

Final Notes

1. Details on PPT

Classes of computers
 personal computers: general purpose, cost/performance tradeoff
 server computers: network based, high capacity, performance
 supercomputers: highest capability
 embedded computers: stringent power/cost constraints

Common number system
 binary
 decimal
 hexadecimal



Components of the computer

Input (data transferred to memory)
 Output (data read from memory)

Memory

Processor
 native data types
 instructions
 registers
 addressing modes
 memory architecture
 interrupt and exception handling
 external I/O

ISA
 native data types
 instructions
 registers
 addressing modes
 memory architecture
 interrupt and exception handling
 external I/O
 → 一种叫 RISC 的, MIPS 和 RISC
 较简单, 轻量级的

2. Digital Logic

{ digital waveforms: discrete (e.g.: > certain threshold →)
 analog waveforms: continuous

logic
 truth table
 logical expression
 Graphic form (k-map)

Universal Gate (单一Gate可搭建所有电路): AND 和 NOR

Gate: NOT → AND → OR → INAND → NOR → XOR →

{ combinational logic circuits: no memory → output can be expressed as logic functions of input

sequential logic circuits: with memory! cannot

Combinational logic

① Multiplexer (selector)

2^n to 1 multiplexor
 n^+ selection input (which to choose) → 这个 n 位输入对应二进制数就是选择输入
 1⁺ output (相当于先在前面的电路里用 not gate 做 n 位 selector
 input 转成 2ⁿ 相位, 然后每相位和一个 input - 起进 and gate
 也就是说必须把指定相位, input 才能连上)

② Decoder (n to 2^n decoder)

n to 2^n decoder
 n input
 2^n outputs (每个 input 选择 In₁ 还是 not In₁, 每种组合对一个输出)
 one-hot vector
 ③ minterm and maxterm
 canonical form SOP (min): 所有结果为 1 的 input 算起来之后的和 Sum (min)
 POS (max): 所有结果为 0 的 input 算起来的和 Prod (max)

④ two-level logic and PLA

⑤ K-map

要注意的是: 表中 entry 都是 minterm
 简化方法可能唯一
 是上一行和下一行, 最左边一行和最右边一行也是 adjacent 的
 相邻的 2ⁿ 才能被圈

Sequential logic

① S-R Latch

 NOR Gate implementation:
 $R = 1, S = 1 \Rightarrow \text{forbidden state} \rightarrow 1 - 0 = 0 - 0 = 0_0 (X)$
 $R = 0, S = 0 \Rightarrow \text{Latch}$
 $R = 0, S = 1 \Rightarrow \text{Set}$
 $R = 1, S = 0 \Rightarrow \text{Reset}$
 NAND Gate implementation:
 $S = 1 \text{ AND } Q = 0 \Rightarrow \text{forbidden state} \rightarrow 0 - 1 = 1 - 0 = 0_0 (X)$
 $R = 0, S = 1 \Rightarrow \text{Latch}$
 $R = 0, S = 0 \Rightarrow \text{Reset}$
 $R = 1, S = 0 \Rightarrow \text{Set}$

② D-Latch

 $WE = 1: Q = D$
 $WE = 0: Q = \text{previous value}$

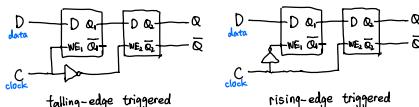
③ Register

a collection of D-latches, controlled by a common WE

④ clock signal

a shared signal that controls all the circuit element about time
 clock speed / clock frequency × cycle time =
 edge-triggering terminology (state change only occur in edge time)
 ↳ flip-flops (latch can do state change any time)

⑤ Master-Slave D Flip-flop
 built by a pair of D-latches



3. Data representation

2's complement

① signed integer

$$(6)_10 = (0110)_2 \xrightarrow{\text{1 bit}} 010 \xrightarrow{\text{1 bit}} 100 \xrightarrow{\text{1 bit}} 1010_10$$

$$(-6)_10 = (1010)_2 \xrightarrow{\text{1 bit}} 100 \xrightarrow{\text{1 bit}} 010 \xrightarrow{\text{1 bit}} 0110_10$$

∴ 对 k-bit 来说, 能表示的是 $[-2^{k-1}, 2^{k-1}]$

Overflow & Underflow



(unsigned int 0 ~ $2^{31}-1$, 逻辑上是 trivial - 1.)

(signed extension: for signed ints
 zero extension: for unsigned ints
 (more in MIPS))

② floating point numbers

finite range, limited precision

(能表示的数不是 equally distributed 的, 前密后疏)

$$\frac{1}{2} \frac{8}{32} \frac{23}{52} \quad (32 \text{ bits, single precision}) \quad 7 位精度$$

$$\frac{1}{2} \frac{11}{52} \quad (64 \text{ bits, double precision}) \quad 16 位精度$$

IEEE 754 Standard

bias 目的: 更容易做 comparison (直接 bit-by-bit 比较)

1° normalized case

拿到一个数, 先化为 $\pm 1 \times 2^{-n} \times \dots$ 的形式
 整数部分: 小数点前部分拆了再合起来
 小数化二进制用“归整整数部分”法

然后把二进制表示写出来, 化为 normalized scientific representation

exponent 是原数的 exp + bias 之后转为二进制源码的值

最后加 sign bit, 0 for 正, 1 for 负

注意: exponent 在这个 case 里不能 = 000...0 或 111...

假设是 1 位的偏移:
 representable range:
 $\begin{aligned} \text{max} &= (-1)^s \times (1.1 \dots 1) \times 2^{(11-10)-127} = -1.1 \dots 1 \times 2^{128} \approx -2^{128} \\ \text{min} &= (-1)^s \times 1.00 \dots 0 \times 2^{(10-10)-127} = 2^{128} \\ \text{min} &= (-1)^s \times 1.00 \dots 0 \times 2^{(10-10)-127} = 2^{128} \\ \text{max} &= (-1)^s \times 1.1 \dots 1 \times 2^{(11-10)-127} = 1.1 \dots 1 \times 2^{128} \approx 2^{128} \end{aligned}$

underflow & overflow (都是 exponent 意义上讲的, significand 控制的更多是精度)

underflow: $-2^{128} \approx -2^{127} \approx 2^{127} \approx 2^{128} + 2^{127}$

说明 significand 有错吗!!!

2° de-normalized case

exponent = 00...0 = the exponent of $(00 \dots 0)_2 = 1 - 127 = -126$

significand 直接读, 但没有 default 1

bit 位数看是 0 000...0 01100...0

value = $(-1)^s \times 2^{-126} \times (0.01100 \dots 0)_2$

if exponent = significand 全是 0, 那就表示 0

3° Special case

exponent = 11...1 significand 是全 0: infinite number, 正负看 sign bit

significand 不是全 0: NaN

一般习惯用 hexdecimal 的时候在前面加“0x”

④ characters ⇒ unsigned bytes!

usually 1 byte (8-bit), follow ASCII standard

4. MIPS

Doing abstraction makes life easier!

Ideally for each ISA there is an assembly language (satisfying/following the ISA standard)

但其实对于一个 ISA 来说, 有九种! 在细微差异, 也挺常见

Assembler 把它转为 machine code

Basic Knowledge

① Register: inside processor, 处理速度 (transfer速度) 比内存快很多
 in MIPS: 32 registers [number tradeoff]

data 从 memory 里拿出来, 到 processor 里才能做操作

② data transfer [lw (load word)]

sw (store word)

memory location: Offset (base address)

↳ must be constant

③ Endianness

data 之间总体都是从地址到地址的

寻址读取永远先读 memory value 小字节, 但写回是

0x00 0x01 0x34 0x12

但都是先读 0x00 再读 0x01, 再根据 endianness 来决定怎么“interpret”数据

④ Immediate operands

(addi, no addi!)

constant zero: \$zero (read-only)

↳ 但当成 register 用 (add \$t2,\$t1,\$zero)

addi \$t2,\$zero,5

Memory register constant → 快

⑤ Logical operation [bit-by-bit, 每个位分别做 logical operation 再合起来]

and, or, nor, andi, ori (with an immediate value)

shift [sll 整体左移/右移 k 位 e.g.: sll destination, source, shift_amount]

得是 constant 它是 register

MIPS Program

① Areas in memory: Reserved, text, data, heap, stack
 小地址 → 大地址

当我们在 data 中 define variable 的时候, 占用的内存都是连续的

→ +4 byte

hi, word 1234 ⇒ 12 34 5 = 可通过 12 + \$t0, 5 来访问 h[3] = 4

hi, word 5 = 5 来访问 h[4] = 5

declare float: .byte (1 byte →)

.word (4 byte →)

.double (8 byte →)

.ascii str (string)

.asciz str (string, 末尾 + '0')

② Making decisions in MIPS

Conditional jump

beq: branch if equal e.g.: beq \$r1, \$r2, L1 go to label L1 if \$r1==\$r2	bne: branch if not equal e.g.: bne \$r1, \$r2, L2 go to label L2 if \$r1!=\$r2
--	--

Unconditional jump

j: jump e.g.: jump L1	jr: jump register e.g.: jr reg [reg存着address] 通常在procedure call/case语句里用
--------------------------	--

这个会起来implement if-else-if-else语句

③ Loops

Basic blocks 就是一串 sequentially 执行的代码块 (不能从中间进去/出去)

One of the first early phases of compilation is breaking the program into basic blocks.
A compiler identifies basic blocks for optimization
An advanced processor can accelerate execution of basic blocks

④ Comparison

slt: set on less than
e.g.: slt \$r1, \$r2, \$r3
\$r2<\$r3: \$r1=1
otherwise: \$r1=0

slti: set on less than immediate
e.g.: slti \$r1, \$s2, 10
\$s2=10: \$r1=1 if \$s2<10
otherwise: \$r1=0

如果重做 if-else jump 操作, 只能组合两条 instruction
如果把 comparison 和 jump 统一到一起, 那单条 instruction 就会变慢 [slower clock]
而且 beq, bne 才是 common case

→ pseudo instructions

$blt(z)$	$bgt(z)$ branch if $\{l(less) + t(than) + e(equal) + z(zero, \text{fromc})\}$
$bge(z)$	$ble(z)$

⑤ Representing instructions

All 32 bits, 4 bytes, 3 formats

R-format (register)	I-format (immediate)	J-format (Jump)
I = rs, rt, rd, shamt, funct	J = opcode, address	decoding 需要跟着 opcode 和功能码看哪个 format!

R: 注意结果(也就是写 code 时候的第一个)是存在 R 里的
rs, rt 的顺序和 code 里一样
shamt: shift amount
funct: 当 opcode 为 (00000) 时, 有很多种指令, 这时候就靠 function code 区分

I: rs: 在着 value register
rt: 在着 address bit register
sw: rs → rt
immediate: $\{ \text{constant} + \text{offset} \}$

J:

⑥ MIPS Arithmetic

1° instruction details about signed/unsigned adding & extension

add, addu, addi 都是直接加, signed 和 unsigned 的加法只是涉及到 signed 的要检查 overflow 而 unsigned 的不用
addu 和 addi 不说明 instruction 中的 register/immediate 中的值是 unsigned 的, 及说明加法之后不检查是否是 overflow

addi 是 sign extension, 因为 pressence value

slt 表示 unsigned comparison (也就是说溢出来 from 值)
无元表示 signed comparison (也就是说 2's complement 表现每一位)
slti 表示 sign extension, 都是 sign extension (slti 表示的是 sign extension, 再 unsignedly 做去比较)

→ sign extension 再跟 \$rs 加
{lb \$rt, imm(\$rs)} grab一个 byte 出来, 再 sign extension 有进位吗
{lbu \$rt, imm(\$rs)} sign extension 再跟 \$rs 加, grab一个 byte 出来, 再 sign extension 有进位吗
→ load half 同理, 取 low, 没有 imm

andi \$rt, \$rs, imm zero extension

2° ALU: arithmetic logic unit

hardware in charge of arithmetic operations
logical operations

1-bit full adder

用 truth table 写出 Cout, Sum
 $Cout = ab+bc+a$ (用 k-map 简化过)
 $Sum = abc+abc+a\bar{b}\bar{c}$

1-bit ALU (and, or, addition)

Carryin control signal (应该是 opcode 变成的)
multiplexer (4 to 1)
L2 to 1 不够用
所以 (do) 实现了 or operation
实际上做了 +, and, or 或者 or
操作先把 b 取反 -b

→ 假如想实现 subtract 功能的话可以在里面这个
multiplexer (4 to 1) 然后把 carryin 设为 1 → 下面有
操作先把 b 取反 -b

1-bit ALU (and, or, addition, subtraction, beg, slt)

最大位
opcode: 11
把 sign bit 的 sub 的结果放到
最小位的 loc 里, 然后其最高
位输出 0, 其余全部 0
最大的输出是 sub 的结果
最小位的值 (通过 sign bit)
SLT
→ 假如想实现 subtract 功能的话可以在里面这个
multiplexer (4 to 1) 然后把 carryin 设为 1 → 下面有
操作先把 b 取反 -b

→ beg: basic idea 是做 subtraction, 用两个 register 中的值相减看是否为 0
→ equal comparison (branch in next step)

4-bit ALU (and, or, +, -)

当 b = 0 时
b inverse
opcode
ALU₀ ALU₁ ALU₂ ALU₃ R₀ R₁ R₂ R₃

Ripple Structure

优点: 32 bit 的 data 只用了 32 个 1 bit ALU 在, 简便
缺点: 下一个 bit 需等上一个 bit 完成, 单线程, 慢
beg 实现: 每位 sub 完, 只要有一个不是 0, 就不 =

Carry look ahead (Optional)

实际的 ALU: 每四位之间用传统的办法, 但是四比特和四
bit 之间用 carry look ahead (8 个并行单元)

define { G_i = A_i · B_i
P_i = A_i + B_i
C_{i+1} = G_i + P_i · C_i
(要么该位全为 1, 要么其中一个为 1, 而且 C_i=1, C_{i+1} 才能是 1)

这个办法在用硬件找效率

⑦ Procedures

General purpose registers for procedure calling:
\$a0 ~ \$a3: arguments (reg 4 ~ 7)
\$v0, \$v1: result values (reg 2 and 3)
\$t0 ~ \$t9: temporary registers (reg 10 ~ 19)
\$s0 ~ \$s7: saved registers (reg 20 ~ 27)
\$gp: stack pointer (reg 29)
\$sp: frame pointer (reg 30)
\$ra: return address (reg 31)

PC: 在着当前在执行的 instruction 的地址
(全局寄存器, 在 instruction fetch 之后就会执行 PC=PC+4)

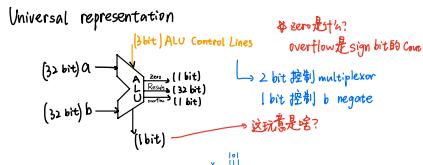
jal ProcedureLabel1 {当前 "jal" 指令的下一行的 address 被 assign 临时
(不是 pseudo) jump to ProcedureLabel1
往堆顶 \$ra 处于回到原处位置}

由于在 nested procedure 里, \$ra 会被 override, 引入 \$sp 来辅助
callc 开头先把 \$ra 的值 push 到 stack 里去; ldi \$sp, \$sp-4
然后在 procedure 结束之前 (jr \$ra, \$ra-4) 再把值从 stack 里 pop 出来 (lw \$ra, 0(\$sp))
(leaf procedure 不用做 push & pop 但是做肯定没意义)
→ 因为 \$ra 不会被 override

⑧ char & string

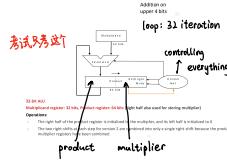
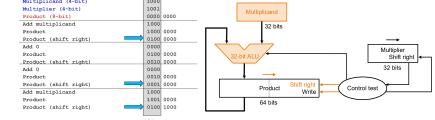
lb & sb 用来 load/store char, 通过 sign extension
pseudo instruction: la { lui: load upper immediate
ori }

la \$s1, int { lui \$s1, int (把高前 16 位给移位的前 16 位, then set the lower 16 bits of \$s1 to be 0)
ori \$s1, \$s1, int (再把后 16 位和 \$s1 的后 16 位, or 一起, \$s1 前 16 位不变)}

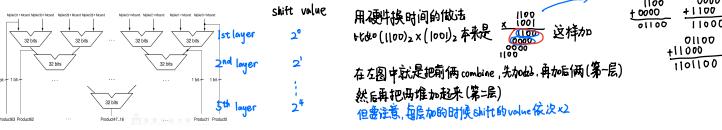


Unsigned multiplication

$A \times B \rightarrow$ multiplier
 $B \times A \rightarrow$ multiplicand
假设 A 有 m bit, 那么 x_B 至多 (m+n) bit



Hardware Speedup



Multiplication in MIPS

Hi. Lo is used to store the product
mul \$rs, \$rt
mult \$rs, \$rt
乘积可能超过32bit → 用两个特殊register存储

Fetch Results
mtlo (move from lo) mtlo \$f1
mthi (move from hi) mthi \$f2

Overflow testing (不能把整个结果看在一个32bit register里)

① Unsigned overflow

mulu \$d, \$a1
mtlo \$a1
mtlo \$a2
bneq \$at, \$0, no_overflow
j overflow
no_overflow:

② Signed overflow

mult \$d, \$a1
mtlo \$a1
mtlo \$a2
xra \$t0, \$a2, \$1
bneq \$at, \$0, no_overflow
j overflow
no_overflow:

Example

Division of a 4-bit unsigned number (0111) by another one (0011)

Iteration	Divisor (D)	Remainder (R)	Remark
0	0000	0011	Initial state
1	0011	1100	$R = D - 1$
2	0011	1110	$R = R - D = 0$
3	0011	0000	$R = R - D = 0$
4	0011	0000	余数为0, 不能再减了
extra			

initial state: 左半边是0, 右半边是 dividend

起步先 shift left 移!!

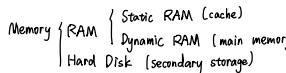
第1轮: 左半边减 divisor, 观察相减后的 sign bit
↓ 1: undo 相减操作并左移一格, 右边填0
↓ 0: 左移一格, 右边填1
↓ i: divisor 的 bit 数

↓ 结束时 shift right 一格!!

只 shift remainder
不 shift quotient

</div

⑭ Memory technology



SRAM Technology

- SRAM is a type of semiconductor memory
 - Each cell is constructed using transistors only (i.e. no capacitors)
 - Density not as high as DRAM, it is more power hungry
 - Does not need to be periodically refreshed; the memory retains its information as long as power is applied
 - Read operation discharges the electrons; read is destructive

Because of this refresh requirement, it is called **dynamic memory**

Hard Disk Calculation

Access to a sector involves

- Queueing delay if other accesses are pending
- Seek: move the heads
- Rotational latency
- Data transfer
- Controller overhead

$$\text{average} = \text{seek time} + \frac{\text{半圈的时间}}{2} + \text{data transfer time} + \text{controller delay}$$

Disk Access Example

Given:

- 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk

Average read time:

- 4ms seek time
 - + $\frac{1}{2} / (15,000/60) = 2\text{ms}$ rotational latency
 - + $512 / 100MB/s = 0.005\text{ms}$ transfer time
 - + 0.2ms controller delay
 - = 6.205ms

If actual average seek time is 1ms

- Average read time = 3.205ms

Locality

- {| temporal locality: 最近用过的很快还会再用
- spacial locality: 最近用过的边上的很快会用上

Memory hierarchy



Cache hit time:

= time to determine miss or hit + time to access the cache

Cache miss penalty:

= time to bring a block from lower level to upper level

Cache Performance

Average memory access latency

$$= \text{hit_time} + \text{miss_rate} * \text{miss penalty}$$

$$(\text{两层}) = \text{hit_time} + \text{miss_rate} * \frac{(\text{hit_time} + \text{miss_rate} * \text{miss penalty})}{\text{miss penalty}}$$

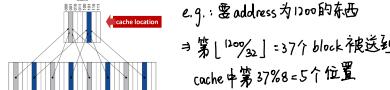
Cache Operations

Cache Organization

↳ array of cache blocks 相当于最小单位 (-般32/64 bits)

- direct-map
- set-associative
- fully-associative

direct-map 一般取模 (缺点: frequent move in & kick out)



标签: 语音该 block 在哪 8↑ block 里



Byte address 1200: 0000 0000 0000 0000 0000 0100 1011 0000

Block address

Byte offset

Valid bit: mark着该 block 是否有有效数字
(initialize 的时候都是垃圾值, valid bit 都是 0)

Bits in Cache

- A direct-mapped cache with 16KB of data and 8-word (2 bytes) blocks
- Assume 32-bit memory address
- What is the actual cache size?

16KB = 4K words = 2^{12} words
Block size of 8 words $\rightarrow 2^3$ blocks

Each block has 8 × 32 = 256 bits of data
A tag has 32 - 9 - 8 = 18 bits
Add 1 for valid bit
Total bits per entry = $(256 + 18 + 1)$

Total cache size: $2^3 \times (256 + 18 + 1) = 2^3 \times 275 = 140800 = 137.5$ Kbytes

Cache miss

C7BC >

Associate Cache

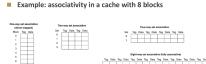
Fully-associative

每个 block 都能放在 cache 里任意位置
(缺点: 找空位麻烦) 像 open addressing

N-way set associative

每个 block 只能放在 cache 里指定的 n 个位置

Example: associativity in a cache with 8 blocks



increase N (associativity)
⇒ 缺处: miss rate ↓
环处: hit time ↑
一般用 2-way / 4-way map

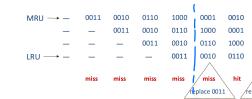
Block replacement

When cache misses:

{ direct map: 没有这个 concern
 n associative: LRU (Least recently used)

① replace LRU

② set it to MRU



- ① 0001 替换 0111
- ② 把 0001 设为 MRU
- ③ 把其他往下移
(0010 自动设为 LRU)

Data modification (dirty data processing)

{ write back: 数据修改后写回 cache, cache 回收 返回 memory 时将对 main memory 修改
 write through: 只要数据发生变动就写回 cache 和 main memory
(④群作, 慢, 需要 write buffer)

Virtual memory

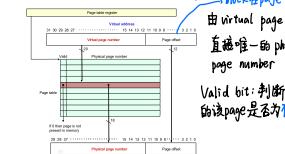
hard disk 把 main memory 当 cache 用

physical address = [F] virtual address

↳ page table
→ 一个大于 block 的单位

Page table (在 memory 里, 起始地址 stored in page table register)

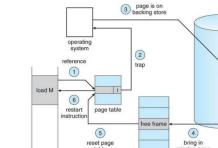
不用 tag 因为是一对一 mapping (每块), 不会有 2 个 virtual page
(cache 是 memory) map 到同一个 physical page ∴ 不用验证
block 在 memory 里的位置



Page fault

valid bit = 0 或 page table 中没有该 mapping ⇒ 不存在 valid 的 memory

∴ 必把 data 从 hard disk 里调出来 慢
流程:



程序会停下来等 page fault resolve

∴ page size 太大 浪费空间, 导致 page fragmentation
太小 总有 page fault

Fast address Translation

在 cache 里存一个 page table 和 copy ⇒ TLB
要不然每次 access 两次 memory (→ page table → 读 data)

Byte address 1200: 0000 0000 0000 0000 0000 0100 1011 0000

Block address

Byte offset

Valid bit: mark着该 block 是否有有效数字

(initialize 的时候都是垃圾值, valid bit 都是 0)

Bits in Cache

- A direct-mapped cache with 16KB of data and 8-word (2 bytes) blocks
- Assume 32-bit memory address
- What is the actual cache size?

16KB = 4K words = 2^{12} words
Block size of 8 words $\rightarrow 2^3$ blocks

Each block has 8 × 32 = 256 bits of data
A tag has 32 - 9 - 8 = 18 bits
Add 1 for valid bit
Total bits per entry = $(256 + 18 + 1)$

Total cache size: $2^3 \times (256 + 18 + 1) = 2^3 \times 275 = 140800 = 137.5$ Kbytes

Cache miss

C7BC >

Associate Cache

Fully-associative

每个 block 都能放在 cache 里任意位置
(缺点: 找空位麻烦) 像 open addressing

N-way set associative

每个 block 只能放在 cache 里指定的 n 个位置

Example: associativity in a cache with 8 blocks



