

**Chapter 1** ISA的七个维度: ① Class of ISA; Memory Addressing; ② Addressing Modes, Types and sizes of operands; Operations; ③ Control flow inst.; Encoding an ISA. 8-24月11号

3 Walls: ILP Wall; Memory Wall; Power Wall  
SISD: ILP; SIMD, MISD; MIMD(TLP)

Dynamic Power:  $Power_{dy} = \frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency Switched}$ ; Energy  $= Cap\_load \times \text{Voltage}$

Static Power:  $Power_{static} = \text{current static} \times \text{Voltage}$

MTTF: To Failure; MTTR: To Repair; FIT = 1/MTTF

MTBF = MTTF + MTTR; Availability =  $\frac{MTTF}{MTBF}$

MIPS =  $\frac{1}{1000000} \left( \frac{\# \text{ of inst.}}{\text{benchmark}} \times \frac{\text{benchmark}}{\text{total runtime}} \right)$  每秒百万条指令

CPU Time =  $IC \times CPI \times \frac{1}{rate}$  IPC: 每时钟周期指令数

Amdahl's Law =  $Speedup = \frac{1}{(1-a) + a/f}$

**Chapter 2** trace cache: 动态指令重定向, 一个ins 保存一小部分多指令

DRAM (Dynamic Random) = 高密度, 便宜, 慢, 定时刷新, 主存

SRAM (Static Random) = 低密度, high power, 贵, cache

Cache原理: 时间, 空间局部性, Temporal/Spatial

Cache 四个问题: Q1: Block Placement (FA, DM, SA)

Q2: Block Identification (Tag/Block)

Q3: Block Replacement (LRU/Random/FIFO)

Q4: Write Strategy (Write Buffer / Write Through)

Write Buffer: 阻塞  $\rightarrow$  L2 Cache; Write Policy: Write Allocate / Write Around (write miss)

CPU Ex Time =  $(CPU\ Ck + Mem\ Stall\ cycle) \times Ck\ time$

Mem Stall Time =  $IC \times Mem\ reqs\ per\ inst. \times MR \times MP$

AMAT = Hit Time +  $MR \times MP$ . 单次访问指令数

提升 Cache 性能: ① Reduce Hit Time: Small and simple first-level caches, way prediction; avoiding address translation, trace cache | ② Increase Bandwidth: Pipelined caches, multibanked caches, non-blocking caches. (L2US 反冲)

③ Reduce Miss Penalty: multilevel caches, read miss prior to writes, critical word first, merging write buffers, victim caches (early restart)

④ Reduce Miss Rate: larger block size, large cache size, higher associativity; compiler optimizations (moving arrays, Loop Interchange, Loop Fusion, Blocking)

⑤ Reduce miss penalty or miss rate via parallelization: Hardware or compiler prefetching

提升 DRAM 性能 (CFM: Fast Page Mode DRAM; Synchronous DRAM; DDR; Rambus)

VIPT: 要求 pg offset  $\geq$  block offset + index. 地址与索引可重叠

多路组相选, cache 大小可扩大 n 倍. 不同组用不同 idx, tag, ofs.

多级:  $AMAT = HT_{L1} + MR_{L1} (HT_{L2} + MR_{L2} \times MP_{L2})$

第①个: Block size  $\uparrow$ , Compulsory  $\downarrow$ ; Cache larger  $\uparrow$ , Capacity  $\downarrow$ ; Higher associativity  $\downarrow$ , conflict  $\downarrow$ .  $MR \rightarrow$  Block size - 块小 compulsory  $\downarrow$ , 后来 conflict  $\uparrow$

**Chapter 3**

Latency: 一个指令产生结果到另一个指令可使用结果的 cycle 差值

Initiation Interval: FU 接受两个指令 cycle 差值. 完全流水: 1; 完全非流水: (latency + 1)

Latency = FU Time - 1 cycle

RAW: True Dependence; WAW: output -; WAR: anti -

ILP: Basic Block 中 ILP 的指令; 软件: Loop LP

两个基本层次: ① exception behavior ② data flow. 程序正确性

软件指令: Loop unrolling; Static Branch Prediction; Static Multiple Issue (VLIW); Advanced Compiler (软件语言或汇编代码调优); 硬件辅助: Conditional or Predicted Inst.; compiler speculation, hw support

硬件指令: dynamic scheduling

Scoreboard: IF, IS: 避免结构冲突, WAW 冲突

RO: 直到所有可用; 动态解决 RAW 冲突

EX, WB: 解决 WAR 冲突

指令级并行: Inst. Status: 指令阶段; RRS (FU 和寄存器的寄存器); FUS (busy, op, Fj 目标, Fj, Fj 源, Qj, Qj 从哪个 FU 得数, Rj, Rk 是否可用)

限制: ILP; issued window 不能是串; FU; WAR 和 WAW stalls

特点: Multiple Funcs, (issue in order, 000, Scoreboard centralized control stall when waw, war)

Tomasulo: Fewer Func, unpipelined, issue in order, complete 000, FP.op. queue, Reservation Station, LD/ST buf, OCB. Renan: no waw, war. OCB: forwarding, decentralized. 不局限于基本

保留站 (busy, op, Vj, Vj 源, Qj, Qj 源操作数寄存器. A read)

(ZF) IS: 保留站空闲, 发射

(EX) (WB): 监听总线. 两个都有命了, 执行

好处: 分布式, 消除 WAW/WAR. 缺点: OCB 瓶颈. 一条线, 非精确控制

Scoreboard 显式寄存器重命名: 分配新的物理寄存器, translation table

发射: IS: 要合的一个 rd. 若有 free rd, 不 IS

Branch Prediction: 1-bit, 2-bit, Correlating, Tournament, Branch Target Buffer, Integrated Inst. fetch Units, Return Addr Predictors

Correlating: (m, n) 预测器: 最近 m 个行为. 在 2^n 个位预测器中选. global / local

Tournament: 用地址索引局部预测器和选择器. 用分支历史索引全局预测器

gshare: 全局分支历史与分支地址异或做索引

Integrated: branch prediction, inst. prefetch, inst. memory access and buffering

Return Addr: 用栈寄存器 ra

Speculation: 按机面顺序执行, 按序提交. Tomasulo 加 ROB. 取指 store buffer, ROB 换名

ROB Entry: op, rd, value ready, exception vector

Commit: 就绪: 刷新 ROB 重新开始. 其他: 更新 reg/mem, 移除 ROB. 会像按序 issue (RS 和 ROB 保持顺序)

	Issue Structure	Hazard Detection	Scheduling	Characteristic
Superscalar (static)	Dynamic	HW	Static	顺序执行
Superscalar (dynamic)	Dynamic	HW	Dynamic	乱序
Superscalar (speculative)	Dynamic	HW	Dynamic	乱序+投机
VLIW/VLW	Static	SW	Static	packet 内无 hazard
EPIC	mostly static	mostly SW	mostly static	明显 hazard 被 compiler 标注

静态多发射: 对指令序列提前, 按序提交. 需要寄存器或通路

VLIW: 空槽 (ILP 开发). 二进制兼容

# Chapter 4 DLP: Vector Processor, GPU.

Vector Processors: memory-memory (通过寄存器, 检查内存地址的依赖. 启动延迟更高, 内存带宽要求更高).

vector-register.

向量处理器 Pos: 向量内存依赖 SIMD. Regular Mem access pattern.

Fewer branches  
Cons: 并行少, 性能较差; 内存带宽瓶颈

三种寻址方式: unit stride (带状), Non-unit (constant) stride.

Indexed (gather-scatter).

向量长度 = MVL, Strip mining: 先取  $N \bmod MVL$  个元素.

Vector chaining: 有 RAW 的依赖. Convoy  $\rightarrow$  Chime  $\rightarrow$  chain.

opt2: Conditional Exec / Masked Vector Inst.

opt3: Sparse Matrix: Indexed load (gather), Indexed store (scatter).

opt4: Multi-Lane.

GPU: Threads are organized into blocks, blocks are organized in a grid.

GPU handles thread management.

与向量处理器: 无寄存器 processor, 用 multi-threading 来降低 memory latency.

All GPU loads are gather inst. All GPU stores are scatter inst.

# Chapter 5 TLP: Multiprocessors.

Flynn 分类: SISD, SIMP, MISD, (MIMD)

硬件下的内存: Shared / Distributed.

软件下的内存: Shared Memory / Message bypassing

Centralized shared mem (UMA/SMP).

Distributed memory { Distributed shared mem (DSM/NUMA).

multiple computer: msg bypassing

UMA: uniform mem access. Small multiprocessor only.

NUMA: more scalable. 种类更多.

不同处理器访问的内存 可不同. 互斥不互斥.

Programming Model: multiprogramming: 无交流; shared addr space (通过内存共享), message passing (消息传递). data parallel (不同处理器同时处理, 每个处理器接收信息).

Communication Abstraction: shared addr space (ld/st/atomic swap). msg passing (send / library calls / receive)

$$CPI = \text{Base CPI} + \text{Remote Request Rate} \times \text{Remote Request Cost}.$$

Cache Coherence

(-) Snooping Protocol.

Write Invalid Protocol

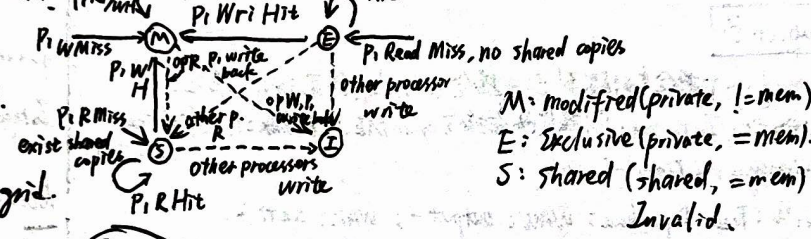
Write Broadcast Protocol.

并行读写策略: BWS.

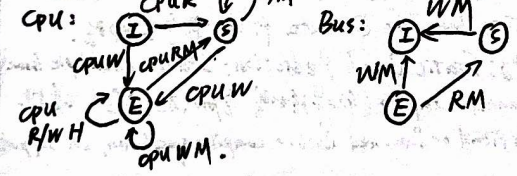
P1 先写 X 后读 X, return 新值  
P1 先写 X, P2 后读, ...  
P1 先写, P2 后写, 后读.

$V=0, D=x \rightarrow$  Invalid  
 $V=1, D=0 \rightarrow$  Shared (not dirty)  
 $V=1, D=1 \rightarrow$  Exclusive dirty.

MESI



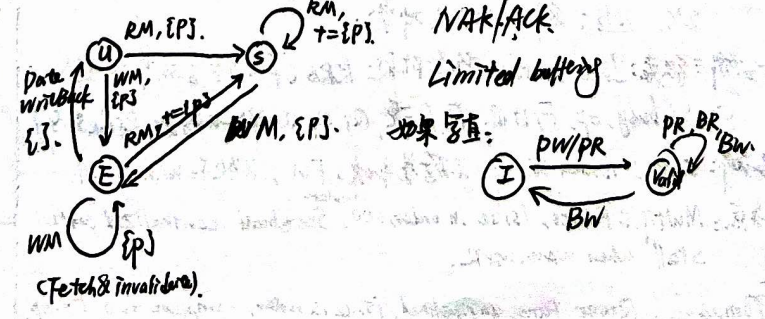
Snoopy



(=) Directory-Based.

Block Status: Shared, Uncached, Exclusive

Local node / Home node / Remote node



同步: Atomic Exchange / Test-and-Set / Fetch and Increment

```

try: mv x3, x4
    lr x2, x1
    sc x3, 0(x1)
    bnez x3, try
    mv x4, x2.
    
```

Barrier 3000  $\rightarrow$  Local sense  
if (count == total) { count++; release = 10-5;

Models of Memory Consistency.

- Sequential Consistency: in order; interleaved. (bad performance).
- Relaxed. 共享变量更新. 其他操作.

$\rightarrow W \rightarrow E$   
 $\rightarrow W \rightarrow W$