

Machine-Level Programming I: Basics

- C/Assembly and Machine Code-

CENG331 - Computer Organization

Middle East Technical University

Instructor:

Murat Manguoglu (Section 1)

Adapted from slides of the textbook: <http://csapp.cs.cmu.edu/>

Today: Machine Programming I: Basics

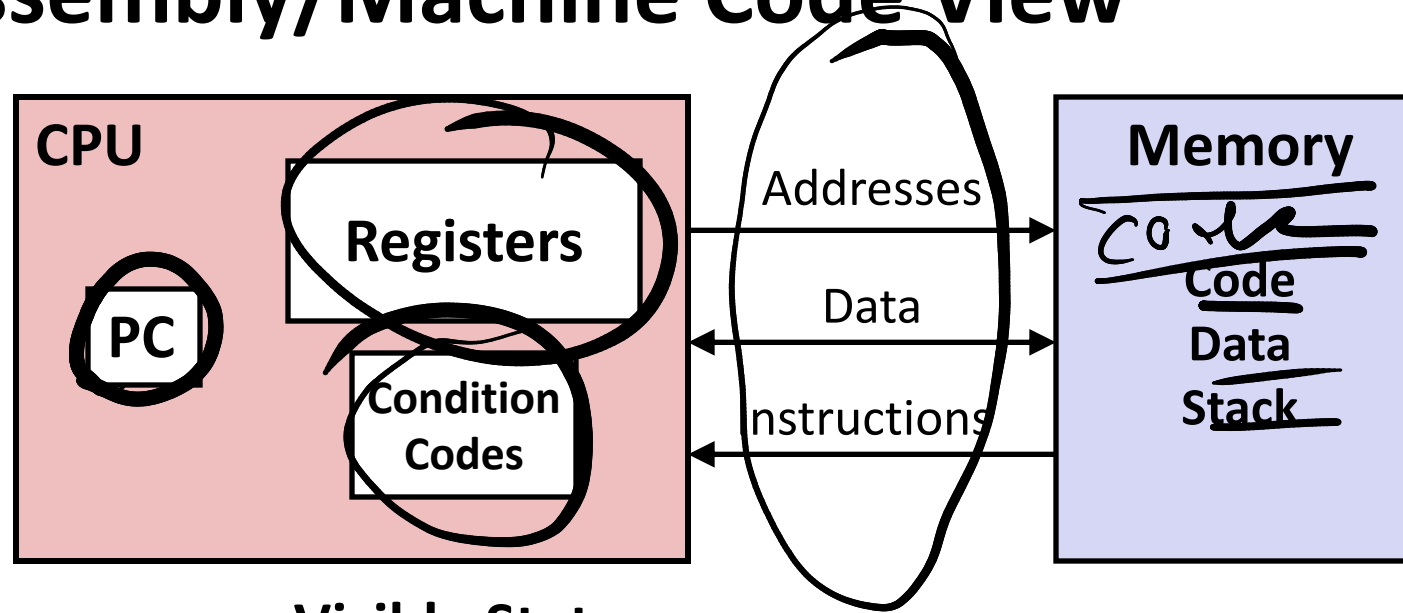
- History of Intel processors and architectures
- **C, assembly, machine code**
- Assembly Basics: Registers, operands, move
- Arithmetic & logical operations

Definitions

- **Architecture**: (also ISA: instruction set architecture) The parts of a processor design that one needs to understand or write assembly/machine code.
 - Examples: instruction set specification, registers.
- **Microarchitecture**: Implementation of the architecture.
 - Examples: cache sizes and core frequency.
- **Code Forms**:
 - **Machine Code**: The byte-level programs that a processor executes
 - **Assembly Code**: A text representation of machine code
- **Example ISAs**:
 - Intel: x86, IA32, Itanium, x86-64
 - ARM: Used in almost all mobile phones and now more common in desktops/laptops

RISC

Assembly/Machine Code View

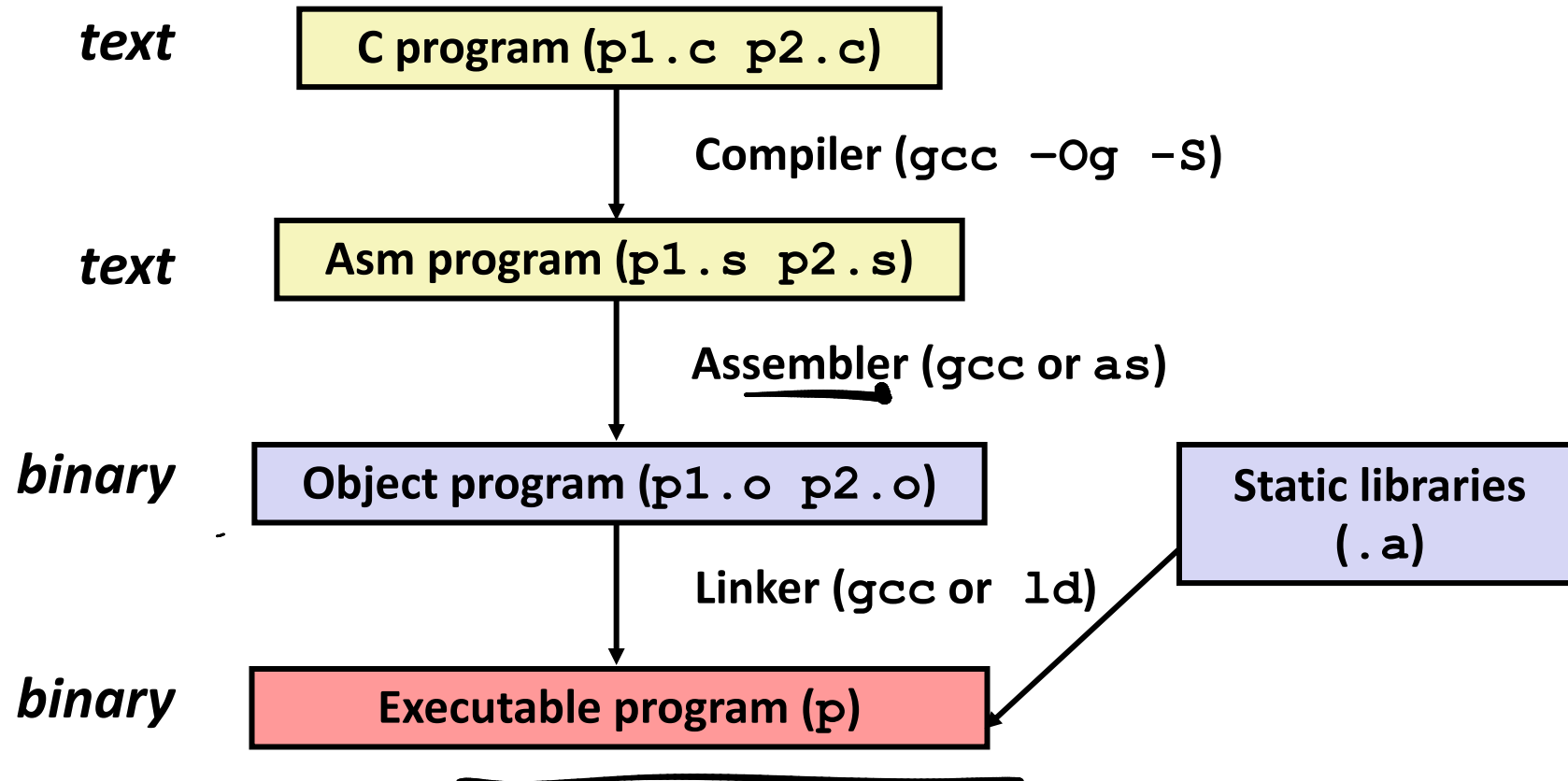


Programmer-Visible State

- **PC: Program counter**
 - Address of next instruction
 - Called the instruction pointer register "RIP" (x86-64)
- **Register file**
 - Heavily used program data
- **Condition codes**
 - Store status information about most recent arithmetic or logical operation
 - Used for conditional branching
- **Memory**
 - Byte addressable array
 - Code and user data
 - Stack to support procedures

Turning C into Object Code

- Code in files `p1.c` `p2.c`
- Compile with command: `gcc -Og p1.c p2.c -o p`
 - Use basic optimizations (`-Og`) [New to recent versions of GCC]
 - Put resulting binary in file `p`



Compiling Into Assembly

C Code (sum.c)

```
long plus(long x, long y);  
  
void sumstore(long x, long y,  
               long *dest)  
{  
    long t = plus(x, y);  
    *dest = t;  
}
```

Generated x86-64 Assembly

```
sumstore:  
→ pushq    %rbx  
→ movq     %rdx, %rbx  
→ call    plus  
→ movq     %rax, (%rbx)  
→ popq     %rbx  
→ ret
```

✓ AT&T or Intel

Obtain with command

```
gcc -O1 -S sum.c
```

```
gcc -O1 -Wa,-aslh -c sum.c > sum.s
```

Produces file `sum.s`

Warning: probably get very different results due to different architectures, different compiler/versions of compiler and different compiler settings.

Assembly Characteristics: Data Types

- “Integer” data of 1, 2, 4, or 8 bytes

- Data values
- Addresses (untyped pointers)

- Floating point data of 4, 8, or 10 bytes

- Code: Byte sequences encoding series of instructions

- No aggregate types such as arrays or structures

- Just contiguously allocated bytes in memory

Assembly Characteristics: Operations

- Perform arithmetic function on register or memory data
- Transfer data between memory and register
 - Load data from memory into register
 - Store register data into memory
- **Transfer control**
 - Unconditional jumps to/from procedures
 - Conditional branches

Object Code

Code for sumstore

0x0400595:

→ 0x53
0x48
0x89
0xd3
0xe8
0xf2
0xff
0xff
0xff
0x48
0x89
0x03
0x5b
↓ 0xc3

- Total of 14 bytes
- Each instruction 1, 3, or 5 bytes
- Starts at address 0x0400595

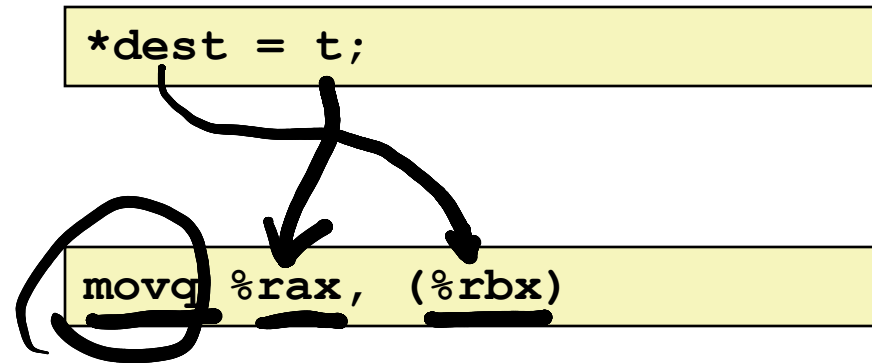
■ Assembler

- Translates .s into .o
- Binary encoding of each instruction
- Nearly-complete image of executable code
- Missing linkages between code in different files

■ Linker

- Resolves references between files
- Combines with static run-time libraries
 - E.g., code for **malloc**, **printf**
- Some libraries are *dynamically linked*
 - Linking occurs when program begins execution

Machine Instruction Example



`0x40059e: 48 89 03`

■ C Code

- Store value `t` where designated by `dest`

■ Assembly

- Move 8-byte value to memory
 - Quad words in x86-64 parlance
- Operands:
 - `t`: Register `%rax`
 - `dest`: Register `%rbx`
 - `*dest`: Memory `M[%rbx]`

■ Object Code

- 3-byte instruction
- Stored at address `0x40059e`

Disassembling Object Code

Disassembled

```
0000000000400595 <sumstore>:
400595: 53                push    %rbx
400596: 48 89 d3          mov     %rdx,%rbx
400599: e8 f2 ff ff ff    callq   400590 <plus>
40059e: 48 89 03          mov     %rax, (%rbx)
4005a1: 5b                pop     %rbx
4005a2: c3                retq
```

■ Disassembler

objdump -d sum

- Useful tool for examining object code
- Analyzes bit pattern of series of instructions
- Produces approximate rendition of assembly code
- Can be run on either a .out (complete executable) or .o file

Alternate Disassembly

Object

0x0400595:

0x53

0x48

0x89

0xd3

0xe8

0xf2

0xff

0xff

0xff

0x48

0x89

0x03

0x5b

0xc3

Disassembled

Dump of assembler code for function sumstore:

0x0000000000400595 <+0>: push %rbx

0x0000000000400596 <+1>: mov %rdx,%rbx

0x0000000000400599 <+4>: callq 0x400590 <plus>

0x000000000040059e <+9>: mov %rax, (%rbx)

0x00000000004005a1 <+12>: pop %rbx

0x00000000004005a2 <+13>: retq

■ Within gdb Debugger

`gdb sum`

`disassemble sumstore`

■ Disassemble procedure

`x/14xb sumstore`

■ Examine the 14 bytes starting at `sumstore`

What Can be Disassembled?

```
% objdump -d WINWORD.EXE

WINWORD.EXE:      file format pei-i386

No symbols in "WINWORD.EXE".
Disassembly of section .text:

30001000 <.text>:
30001000:
30001001:
30001003:
30001005:
3000100a:
```

**Reverse engineering forbidden by
Microsoft End User License Agreement**

- Anything that can be interpreted as executable code
- Disassembler examines bytes and reconstructs assembly source

Today: Machine Programming I: Basics

- History of Intel processors and architectures
- C, assembly, machine code
- **Assembly Basics: Registers, operands, move**
- Arithmetic & logical operations

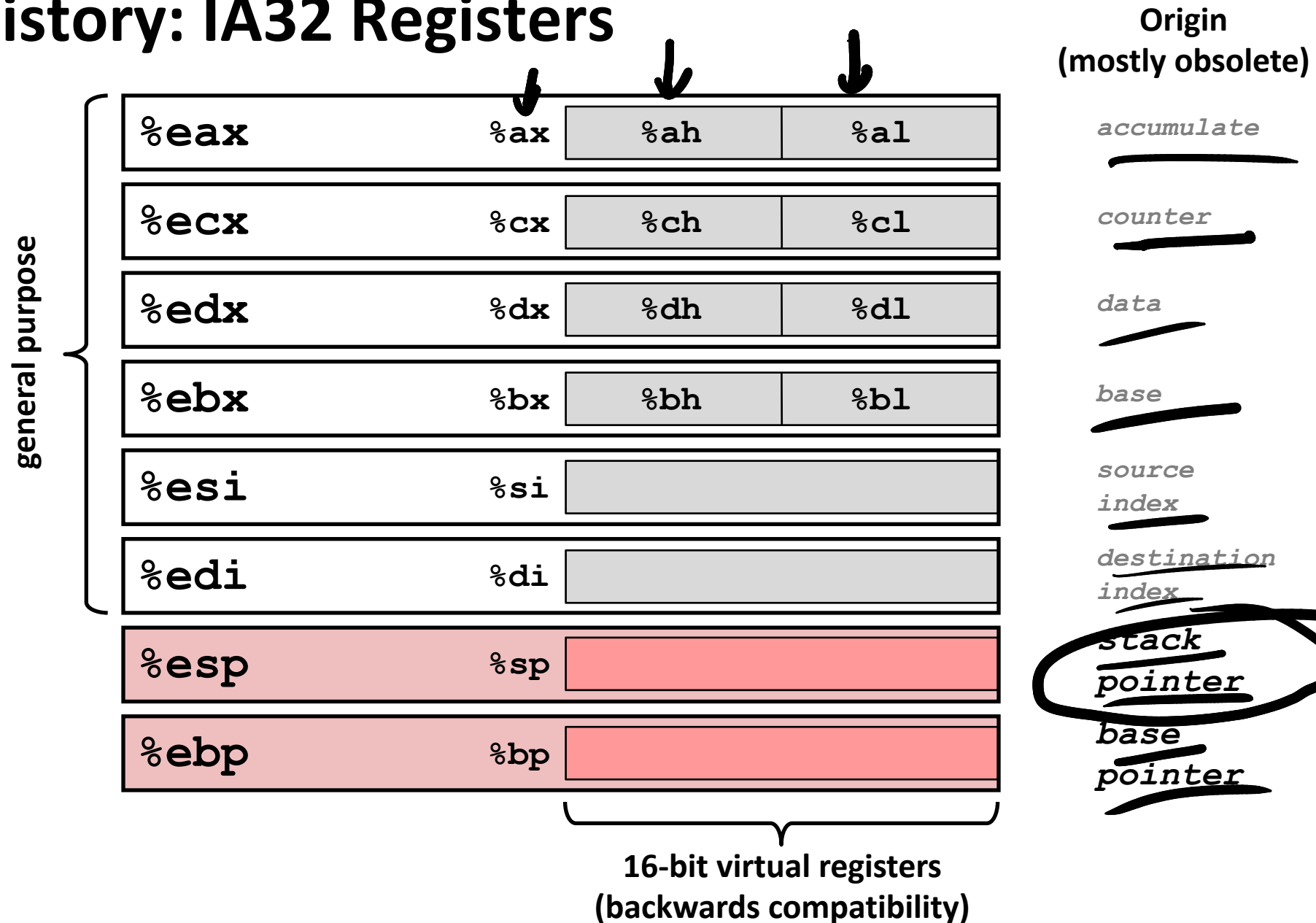
x86-64 Integer Registers

| | |
|-------------|-------------|
| %rax | %eax |
| %rbx | %ebx |
| %rcx | %ecx |
| %rdx | %edx |
| %rsi | %esi |
| %rdi | %edi |
| %rsp | %esp |
| %rbp | %ebp |

| | |
|-------------|--------------|
| %r8 | %r8d |
| %r9 | %r9d |
| %r10 | %r10d |
| %r11 | %r11d |
| %r12 | %r12d |
| %r13 | %r13d |
| %r14 | %r14d |
| %r15 | %r15d |

- Can reference low-order 4 bytes (also low-order 1 & 2 bytes)

Some History: IA32 Registers



X86-64 Registers

| 63 | 31 | 15 | 8 | 7 | 0 |
|------|-------|-------|-----|-------|---|
| %rax | %eax | %ax | %ah | %al | |
| %rbx | %ebx | %bx | %bh | %bl | |
| %rcx | %ecx | %cx | %ch | %cl | |
| %rdx | %edx | %dx | %dh | %dl | |
| %rsi | %esi | %si | ✗ | %sil | |
| %rdi | %edi | %di | ✗ | %dil | |
| %rbp | %ebp | %bp | ✗ | %bpl | |
| %rsp | %esp | %sp | | %spl | |
| %r8 | %r8d | %r8w | | %r8b | |
| %r9 | %r9d | %r9w | | %r9b | |
| %r10 | %r10d | %r10w | | %r10b | |
| %r11 | %r11d | %r11w | | %r11b | |
| %r12 | %r12d | %r12w | | %r12b | |
| %r13 | %r13d | %r13w | | %r13b | |
| %r14 | %r14d | %r14w | | %r14b | |
| %r15 | %r15d | %r15w | | %r15b | |

Moving Data

■ Moving Data

movq Source, Dest:

■ Operand Types

- **Immediate**: Constant integer data
 - Example: \$0x400, \$-533
 - Like C constant, but prefixed with '\$'
 - Encoded with 1, 2, or 4 bytes
- **Register**: One of 16 integer registers
 - Example: %rax, %r13
 - But %rsp reserved for special use
 - Others have special uses for particular instructions
- **Memory**: 8 consecutive bytes of memory at address given by register
 - Simplest example: (%rax)
 - Various other “address modes”

| |
|------|
| %rax |
| %rcx |
| %rdx |
| %rbx |
| %rsi |
| %rdi |
| %rsp |
| %rbp |
| |
| %rN |

movx Operand Combinations

| | Source | Dest | Src, Dest | C Analog |
|-------------|------------|------------|--------------------------------------|-----------------------|
| <u>movq</u> | <u>Imm</u> | <u>Reg</u> | movq <u>\$0x4</u> , <u>%rax</u> | <u>temp = 0x4;</u> |
| | | <u>Mem</u> | movq <u>\$-147</u> , (<u>%rax</u>) | <u>*p = -147;</u> |
| | <u>Reg</u> | <u>Reg</u> | movq <u>%rax</u> , <u>%rdx</u> | <u>temp2 = temp1;</u> |
| | | <u>Mem</u> | movq <u>%rax</u> , (<u>%rdx</u>) | <u>*p = temp;</u> |
| | <u>Mem</u> | <u>Reg</u> | movq (<u>%rax</u>), <u>%rdx</u> | <u>temp = *p;</u> |
| | | | | |

Cannot do memory-memory transfer with a single instruction

Simple Memory Addressing Modes

■ Normal (R) Mem[Reg[R]]

- Register R specifies memory address
- Aha! Pointer dereferencing in C

```
movq (%rcx), %rax
```

■ Displacement D(R) Mem[Reg[R]+D]

- Register R specifies start of memory region
- Constant displacement D specifies offset

```
movq 8(%rbp), %rdx
```



Example of Simple Addressing Modes

```
void swap
(long *xp, long *yp)
{
    long t0 = *xp;
    long t1 = *yp;
    *xp = t1;
    *yp = t0;
}
```

```
swap:
— movq    (%rdi), %rax
— movq    (%rsi), %rdx
— movq    %rdx, (%rdi)
— movq    %rax, (%rsi)
ret
```

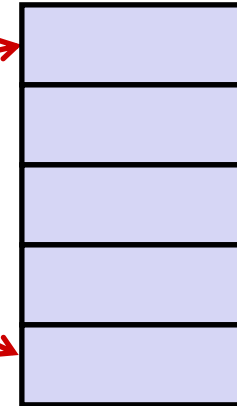
Understanding Swap()

```
void swap
(long *xp, long *yp)
{
    long t0 = *xp;
    long t1 = *yp;
    *xp = t1;
    *yp = t0;
}
```

Registers

| | |
|------|---|
| %rdi | - |
| %rsi | |
| %rax | |
| %rdx | |

Memory



| Register | Value |
|-------------|-----------|
| <u>%rdi</u> | <u>xp</u> |
| <u>%rsi</u> | <u>yp</u> |
| <u>%rax</u> | <u>t0</u> |
| <u>%rdx</u> | <u>t1</u> |

swap:

```
movq    (%rdi), %rax    # t0 = *xp
movq    (%rsi), %rdx    # t1 = *yp
movq    %rdx, (%rdi)    # *xp = t1
movq    %rax, (%rsi)    # *yp = t0
ret
```

Understanding Swap()

Registers

| | |
|------|-------|
| %rdi | 0x120 |
| %rsi | 0x100 |
| %rax | |
| %rdx | |

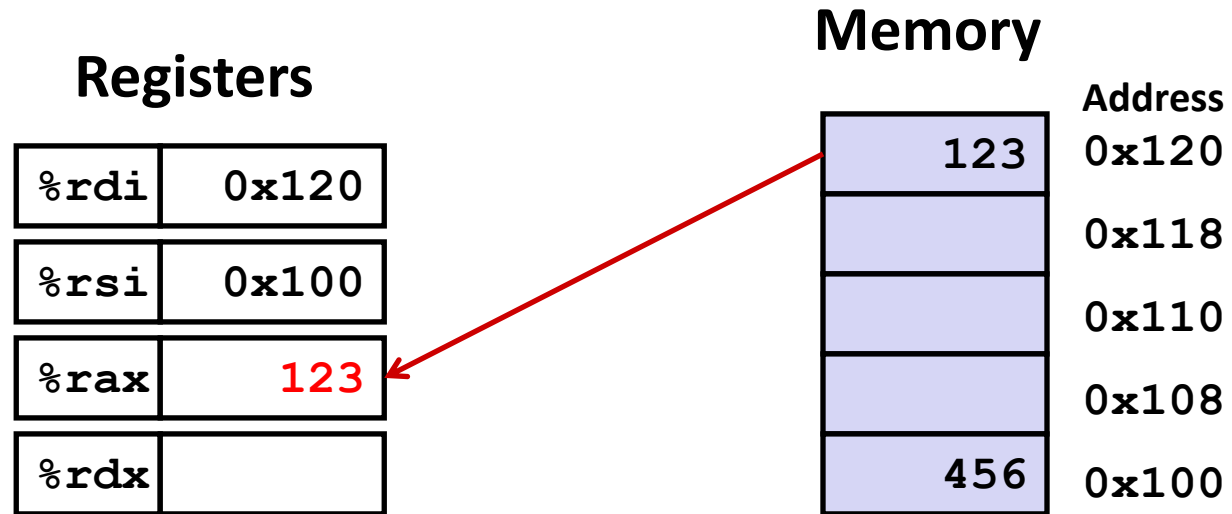
Memory

| Address |
|---------|
| 0x120 |
| 123 |
| 0x118 |
| 0x110 |
| 0x108 |
| 0x100 |
| 456 |

swap:

```
    movq    (%rdi), %rax    # t0 = *xp
    movq    (%rsi), %rdx    # t1 = *yp
    movq    %rdx, (%rdi)    # *xp = t1
    movq    %rax, (%rsi)    # *yp = t0
    ret
```

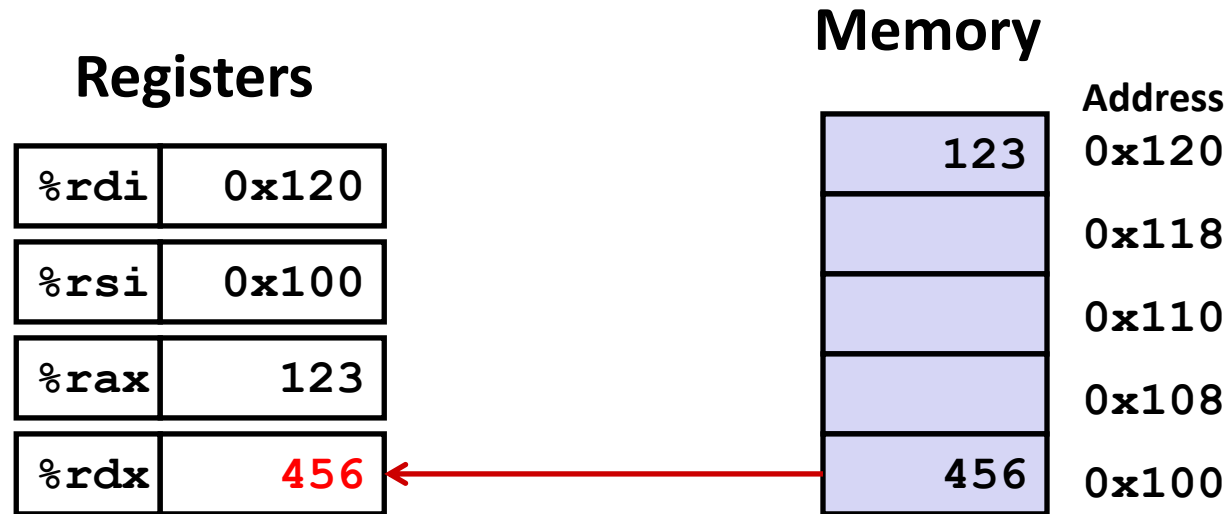
Understanding Swap()



swap:

```
movq    (%rdi), %rax    # t0 = *xp
movq    (%rsi), %rdx    # t1 = *yp
movq    %rdx, (%rdi)    # *xp = t1
movq    %rax, (%rsi)    # *yp = t0
ret
```

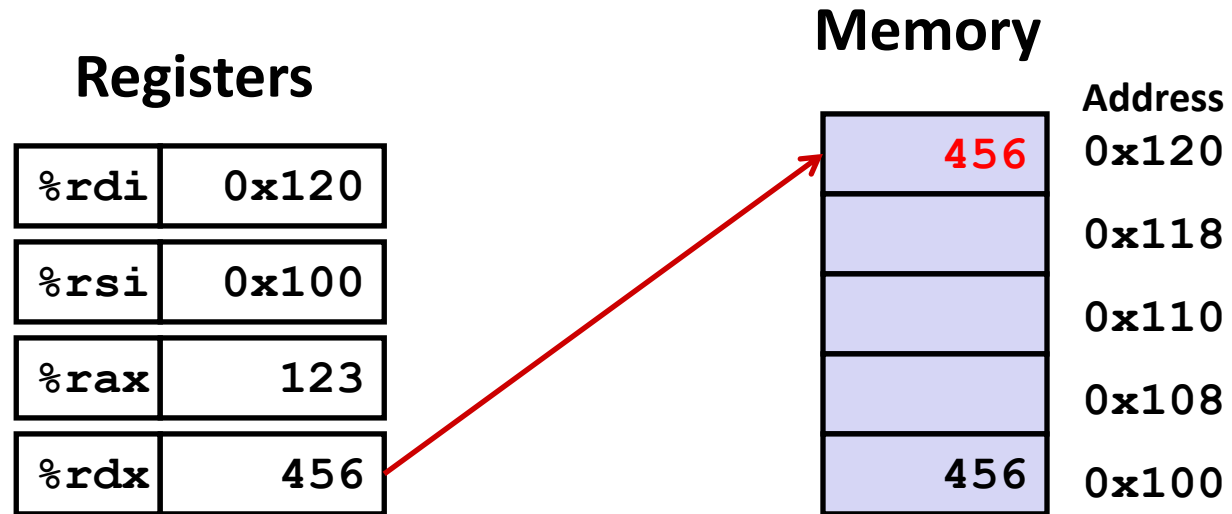

Understanding Swap()



swap:

```
movq    (%rdi), %rax    # t0 = *xp
movq    (%rsi), %rdx    # t1 = *yp
movq    %rdx, (%rdi)    # *xp = t1
movq    %rax, (%rsi)    # *yp = t0
ret
```

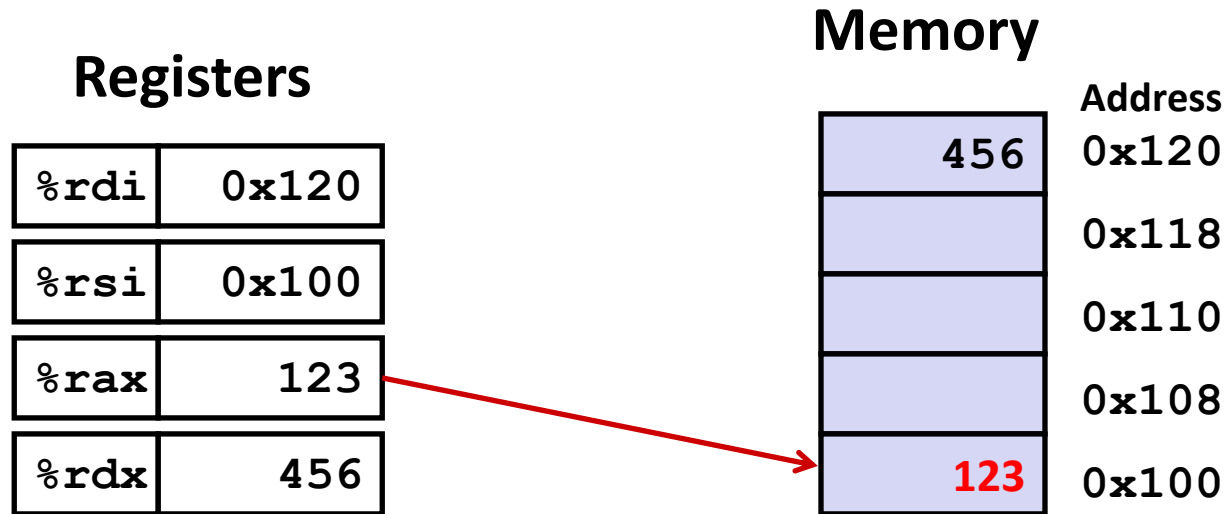
Understanding Swap()



swap:

```
movq    (%rdi), %rax    # t0 = *xp
movq    (%rsi), %rdx    # t1 = *yp
movq    %rdx, (%rdi)    # *xp = t1
movq    %rax, (%rsi)    # *yp = t0
ret
```

Understanding Swap()



```
swap:
    movq    (%rdi), %rax    # t0 = *xp
    movq    (%rsi), %rdx    # t1 = *yp
    movq    %rdx, (%rdi)    # *xp = t1
    movq    %rax, (%rsi)    # *yp = t0
    ret
```

Example of Simple Addressing Modes

```
void swap
(long *xp, long *yp)
{
    long t0 = *xp;
    long t1 = *yp;
    *xp = t1;
    *yp = t0;
}
```

```
swap:
    movq    (%rdi), %rax
    movq    (%rsi), %rdx
    movq    %rdx, (%rdi)
    movq    %rax, (%rsi)
    ret
```

4 moves

```
void swap2
(long *xp, long *yp)
{
    *xp = *xp ^ *yp;
    *yp = *xp ^ *yp;
    *xp = *xp ^ *yp;
}
```

```
swap2:
    movq    (%rsi), %rax
    xorq    (%rdi), %rax
    movq    %rax, (%rdi)
    xorq    (%rsi), %rax
    movq    %rax, (%rsi)
    xorq    %rax, (%rdi)
    ret
```

6 moves

Simple Memory Addressing Modes

■ Normal (R) Mem[Reg[R]]

- Register R specifies memory address
- Aha! Pointer dereferencing in C

```
movq (%rcx), %rax
```

■ Displacement D(R) Mem[Reg[R]+D]

- Register R specifies start of memory region
- Constant displacement D specifies offset

```
movq 8(%rbp), %rdx
```

Complete Memory Addressing Modes

■ Most General Form

D(Rb, Ri, S) Mem[Reg[Rb]+S*Reg[Ri]+ D]

- D: Constant “displacement” 1, 2, or 4 bytes
- Rb: Base register: Any of 16 integer registers
- Ri: Index register: Any, except for %rsp
- S: Scale: 1, 2, 4, or 8 (*why these numbers?*)

■ Special Cases

(Rb, Ri)

^{/*}
Mem[Reg[Rb]+Reg[Ri]]

D(Rb, Ri)

Mem[Reg[Rb]+Reg[Ri]+D]

(Rb, Ri, S)

Mem[Reg[Rb]+S*Reg[Ri]]

Address Computation Examples

| | |
|-------------|---------------|
| <u>%rdx</u> | <u>0xf000</u> |
| <u>%rcx</u> | <u>0x0100</u> |

| Expression | Address Computation | Address |
|---|--------------------------------|----------------|
| <u>0x8</u> (<u>%rdx</u>) | <u>0xf000</u> + 0x8 | <u>0xf008</u> |
| (<u>%rdx</u> , <u>%rcx</u>) | <u>0xf000</u> + <u>0x100</u> | <u>0xf100</u> |
| (<u>%rdx</u> , <u>%rcx</u> , 4) | <u>0xf000</u> + <u>4*0x100</u> | <u>0xf400</u> |
| <u>0x80</u> (<u>,</u> <u>%rdx</u> , 2) | <u>2*0xf000</u> + 0x80 | <u>0x1e080</u> |

Today: Machine Programming I: Basics

- History of Intel processors and architectures
- C, assembly, machine code
- Assembly Basics: Registers, operands, move
- **Arithmetic & logical operations**

Address Computation Instruction

■ leaq Src, Dst

- Src is address mode expression
- Set Dst to address denoted by expression

■ Uses

- Computing addresses without a memory reference
 - E.g., translation of `p = &x[i];`
- Computing arithmetic expressions of the form $x + k*y$
 - $k = 1, 2, 4, \text{ or } 8$

■ Example

```
long m12(long x)
{
    return x*12;
}
```

$x * 12$

Converted to ASM by compiler:

```
leaq (%rdi,%rdi,2), %rax # t <-  $x+x*2$ 
salq $2, %rax           # return  $t \ll 2$ 
```

\rightarrow $\text{incrg} \cdot \text{rax}$

$3x$
 $4 \times \text{rax}$

Some Arithmetic Operations

■ Two Operand Instructions:

Format

Computation

| | | |
|--------------------|------------------------|---------------------------------------|
| <code>addq</code> | <code>Src, Dest</code> | <code>Dest = Dest + Src</code> |
| <code>subq</code> | <code>Src, Dest</code> | <code>Dest = Dest - Src</code> |
| <code>imulq</code> | <code>Src, Dest</code> | <code>Dest = Dest * Src</code> |
| <code>salq</code> | <code>Src, Dest</code> | <code>Dest = Dest << Src</code> |
| <code>sarq</code> | <code>Src, Dest</code> | <code>Dest = Dest >> Src</code> |
| <code>shrq</code> | <code>Src, Dest</code> | <code>Dest = Dest >> Src</code> |
| <code>xorq</code> | <code>Src, Dest</code> | <code>Dest = Dest ^ Src</code> |
| <code>andq</code> | <code>Src, Dest</code> | <code>Dest = Dest & Src</code> |
| <code>orq</code> | <code>Src, Dest</code> | <code>Dest = Dest Src</code> |

Also called `shlq`

Arithmetic

Logical

■ Watch out for argument order!

■ No distinction between signed and unsigned int (why?)

Some Arithmetic Operations

■ One Operand Instructions

| | | |
|-------------------|-------------------|---------------------------------|
| <code>incq</code> | <code>Dest</code> | $\text{Dest} = \text{Dest} + 1$ |
| <code>decq</code> | <code>Dest</code> | $\text{Dest} = \text{Dest} - 1$ |
| <code>negq</code> | <code>Dest</code> | $\text{Dest} = -\text{Dest}$ |
| <code>notq</code> | <code>Dest</code> | $\text{Dest} = \sim\text{Dest}$ |

■ See the book or Intel x86-64 manual for more instructions

Arithmetic Expression Example

```
long arith
(long x, long y, long z)
{
    long t1 = x+y;
    long t2 = z+t1;
    long t3 = x+4;
    long t4 = y * 48;
    long t5 = t3 + t4;
    long rval = t2 * t5;
    return rval;
}
```

rax

arith:

```
→ leaq
→ addq
→ leaq
→ salq
→ leaq
→ imulq
→ ret
```

```
(%rdi,%rsi), %rax
%rdx, %rax ← t1
(%rsi,%rsi,2), %rdx
$4, %rdx
4(%rdi,%rdx), %rcx
%rcx, %rax
```

$t_1 \leftarrow t_1 + t_2$
 $3 \times y \rightarrow z$
 $4 + 1 \times di + rdx$

Interesting Instructions

- `leaq`: address computation
- `salq`: shift
- `imulq`: multiplication
 - But, only used once

Understanding Arithmetic Expression Example

```
long arith
(long x, long y, long z)
{
    long t1 = x+y;
    long t2 = z+t1;
    long t3 = x+4;
    long t4 = y * 48;
    long t5 = t3 + t4;
    long rval = t2 * t5;
    return rval;
}
```

arith:

```
leaq    (%rdi,%rsi), %rax    # t1
addq    %rdx, %rax          # t2
leaq    (%rsi,%rsi,2), %rdx
salq    $4, %rdx            # t4
leaq    4(%rdi,%rdx), %rcx   # t5
imulq    %rcx, %rax          # rval
ret
```

| Register | Use(s) |
|-------------|-------------------|
| <u>%rdi</u> | Argument x |
| <u>%rsi</u> | Argument y |
| <u>%rdx</u> | Argument z |
| %rax | t1, t2, rval |
| %rdx | t4 |
| %rcx | t5 |

Machine Programming I: Summary

■ History of Intel processors and architectures

- Evolutionary design leads to many quirks and artifacts

■ C, assembly, machine code

- New forms of visible state: program counter, registers, ...
- Compiler must transform statements, expressions, procedures into low-level instruction sequences

■ Assembly Basics: Registers, operands, move

- The x86-64 move instructions cover wide range of data movement forms

■ Arithmetic

- C compiler will figure out different instruction combinations to carry out computation