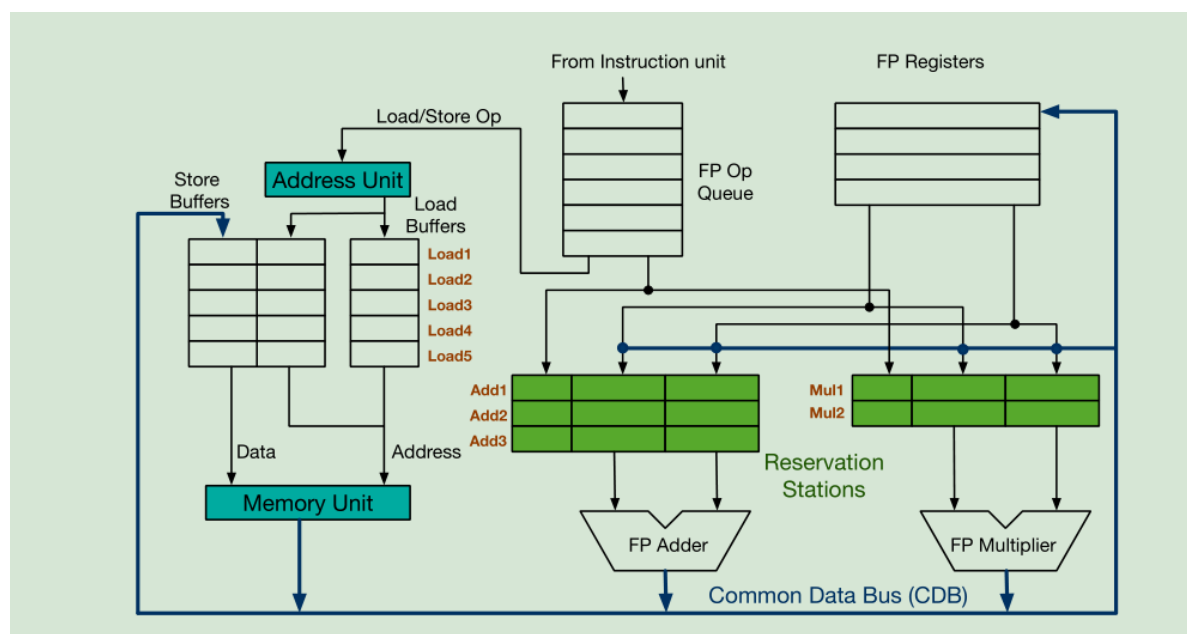


# 2022 Spring 体系结构回忆版

## 1. 乱序执行

考察Tomasulo with ROB, 作业原题

试卷图参考



1. ROB有什么作用？解释ROB entry各参数的意义
2. 指令在什么阶段广播？会发送什么到CDB上？CDB为什么不直接与浮点寄存器堆相连？
3. 参考乱序执行作业，根据ROB、RS和FU的当前状态推导五条指令  
(题目描述几乎一致，第一条指令已经commit并且清空了RS、ROB entry的busy位也置0，第二条指令尚未广播)
4. 填写此时的ROB状态

## 2. VLIW

几乎也是作业原题

计算 $((a4 * x + a3) * x + a2) * x + a1) + a0$

1. 给出了Horner's method的一串指令（共8条）
  1. 画数据流图
  2. 找到数据冒险，并且判断类型
  3. 一条VLIW可以调度三条指令，所有指令都可以一周期内完成，求8条指令的总执行周期数
2. 优化了指令序列（共10条）
  1. 找到数据冒险，如果有WAW和WAR类型的冒险，则修改避免此冒险
  2. 求此时用VLIW需要的执行周期数，并说明为什么优化后指令条数变多，执行效率却更高

## 3. 分支预测

1. 解释BTB entry中，索引、tag的意义和具体的位数  
机器为32位，BTB采用2位饱和和分支预测器，表格大致如下

| BTB index | valid | tag   | target PC | state |
|-----------|-------|-------|-----------|-------|
| 0         | 0     | ...   | ...       | 00    |
| 1         | 1     | 0CFFD | 0CFFD00C  | 01    |
| ...       | ...   |       |           |       |
| 1023      | 1     | 0CFFC | 0CFFD008  | 11    |

## 2. 给定指令序列如下（作业原题）

初始时R3 = 5, R5 = R6 = 0, R2 = 0x80001000, 此处对应的整数数组为[5, 2, 7, 2, 6, 4, 3, 5, 9, 2]

指令op r1, r2, r3中, r1为目的寄存器, r2、r3为源寄存器, l-type指令类似处理

beqz当标志位ZF=0时跳转, bnez为非零跳转, j为无条件跳转

```

loop:
    lw R1, 0(R2)
    andi R4, R1, 1
b_1:
    beqz R4, l_3
    add R5, R5, R1
    j l_4
l_3:
    add R6, R6, R1
l_4:
    subi R3, R3, 1
    addi R2, R2, 4
b_2:
    bnez R3, loop

```

1. 分析程序功能, 求出完成程序后R5、R6的值
2. 填写分支预测表（和作业类似）
3. 计算全局分支预测准确率
4. 如果想优化b\_2分支, 应该如何进行? 如果采用全局预测, 有什么缺点
5. 如果机器支持推测执行（Predicated Execution）, ADDEQ表示当标志位表示相等时执行此步, ADDNE表示不相等时执行此步, 请优化上述指令, 并且说明为什么采用推测执行有利于优化程序

## 4. DRAM

DRAM读取行用时50ns, 采用row buffer以后, 在Row buffer内读取时延较小, 为20ns

在某个时刻t=0, A、B、C三个线程同时发送以下访存请求给DRAM, 到达DRAM的顺序如下（均在t=0到达）, 其中AR1表示A线程请求访问第一行, AR2表示A线程请求访问第二行

| 请求顺序 | 请求行 |
|------|-----|
| A    | 3   |
| A    | 2   |
| A    | 1   |
| B    | 7   |
| A    | 4   |
| C    | 7   |
| B    | 1   |
| A    | 1   |

1. 采用FCFS的排序，不使用row buffer，求三个线程的时延以及平均时延  
时延指从开始请求到线程的最后一个请求完成  
表格如下

| 请求名字 | 访存时延 | 累计时延 |
|------|------|------|
| AR3  | 50   | 50   |
| AR2  | 50   | 100  |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |

2. 采用Row buffer的调度模式

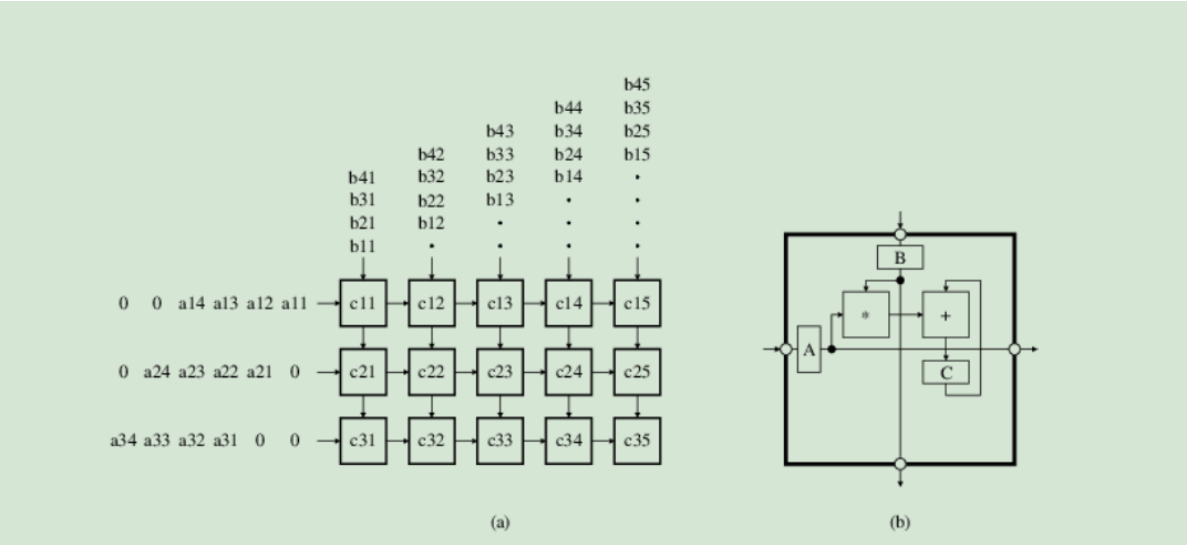
| 请求名字 | 访存时延 | 累计时延 |
|------|------|------|
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |

3. 采用Row buffer，最小化平均时延

| 请求名字 | 访存时延 | 累计时延 |
|------|------|------|
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |
|      |      |      |

## 5. 脉动阵列

脉动阵列用于计算两个矩阵的相乘



脉动阵列的解释自行参考wiki

题目中给出的初始矩阵是3x3的矩阵，即上图去掉最后两列

Part1问题：

1. 为什么硬件加速器使用脉动阵列？和传统的并行计算单元各自独立地读取和写回寄存器相比，有什么优势
2. 求第三周期时脉动阵列各单元正在计算的数值（只用写 $a_{ik} * b_{kj}$ 的形式，无需考虑累加器）
3. 证明第三周期结束，C11输出了正确的矩阵结果
4. 完成3x3矩阵计算需要多少周期？推广到nxn矩阵呢？

Part2问题：

采用循环移位的脉动阵列,移动方向为 $a_{i,n-1} \rightarrow a_{i,0} \quad b_{n,j-1} \rightarrow b_{0,j}$

1. 简单说明为什么循环移位可以计算出正确的矩阵相乘
2. 求第二第三周期时每一个PE正在计算的数值（和Part1 2一样）
3. 完成3x3矩阵计算需要多少周期？推广到nxn矩阵呢？