# Machine Prog (Assembly Code) ——Basics & Control section

王善上 贾博暄 倪嘉怡 许珈铭

2023.9.27

# Preface

许珈铭

# C/C++编译

- Pre-processor
- Compiler
- Assembler
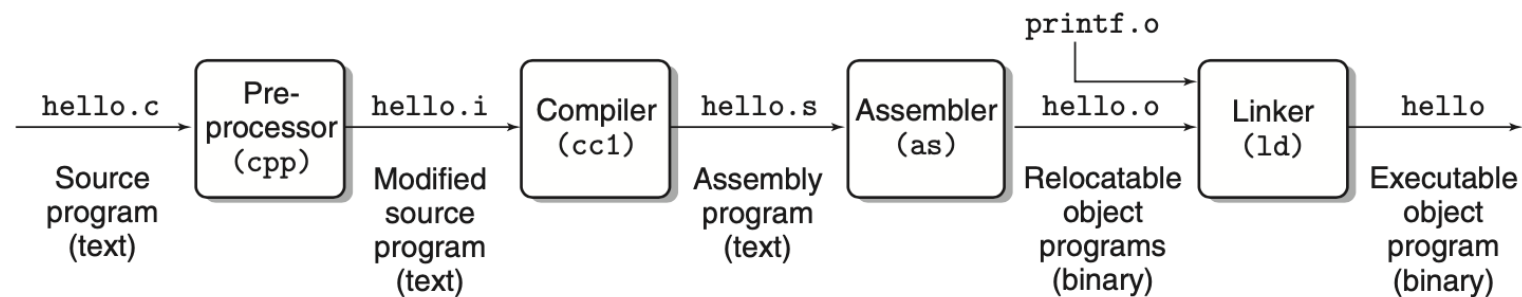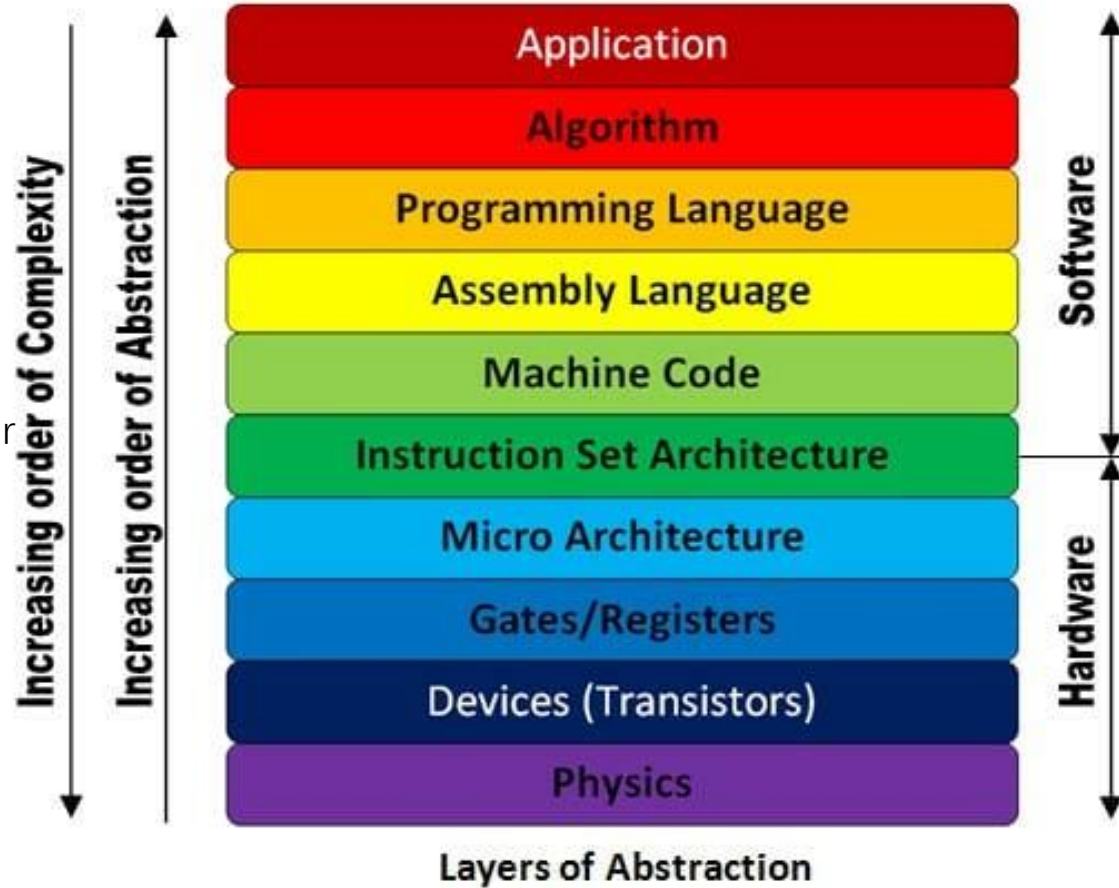- Linker&loader



**Figure 1.3   The compilation system.**

# Recall

- C/C++ (.c/.cpp, .h/.hpp)

GCC
↓ Compiler
- Assembly language (.s)

↓ Assembler
- Machine Code (.o, .so, .a)

↓ Linker

CLANG (LLVM)
↓ Compiler
- LLVM IR

↓ LLVM Optimizer
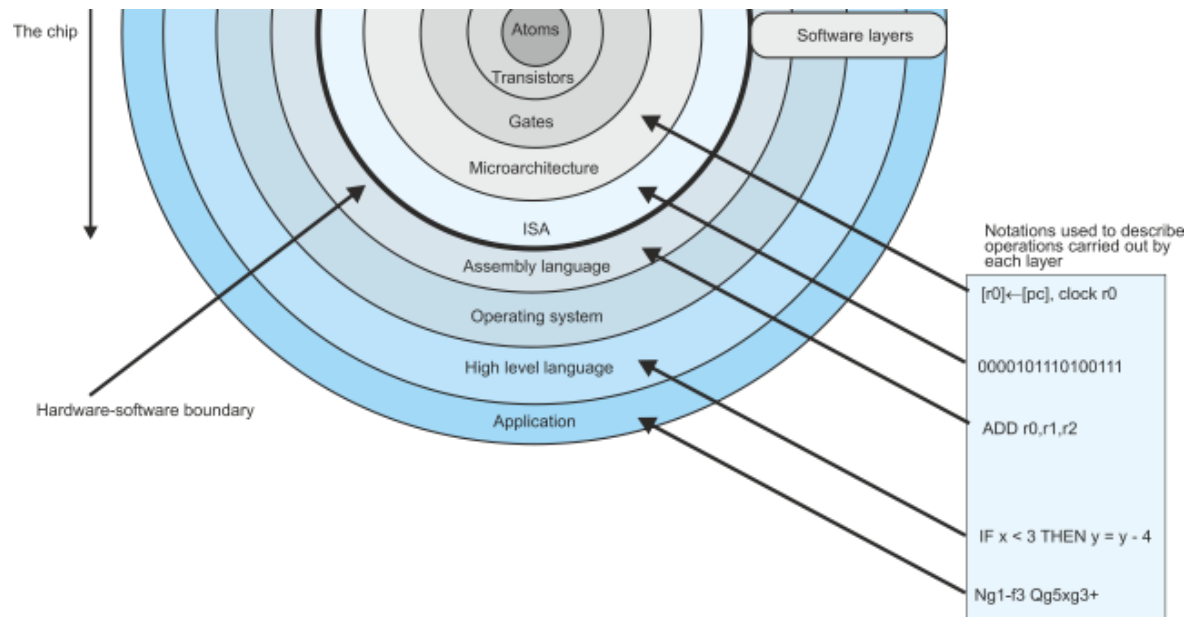- LLVM Backend

↓ Assembler

↓ Linker

- Java

JVM
↓ Compiler
- Java Bytecode

(Java Virtual Machine)
↓ Assembler

↓ Linker

- Machine Code

JIT
↓ Compiler
- Java Bytecode

↓ Just In Time compiler
- Machine code
- (Java Virtual Machine)

↓ Linker

- Machine Code

| | | |
|---|---|---|
| **Application** | | |
| **Algorithm** | | Software |
| **Programming Language** | | |
| **Assembly Language** | | |
| **Machine Code** | | |
| **Instruction Set Architecture** | | |
| **Micro Architecture** | | |
| **Gates/Registers** | | Hardware |
| **Devices (Transistors)** | | |
| **Physics** | | |

Increasing order of Complexity

Increasing order of Abstraction

**Layers of Abstraction**

# Two abstractions: ISA
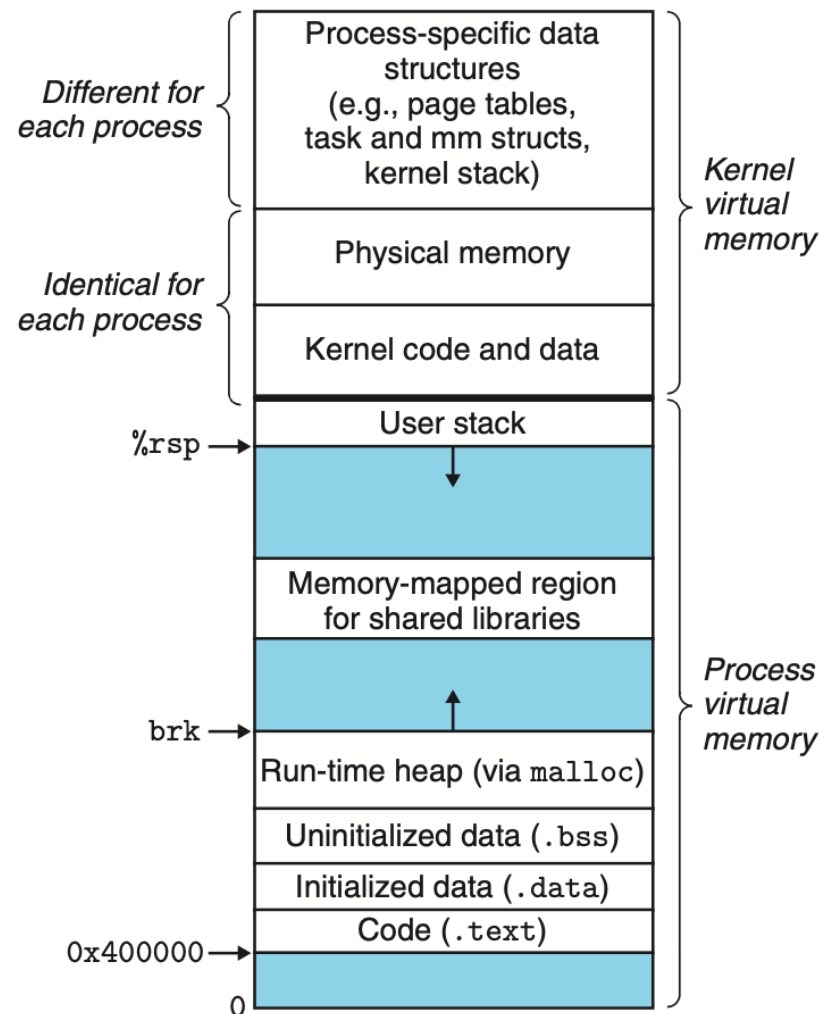
- 定义处理器状态、指令格式和指令对状态的影响
- The hardware/software **interface**
- Examples
  - IA32, AMD64 (x86-64), ARM64, RISC-V, etc.

- ISA ≠ 汇编语言

# Two abstractions: Virtual Address

- 每个 **进程** 享有同样的虚拟地址空间
- 操作系统的四大抽象之一，随着ICS的学习会逐渐补全拼图

# 工具链
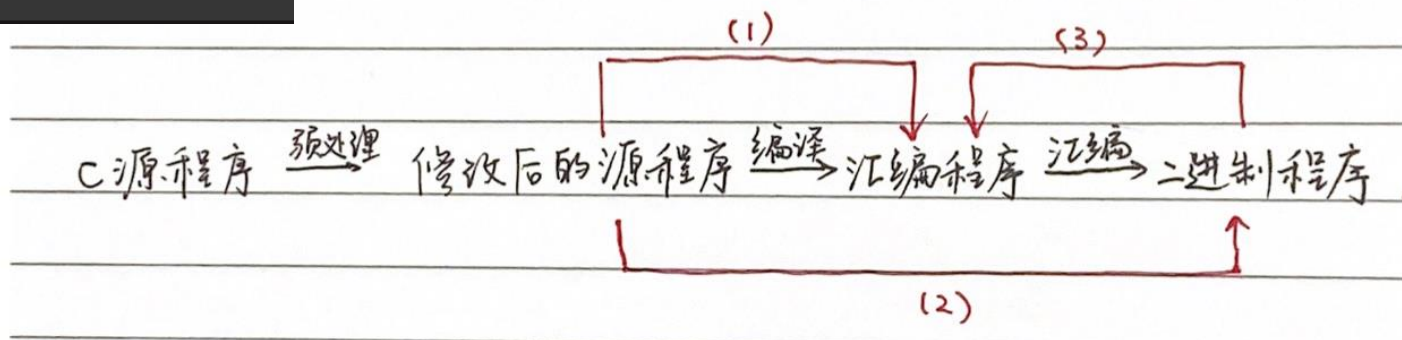
(1) gcc [-O] -S <file> [-o]

(2) gcc [-O] -c <file> [-o]

(3) objdump -d <file> (> ...)

·**gdb** and bomblab

(c->s) Compiler Explorer 

(o->c) ByteNinja 

# 编译器优化

助教 李翰禹

再具体解释一下优化等级的问题
-O0是几乎所有优化都不做
-O 或者-O1是做一部分优化
-Og是-O1优化去除会影响调试器正常工作的优化
-O2是-O1不考虑时间和空间的权衡的情况下做更多优化
-O3是-O2基础上做更多优化，比如一些向量化和循环展开的优化
-Os是-O2优化去除会增加目标代码长度的优化
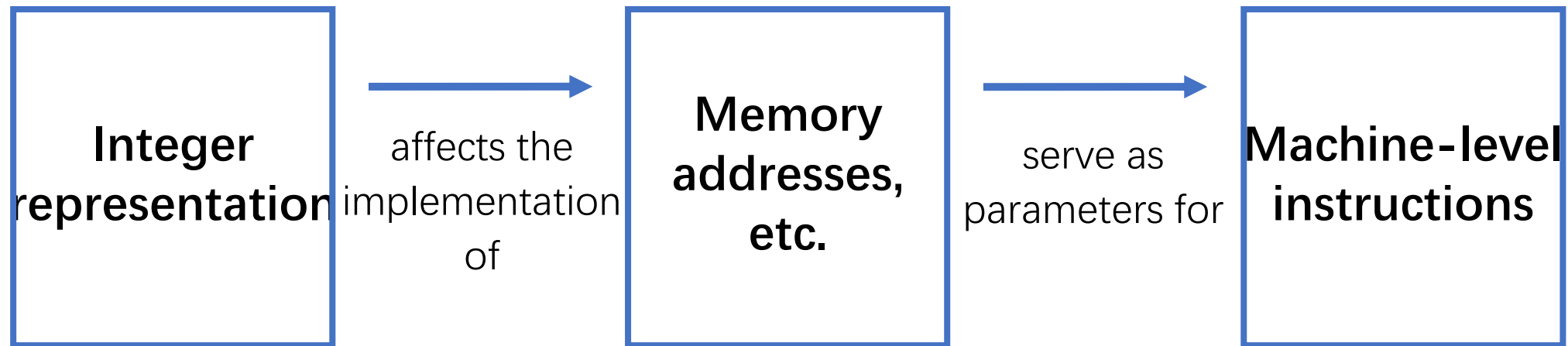-Ofast是无视严格标准的优化

# Basics (CS:APP Ch. 3.1 - 3.5)

贾博暄

# Relationships between topics

- (courtesy of TA)

# Registers

Note on naming
**b,w,l,q**

- Correspondence
- Memory reference in parens:  ( )
  - **%rax**     -> register
  - **(%rax)** -> memory
    - Compare: pointers in C

| C declaration | Intel data type | Assembly-code suffix | Size (bytes) |
|---|---|---|---|
| char | Byte | b | 1 |
| short | Word | w | 2 |
| int | Double word | l | 4 |
| long | Quad word | q | 8 |
| char * | Quad word | q | 8 |
| float | Single precision | s | 4 |
| double | Double precision | l | 8 |

**Figure 3.1  Sizes of C data types in x86-64.** With a 64-bit machine, pointers are 8 bytes long.

| 63 | 31 | 15 | 7 | 0 | |
|---|---|---|---|---|---|
| %rax | %eax | %ax | %al | | Return value |
| %rbx | %ebx | %bx | %bl | | Callee saved |
| %rcx | %ecx | %cx | %cl | | 4th argument |
| %rdx | %edx | %dx | %dl | | 3rd argument |
| %rsi | %esi | %si | %sil | | 2nd argument |
| %rdi | %edi | %di | %dil | | 1st argument |
| %rbp | %ebp | %bp | %bpl | | Callee saved |
| %rsp | %esp | %sp | %spl | | Stack pointer |
| %r8 | %r8d | %r8w | %r8b | | 5th argument |
| %r9 | %r9d | %r9w | %r9b | | 6th argument |
| %r10 | %r10d | %r10w | %r10b | | Caller saved |
| %r11 | %r11d | %r11w | %r11b | | Caller saved |
| %r12 | %r12d | %r12w | %r12b | | Callee saved |
| %r13 | %r13d | %r13w | %r13b | | Callee saved |
| %r14 | %r14d | %r14w | %r14b | | Callee saved |
| %r15 | %r15d | %r15w | %r15b | | Callee saved |

# Conventions...

- Instructions that generate 4-byte quantities for a register set the upper 4 bytes of the register to zero

```
1    movabsq $0x0011223344556677, %rax    %rax = 0011223344556677
2    movb    $-1, %al                     %rax = 00112233445566FF
3    movw    $-1, %ax                     %rax = 001122334455FFFF
4    movl    $-1, %eax                    %rax = 00000000FFFFFFFF
5    movq    $-1, %rax                    %rax = FFFFFFFFFFFFFFFF
```

- Consequences
  - **MOV** instructions get a little weird (e.g. there is no **movzlq**)
  - You might see (for example):
  - **movzbl** when your intuition says **movzbq**
      Textbook: Problem 3.4
  - Writing to **%eax** on one line and quoting **%rax** on the next

# Conventions...

- In a memory reference, the scaling factor must be either _1, 2, 4, or 8_

$$Imm(rb, r_i, s)$$

**9(%rax,%rdx,4)**

- Why?
  - Useful when referencing array and structure elements

1、某 C 语言程序中对数组变量 a 的声明为"int a[10][10];"，有如下一段代码：

```
for (i=0; i<10; i++)
    for (j=0; j<10; j++)
        sum+= a[i][j];
```

假设执行到"sum+= a[i][j];"时，sum 的值在%rax 中，a[i][0]所在的地址在%rdx 中，j 在%rsi 中，则"sum+= a[i][j];"所对应的指令是（      ）。

A. addl 0 (%rdx, %rsi, 4), %rax

B. addl 0 (%rsi, %rdx, 4) , %rax

C. addl 0 (%rdx, %rsi, 2) , %rax

D. addl 0 (%rsi, %rdx, 2) , %rax

# Conventions…

- x86-64 不允许将除 64 位寄存器以外的寄存器作为寻址模式基地址

1. 判断下列 x86-64 ATT 操作数格式是否合法。

(7) ( ) (%ecx)

- Textbook: Problem 3.3

Here is the code with explanations of the errors:

movb $0xF, (%ebx)        *Cannot use %ebx as address register*

# Conventions…

- Binary instructions' operands are, in order, _Source_ and _Destination_
  - Intuitive (eg. **MOV**)
  - Counter-intuitive (eg. **SUB**, **CMP**)

- the operand pairs CANNOT be of what type? _Any -> Imm, Mem -> Mem_
  - this applies to **MOV** and Arithmetic

# Instruction Classes – Move

- **MOV**
  - Operand types:    $Imm,\ Reg,\ Mem$
  - Operand pairs:    Some are forbidden
- Memory reference
  - Full form:

Offset (Base, Index, scale factor)

$$Imm + R[rb] + R[ri] \cdot s$$

  - Selected omission yields the other formats
  - $s$ limited to 1, 2, 4, 8
  - No $ before Imm here

| From\To | Imm | Reg | Mem |
|---------|-----|-----|-----|
| Imm | | √ | √ |
| Reg | | √ | √ |
| Mem | | √ | |

# Instruction Classes – Move

- **MOV: `movl` vs `movq` vs `movabsq`**
  - **`movl`**
    - Moves double word
    - Higher bits automatically set to 0
  - **`movq`**
    - Takes 32-bit value, sign-extends to 64-bit (i.e. quad word)
    - **THEN** moves to destination
  - **`movabsq`**
    - Takes 64-bit immediate value and moves to destination
    - **Only allows Imm -> Reg**

# Instruction Classes - Move

- "Corner cases"
  - **MOVZ**
    - there is no `movzlq`

  - **MOVS**
    - `cltq`
      - effectively a compact encoding of
      `movslq %eax, %rax`

MOVZ

movl
suffices

b w l q

MOVS

Special case:
cltq

b w l q

**z**: zero-extend
**s**: sign-extend

# Instruction Classes – Arithmetic

- **leaq**
  - (Important!) **lea** vs **mov**
  - Same format
  - **lea** DOES NOT reference memory; purely numerical

```
B.  leal (%ecx,%ebx,4), %eax
C.  movl (%ecx,%ebx,4), %eax
```

- Unary: **INC**, **DEC**, **NEG**, **NOT**
  - Somewhat self-evident

- Binary: **ADD**, **SUB**, **IMUL**, **XOR**, **OR**, **AND**
  - 2nd operand ($D$) = Destination
    - Mnemonic: Add       $S$ to      $D$
                Subtract $S$ from $D$

|      |       |                     |
|------|-------|---------------------|
| ADD  | $S, D$ | $D \leftarrow D + S$ |
| SUB  | $S, D$ | $D \leftarrow D - S$ |

# Instruction Classes – Shifts

- Shifts
  - **SAL(SHL)**, **SAR**, **SHR**
  - **sh**ifts (logical) vs **s**hifts (**a**rithmetic)
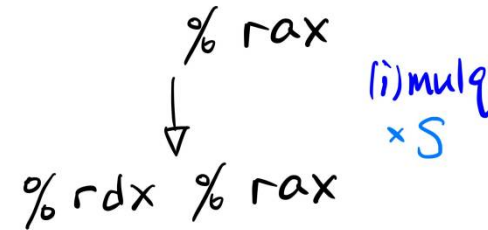
**sh**: logical shift

**sa**: arithmetic shift

- 1st operand ($k$) = shift amount
  - **must be** $Imm or %cl
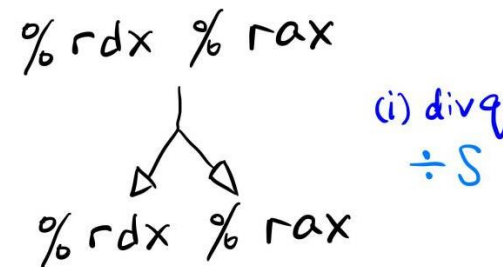  - actual shift amount (also) depends on type of value to be shifted

| SAL | $k, D$ | $D \leftarrow D << k$ |
|-----|--------|-----------------------|
| SHL | $k, D$ | $D \leftarrow D << k$ |
| SAR | $k, D$ | $D \leftarrow D >>_A k$ |
| SHR | $k, D$ | $D \leftarrow D >>_L k$ |

# Instruction Classes - 128-bit

- Single operand
- **imulq**, **mulq**
  - *S* is multiplicand
  - **imulq** *S* vs **IMUL** *S, D*
- **idivq**, **divq**
  - *S* is divisor
  - Quotient at **%rax**, remainder at **%rdx**
- **cqto**
  - Sign extend, presets **%rdx**
  - Compare: **cltq**

% rax
↓
% rdx % rax

(i)mulq
×S

**i**: integer

% rdx % rax
↓
% rdx % rax

(i) divq
÷S

**cltq**:
   convert **l** to **q**
**cqto**:
   convert **q** to **o**

On to the next section!

# Control (CS:APP Ch 3.6)

倪嘉怡

# Contents

- 条件码，跳转指令
- 逆向工程和goto代码
- **if-else, do-while, while, for, switch, …**
- 条件传送

# 条件码

- CF Carry Flag：进位；无符号溢出。1进位0不进
- ZF Zero Flag ： 零。1零0非零
- SF Sign Flag：符号。1负0正
- OF Overflow Flag：溢出；正/负有符号溢出。1溢出0不溢出

- 逻辑运算：进位 CF、溢出 OF 置 0
- 移位运算 ： 溢出 OF 置 0；进位 CF 设置为最后一个被移出的位
- inc, dec ： 不改变进位标志
- lea：不改变条件码

- cmp = sub -
- test = and &

- https://www.felixcloutier.com/x86/
- x86 and amd64 instruction reference

| 指令 | 效果 | 描述 |
|---|---|---|
| leal $S,D$ | $D \leftarrow \&S$ | 加载有效地址 |
| INC $D$ | $D \leftarrow D+1$ | 加1 |
| DEC $D$ | $D \leftarrow D-1$ | 减1 |
| NEG $D$ | $D \leftarrow -D$ | 取负 |
| NOT $D$ | $D \leftarrow \sim D$ | 取补 |
| ADD $S,D$ | $D \leftarrow D+S$ | 加 |
| SUB $S,D$ | $D \leftarrow D-S$ | 减 |
| IMUL $S,D$ | $D \leftarrow D*S$ | 乘 |
| XOR $S,D$ | $D \leftarrow D \,\hat{}\, S$ | 异或 |
| OR $S,D$ | $D \leftarrow D \mid S$ | 或 |
| AND $S,D$ | $D \leftarrow D \,\&\, S$ | 与 |
| SAL $k,D$ | $D \leftarrow D << k$ | 左移 |
| SHL $k,D$ | $D \leftarrow D << k$ | 左移（等同于SAL） |
| SAR $k,D$ | $D \leftarrow D >>_A k$ | 算术右移 |
| SHR $k,D$ | $D \leftarrow D >>_L k$ | 逻辑右移 |

| | CF | ZF | SF | OF | |
|---|---|---|---|---|---|
| LEA | ✗ | ✗ | ✗ | ✗ | |
| INC | ✗ | ✓ | ✓ | ✓ | INC, DEC |
| DEC | ✗ | ✓ | ✓ | ✓ | |
| NEG | ✓ | ✓ | ✓ | ✓ | |
| NOT | ✗ | ✗ | ✗ | ✗ | |
| ADD | ✓ | ✓ | ✓ | ✓ | 算术运算 |
| SUB | ✓ | ✓ | ✓ | ✓ | |
| IMUL | ✓ | ✓ | ✓ | ✓ | |
| XOR | 0 | ✓ | ✓ | 0 | 逻辑运算 |
| OR | 0 | ✓ | ✓ | 0 | |
| AND | 0 | ✓ | ✓ | 0 | |
| SAL | ✓ | ✓ | ✓ | 0 | 移位运算 |
| SHL | ✓ | ✓ | ✓ | 0 | |
| SAR | ✓ | ✓ | ✓ | 0 | |
| SHR | ✓ | ✓ | ✓ | 0 | |

# 跳转指令

%rax 寄存器
(%rax) 内存
寄存器+括号=内存

| | 指令 | 同义名 | 跳转条件 | 描述 |
|---|---|---|---|---|
| Jump | jmp Label | | 1 | 直接跳转 |
| | jmp *Operand | | 1 | 间接跳转 |
| if Equal | je Label | jz | ZF | 相等/零 |
| if Not Equal | jne Label | jnz | ~ZF | 不相等/非零 |
| if Sign | js Label | | SF | 负数 |
| if Not Sign | jns Label | | ~SF | 非负数 |
| if Greater | jg Label | jnle | ~(SF ^ OF) & ~ZF | 大于（有符号>) |
| if Not Greater | jge Label | jnl | ~(SF ^ OF) | 大于或等于（有符号>=) |
| if Less | jl Label | jnge | SF ^ OF | 小于（有符号<) |
| if Not Less | jle Label | jng | (SF ^ OF) \| ZF | 小于或等于（有符号<=) |
| if Above | ja Label | jnbe | ~CF & ~ZF | 超过（无符号>) |
| if Not Above | jae Label | jnb | ~CF | 超过或相等（无符号>=) |
| if Below | jb Label | jnae | CF | 低于（无符号<) |
| if Not Below | jbe Label | jna | CF \| ZF | 低于或相等（无符号<=) |

# 逆向工程和goto

- goto代码：描述汇编代码程序控制流的C程序
- 逆向工程：汇编代码->C语言代码

```
    x at %ebp+8, y at %ebp+12
1       movl      8(%ebp), %edx      Get x
2       movl      12(%ebp), %eax     Get y
3       cmpl      %eax, %edx         Compare x:y
4       jge       .L2                if >= goto x_ge_y
5       subl      %edx, %eax         Compute result = y-x
6       jmp       .L3                Goto done
7    .L2:                            x_ge_y:
8       subl      %eax, %edx         Compute result = x-y
9       movl      %edx, %eax         Set result as return value
10   .L3:                            done: Begin completion code
```

c). 产生的汇编代码

```
1  int absdiff(int x, int y) {
2      if (x < y)
3          return y - x;
4      else
5          return x - y;
6  }
```

a）原始的 C 语言代码

# goto：正向视角



a）原始的 C 语言代码



b）与之等价的 goto 版本

Basic：mov, add, sub…
专用于设置条件码：cmp, test…
跳转：jmp, je, jg…

如何把else转化为汇编?



c）产生的汇编代码

# goto：逆向视角

```
x at %ebp+8, y at %ebp+12
1       movl    8(%ebp), %edx
2       movl    12(%ebp), %eax
3       cmpl    %eax, %edx
4       jge     .L2
5       subl    %edx, %eax
6       jmp     .L3
7  .L2:
8       subl    %eax, %edx
9       movl    %edx, %eax
10 .L3:
```

c）产生的汇编

```
1  int gotodiff(int x, int y) {
2      int result;
3      if (x >= y)
4          goto x_ge_y;
5      result =  y - x;
6      goto done;
7  x_ge_y:
8      result = x - y;
9  done:
10     return result;
11 }
```

b）与之等价的 goto 版本

```
1  int absdiff(int x, int y) {
2      if (x < y)
3          return y - x;
4      else
5          return x - y;
6  }
```

a）原始的 C 语言代码

# 循环：do-while

```
do
    body-statement
while (test-expr);
```

```
loop:
    body-statement
    t = test-expr;
    if (t)
        goto loop;
```

```
1   int fact_do(int n)
2   {
3       int result = 1;
4       do {
5           result *= n;
6           n = n-1;
7       } while (n > 1);
8       return result;
9   }
```

a）C 代码

```
    Argument: n at %ebp+8
    Registers: n in %edx, result in %eax
1   movl    8(%ebp), %edx       Get n
2   movl    $1, %eax            Set result = 1
3   .L2:                        loop:
4   imull   %edx, %eax          Compute result *= n
5   subl    $1, %edx            Decrement n
6   cmpl    $1, %edx            Compare n:1
7   jg      .L2                 If >, goto loop
    Return result
```

c）对应的汇编代码

# 循环：while

```
while (test-expr)
    body-statement
```

```
    t = test-expr;
    if (!t)
        goto done;
loop:
    body-statement
    t = test-expr;
    if (t)
        goto loop;
done:
```

```
1   int fact_while(int n)
2   {
3       int result = 1;
4       while (n > 1) {
5           result *= n;
6           n = n-1;
7       }
8       return result;
9   }
```

a) C 代码

```
1   int fact_while_goto(int n)
2   {
3       int result = 1;
4       if (n <= 1)
5           goto done;
6   loop:
7       result *= n;
8       n = n-1;
9       if (n > 1)
10          goto loop;
11  done:
12      return result;
13  }
```

b) 等价的 goto 版本

```
Argument: n at %ebp+8
Registers: n in %edx, result in %e

1       movl    8(%ebp), %edx       Get n
2       movl    $1, %eax            Set result = 1
3       cmpl    $1, %edx            Compare n:1
4       jle     .L7                 If <=, goto done
5   .L10:                           loop:
6       imull   %edx, %eax          Compute result *= n
7       subl    $1, %edx            Decrement n
8       cmpl    $1, %edx            Compare n:1
9       jg      .L10                If >, goto loop
10  .L7:                            done:
    Return result
```

c) 对应的汇编代码

# 循环：for

```
1    int fact_for(int n)
2    {
3        int i;
4        int result = 1;
5        for (i = 2; i <= n; i++)
6            result *= i;
7        return result;
8    }
```

```
for (init-expr; test-expr; update-expr)
    body-statement
```

```
init-expr;
while (test-expr) {
    body-statement
    update-expr;
}
```

```
init-expr;
t = test-expr;
if (!t)
    goto done;
loop:
    body-statement
    update-expr;
    t = test-expr;
    if (t)
        goto loop;
done:
```

```
1    int fact_for_goto(int n)
2    {
3        int i = 2;
4        int result = 1;
5        if (!(i <= n))
6            goto done;
7    loop:
8        result *= i;
9        i++;
10       if (i <= n)
11           goto loop;
12   done:
13       return result;
14   }
```

```
Argument: n at %ebp+8
Registers: n in %ecx, i in %edx, result in %eax

1     movl    8(%ebp), %ecx      Get n
2     movl    $2, %edx           Set i to 2         (init)
3     movl    $1, %eax           Set result to 1
4     cmpl    $1, %ecx           Compare n:1        (!test)
5     jle     .L14               If <=, goto done
6   .L17:                        loop:
7     imull   %edx, %eax         Compute result *= i (body)
8     addl    $1, %edx           Increment i        (update)
9     cmpl    %edx, %ecx         Compare n:i        (test)
10    jge     .L17               If >=, goto loop
11  .L14:                        done:
```

# switch

- 跳转表内的数组引用

```
1   int switch_eg(int x, int n) {
2       int result = x;
3
4       switch (n) {
5
6       case 100:
7           result *= 13;
8           break;
9
10      case 102:
11          result += 10;
12          /* Fall through */
13
14      case 103:
15          result += 11;
16          break;
17
18      case 104:
19      case 106:
20          result *= result;
21          break;
22
23      default:
24          result = 0;
25      }
26
27      return result;
28  }
```
a) switch 语句

```
1   int switch_eg_impl(int x, int n) {
2       /* Table of code pointers */
3       static void *jt[7] = {
4           &&loc_A, &&loc_def, &&loc_B,
5           &&loc_C, &&loc_D, &&loc_def,
6           &&loc_D
7       };
8
9       unsigned index = n - 100;
10      int result;
11
12      if (index > 6)
13          goto loc_def;
14
15      /* Multiway branch */
16      goto *jt[index];
17
18  loc_def:   /* Default case*/
19      result = 0;
20      goto done;
21
22  loc_C:     /* Case 103 */
23      result = x;
24      goto rest;
25
26  loc_A:     /* Case 100 */
27      result = x * 13;
28      goto done;
29
30  loc_B:     /* Case 102 */
31      result = x + 10;
32      /* Fall through */
33
34  rest:      /* Finish case 103 */
35      result += 11;
36      goto done;
37
38  loc_D:     /* Cases 104, 106 */
39      result = x * x;
40      /* Fall through */
41
42  done:
43      return result;
44  }
```
b) 翻译到扩展的 C 语言

```
    x at %ebp+8, n at %ebp+12
1       movl    8(%ebp), %edx        Get x
2       movl    12(%ebp), %eax       Get n
    Set up jump table access
3       subl    $100, %eax           Compute index = n-100
4       cmpl    $6, %eax             Compare index:6
5       ja      .L2                  If >, goto loc_def
6       jmp     *.L7(,%eax,4)        Goto *jt[index]
    Default case
7   .L2:                             loc_def:
8       movl    $0, %eax             result = 0;
9       jmp     .L8                  Goto done
    Case 103
10  .L5:                             loc_C:
11      movl    %edx, %eax           result = x;
12      jmp     .L9                  Goto rest
    Case 100
13  .L3:                             loc_A:
14      leal    (%edx,%edx,2), %eax  result = x*3;
15      leal    (%edx,%eax,4), %eax  result = x+4*result
16      jmp     .L8                  Goto done
    Case 102
17  .L4:                             loc_B:
18      leal    10(%edx), %eax       result = x+10
    Fall through
19  .L9:                             rest:
20      addl    $11, %eax            result += 11;
21      jmp     .L8                  Goto done
    Cases 104, 106
22  .L6:                             loc_D:
23      movl    %edx, %eax           result = x
24      imull   %edx, %eax           result *= x
    Fall through
25  .L8:                             done:
    Return result
```
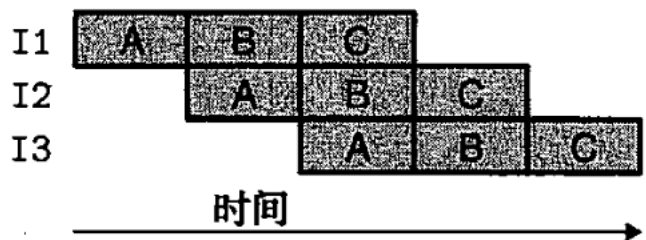
图 3-19    图 3-18 中 switch 语句示例的汇编代码

# 条件传送

- 条件**控制**转移->条件**数据**传送：可能提高效率
- 计算流水线（p264）：取指—译码—执行—访存—写回



a）原始的 C 语言代码

```
int absdiff(int x, int y) {
    return x < y ? y-x : x-y;
}
```



时间

b）使用条件赋值的实现

```
int cmovdiff(int x, int y) {
    int tval = y-x;
    int rval = x-y;
    int test = x < y;
    /* Line below requires
        single instruction: */
    if (test) rval = tval;
    return rval;
}
```

# Bad Cases for Conditional Move

## Expensive Computations

```
val = Test(x) ? Hard1(x) : Hard2(x);
```

Bad Performance

- **Both values get computed**
- **Only makes sense when computations are very simple**

## Risky Computations

```
val = p ? *p : 0;
```

Unsafe

当test-expr p为False，*p间接引用空指针

- Both values get computed
- May have undesirable effects

## Computations with side effects

```
val = x > 0 ? x*=7 : x+=3;
```

Illegal

两个分支同时改变全局变量

- Both values get computed
- Must be side-effect free

# Practice

王善上

# The End