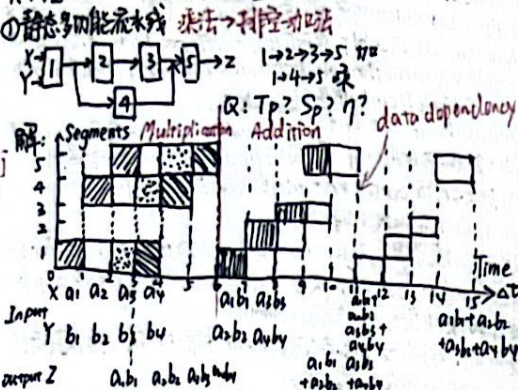


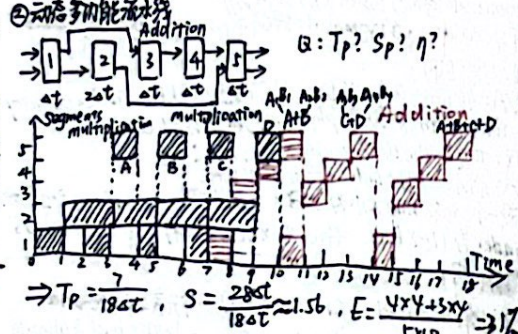
公式

- ① $CPU\ Time = \frac{Instruction\ Count \times CPI}{Clock\ Rate}$
- ② Amdahl's law: $T_{improved} = \frac{T_{unaffected}}{improvement\ factor} + T_{unaffected}$
- ③ $S_{p\ overall} = \frac{T_{old}}{T_{new}} = \frac{f}{(1-f) + \frac{f}{S_{p\ optimized}}}$
- ④ 吞吐率 $Throughput = \frac{指令数量}{时间}$
- ⑤ 加速比 $Speedup = \frac{1 \times m}{m+n-1}$
- ⑥ 效率 $Efficiency = \frac{m}{m+n-1}$
- ⑦ 内存性能度量
- ⑧ 从CPU的执行和 stall 角度分类讨论: $CPU\ time = IC \times (CPI_{execute} + \frac{Mem\ Access}{Inst} \times Miss\ Rate \times Miss\ Pen\ alty) \times Cycle\ Time$
- ⑨ 从CPU的指令类型 (ALU指令还是Memory指令) 分类讨论
- AMAT = $\frac{Whole\ access\ time}{All\ memory\ accesses\ in\ program}$
- $CPU\ time = (ALU\ Ops \times CPI_{ALU} + Mem\ Access \times AMAT) \times Cycle\ time$

1. 向量 A (a1, a2, a3, a4), 向量 B (b1, b2, b3, b4), 算 A.B.



$T_p = \frac{7}{15} \text{ sat}, S_p = \frac{12+8+4}{15} = 1.6, \eta = 32\%$

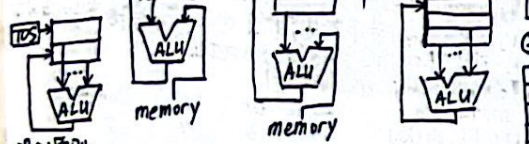


$T_p = 18 \text{ sat}, S = \frac{28 \text{ sat}}{18 \text{ sat}} \approx 1.56, E = \frac{4 \times 4 + 5 \times 4}{5 \times 18} \Rightarrow 1\%$

ISA boundary between software and hardware

GPR classification: W, A, B, R

Stack	Accumulator	Register RM	Register RR
Push A	Load A	Load R1, A	Load R1, A
Push B	Add B	Add R1, B	Load R2, B
Add	Store C	Store C, R1	Add R3, R1, R2
Pop C			Store C, R3

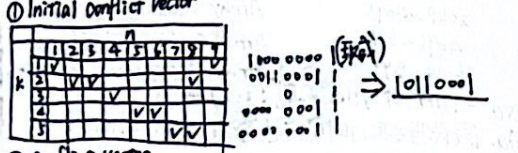


- ① Simplicity favors regularity
- ② Smaller is faster
- ③ make the common case fast
- ④ Good design demand good compromises

Pipeline

Instruction Dependences	Pipeline Hazards
- Data Dependences	- Data Hazards
- Name Dependences	- RAW
- Anti-dependence	- WAR
- Output-dependence	- WAW
- Control Dependences	- Branch Hazards
	- Structural Hazards

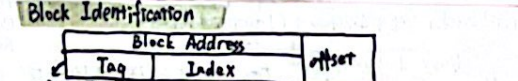
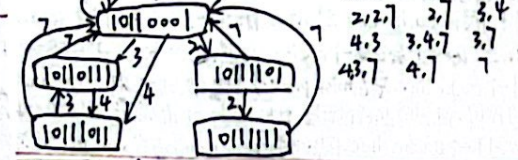
2. 非并行流水线的扩展



② Conflict vector

Interval	Initial	1	2	3	4	5	6	7
1-2	1011001	0010110	0000111	0000000	0000000	0000000	0000000	0000000
2-3		1011001	0010110	0000111	0000000	0000000	0000000	0000000
3-4			1011001	0010110	0000111	0000000	0000000	0000000
4-5				1011001	0010110	0000111	0000000	0000000

③ State transition graph



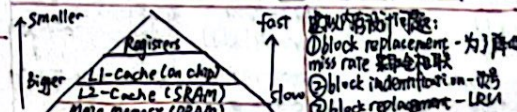
性能指标

- 1. Clock Rate = 200MHz, Ideal CPI=1, 50% arith/logic, 30% id/st, 20% control, 10% of memory operations get 50 cycle miss penalty, 1% of instructions same
- Q: CPU time and the AMAT?
- A: $CPI = ideal\ CPI + average\ stalls\ per\ instruction = 1.1 + 0.3 \times 0.1 \times 50 + 1 \times 0.01 \times 50 = 3.1$
- $AMAT = \frac{1}{1.3} \times (1 + 0.01 \times 50) + \frac{0.1}{1.3} \times (10.01 \times 50) = 2.54\ cycles$
- 2. Ideal CPI=1, 50% data access, miss penalty: 25 cycles, Miss rate: 2%, Q: How faster if all inst cache hits?
- A: Memory stall cycles = $IC \times \frac{Memory\ accesses}{Instruction} \times Miss\ rate \times Miss\ penalty = IC \times 0.5 \times 0.02 \times 25 = 10 \times 0.5$
- $\frac{IC \times 0.5 + IC}{IC} = 1.75\ faster$
- 3. Ideal CPI=2, clock cycle Time = 1ns, MPI=1.5, cache size = 4KB, block size = 64 bytes.
- ① direct mapped miss rate = 1.4%, ② two-way set associative miss rate = 1.0%, cycle time = 1st x cycle time, miss penalty: 75ns, hit time: 1 cycle
- Q: 两种情况的AMAT和CPU time.
- A: AMAT: ① $1 + 0.014 \times 75 = 2.05ns$, ② $1 \times 1.25 + 0.017 \times 75 = 2ms$
- CPU time: ① $IC \times (2 \times 1.0 + (1.5 \times 0.014 \times 75)) = 3.58 \times IC$, ② $IC \times (2 \times 1.0 + (1.5 \times 0.017 \times 75)) = 3.63 \times IC$
- 用公式: $IC \times [CPI_{execute} \times Clock\ Cycle\ time + (Miss\ rate \times \frac{Memory\ accesses}{Instruction} \times Miss\ penalty \times clock\ cycle\ time)]$

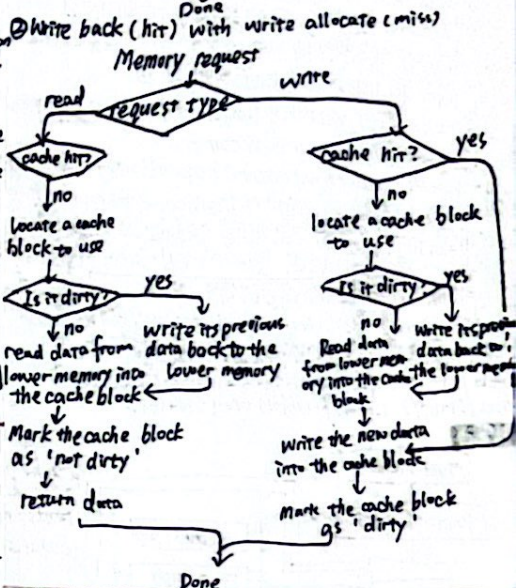
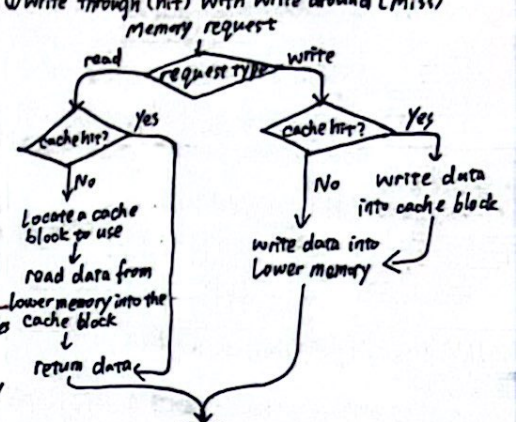
提高AMAT的优化

- 1. 减少miss penalty: multilevel caches, critical word first, read miss before write miss, merging write buffers, victim caches.
- 2. 减少miss rate: larger block size, larger cache size, higher associativity, way predication, pseudo-associativity, compiler optimizations
- 3. 减少hit time: small and simple caches, avoiding address translation, pipelined cache access, trace caches.
- 4. 通过并行减少miss rate和miss penalty, non-blocking caches, hardware prefetching, compiler prefetching.

Stack replacement algorithm



Write Strategy



TLB是缓存的缓存, 加速的是物理地址的对比, 地址字长

