

Chapter 1:

一、性能衡量

- Latency/Response Time: 同一项任务从开始到结束的时间
- Bandwidth/Throughput: 一定时间内进行的工作总量
- Performance = $1 / \text{Exe Time}$
- CPU Time: CPU的实际执行时间, 不含等待I/O等时间, 与之对应的是 Elapsed Time

$$\text{CPU Time} = \text{Clock Cycles} \times \text{Clock Cycle Time}$$

$$\text{Clock Cycles} = \text{Inst. Count} \times \text{CPI}$$

5. Amdahl's Law

$$\begin{aligned} \textcircled{1} \text{ TImproved} &= \frac{\text{Taffected}}{\text{Tunaffected}} / \text{加速比} + \text{Tunaffected} \\ \textcircled{2} \text{ Speedup} &= \frac{1}{(1-F) + F / \text{局部加速比}} < \frac{1}{1-F} \end{aligned}$$

二、八条设计思想

- Design for Moore's Law
- Use abstraction (抽象) to simplify design
- Make common case fast
- Improve performance via parallelism
- ... via pipelining
- Via prediction
- Use Memory Hierarchy
- Improve dependability via redundancy

三、ISA

- Class of ISA
- Stack Architecture: 隐式操作数在栈顶(TOS), 操作数在栈中弹出, 结果压入栈中
- Accumulator Architecture: 1个隐式操作数(累加器), 1个显式操作数(内存地址), 累加器同时也是结果
- General-Purpose Register Arch: 只支持显式操作数
- Reg-Mem Arch: 任何指令都能访问内存
- Load-Store Arch: 仅 load, store 可访 (80x86 是 Reg-Mem, ARM-RISC-V 是 Load-store)

- 地址对齐 (Addr. Alignment): 设对象大小 s 字节, 其内存地址为 A, 若 $A \% s = 0$ 则称地址对齐

四、Flynn'分类

- SISD: 单处理器, ILP
- SIMD: 多处理器, ILP, DLP
- MIMD: ILP, DLP, TLP

五、功耗公式

- 单位: Power (W, $1W=1J/s$), Energy (J)

- CPU 耗能 { Dynamic Energy (主要来自开关晶体管) }
- Static Energy (来自泄露电流)

$$0 \rightarrow 1 \text{ 或 } 1 \rightarrow 0: \text{Energy-Dynamic} \propto \frac{1}{2} \text{Capacitive Load} \times$$

$$\text{电压}^2, \text{Power-Dynamic} \propto \frac{1}{2} \text{Capacitive Load} \times \text{电压}^2 \times t/\Delta t$$

换频率 (对同一任务, 降低时钟频率可减小 Power, 但无法减小 Energy), Power-Static = I-static $\times V$

(与器件数目成正比)

六、Dependability

- Fault 是设计的瑕疵, Fault 存在时会产生潜在的

Error, Error 影响实际交付的服务时, 就发生 Failure

4. VIPT: 使用 Page Offset (虚拟地址和物理

地址的共同部分索引 L1 Cache, 从而可在访问 TLB

2. Reliability ① MTTF, MTTR ② MTBF (Mean Time Between Failure) = $MTTF + MTTR$

$$3. Availability = \frac{MTTF}{MTTF + MTTR}$$

4. MTTF 的倒数表示故障率 (FIT)

$$5. 10 \text{ 个 MTTF 为 } 10^6 \text{ h 的 Disk, 和 1 个 MTTF 为 } 10^4 \text{ h 的 ATA, 求总 MTTF: 先求总故障率 } 10 \times \frac{1}{10^6} + \frac{1}{10^4} \text{ 再倒数}$$

6. 集成电路成本

- 晶片 (Dies) 由晶圆 (Wafer) 切割而成
- Dies per wafer = $\frac{\pi \times (\text{Wafer Diameter}/2)^2}{\text{Die Area}}$

$$\frac{\sqrt{2 \times \text{Die Area}}}{\text{Cost of Wafer}}$$

$$3. Cost of Die = \frac{\text{Cost of Wafer}}{\text{Dies per wafer} \times \text{Die Yield (良率)}}$$

$$4. Die Yield = \frac{\text{Wafer Yield} \times \frac{1}{(1 + \text{Defects per unit area} \times \text{Die area})}}{N}$$

八、补充 Normalized Geometric Mean 具有统一结果, 不论机器

Chapter 2: Memory Hierarchy

一、定量计算

- Latency: 从内存块读出第一个块的时间 (难以量化)
- Bandwidth: 读出内存块剩余部分的时间

$$2. AMAT = \text{Hit Rate} \times \text{Hit Time} + \text{Miss Rate} \times \text{Miss Time}$$

$$= (1 - \text{Miss Rate}) \times \text{Hit Time} + \text{Miss Rate} \times (\text{Hit Time} + \text{Miss Penalty}) = \text{Hit Time} + \text{Miss Rate} \times \text{Miss Penalty}$$

$$= \text{Hit Time}_{L1} + \text{Miss Rate}_{L1} \times (\text{Hit Time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2})$$

$$3. Local Miss Rate = \frac{\text{该级 Cache Miss 次数}}{\text{该级 Cache 访问次数}} \quad Global Miss Rate = \dots$$

$$4. CPU Time = (\text{CPU Clock Cycles} + \text{Mem. Stall Cycles}) \times$$

$$\text{clock cycle time}, \text{其中 Mem. Stall Cycles} = \text{Miss 数} \times \text{Miss Penalty} = IC \times \text{Mem. 指令占比} \times \text{Miss Rate} \times \text{Miss}$$

$$Penalty$$

$$e.g. 假设理想 CPI 1.1, 30\% \text{ 指令为 MEM, 有 } 10\% \text{ 的 MEM}$$

指令发生 Miss, Miss 代价为 50 cycles; 有 1% 的指令发

生 Miss, Miss 代价为 50 cycles, 求 CPI, AMAT

$$解: CPI = 1.1 + 30\% \times 10\% \times 50 + 1\% \times 50 = 3.1$$

$$AMAT = \frac{1}{1.1} + [1 + 1\% \times 50] + \frac{0.3}{1.1} + [1 + 10\% \times 50] = 2.56$$

二、设计思想

- 局部性原理: Temporal / Spatial Locality

- Unified Cache: 命中率低

Split Cache: 指令与数据一级 Cache 分开 (二级不分)

3. Miss 的种类 ① Compulsory: 冷启动/数据第一次被访问的 Miss (可通过预取降低) ② Capacity: Cache

已满时, 一些数据被替换出去, 再访问这些数据的

Miss ③ Conflict: 在组相联/直接映射中多个块

被映射到同一个 Cache 组中, 导致频繁替换

4. VIPT: 使用 Page Offset (虚拟地址和物理

地址的共同部分索引 L1 Cache, 从而可在访问 TLB

则完成对 L1 Cache 的访问

- 多级 Cache: L1 足够小, 以减小 Hit Time, L2 足够大, 以减小 L1 Miss Penalty
- 写策略 ① Write-Through: 同时写进 Cache 和 Mem, Mem 总是最新的, Cache 无 dirty 位 ② Write-Back: 只写 Cache, 该块要被替换时才写回 Mem

处理 Write Miss: ① Write Allocate: Miss 的块先拿到 Cache 中, 再正常写入 (Write-Back 只能用此方法)

- No Write Allocate (Write Around): 不经过 Cache, 直接写入 Memory (常搭配 Write-Through)
- Interleave: 在读写大量/分散性数据时, 需对不同 bank 连续操作, 通过合理控制, 保证各 bank 传输周期相连贯 (U-Shape)

三、Cache 优化

初期: Compulsory, 后期: Conf/Cap?

Chapter 2: Memory Hierarchy

一、定量计算

- Latency: 从内存块读出第一个块的时间 (难以量化)
- Bandwidth: 读出内存块剩余部分的时间

$$2. AMAT = \text{Hit Rate} \times \text{Hit Time} + \text{Miss Rate} \times \text{Miss Time}$$

$$= (1 - \text{Miss Rate}) \times \text{Hit Time} + \text{Miss Rate} \times (\text{Hit Time} + \text{Miss Penalty}) = \text{Hit Time} + \text{Miss Rate} \times \text{Miss Penalty}$$

$$= \text{Hit Time}_{L1} + \text{Miss Rate}_{L1} \times (\text{Hit Time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2})$$

$$3. Local Miss Rate = \frac{\text{该级 Cache Miss 次数}}{\text{该级 Cache 访问次数}} \quad Global Miss Rate = \dots$$

$$4. CPU Time = (\text{CPU Clock Cycles} + \text{Mem. Stall Cycles}) \times$$

$$\text{clock cycle time}, \text{其中 Mem. Stall Cycles} = \text{Miss 数} \times \text{Miss Penalty} = IC \times \text{Mem. 指令占比} \times \text{Miss Rate} \times \text{Miss}$$

$$Penalty$$

$$e.g. 假设理想 CPI 1.1, 30\% \text{ 指令为 MEM, 有 } 10\% \text{ 的 MEM}$$

指令发生 Miss, Miss 代价为 50 cycles; 有 1% 的指令发

生 Miss, Miss 代价为 50 cycles, 求 CPI, AMAT

$$解: CPI = 1.1 + 30\% \times 10\% \times 50 + 1\% \times 50 = 3.1$$

$$AMAT = \frac{1}{1.1} + [1 + 1\% \times 50] + \frac{0.3}{1.1} + [1 + 10\% \times 50] = 2.56$$

二、冒险情形与处理

- Structural Hazard: 多条指令竞争 FU, 数据读写口

- Data Hazard

① Data dependence: RAW, 真依赖

② Name dependence: WAR, WAW, 假依赖

I. Antidependence: WAR, 指令 A 读取数据后, 指令 B

向同一寄存器写入数据, A 可能会错误读取到 B 写入

的值 II. Output Dependence: WAW, 两指令都向同一

寄存器写入数据, 可能前一条指令覆盖后一条指令写的值

3. 寄存器重命名: 用于解决 Name Dependence, WAR, WAW 分别对后面/前面的寄存器重命名

e.g. mul x2, x2, x2 add x1, x1, x2 \Rightarrow mul p7, p2, p2 add p8, p1, p7

mul x2, x4, x4 mul p9, p4, p4

四、多发射

- 多发射可使 CPI < 1
- Superscalar

静态、动态都可

3. VLIW

一级流水线用时少

五、精确异常

1. 确保异常指令前的所有指令完成执行, 异常指令及之后的指令尚未执行

2. mepc: 异常指令的地址, mtvec: 异常处理程序的入口地址, mtral: 存储异常相关信息

六、ScoreBoard

若不允许发射, 则它及它的后续指令都无法进入 IS, Tomasulo 也这样

I. 算法流程

① IS: 允许发射的条件: I. FU 空闲 (避免结构冲突)

II. 正在执行的其它指令与该指令的目标寄存器都不

相同 (避免 WAW) ② RO: 两个源寄存器都准备好了才读数, 进入 EX (避免 RAW 真依赖) ③ WB: 写回时看目标寄存器是否在某个 FU 的 Ready List 中为 yes

若是说明有指令在读寄存器, 等读完再写回 (避

5. 实现前递时, 若 EX Hazard 与 MEM Hazard 同时满足, 只执行 EX 前递

6. 静态 - Compiler, 动态 - Hardware

三、控制冒险与分支预测 (动态解决)

1. 2-bit predictor

① 使用一个分支预测缓冲区 (BHT). 在 BHT 中, 每条分支指令使用 2 位记录状态, 还记录对应的目标 Address

到 BTB Strongly taken $11 \leftrightarrow 10$ Weakly taken

Strongly Untaken $00 \leftrightarrow 01$ Weakly Untaken

② BHT 形似 Cache, 使用分支指令地址低位作 Index, 高位作 tag.

③ 引申 N-bit predictor: 若 BHT 记录值 $\geq (2^n - 1)/2$, 则预测为 Taken

Is Inst. in BTB? Predict Reality Delay Cycle

✓	Taken	Taken
✓	Taken	Not Taken
✗	Not Taken	Taken
✗	Not Taken	Not Taken

<p

免WAR)。
2. 由于 ScoreBoard 通过停顿的方式处理 WAR/RAW, 所以说并未解决这些冒险; 升级版可显式 Reg. 重命名
3. 按序发射, 乱序完成, 不能实现精确中断
4. 三张状态表
① Inst. Status Table: 记录指令所处阶段 IS/RO/EX/WB
② FU Status: 每个 FU 对应如下记录:
Busy, OP: 正在执行的指令类型
Fi, Fj, Fk: 目标(i), 源(j,k)寄存器号
Qi, Qk: Fj, Fk 若还没准备好, 应该从哪个 FU 的目标寄存器读取,添写 FU 的名字
Rj, Rk: 指示源寄存器的读取情况, 刚准备就绪且尚未读取时为 Yes, 若还未准备就绪或已读取则为 No, 要求两个操作数都就绪后一起读取
③ Register Status Table: 若某个 FU 以该寄存器为目标, 则填 FU 的名字, 否则留空

	IS	EX	WB	COMMIT
FLD F6,34(R2)	1	2	3	4
FLD F2,45(R3)	5	6	7	8
FMUL F0,F2,F4	6	9	10-19	20
FSUB F8,F6,F2	7	9	10-11	12
FDIV F10,F0,F6	8	21	22-61	62
FADD F6,F8,F2	13	14	15-16	22

	IS	RO	EX	WB	MUL:10
FLD F6,34(R2)	1	2	3	4	ADD:2
FLD F2,45(R3)	5	6	7	8	SUB:2
FMUL F0,F2,F4	6	9	10-19	20	LD:1
FSUB F8,F6,F2	7	9	10-11	12	DIV:40
FDIV F10,F0,F6	8	21	22-61	62	
FADD F6,F8,F2	13	14	15-16	22	

t. Tomasulo's Algorithm

1. 算法流程 (分为 IS, EX, WB)
① IS: 保留站有空余时, 从指令序列头部取指令发送到保留站 (Reservation Center) ② EX: 全部操作数可用时, 将指令送入功能单元执行; 若遇分支指令, 在分支指令执行完前, 不准执行任何指令
③ WB: 将结果放到 CDB 上广播

2. 三张状态表 (仅第二张不同):

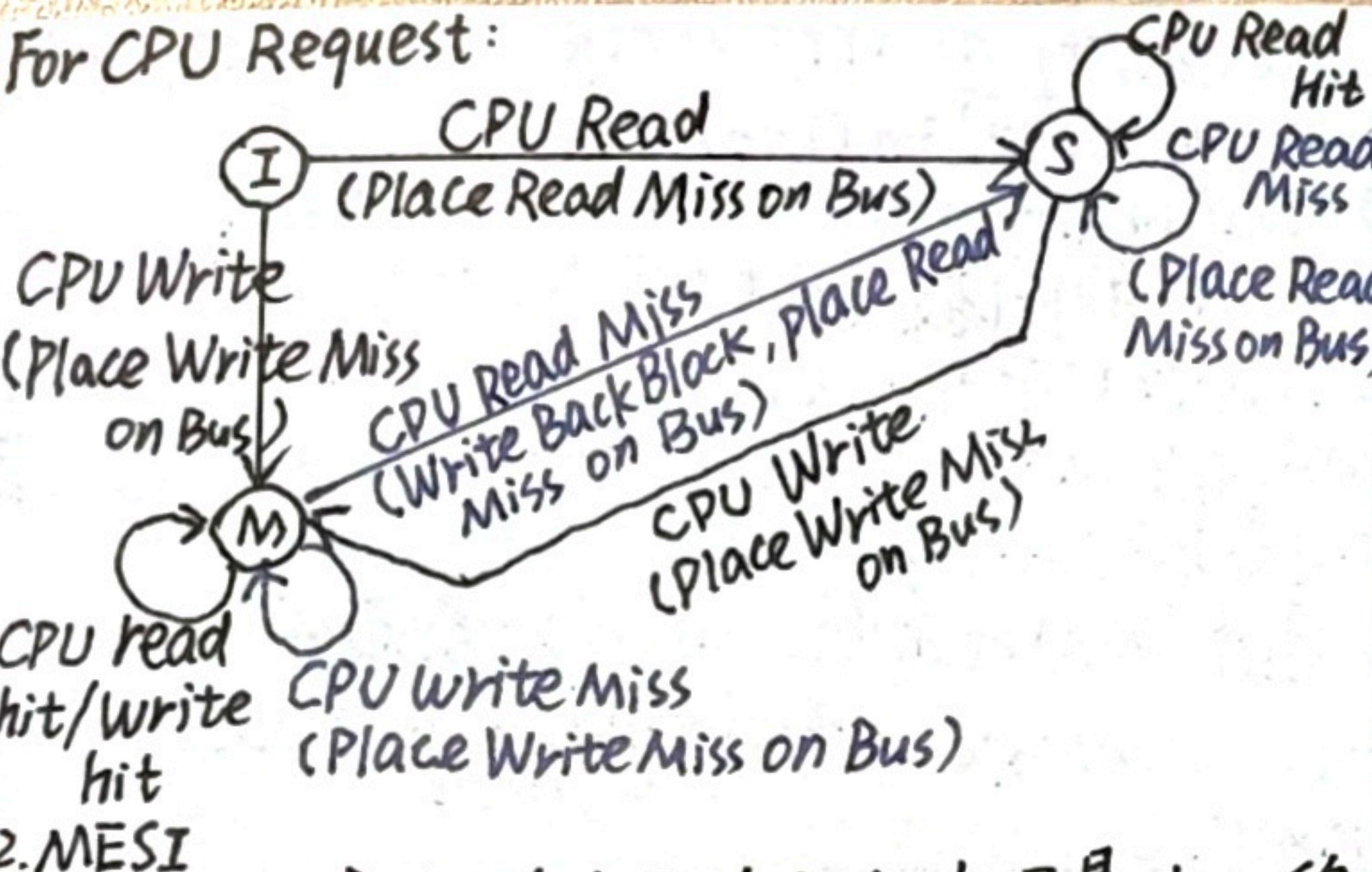
ROB Status: 每个 ROB 对应如下记录:
busy, Op, Vi, Vj: 操作数的真实值
Qi, Qk: 将产生 Vi, Vj 结果的保留站, 若操作数已产生并送入保留站, 则置空 (即 Qi, Vi 只能有一个具有有效值)
A: 对于 load, store 指令, 存放 Mem. Addr. 计算的信息
初始时存放立即数, 地址计算完后存放实际的 Effective Addr. (使用 Vj 存放计算地址所用的寄存器值)
3. 实现了隐式的 Reg. 重命名, 不能实现精确中断

	IS	EX	WB	在Tomasulo
FLD F6,34(R2)	1	3	4	中, Load 指令
FLD F2,45(R3)	2	4	5	的 IS 和 EX 之间要停一拍
FMUL F0,F2,F4	3	6-15	16	
FSUB F8,F6,F2	4	6-7	8	
FDIV F10,F0,F6	5	17-56	57	
FADD F6,F8,F2	6	9-10	11	

1. Hardware Speculation (Tomasulo with ROB)

- ROB 的条目字段
- Inst. Type: 存储指令类型: I. 分支指令(无目的地)
- Store 指令(目的地为存储器地址) II. 寄存器指令 (Load/ALU)(目的地是指令的目标寄存器号)
- Destination(目的地) ③ Value: 结果值
- Ready: 表明指令完成执行, 结果值就绪但还未提交
- 算法流程
Commit: 提交预测错误的分支指令时, 清空 ROB, 执行过程从后续正常指令处重新开始
- 顺序提交实现了精确中断

For CPU Request:



2. MESI

Exclusive: 该块只存在于一个 Cache 中, 且是 clean 的, 该状态设置于某个块刚被读进来时
E → (Read by Others) → S
E → (Local Write) → M, 且无需在总线上发送 invalidate 信号

3. MOESI

Owned: 表明该块被 Cache 所有且为 dirty

四. True/False Sharing Miss

1. True Sharing Miss

① Centralized Shared Memory (SMP/UMA): 所有处理器统一共享物理内存, 所有处理器访问任何内存字的时间相等 ③ Distributed Shared Memory (DSM/NUMA): 每个核有独立的内存, 但地址空间是统一的, 访问远程内存较慢.

二. 缓存一致性问题 Cache Coherence Problem

1. Coherence: 规定读操作应当返回何值

① P.read(x) → P. P.write(x), 总返回 P 写入的值

② P1.read(x) → P2. write(x), 返回 P2 写入的值

③ Write Serialization

④ 不同内存位置

2. Consistency: 针对内存决定写操作的值何时能被读操作返回. 内存一致性模型:

① Sequential Consistency: 读、写顺序都与程序一致

② Total Store Order: 读乱序, 写顺序

③ Weak Order: 读、写都乱序, 使用 Mem. Barrier / Fence 来强制顺序.

⑤ 不同内存位置

3. UMA 使用 Snoopy Protocol: 所有处理器连总线, 写 Cache 时在总线上发送 Invalid 信号; NUMA 使用 Directory Protocol, 使用 Directory 记录哪些处理器在 Cache 中有特定 block 的拷贝. 修改时, 通过 Directory 以点对点方式通知

4. Write-Invalidate / Write-Broadcast

三. Write Invalidate Protocol

1. MSI: Invalid/Shared/Modified

For Bus Request:

Write Miss

Read Miss

Write back block

Read back block

Write back block

Read Miss

Write back block

Read Miss