之后, 预测器全局历史都进入“强发生”状态, 它总是预测发生, 只有最后一次循环迭代 (循环分支不发生) 的结果预测错误。

分支 2:

75%. 分析分支2和分支3的相关性很有帮助。全局分支历史包括分支3的最后一次结果以及总是发生的循环分支结果, 分别为偶数次和奇数次循环迭代设立独立的饱和计数器很有效。当最后一个分支3 的分支不发生, 分支2的模式是 **N T N T** ... 因为永远不会有连续两个 **T** 出现所以饱和计数器会在强不发生和弱不发生之间振荡, 导致全部预测都是不发生以及 50% 准确率。当最后一个分支3的分支发生, 分支 2的模式是 **N N N N** ... 将导致100%的精度, 因此总的准确率为75%。

分支 3:

75%. 与上述情况类似, 分支 2的分支历史与分支3相关, 并且分支 3 使用两个饱和计数器, 由If条件分支2的结果决定用哪一个。当分支 2不发生, 分支 3的模式是 **N T N T N T N T** ... (这些分支结果来自于第 1, 2, 3, 5, 6, 7, 9, 10, 11, ...次迭代), 预测器将会在弱不发生和强不发生之间振荡, 但总是预测 **N**; 当分支 2发生, 分支3 同样总是发生, 所以预测器会100%准确。所以, 每4次迭代只有一次预测错误, 总预测精度 75%。

### 4 两层分支预测器

假设一个两层全局预测器由全局历史寄存器和**所有分支共享**的一张模式历史表组成(称为预测器 A)

1) 我们把分支预测器中不同的分支映射到相同位置的情况称为“分支干扰”。预测器A的结构中, 不同分支会在哪里发生这种干扰?

参考答案:

全局历史寄存器 (GHR), 模式历史表 (PHT)

2) 另一个两层全局预测器由全局历史寄存器和**每个分支一张**模式历史表组成 (称为预测器 B”)

(a) 什么情况下预测器A的预测准确率低于预测器B? 请解释理由, 并举例说明。可以通过代码来说明。

参考答案:

当映射到同一个PHT条目的两个分支转向相反方向时预测器A的预测精度会低。考虑一个分支B1, 它对于给定的全局历史总是发生, 如果B1有自己的PHT, 它永远能预测正确。现在, 考虑另一个分支B2, 对于同样的历史总是不发生, 如果B2有自己的PHT, 它也总是能预测正确。但是, 如果B1和B2共享一个PHT, 它们映射到同一个PHT条目, 因此互相影响并降低了彼此的预测精度。

## 4 两层分支预测器

假设一个两层全局预测器由全局历史寄存器和**所有分支共享**的一张模式历史表组成(称为预测器 A)

2) 另一个两层全局预测器由全局历史寄存器和**每个分支一张**模式历史表组成 (称为预测器 B")

(b) 预测器A能获得比预测器B更高的预测准确率吗? 请解释理由, 并举例说明。可以通过代码来说明。

参考答案:

如果一个分支B1在有另一个分支B2共享同一个PHT条目时比独有自己的PHT预测结果更准确, 就有这种可能。考虑这种场景, 分支B1对于给定的全局历史(当B1有自己的PHT)总是预测错误, 原因是这段历史中它总是在发生和不发生之间振荡, 现在有一个总是发生的分支B2映射到相同的PHT条目, 这将会改进B1的预测精度, 因为这个时候分支B1会由于B2总是发生而总是被预测发生。如果B2在相同的历史中比B1执行的更频繁的话, 这也可能不会降低B2的预测精度。因此, 总的预测精度将会得到改善。

## 4 两层分支预测器

假设一个两层全局预测器由全局历史寄存器和**所有分支共享**的一张模式历史表组成(称为预测器 A)

2) 另一个两层全局预测器由全局历史寄存器和**每个分支一张**模式历史表组成 (称为预测器 B")

(c) 分支干扰是否总是影响预测器的预测准确率? 请解释理由, 并举例说明。可以通过代码来说明。

参考答案:

如果映射到同一个PHT条目的分支转向相同, 那么预测器A和B会获得同样的预测精度。在这种情况下, 分支之间的干扰不会影响预测精度。考虑两个分支B1和B2在某段固定的全局历史中总是发生, 那么不管是有自己的PHT还是共享PHT条目, 预测精度是相同的。

## 5 分支预测和推断

考察两台具有15段流水线的机器A和B, 流水段分布如下:

取指 (1个阶段)

译码 (8个阶段)

执行 (5个阶段)

写回 (1个阶段)

两台机器都会在发生流相关时采用数据转发。在译码的最后一个阶段检测是否有流相关, 指令会在停顿在这个阶段等待检测结果。

机器A有一个预测准确率为P%的分支预测器, 分支方向和目标在执行的最后一个阶段产生。

机器B采用推断执行, 类似咱们在课堂上讲的方式。

## 6 分支预测和推断

1) 考察以下在机器A上执行的代码段: 我们把这段代码转化为在机器B上执行的推断代码, 大概是下面这个样子:

Add r3 r1, r2	add r3 r1, r2
sub r5 r6, r7	sub r5 r6, r7
beq r3, r5, X	cmp r3, r5
addi r10 r1, 5	addi.ne r10 r1, 5
add r12 r7, r2	addi.ne r12 r7, r2
add r1 r11, r9	addi.ne r14 r11, r9
X: addi r15 r2, 10	addi r15 r2, 10
.....	.....

(假设条件结果由'cmp'指令计算, 推断由'.ne'指令根据条件结果执行。条件结果在执行的最后一个阶段计算并且可以像其他值一样被转发)

这一段代码会重复执行上百次, 分支40%可能性发生, 60%可能性不发生, 平均而言, P取什么样的值会使机器A比机器B具有更高的指令吞吐量?

## 6 分支预测和推断

参考答案:

这个问题显示了在分支预测机器上预测错误的惩罚和在推断执行机器上执行无用指令浪费的时钟周期之间的折衷。

我们在这里假设:

- 机器A和B具有分离的(流水)分支/比较和加法执行单元,即,在分支/比较指令停顿时可以执行加法指令
- 按序写回
- 当推断执行的指令被发现是无用的(在比较指令结果之后),它仍然会继续剩余的流水段,相当于插入空操作。

对于不同的假设,这个问题的解答会不同,都是可能的正确答案。

在机器A,当 beq r3, r5, X 分支不发生并且预测正确,执行的过程如下:

```
add r3 <- r1, r2      F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
sub r5 <- r6, r7        F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
beq r3, r5, X           F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|E1|E2|E3|E4|E5|WB|
addi r10 <- r1, 5        F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
add r12 <- r7, r2        F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
add r1 <- r11, r9        F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
X: addi r15 <- r2, 10      F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
.....
```

北京航空航天大学

93

93

## 6 分支预测和推断

参考答案:

这个问题显示了在分支预测机器上预测错误的惩罚和在推断执行机器上执行无用指令浪费的时钟周期之间的折衷。

我们在这里假设:

- 机器A和B具有分离的(流水)分支/比较和加法执行单元,即,在分支/比较指令停顿时可以执行加法指令
- 按序写回
- 当推断执行的指令被发现是无用的(在比较指令结果之后),它仍然会继续剩余的流水段,相当于插入空操作。

对于不同的假设,这个问题的解答会不同,都是可能的正确答案。

在机器A,当 beq r3, r5, X 分支发生并且预测正确,执行过程如下:

```
add r3 <- r1, r2      F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
sub r5 <- r6, r7        F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
beq r3, r5, X           F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|E1|E2|E3|E4|E5|WB|
X: addi r15 <- r2, 10      F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
.....
```

北京航空航天大学

94

94

## 6 分支预测和推断

参考答案:

机器A因为分支预测错误(不管是发生还是未发生)受到17个时钟周期(8个译码+5个执行+4个停顿)的惩罚。

机器B在分支不发生且预测正确时的执行跟机器A完全一样,但是,当分支发生(比较结果相等)时机器B浪费3个时钟周期,如下图所示。

```
add r3 <- r1, r2      F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
sub r5 <- r6, r7        F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
cmp r3, r5             F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|E1|E2|E3|E4|E5|WB|
addi.ne r10 <- r1, 5      F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
add.ne r12 <- r7, r2      F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
add.ne r14 <- r11, r9      F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
addi r15 <- r2, 10        F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
.....
```

所以,如果预测错误的代价低于执行无用指令浪费的周期,机器A比机器B的指令吞吐量高。

$$(1-P) \times 17 < 3 \times 0.4$$

因此,  $P > 0.9294$  时, 机器A 比机器B 有更高的指令吞吐量。

北京航空航天大学

95

95

## 6 分支预测和推断

2) 考察在机器A上执行的另一段代码: 我们把这段代码转化为在机器B上执行的推断代码,大概是下面这个样子:

add r3 r1, r2	add r3 r1, r2
sub r5 r6, r7	sub r5 r6, r7
beq r3, r5, X	beq r3, r5, X
addi r10 r1, 5	cmp r3, r5
add r12 r10, r2	addi.ne r10 r1, 5
add r14 r12, r9	add.ne r12 r10, r2
X: addi r15 r14, 10	add.ne r14 r12, r9
.....	addi r15 r14, 10
	.....

(假设条件结果由'cmp'指令计算,推断由'.ne'指令根据条件结果执行。条件结果在执行的最后一个阶段计算并且可以像其他值一样被转发)  
这一段代码会重复执行上百次,分支40%可能性发生,60%可能性不发生,平均而言,P取什么样的值会使机器A比机器B具有更高的指令吞吐量?

北京航空航天大学

96

96



## 6 分支预测和推断

参考答案:

在机器A上, 当 beq r3, r5, X 分支未发生且预测正确, 执行过程如下:

```
add r3 <- r1, r2  F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
sub r5 <- r6, r7   F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
beq r3, r5, X      F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|E1|E2|E3|E4|E5|WB|
addi r10 <- r1, 5  F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
add r12 <- r10, r2 F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|E1|E2|E3|E4|E5|WB|
add r14 <- r12, r9 F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|-|-|-|-|E1|E2|E3|E4|E5|WB|
X: addi r15 <- r14, 10 F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|-|-|-|-|E1|E2|E3|E4|E5|WB|
```

当分支发生且预测正确, 执行过程如下:

```
add r3 <- r1, r2  F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
sub r5 <- r6, r7   F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
cmp r3, r5         F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|E1|E2|E3|E4|E5|WB|
addi r15 <- r14, 10 F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
.....
```

## 6 分支预测和推断

参考答案:

机器A因为分支预测错误(不管是发生还是未发生)受到17个时钟周期(8个译码+5个执行+4个停顿)的惩罚。

机器B在分支不发生且预测正确时的执行跟机器A完全一样, 但是, 当分支发生(比较结果相等)时机器B浪费11个时钟周期, 如下图所示。

```
add r3 <- r1, r2  F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
sub r5 <- r6, r7   F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|WB|
cmp r3, r5         F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|E1|E2|E3|E4|E5|WB|
addi.ne r10 <- r1, 5 F|D1|D2|D3|D4|D5|D6|D7|D8|E1|E2|E3|E4|E5|-|-|-|WB|
addi.ne r12 <- r10, r2 F|D1|D2|D3|D4|D5|D6|D7|D8|-|-|-|E1|E2|E3|E4|E5|WB|
addi.ne r14 <- r12, r9 F|D1|D2|D3|D4|D5|D6|D7|-|-|-|D8|-|-|-|E1|E2|E3|E4|E5|WB|
addi r15 <- r14, 10  F|D1|D2|D3|D4|D5|D6|-|-|-|D7|-|-|-|D8|-|-|-|E1|E2|E3|E4|E5|WB|
```

所以, 如果预测错误的代价低于执行无用指令浪费的周期, 机器A比机器B的指令吞吐高。

$$(1-P) \times 17 < 11 \times 0.4$$

因此,  $P > 0.7411$  时, 机器A 比机器B 有更高的指令吞吐量。

## 作业5——Cache和Memory 参考答案

## 1 Cache

下面给出了运行在带数据cache的处理器上的程序所生成的四种不同的地址序列, 同时给出了每种序列的cache命中率。假设cache在每个序列开始时是空的, 请回答该处理器数据cache的下述参数分别是多少:

(a) 相联度(1, 2 还是4路)

假设: 所有的访存都是单字节的访问, 所有的地址都是字节地址。

序列	地址序列	命中率
1	0, 2, 4, 8, 16, 32	0.33
2	0, 512, 1024, 1536, 2048, 1536, 1024, 512, 0	0.33
3	0, 64, 128, 256, 512, 256, 128, 64, 0	0.33
4	0, 512, 1024, 0, 1536, 0, 2048, 512	0.25

参考答案:

4

对于序列 2, 块 0, 512, 1024 和 1536 是仅有的重用块, 也就是会第二次访问并可能会导致cache命中的块, 其中3块应该在第二次被访问时命中, 因此命中率才会是 0.33 (3/9)。

块大小是 8 字节(见下一问), 对于任何的 cache 大小 (256B 或 512B), 这些块都映射到set 0。因此, 相联度是1或者2会造成四块中最多1或2块在第二次访问时在cache中, 使得最大的可能命中率小于3/9, 而这个序列的命中率是3/9, 说明相联度只能是4。

## 1 Cache

下面给出了运行在带数据cache的处理器上的程序所生成的四种不同的地址序列，同时给出了每种序列的cache命中率。假设cache在每个序列开始时是空的，请回答该处理器数据cache的下述参数分别是多少：

(b) 块大小(1, 2, 4, 8, 16 还是32 字节)

假设：所有的访存都是单字节的访问，所有的地址都是字节地址。

序列	地址序列	命中率
1	0, 2, 4, 8, 16, 32	0.33
2	0, 512, 1024, 1536, 2048, 1536, 1024, 512, 0	0.33
3	0, 64, 128, 256, 512, 256, 128, 64, 0	0.33
4	0, 512, 1024, 0, 1536, 0, 2048, 512	0.25

参考答案：

8 字节

对于序列 1，6次访问中只有2次(地址2和 4) 能够cache命中，命中率是0.33。除了8字节外其他的块大小都不能满足命中率为0.33，要么大要么小。

## 1 Cache

下面给出了运行在带数据cache的处理器上的程序所生成的四种不同的地址序列，同时给出了每种序列的cache命中率。假设cache在每个序列开始时是空的，请回答该处理器数据cache的下述参数分别是多少：

(c) cache总容量(256还是 512 字节)

假设：所有的访存都是单字节的访问，所有的地址都是字节地址。

序列	地址序列	命中率
1	0, 2, 4, 8, 16, 32	0.33
2	0, 512, 1024, 1536, 2048, 1536, 1024, 512, 0	0.33
3	0, 64, 128, 256, 512, 256, 128, 64, 0	0.33
4	0, 512, 1024, 0, 1536, 0, 2048, 512	0.25

参考答案：

256 字节

对于序列3，512字节的容量会使命中率达到4/9(4路组相联，8字节cache块，无论什么替换策略)，高于 0.33，所以 cache 总容量是 256 字节。

## 1 Cache

下面给出了运行在带数据cache的处理器上的程序所生成的四种不同的地址序列，同时给出了每种序列的cache命中率。假设cache在每个序列开始时是空的，请回答该处理器数据cache的下述参数分别是多少：

(d) 替换策略(LRU 还是 FIFO)

假设：所有的访存都是单字节的访问，所有的地址都是字节地址。

序列	地址序列	命中率
1	0, 2, 4, 8, 16, 32	0.33
2	0, 512, 1024, 1536, 2048, 1536, 1024, 512, 0	0.33
3	0, 64, 128, 256, 512, 256, 128, 64, 0	0.33
4	0, 512, 1024, 0, 1536, 0, 2048, 512	0.25

参考答案：

LRU

对于上述的参数，所有序列4中的cache块被映射到组0，如果使用 FIFO 替换策略，命中率是 3/8，而采用LRU 替换策略命中率是 1/4，所以替换策略是 LRU。

## 2 内存的交叉存取

2.1 一台机器有4 KB的主存，由1个通道、1个rank和N(N>1)个bank构成。系统没有虚拟存储。

1) 数据采用cache交叉存取策略，即连续的cache块对应到连续的bank上；

2) cache块大小为32字节，bank的1行有128字节；

3) 采用打开行策略，即行缓冲中的行在被访问后继续保持在行缓冲中，直到有别的行被访问；

4) 行缓冲命中指访问的行存在于行缓冲，行缓冲缺失指访问的行不在行缓冲。

(a) 某个程序在这台机器上执行，访问以下字节时(数字表示字节的位置，比如320表示第320个字节)发生片上cache缺失而需要访存：0, 32, 320, 480, 4, 36, 324, 484, 8, 40, 328, 488, 12, 44, 332, 492，若行缓冲命中率为0，即所有访问的行都不在行缓冲中，请问bank数N的最小值是多少？

参考答案：

2 个

Cache块大小是32字节，所以，对于给定的访存序列相应的cache块访问序列是 0, 1, 10, 15, 0, 1, 10, 15, 0, 1, 10, 15, 0, 1, 10, 15。

当bank数是1时，所有cache块映射到同一个bank，块 0、1、10和 15 映射到行0、0、2和3。所以，当块1紧接着块0被访问时，会产生行缓冲命中，即行缓冲命中率为0。

当 bank数是 2时，块 0、1、10和15映射到不同的行和bank (bank 0, 行 0; bank 1, 行 0; bank 0, 行 1; bank 1, 行 1)。这样，访问序列就是(bank 0, 行 0), (bank 1, 行 0), (bank 0, 行 1), (bank 1, 行 1) (重复四次)。

因此，每个bank上的行 0和1被交替访问，导致行缓冲命中率为0。

## 2 内存的交叉存取

2.1 一台机器有4 KB的主存，由1个通道、1个rank和N(N>1)个bank构成。系统没有虚拟存储。

- 1) 数据采用cache块交叉存取策略，即连续的cache块对应到连续的bank上；
  - 2) cache块大小为32字节，bank的1行有128字节；
  - 3) 采用打开行策略，即行缓冲中的行在被访问后继续保持在行缓冲中，直到有别的行被访问；
  - 4) 行缓冲命中指访问的行存在于行缓冲，行缓冲缺失指访问的行不在行缓冲。
- (b) 如果对于同一个序列，行缓冲命中率是75%，请问bank数N的最小值是多少？

参考答案：

4 个

当 bank 数是1时，cache 块0、1、10和15映射到行0、0、2和3 (与a中相同)。这时的访问序列为 0, 0, 2, 3 (重复四次)，其中3次访问行缓冲不命中，命中率是 25%。

对于其它数量的 bank，块0、1、10和15映射到不同的行。

假设有这四个cache块的行没有行在行缓冲中打开，那么最大的命中率只能是75%，因为对每个块的第一次访问不命中(强制缺失)。这个最大命中率(75%)意味着每个块除了第一次访问后续访问不能有命中，因此，含有四个块的行必须映射到不同的bank，即最少有4个bank。

## 2 内存的交叉存取

2.1 一台机器有4 KB的主存，由1个通道、1个rank和N(N>1)个bank构成。系统没有虚拟存储。

- 1) 数据采用cache块交叉存取策略，即连续的cache块对应到连续的bank上；
  - 2) cache块大小为32字节，bank的1行有128字节；
  - 3) 采用打开行策略，即行缓冲中的行在被访问后继续保持在行缓冲中，直到有别的行被访问；
  - 4) 行缓冲命中指访问的行存在于行缓冲，行缓冲缺失指访问的行不在行缓冲。
- (c) i) 对于同一序列，行缓冲的命中率能达到100%吗？请解释原因

参考答案：

四个cache块映射到不同的行，因此行缓冲命中率达到100%的唯一可能就是包含每个块的行都已经在行缓冲中打开。

## 2 内存的交叉存取

2.1 一台机器有4 KB的主存，由1个通道、1个rank和N(N>1)个bank构成。系统没有虚拟存储。

- 1) 数据采用cache块交叉存取策略，即连续的cache块对应到连续的bank上；
  - 2) cache块大小为32字节，bank的1行有128字节；
  - 3) 采用打开行策略，即行缓冲中的行在被访问后继续保持在行缓冲中，直到有别的行被访问；
  - 4) 行缓冲命中指访问的行存在于行缓冲，行缓冲缺失指访问的行不在行缓冲。
- (c) ii) 如果能达到，最少需要多少bank才能够获得100%的行缓冲命中率？

参考答案：

4 个bank就足以实现，只要4个分别包含4个cache块的行都已经打开(分别在4个bank)。

## 2 内存的交叉存取

2.2 一个DRAM主存储系统由1个通道、1个rank和N个bank构成。Bank一行256字节，一个cache块64字节。数据采用跨bank的行交叉存取方式组织，物理地址的分配方案如下：

行	Bank	列	BiB(Bytes in Bus)
---	------	---	-------------------

采用打开行策略，即行缓冲中的行在被访问后继续保持在行缓冲中，直到有别的行被访问。初始时，所有bank的第1024行打开。

(a) 当有如下的cache块访问序列时，如果系统的行缓冲命中率为33.3% (即1/3)，请问系统中共有多少个bank：

0, 4, 8, 16, 32, 64, 128, 256, 128, 64, 32, 16, 8, 4, 0

参考答案：

$2^5 = 32$  bank

(b) 如果行缓冲命中率是7/15，请问系统中共有多少个bank？

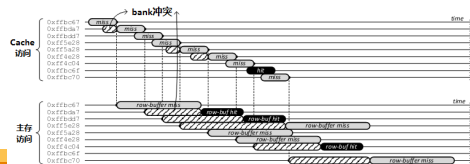
参考答案：

$2^7 = 128$  bank

### 3 Bank

一个处理器的分层存储结构由一个小的SRAM L1-cache和一个大的DRAM主存储器组成，SRAM和DRAM都被划分成bank。处理器有24位物理地址空间，并且不支持虚拟存储(即所有地址都是物理地址)。某个应用开始在这个处理器上运行，下图显示了在时间尺度上应用对存储系统引用的过程(包括在L1-cache和主存中)。

例如，应用对存储第一次引用的字节地址是0xffbc67(假设所有的引用都是对按字节编址的内存地址的按字节读取)，但是这次引用在L1-cache中不命中(假设L1-cache初始时为空)。紧接着，应用访问主存，这会经历一次行缓冲的不命中(初始时，假设主存的所有bank都打开一个永远不会被任何应用访问的行)。最后，包含字节地址0xffbc67的cache块从主存取到cache中。随后的内存引用可能会经历L1-cache和/或主存的bank冲突(当某个特定的bank还在提供之前的某个引用时)。



109

### 3 Bank

下表用16进制和2进制分别给出了该应用对存储系统引用的地址序列。

16进制	2进制
ffbc67	1111 1111 1011 1100 0110 0111
ffbda7	1111 1111 1011 1101 1010 0111
ffbdd7	1111 1111 1011 1101 1101 0111
ff5e28	1111 1111 0101 1110 0010 1000
ff5a28	1111 1111 0101 1010 0010 1000
ff4e28	1111 1111 0100 1110 0010 1000
ff4c04	1111 1111 0100 1100 0000 0100
ffbc6f	1111 1111 1011 1100 0110 1111
ffbc70	1111 1111 1011 1100 0111 0000

请分析上面的图和表，回答下列有关处理器上cache和主存的组织相关的问题，以下是一些假设：

1) L1-cache的假设  
块大小: ? (2的幂, 大于2)  
相联度: ? (2的幂, 大于2)  
数据存储的大小: ? (2的幂, 大于2)  
Bank数: ? (2的幂, 大于2)

注意: 对于以下问题，假设所有的偏移量和索引来自连续的地址位

(a) L1-cache的块大小是多少字节? 24-bit物理地址中哪些位是cache块偏移量? (物理地址的最低位为0位)

参考答案:  
块大小: 16字节  
块偏移量所在位置: 0-3位

初始时为0  
2) 主存的假设  
通道数: 1  
每通道rank数: 1  
每rank的bank数: ? (2的幂, 大于2)  
每bank的行数: ? (2的幂, 大于2)  
每行的cache块数: ? (2的幂, 大于2)  
包含应用的整个工作集  
初始时，所有的bank打开第0行，应用永远不会访问该行

110

### 3 Bank

下表用16进制和2进制分别给出了该应用对存储系统引用的地址序列。

16进制	2进制
ffbc67	1111 1111 1011 1100 0110 0111
ffbda7	1111 1111 1011 1101 1010 0111
ffbdd7	1111 1111 1011 1101 1101 0111
ff5e28	1111 1111 0101 1110 0010 1000
ff5a28	1111 1111 0101 1010 0010 1000
ff4e28	1111 1111 0100 1110 0010 1000
ff4c04	1111 1111 0100 1100 0000 0100
ffbc6f	1111 1111 1011 1100 0110 1111
ffbc70	1111 1111 1011 1100 0111 0000

注意: 对于以下问题，假设所有的偏移量和索引来自连续的地址位

(b) L1-cache有多少bank? 24-bit物理地址中哪些位是L1-cache的bank索引? (物理地址的最低位为0位)

参考答案:  
L1-cache的 bank数: 4  
L1-cache bank 索引位的位置: 4-5

请分析上面的图和表，回答下列有关处理器上cache和主存的组织相关的问题，以下是一些假设：

1) L1-cache的假设  
块大小: ? (2的幂, 大于2)  
相联度: ? (2的幂, 大于2)  
数据存储的大小: ? (2的幂, 大于2)  
Bank数: ? (2的幂, 大于2)

初始时为0  
2) 主存的假设  
通道数: 1  
每通道rank数: 1  
每rank的bank数: ? (2的幂, 大于2)  
没bank的行数: ? (2的幂, 大于2)  
每行的cache块数: ? (2的幂, 大于2)  
包含应用的整个工作集  
初始时，所有的bank打开第0行，应用永远不会访问该行

111

### 3 Bank

下表用16进制和2进制分别给出了该应用对存储系统引用的地址序列。

16进制	2进制
ffbc67	1111 1111 1011 1100 0110 0111
ffbda7	1111 1111 1011 1101 1010 0111
ffbdd7	1111 1111 1011 1101 1101 0111
ff5e28	1111 1111 0101 1110 0010 1000
ff5a28	1111 1111 0101 1010 0010 1000
ff4e28	1111 1111 0100 1110 0010 1000
ff4c04	1111 1111 0100 1100 0000 0100
ffbc6f	1111 1111 1011 1100 0110 1111
ffbc70	1111 1111 1011 1100 0111 0000

请分析上面的图和表，回答下列有关处理器上cache和主存的组织相关的问题，以下是一些假设：

1) L1-cache的假设  
块大小: ? (2的幂, 大于2)  
相联度: ? (2的幂, 大于2)  
数据存储的大小: ? (2的幂, 大于2)  
Bank数: ? (2的幂, 大于2)

注意: 对于以下问题，假设所有的偏移量和索引来自连续的地址位

(c) 主存中有多少bank? 24-bit物理地址中哪些位是主存的bank索引? (物理地址的最低位为0位)

参考答案:  
主存bank数: 8  
主存bank索引位的位置: 10-12

初始时为0  
2) 主存的假设  
通道数: 1  
每通道rank数: 1  
每rank的bank数: ? (2的幂, 大于2)  
没bank的行数: ? (2的幂, 大于2)  
每行的cache块数: ? (2的幂, 大于2)  
包含应用的整个工作集  
初始时，所有的bank打开第0行，应用永远不会访问该行

112

### 3 Bank

下表用16进制和2进制分别给出了该应用对存储系统引用的地址序列。

16进制	2进制	
ffbc67	1111 1111 1011 1100 0110 0111	请分析上面的图和表, 回答下列有关
ffbda7	1111 1111 1011 1101 1010 0111	处理器上 cache 和主存的组织相关的问
ffbdd7	1111 1111 1011 1101 1101 0111	题, 以下是一些假设:
ff5e28	1111 1111 0101 1110 0010 1000	1) L1-cache 的假设
ff5a28	1111 1111 0101 1010 0010 1000	块大小: ? (2的幂, 大于2)
ff4e28	1111 1111 0100 1110 0010 1000	相联度: ? (2的幂, 大于2)
ff4c04	1111 1111 0100 1100 0000 0100	数据存储的大小: ? (2的幂, 大于2)
ffbc6f	1111 1111 1011 1100 0110 1111	Bank数: ? (2的幂, 大于2)
ffbc70	1111 1111 1011 1100 0111 0000	初始时为空

注意: 对于以下问题, 假设所有的偏移量和索引来自连续的地址位

(d) 物理地址向主存映射时用了什么样的交叉存取方案?

参考答案:  
行交叉存取

2) 主存的假设  
通道数: 1  
每通道rank数: 1  
每rank的bank数: ? (2的幂, 大于2)  
没bank的行数: ? (2的幂, 大于2)  
每行的cache块数: ? (2的幂, 大于2)  
包含应用的整个工作集  
初始时, 所有的bank打开第0行, 应用永远不会访问该行

113

### 3 Bank

下表用16进制和2进制分别给出了该应用对存储系统引用的地址序列。

16进制	2进制	
ffbc67	1111 1111 1011 1100 0110 0111	请分析上面的图和表, 回答下列有关
ffbda7	1111 1111 1011 1101 1010 0111	处理器上 cache 和主存的组织相关的问
ffbdd7	1111 1111 1011 1101 1101 0111	题, 以下是一些假设:
ff5e28	1111 1111 0101 1110 0010 1000	1) L1-cache 的假设
ff5a28	1111 1111 0101 1010 0010 1000	块大小: ? (2的幂, 大于2)
ff4e28	1111 1111 0100 1110 0010 1000	相联度: ? (2的幂, 大于2)
ff4c04	1111 1111 0100 1100 0000 0100	数据存储的大小: ? (2的幂, 大于2)
ffbc6f	1111 1111 1011 1100 0110 1111	Bank数: ? (2的幂, 大于2)
ffbc70	1111 1111 1011 1100 0111 0000	初始时为空

注意: 对于以下问题, 假设所有的偏移量和索引来自连续的地址位

(e) 为了支持24-bit的物理地址空间, 主存的每个bank需要多少行? 24-bit物理地址中哪些位是主存的行索引? (物理地址的最低位为0位)

参考答案:  
每个主存bank的行数: 2048  
行索引位的位置: 13-23

2) 主存的假设  
通道数: 1  
每通道rank数: 1  
每rank的bank数: ? (2的幂, 大于2)  
没bank的行数: ? (2的幂, 大于2)  
每行的cache块数: ? (2的幂, 大于2)  
包含应用的整个工作集  
初始时, 所有的bank打开第0行, 应用永远不会访问该行

114

### 3 Bank

下表用16进制和2进制分别给出了该应用对存储系统引用的地址序列。

16进制	2进制	
ffbc67	1111 1111 1011 1100 0110 0111	请分析上面的图和表, 回答下列有关
ffbda7	1111 1111 1011 1101 1010 0111	处理器上 cache 和主存的组织相关的问
ffbdd7	1111 1111 1011 1101 1101 0111	题, 以下是一些假设:
ff5e28	1111 1111 0101 1110 0010 1000	1) L1-cache 的假设
ff5a28	1111 1111 0101 1010 0010 1000	块大小: ? (2的幂, 大于2)
ff4e28	1111 1111 0100 1110 0010 1000	相联度: ? (2的幂, 大于2)
ff4c04	1111 1111 0100 1100 0000 0100	数据存储的大小: ? (2的幂, 大于2)
ffbc6f	1111 1111 1011 1100 0110 1111	Bank数: ? (2的幂, 大于2)
ffbc70	1111 1111 1011 1100 0111 0000	初始时为空

注意: 对于以下问题, 假设所有的偏移量和索引来自连续的地址位

(f) 在一行中的每个cache块被称为列, 一行中有多少列? 24-bit物理地址中哪些位是主存的列索引? (物理地址的最低位为0位)

参考答案:  
每行的列数: 64  
列索引位的位置: 4

2) 主存的假设  
通道数: 1  
每通道rank数: 1  
每rank的bank数: ? (2的幂, 大于2)  
没bank的行数: ? (2的幂, 大于2)  
每行的cache块数: ? (2的幂, 大于2)  
包含应用的整个工作集  
初始时, 所有的bank打开第0行, 应用永远不会访问该行

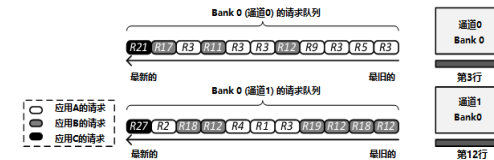
115

### 4 内存调度

为了响应访问请求, 内存控制器会发射1条或多条DRAM命令以从bank访问数据。有4种不同的DRAM命令。

- 1) 激活(ACTIVATE): 取被访问的行装入bank的行缓冲。这一操作也被称为打开行(延迟: 15ns)
- 2) 预充电(PRECHARGE): 将bank的行缓冲中的内容存回行。这一操作也被称为关闭行(延迟: 15ns)
- 3) 读/写: 从行缓冲中访问数据(延迟: 15ns)

下图显示了在时刻0时内存控制器中的内存请求缓冲的快照。每一个请求按照颜色的不同代表了其所属的不同应用(假设所有的应用运行在独立的核上)。同时, 每个请求标注了它要访问的行地址(或索引), 例如R3表示请求的是第3行。另外, 假设所有的请求都是读请求。



116

## 4 内存调度

访存请求在读命令完成后被响应(即读命令被发射15ns之后), 每个应用(A、B或C)停顿直到它所有访存请求被响应为止。  
假设初始时( $t_0$ 时), 每个bank的第3行和第12行分别被取出并存入行缓冲, 没有任何其他应用的请求到达内存控制器。

### 4.1 非应用感知的调度策略

(a) 使用先来先服务调度策略(FCFS), 每个应用的停顿时间是多少?

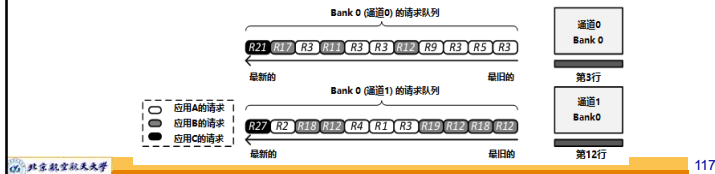
参考答案:

H-命中的延迟, M-缺失的延迟

应用A:  $\text{MAX}(2H+7M, H+9M) = H+9M = 15+405 = 420\text{ns}$

应用B:  $\text{MAX}(2H+8M, H+8M) = 2H+8M = 30+360 = 390\text{ns}$

应用C:  $\text{MAX}(2H+9M, H+10M) = H+10M = 15+450 = 465\text{ns}$



117

## 4 内存调度

访存请求在读命令完成后被响应(即读命令被发射15ns之后), 每个应用(A、B或C)停顿直到它所有访存请求被响应为止。  
假设初始时( $t_0$ 时), 每个bank的第3行和第12行分别被取出并存入行缓冲, 没有任何其他应用的请求到达内存控制器。

### 4.1 非应用感知的调度策略

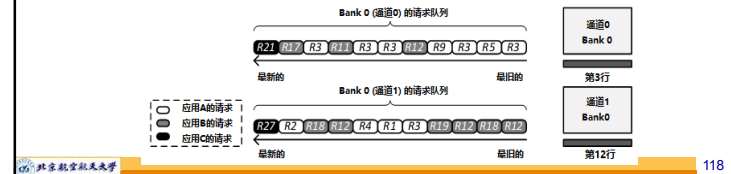
(b) 使用行缓冲优先加先来先服务的调度策略(FR-FCFS), 每个应用的停顿时间是多少?

参考答案:

应用A:  $\text{MAX}(5H+2M, (4H+2M)+4M) = 4H+6M = 60+270 = 330\text{ns}$

应用B:  $\text{MAX}((5H+2M)+3M, 4H+2M) = 5H+5M = 75+225 = 300\text{ns}$

应用C:  $\text{MAX}(((5H+2M)+3M)+M, ((4H+2M)+4M)+M) = 4H+7M = 60 + 315 = 375\text{ns}$



118

## 4 内存调度

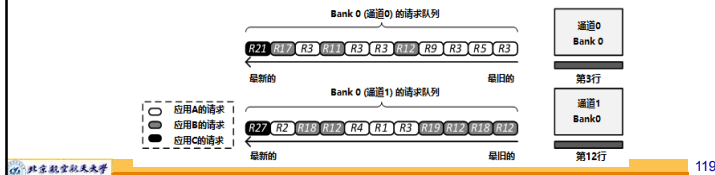
访存请求在读命令完成后被响应(即读命令被发射15ns之后), 每个应用(A、B或C)停顿直到它所有访存请求被响应为止。  
假设初始时( $t_0$ 时), 每个bank的第3行和第12行分别被取出并存入行缓冲, 没有任何其他应用的请求到达内存控制器。

### 4.1 非应用感知的调度策略

(c) FR-FCFS利用的是内存引用行为的什么特征? (6个字 ☺)

参考答案:

行缓冲局部性



119

## 4 内存调度

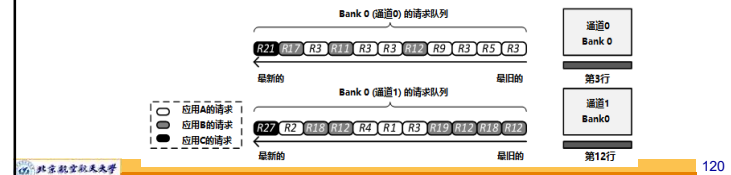
访存请求在读命令完成后被响应(即读命令被发射15ns之后), 每个应用(A、B或C)停顿直到它所有访存请求被响应为止。  
假设初始时( $t_0$ 时), 每个bank的第3行和第12行分别被取出并存入行缓冲, 没有任何其他应用的请求到达内存控制器。

### 4.1 非应用感知的调度策略

(d) 请简要描述可以最大化请求吞吐量的调度策略, 请求吞吐量的意思是每单位时间响应的请求数。(十个字左右☺)

参考答案:

行缓冲优先的先来先服务(FR-FCFS)



120



## 4 内存调度

访存请求在读命令完成后被响应(即读命令被发射15ns之后), 每个应用(A、B或C)停顿直到它所有访存请求被响应为止。  
假设初始时( $t_0$ 时), 每个bank的第3行和第12行分别被取出并存入行缓冲, 没有任何其他应用的请求到达内存控制器。

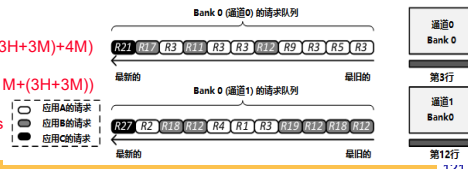
### 4.2 应用感知的调度策略

图中的3个应用, 应用C是内存密集程度最低的(即有最少的请求数)。然而, 它经历了最长的停顿时间, 因为它的请求响应晚于其它多个被优先服务的应用的请求。为了保证应用C的停顿时间最短, 可以为它的请求分配最高优先级, 而给应用A和B的请求分配同样的低优先级。

(a) 调度策略 X: 当应用C分配高优先级并且应用A和B分配相同的低优先级, 每个应用的停顿时间是多少? (对于相同优先级的请求, 假设使用FR-FCFS策略)

参考答案:

应用 A:  $\text{MAX}(M+(4H+3M), M+(3H+3M)+4M)$   
 $= 3H+8M = 45+360 = 405\text{ns}$   
应用 B:  $\text{MAX}(M+(4H+3M)+3M, M+(3H+3M))$   
 $= 4H+7M = 60+315 = 375\text{ns}$   
应用 C:  $\text{MAX}(M, M) = M = 45\text{ns}$



121

## 4 内存调度

访存请求在读命令完成后被响应(即读命令被发射15ns之后), 每个应用(A、B或C)停顿直到它所有访存请求被响应为止。  
假设初始时( $t_0$ 时), 每个bank的第3行和第12行分别被取出并存入行缓冲, 没有任何其他应用的请求到达内存控制器。

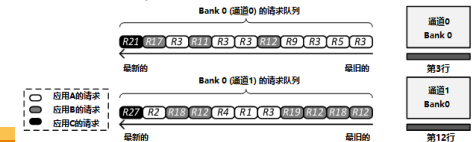
### 4.2 应用感知的调度策略

图中的3个应用, 应用C是内存密集程度最低的(即有最少的请求数)。然而, 它经历了最长的停顿时间, 因为它的请求响应晚于其它多个被优先服务的应用的请求。为了保证应用C的停顿时间最短, 可以为它的请求分配最高优先级, 而给应用A和B的请求分配同样的低优先级。

你能否设计一个更好的调度策略? 虽然应用C的停顿时间小了, 但是应用A和B之间还是会互相影响。

(b) 为其它两个应用分配优先级, 这样你可以最小化所有应用的平均停顿时间。请具体从大到小列出三个应用的优先级(对于相同优先级的请求, 假设使用FR-FCFS策略)

参考答案:  
 $C > B > A$



122

## 4 内存调度

访存请求在读命令完成后被响应(即读命令被发射15ns之后), 每个应用(A、B或C)停顿直到它所有访存请求被响应为止。  
假设初始时( $t_0$ 时), 每个bank的第3行和第12行分别被取出并存入行缓冲, 没有任何其他应用的请求到达内存控制器。

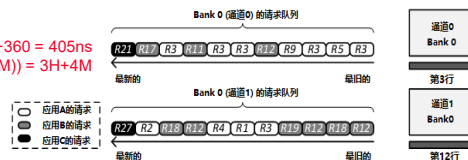
### 4.2 应用感知的调度策略

图中的3个应用, 应用C是内存密集程度最低的(即有最少的请求数)。然而, 它经历了最长的停顿时间, 因为它的请求响应晚于其它多个被优先服务的应用的请求。为了保证应用C的停顿时间最短, 可以为它的请求分配最高优先级, 而给应用A和B的请求分配同样的低优先级。

(c) 调度策略 Y: 使用你的新调度策略, 每个应用的停顿时间分别是多少? (对于相同优先级的请求, 假设使用FR-FCFS策略)

参考答案:

应用 A:  $\text{MAX}(M+(3M)+(4H+3M), M+(3H+3M)+4M) = 3H+8M = 45+360 = 405\text{ns}$   
应用 B:  $\text{MAX}(M+(3M), M+(3H+3M)) = 3H+4M = 45+180 = 225\text{ns}$   
应用 C:  $\text{MAX}(M, M) = M = 45\text{ns}$



123

## 4 内存调度

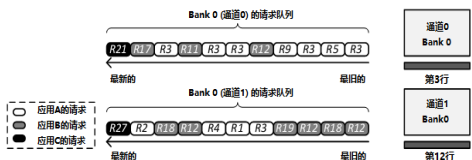
访存请求在读命令完成后被响应(即读命令被发射15ns之后), 每个应用(A、B或C)停顿直到它所有访存请求被响应为止。  
假设初始时( $t_0$ 时), 每个bank的第3行和第12行分别被取出并存入行缓冲, 没有任何其他应用的请求到达内存控制器。

### 4.2 应用感知的调度策略

图中的3个应用, 应用C是内存密集程度最低的(即有最少的请求数)。然而, 它经历了最长的停顿时间, 因为它的请求响应晚于其它多个被优先服务的应用的请求。为了保证应用C的停顿时间最短, 可以为它的请求分配最高优先级, 而给应用A和B的请求分配同样的低优先级。

(d) 请将四种调度策略 (FCFS, FR-FCFS, X, Y) 的平均停顿时间从大到小排列

参考答案:  
 $\text{FCFS} > \text{FR-FCFS} > \text{X} > \text{Y}$



124

## 5 分层存储体系结构

假设你研究出了下一代的存储技术:“魔法RAM”。魔法RAM的位元是非易失性的;它的访问延迟是SRAM的2倍,与DRAM相同;读/写时的能耗和成本与DRAM相当;比DRAM的密度更高。然而,魔法RAM有一个缺点:每个位元在执行2000次写操作之后会停止运转。

(a) 相比DRAM,魔法RAM除了密度高之外,还有什么优势么?请解释。

参考答案:

是的。

魔法RAM不需要刷新,因为它非易失性。这可以降低动态功耗,总线的占用和bank的竞争。魔法RAM的非易失性还可能有利于新的使用模式或编程模型。

(b) 相比SRAM,魔法RAM有什么优势吗?请解释。

参考答案:

是的。

魔法RAM有更高的密度和更低的成本。

北京航空航天大学

125

125

## 5 分层存储体系结构

假设你研究出了下一代的存储技术:“魔法RAM”。魔法RAM的位元是非易失性的;它的访问延迟是SRAM的2倍,与DRAM相同;读/写时的能耗和成本与DRAM相当;比DRAM的密度更高。然而,魔法RAM有一个缺点:每个位元在执行2000次写操作之后会停止运转。

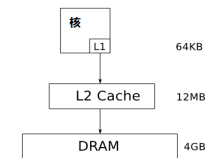
(c) 假设一个系统有64KB SRAM的L1 cache、12MB SRAM的L2 cache和4GB DRAM的主存。

假设你可以利用这个分层存储结构,经过自由地设计和增加任何结构以克服魔法RAM的缺陷(除了修改魔法RAM本身)

(i) 可能将魔法RAM加入这个分层存储结构以减小它的缺陷吗?

参考答案:

是的。



北京航空航天大学

126

126

## 5 分层存储体系结构

假设你研究出了下一代的存储技术:“魔法RAM”。魔法RAM的位元是非易失性的;它的访问延迟是SRAM的2倍,与DRAM相同;读/写时的能耗和成本与DRAM相当;比DRAM的密度更高。然而,魔法RAM有一个缺点:每个位元在执行2000次写操作之后会停止运转。

(c) 假设一个系统有64KB SRAM的L1 cache、12MB SRAM的L2 cache和4GB DRAM的主存。

假设你可以利用这个分层存储结构,经过自由地设计和增加任何结构以克服魔法RAM的缺陷(除了修改魔法RAM本身)

(ii) 如果可能,魔法RAM该放到哪里?根据上面的图来说明,并说明为什么选择放在这个位置。

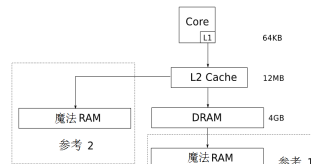
如果不可能,为什么?请解释。

参考答案:

可能的正确答案不止一种。

其中一种:在存储的层次结构中,将魔法RAM放置于DRAM之下,利用DRAM作为魔法RAM的cache。这样,由DRAM执行更多的写操作,使魔法RAM不会过快的磨损。

另一种是把魔法RAM与DRAM并排放置(相同或不同的通道上),利用魔法RAM显式地处理只读数据。



北京航空航天大学

127

127

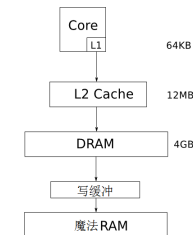
## 5 分层存储体系结构

假设你研究出了下一代的存储技术:“魔法RAM”。魔法RAM的位元是非易失性的;它的访问延迟是SRAM的2倍,与DRAM相同;读/写时的能耗和成本与DRAM相当;比DRAM的密度更高。然而,魔法RAM有一个缺点:每个位元在执行2000次写操作之后会停止运转。

(d) 请给出一种通过修改这个分层存储结构以减少或克服魔法RAM缺陷的方法。请简单清晰地说明你的方法,可以利用图示来说明问题。

参考答案:

采用上述参考1的方案,需要增加一个联合写缓冲的组件以减少写操作对魔法RAM的损耗。同时,存储层次结构中也应该提供一些损耗均衡的机制,或者预测哪些数据被修改的可能性低,将这些数据存入魔法RAM。



北京航空航天大学

128

128



## 6 虚存和cache

一个2路组相联cache，采用写回策略和LRU替换策略，需要 $15 \times 2^9$  bit的标签存储来存储包括有效位、脏位以及LRU等标签信息。Cache虚拟索引，物理标签。虚拟地址空间1MB，页大小是2 KB，cache块大小是8字节。  
(a)该cache的数据存储是多少字节？

参考答案：  
8 KB

Cache是2路组相联，所以每个组有两个标签，每个标签1位，2位有效位，2位脏位，1位LRU位。  
如果我们用 $i$ 表示索引位的位数，  
标签存储大小 =  $2^i \times (2 \times 1 + 2 + 2 + 1) = 15 \times 2^9$   
所以， $2i = 10$ ， $i = 5$   
 $i = 9$   
数据存储大小 =  $2^i \times (2 \times 8) = 2^9 \times (2 \times 8) = 8 \text{ KB}$

## 6 虚存和cache

一个2路组相联cache，采用写回策略和LRU替换策略，需要 $15 \times 2^9$  bit的标签存储来存储包括有效位、脏位以及LRU等标签信息。Cache虚拟索引，物理标签。虚拟地址空间1MB，页大小是2 KB，cache块大小是8字节。  
(b)虚拟索引中有多少位来自虚页号？

参考答案：  
1 位

页大小为2 KB，因此页偏移量是11 位 (10:0)；  
cache 块偏移量是 3 位 (2:0)，虚拟索引 9 位 (11:3)；  
所以虚拟索引中有1位 (11) 来自虚页号。

## 6 虚存和cache

一个2路组相联cache，采用写回策略和LRU替换策略，需要 $15 \times 2^9$  bit的标签存储来存储包括有效位、脏位以及LRU等标签信息。Cache虚拟索引，物理标签。虚拟地址空间1MB，页大小是2 KB，cache块大小是8字节。  
(c)这个存储系统的物理地址空间有多大？

参考答案：  
64 KB

页偏移量11位；物理帧号(物理标签)是 5位；  
所以，物理地址空间是  $2^{(11+5)} = 2^{16} = 64 \text{ KB}$ 。

作业6——预取和并行  
参考答案

## 1 预取 I

假如你是一位架构师，正在为你的机器设计预取引擎。你先在机器上使用跨度预取器执行了A和B两个应用。

应用A:

```
uint8_t a[1000];
sum = 0;
for (i = 0; i < 1000; i += 4)
{
    sum += a[i];
}
```

应用B:

```
uint8_t a[1000];
sum = 0;
for (i = 1; i < 1000; i *= 4)
{
    sum += a[i];
}
```

i 和 sum 在寄存器中，数组a在内存中，一个cache块大小为4个字节。

(a) 使用跨度预取器，应用A和B的预取精度和覆盖率分别是多少？这个跨度预取器检测两次连续访问的跨度，从当前访问的cache块按照这个跨度预取下一个cache块。

参考答案:

应用A的预取精度是248/249，覆盖率是248/250。

应用A访问a[0], a[4], a[8], ... a[996]，有1000/4 = 250次访问。前两次访问a[0]和a[4]不命中，之后，预取器学习到跨度是4并开始预取a[8], a[12], a[16]等等直到a[1000] (访问a[996]时a[1000]被预取，虽然并没有被用到)。统计结果，249个cache块被预取，248个被使用。

因此，预取精度为248/249，覆盖率为248/250。

应用B的预取精度为0，覆盖率为0。

应用B访问a[1], a[4], a[16], a[64]和a[256]，有5次访问。然而，由于这些访问的跨度不是常数，因为数组索引是4的倍数，而不是增或者减一个常数。因此，跨度预取器无法预取到被访问的cache块，使得预测精度和覆盖率均为0。

## 1 预取 I

假如你是一位架构师，正在为你的机器设计预取引擎。你先在机器上使用跨度预取器执行了A和B两个应用。

应用A:

```
uint8_t a[1000];
sum = 0;
for (i = 0; i < 1000; i += 4)
{
    sum += a[i];
}
```

应用B:

```
uint8_t a[1000];
sum = 0;
for (i = 1; i < 1000; i *= 4)
{
    sum += a[i];
}
```

i 和 sum 在寄存器中，数组a在内存中，一个cache块大小为4个字节。

(b) 请分别为应用A和B建议能获得更好的精度和覆盖率的预取器

i) 应用A

参考答案:

下一块预取器总是预取下一个cache块，因此，有a[4]的cache块也会被预取，则预取精度提高到249/250并且覆盖率仍保持249/250。

ii) 应用B

参考答案:

大多数普通的预取器比如跨度、流、下一块等都无法提升应用B的预取精度，因为无法为这些预取器提供一个合适的访问模式。某些采用预执行方法的预取，比如双核执行可能能够改进应用B的预取精度。

## 2 预取II

你跟你的同学一起设计一个预取器，这台机器使用单核、L1和L2 cache以及DRAM内存。我们需要分析不同的预取器和可能的tradeoff。

在本题中，我们要计算预取器在达到稳定状态后的预取精度、覆盖率和带宽开销，所以，所有计算都不包括最开始的6次请求，这6次请求作为预取器的训练集。

(a) 你首先设计一个跨度预取器，观察最后三次cache块请求，如果最后三次请求的跨度是常数，预取器将会使用这一跨度预取下一个cache块。

你执行了一个应用，它具有如下的访问模式 (这些是cache块地址):

A A+1 A+2 A+7 A+8 A+9 A+14 A+15 A+16 A+21 A+22 A+23 A+28 A+29 A+30...

假设这个模式持续了很长时间。

计算你的跨度预取器对于这个应用的精度和覆盖率

参考答案:

0%, 0%。

每三个一组的请求之后，预取按照检测到的跨度触发，但是预取到的块总是无用的；需要的请求不会被这个预取器的预取覆盖。

## 2 预取II

你跟你的同学一起设计一个预取器，这台机器使用单核、L1和L2 cache以及DRAM内存。我们需要分析不同的预取器和可能的tradeoff。

在本题中，我们要计算预取器在达到稳定状态后的预取精度、覆盖率和带宽开销，所以，所有计算都不包括最开始的6次请求，这6次请求作为预取器的训练集。

(b) 你的同学设计了一个新的预取器，当有一个cache块访问时，预取紧接着的N个cache块

(i) 如果用他的预取器执行你刚刚执行过的应用，预取覆盖率和精度分别是66.67%和50%，N是多少？

参考答案:

N = 2。

比如在访问块14之后，预取器预取块15和16，访问15之后预取16 (与已经发射的预取整合) 和17，访问16之后预取17和18。因此，每三个需要的访问中的两个被覆盖 (66.7%)，并且预取的数据一半是有用的 (50%)。

(ii) 假如我们将带宽开销定义为：有预取器时所有cache块的请求数/没有预取器时所有cache块的请求数，那么你同学的预取器在执行刚才那个应用时的带宽开销是多少？

参考答案:

5/3。

对于每一组连续三个访问的cache块，有两个额外的块被预取。比如，取cache块14, 15和16，块17和18也会被预取。

## 2 预取II

你跟你的同学一起设计一个预取器，这台机器使用单核、L1 和 L2 cache 以及 DRAM 内存。我们需要分析不同的预取器和可能的tradeoff。

在本题中，我们要计算预取器在达到稳定状态后的预取精度、覆盖率和带宽开销，所以，**所有计算都不包括最开始的6次请求，这6次请求作为预取器的训练集。**

(c) 你的同学希望改进他的预取器对于刚才那个应用的覆盖率，他可以容忍带宽开销最多两倍。请问他能做到吗？为什么可以/不可以？

参考答案：

不可以。

要获得更好的覆盖率，预取器必须能够跨过上一组3个请求取到下一组中的请求，因为上一组的三个请求都已经在前一次预取到了。比如，访问A+14, A+15 和A+16，它们都已经被预取了，要想提高覆盖率，需要预取到A+21 (下一组三个跨度请求中的第一个)，但是，这需要预取 A+16和 A+21之间的四个块(A+17, A+18, A+19, A+20)，增加的带宽开销超过两倍。

(d) 对于上面的应用，如果想获得100%的覆盖率，N最小得是多少？这个时候的带宽开销是多少？

参考答案：

N= 5 (这样，A+16 预取 A+21，A+21预取A+22, A+23等等)；

带宽开销是7/3。

北京航空航天大学

137

137

## 3 Cache 一致性

(a) MESI cache一致性协议比MSI 协议好在哪里？

参考答案：

允许cache/处理器写一个位置(唯一的干净拷贝)而不需要通知其它的处理器/cache。

(b) 你想要利用MESI置无效协议设计一个基于目录的cache一致性系统，在特定的工作负载下，系统表现得很糟糕，经过仔细的分析，你发现有4个节点持续的发出对某个cache块的置无效请求，这什么情况？

参考答案：

4个节点间发生cache 块乒乓现象。

(c) 如何解决这一问题？

参考答案：

如果现象是由真共享导致的，重写代码减少共享或者使用同步原语减少通信。如果搞不定，可以考虑使用基于更新的一致性协议。

如果是由伪共享导致的，通过编译器或重写代码改变数据的布局，以消除伪共享。

北京航空航天大学

138

138

## 4 一致性协议

假设有一个多处理器系统，系统有512个处理器，每个处理器有1MB的私有写回方式的cache，每个cache块64字节，主存大小为1GB。

(a)如果我们基于MESI cache 一致性协议设计了监听总线，需要多少状态位才能够实现这个一致性协议？这些状态位放在哪？

参考答案：

$2^{24}$ 位。

总共有 $2^{23}$  个cache块 (cache有 $2^{20}$  字节， $2^9$ 个 cache，因此有 $2^{29}$ 字节 在私有 cache中：用每个cache块 $2^6$ 字节去除)，每块需要2位表示状态 (M, E, S或 I)，因此需要 $2^{24}$ 位。

这些位存在于私有cache的标签存储中。

北京航空航天大学

139

139

## 4 一致性协议

假设有一个多处理器系统，系统有512个处理器，每个处理器有1MB的私有写回方式的cache，每个cache块64字节，主存大小为1GB。

(b) 如果用基于目录的cache一致性协议（像我们课堂上讲的例子那样）替换，需要多少状态位？这些状态位在哪？

参考答案：

$2^{24} \times 513 + 2^{24}$  位

仍然需要在私有cache中的 $2^{24}$ 个MESI状态位，然后，必须计算目录存储空间。共有 $2^{24}$  cache块在主存中，每个块需要对应每个处理器1位外加1个独占位(每个块513位)。因此目录中共需要 $2^{24} \times 513$ 位。

目录位：存在于cache 目录；MESI状态位：存在于私有cache。

北京航空航天大学

140

140

## 4 一致性协议

假设有一个多处理器系统，系统有512个处理器，每个处理器有1MB的私有写回方式的cache，每个cache块64字节，主存大小为1GB。  
(c) 对于这个系统，你会选择哪一个协议？为什么？

参考答案：

目录。

总线无法扩展到512个处理器，目录可以。

虽然基于总线的监听系统对于存储的需求要低很多，但是总线无法提供足够的带宽来维持512个处理器的需要。因此，基于目录的系统(用可扩展的互连网络构建)更合适。

北京航空航天大学

141

141

## 5 并行加速比

假如你是一家公司的程序猿，你被要求并行化一个老程序以使它能够在现代多核处理器上跑得更快。

(a) 你并行化了这个程序，然后发现它对于单线程版本的加速比相比于处理器个数的增加而言相差很多。你发现在每个核的数据cache中有大量的cache无效存在，什么样的程序行为导致了这种现象？(请用10个字左右简要说明)

参考答案：

由数据共享导致的cache乒乓

(b) 你修改了程序以解决这个性能问题，然后你发现 程序在每个并行计算之后的一个单线程都会更新一个全局状态，因此导致性能的下降。你的程序有90%的工作是并行的(按照处理器个数x秒计算得出)，另有10%的工作是串行的，并行部分是完美的并行。如果多核处理器核数无限，程序的最大加速比可以到多少？

参考答案：

10。

根据 Amdahl定律: 对于n个处理器，加速比(n) =  $1/(0.1 + 0.9/n)$

由于  $n \rightarrow \infty$ ，加速比(n)  $\rightarrow 10$ 。

(c) 如果要获得4倍的加速比，应该有多少处理器？

参考答案：

6。

由 加速比(n) = 4:

$4 = 1/(0.1 + 0.9/n) \rightarrow 0.25 = 0.1 + 0.9/n \rightarrow 0.15 = 0.9/n, \therefore n = 6$

北京航空航天大学

142

142

## 5 并行加速比

(d) 为了使你改写的程序更高效，公司决定设计一款专门的异构处理器。这款处理器由一个大核(执行代码更快，但是占据的片上面积更大)和多个小核(执行代码更慢，但是消耗面积更小)共享处理器的片上空间。

你的程序并行部分的所有线程将只会在小核上执行：程序的串行部分将会有一个线程执行在大核上。核的性能(执行速度)与它的面积的平方根成正比。

假设芯片面积有16个单元可用，一个小核至少占用1个单元，大核可以占用任意数量的单元。同时假设没有被大核使用的面积会被小核填满。

(i) 如果想让你的程序获得可能的最快执行速度，大核需要多大？

参考答案：

4个单元。

如果给定大核的尺寸是 $n^2$ ，则大核在串行段的加速比是 $n$ ， $16-n^2$ 个小核实现并行段的并行化。这样加速比 =  $1/(0.1/n + 0.9/(16-n^2))$ 。为了最大化加速比，需要最小化分母。对于 $n=1$ ，分母是0.16； $n=2$ ，分母是0.125； $n=3$ ，分母是0.1619。因此 $n=2$ 是最佳的，所以大核占据 $n^2 = 4$ 个单元。

(ii) 如果所有16个单元全部拿来用做小核，这个处理器就变成了同构的多核处理器，对于你的程序而言，它的加速比是多少？假设串行部分跑在一个小核上，并行部分跑在所有16个小核上。

参考答案：

6.4。... 加速比 =  $1/(0.10 + 0.90/16) = 6.4$

(iii) 在串行部分是10%的情况下，使用异构多核(大小核)处理器是有意义的吗？为什么是/不是？

参考答案：

是。

因为串行部分足够大，使得大核在串行部分获得的加速比超过了由于大核带来的并行吞吐的下降

北京航空航天大学

143

143

## 5 并行加速比

(e) 现在你继续优化了你的程序，使得串行部分仅占4%(剩下96%是并行部分)。

(i) 这个时候大核应该有多大(占多少单元)？

参考答案：

2个单元

跟之前的题类似，加速比 =  $1/(0.04/n + 0.96/(16-n^2))$ 。最小化分母以最大化加速比。

$n=\sqrt{2}$ 时最大，因此大核占2个单元。

(ii) 大核这么大的时候加速比是多少？

参考答案：

10。... 加速比 =  $1/(0.04/\sqrt{2} + 0.96/14) = 10.32$

(iii) 假如此时我们采用16个小核的同构多核处理器，你的程序的加速比是多少(假设串行部分跑在一个小核上，并行部分跑在所有16个小核上)？

参考答案：

10。... 加速比 =  $1/(0.04/1 + 0.96/16) = 10$

(iv) 在串行部分是4%的情况下，使用异构多核(大小核)处理器还是有意义的吗？为什么是/不是？

参考答案：

是 and 不是。

虽然仍能获得比同构系统略高的性能，但是异构系统如果无法提供比同构系统显著的性能收益，由于它比同构系统设计复杂得多，也将意义不大。

北京航空航天大学

144

144