

Machine-Level Programming I: Basics

CSci 2021: Machine Architecture and Organization
February 12th-14th, 2018

Your instructor: Stephen McCamant

Based on slides originally by:
Randy Bryant, Dave O'Hallaron

Today: Machine Programming I: Basics

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

Intel x86 Processors

- Dominate laptop/desktop/server market
- Evolutionary design
 - Backwards compatible up until 8086, introduced in 1978
 - Added more features as time goes on
- Complex instruction set computer (CISC)
 - Many different instructions with many different formats
 - But, only a subset encountered with Linux programs
 - Matches performance of more modern Reduced Instruction Set Computers (RISC)
 - In terms of speed. Less so for low power consumption.

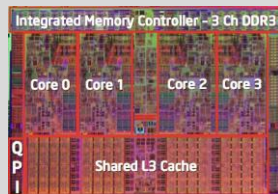
Intel x86 Evolution: Milestones

Name	Date	Transistors	MHz
■ 8086	1978	29K	5-10
<ul style="list-style-type: none"> ▪ First 16-bit Intel processor. Basis for IBM PC & DOS ▪ 1MB address space 			
■ 386	1985	275K	16-33
<ul style="list-style-type: none"> ▪ First 32 bit Intel processor, referred to as IA32 ▪ Added "flat addressing", capable of running Unix 			
■ Pentium 4E	2004	125M	2800-3800
<ul style="list-style-type: none"> ▪ First 64-bit Intel x86 processor, referred to as x86-64 			
■ Core 2	2006	291M	1060-3500
<ul style="list-style-type: none"> ▪ First multi-core Intel processor 			
■ Core i7	2008	731M	1700-3900
<ul style="list-style-type: none"> ▪ Four cores 			

Intel x86 Processors, cont.

- Machine Evolution

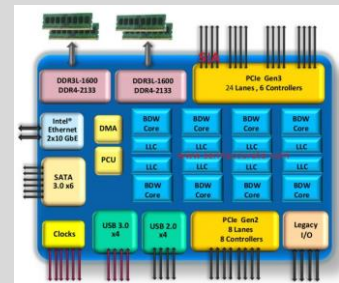
■ 386	1985	0.3M
■ Pentium	1993	3.1M
■ Pentium/MMX	1997	4.5M
■ PentiumPro	1995	6.5M
■ Pentium III	1999	8.2M
■ Pentium 4	2001	42M
■ Core 2 Duo	2006	291M
■ Core i7	2008	731M
- Added Features
 - Instructions to support multimedia operations
 - Instructions to enable more efficient conditional operations
 - Transition from 32 bits to 64 bits
 - More cores



2015 State of the Art

- Core i7 Broadwell 2015

- Desktop Model
 - 4 cores
 - Integrated graphics
 - 3.3-3.8 GHz
 - 65W
- Server Model
 - 8 cores
 - Integrated I/O
 - 2-2.6 GHz
 - 45W



x86 Clones: Advanced Micro Devices (AMD)

- **Historically**
 - AMD has followed just behind Intel
 - A little bit slower, a lot cheaper
- **Then**
 - Recruited top circuit designers from Digital Equipment Corp. and other downward trending companies
 - Built Opteron: tough competitor to Pentium 4
 - Developed x86-64, their own extension to 64 bits
- **Recent Years**
 - Intel got its act together
 - Leads the world in semiconductor technology
 - AMD has fallen behind
 - Spun off its semiconductor factories

Intel's 64-Bit History

- **2001: Intel Attempts Radical Shift from IA32 to IA64**
 - Totally different architecture (Itanium)
 - Executes IA32 code only as legacy
 - Performance disappointing
- **2003: AMD Steps in with Evolutionary Solution**
 - x86-64 (now called "AMD64")
- **Intel Felt Obligated to Focus on IA64**
 - Hard to admit mistake or that AMD is better
- **2004: Intel Announces EM64T extension to IA32**
 - Extended Memory 64-bit Technology (now called "Intel 64")
 - Almost identical to x86-64!
- **All but lowest-end x86 processors support x86-64**
 - But, lots of code still runs in 32-bit mode

Our Coverage

- **IA32**
 - The traditional x86
 - For 2021: RIP, Summer 2015
- **x86-64**
 - The standard
 - `cseclabs> gcc hello.c`
 - `cseclabs> gcc -m64 hello.c`
- **Presentation**
 - Book covers x86-64
 - Web aside on IA32
 - We will only cover x86-64

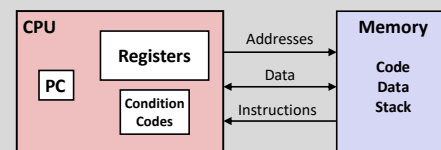
Today: Machine Programming I: Basics

- 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 smartphones

Assembly/Machine Code View

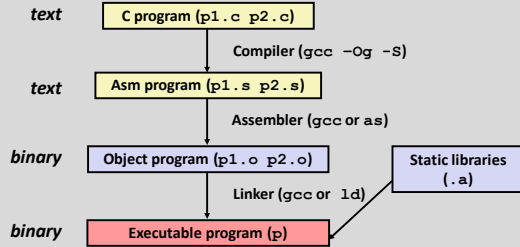


Programmer-Visible State

- **PC: Program counter**
 - Address of next instruction
 - On x86-64, called "RIP"
- **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 since GCC 4.8]
 - Put resulting binary in file `p`



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

13

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
  
```

Obtain (on Ubuntu 14.04 machine) with command

```
gcc -Og -S sum.c
```

Produces file `sum.s`

Note: You may get different results on different machines (older Linux, Mac OS X, ...) due to different versions of gcc and different compiler settings.

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

14

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

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

15

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

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

16

Object Code

Code for `sumstore`

```

0x0400595:
0x53
0x48
0x89
0xd3
0xe8
0xf2
0xff
0xff
0xff
0x48
0x89
0x03
0x5b
0xc3
  
```

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

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

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

17

Machine Instruction Example

```
*dest = t;
```

```
movq %rax, (%rbx)
```

```
0x40059e: 48 89 03
```

C Code

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

Assembly

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

Object Code

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

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

18

Disassembling Object Code

Disassembled

```
000000000400595 <sumstore>:
400595: 53          push    %rbx
400596: 48 89 d3    mov     %rdx,%rbx
400599: e8 f2 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
0x48
0x89
0x03
0x5b
0xc3
```

Disassembled

```
Dump of assembler code for function sumstore:
0x000000000400595 <+0>: push    %rbx
0x000000000400596 <+1>: mov     %rdx,%rbx
0x000000000400599 <+4>: callq   0x400590 <plus>
0x00000000040059e <+9>: mov     %rax, (%rbx)
0x0000000004005a1 <+12>: pop     %rbx
0x0000000004005a2 <+13>: retq
```

■ Within gdb Debugger

- `% gdb sum`
- `(gdb) disassemble sumstore`
- Disassemble procedure
- `(gdb) x/14xb sumstore`
- Examine the 14 bytes starting at `sumstore`

What Can be Disassembled?

```
% objdump -d WINWORD.EXE
```

```
WINWORD.EXE: 0100000000000000
```

```
No symbols in
Disassembly of
```

```
30001000 <.text>:
30001000: 55
30001001: 8b ec      mov     %esp,%ebp
30001003: 6a ff      push    $0xffffffff
30001005: 68 90 10 00 push    $0x30001090
3000100a: 68 91 dc 4c push    $0x304cdc91
```

Legal note: reverse engineering of commercial software is often forbidden by license agreements, and its status under statute varies by jurisdiction

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

Aside: x86 Assembly Formats

- This class uses "AT&T" format, which is standard for Unix/Linux x86(-64) systems
 - Similar to historic Unix all the way back to PDP-11
- Intel's own documentation, and Windows, use a different "Intel" syntax
 - Many arbitrary differences, but more internally consistent

AT&T syntax	Intel syntax
Destination is last operand	Destination is first operand
Size suffixes like "l" in <code>movl</code>	Size on memory operands ("DWORD PTR")
"%" on register names	Just letters in register names
"\$" on immediate values	Just digits in immediates
Addressing modes with <code>()</code>	Addressing modes with <code>[+ *]</code>

Today: Machine Programming I: Basics

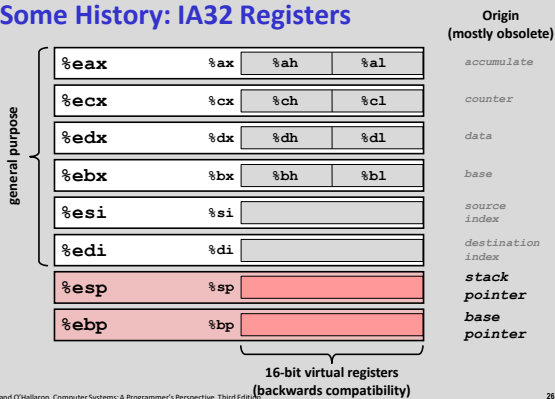
- History of Intel processors and architectures
- C, assembly, machine code
- Assembly Basics: Registers, operands, move
- Arithmetic & logical operations

x86-64 Integer Registers

%rax	%eax	%r8	%r8d
%rbx	%ebx	%r9	%r9d
%rcx	%ecx	%r10	%r10d
%rdx	%edx	%r11	%r11d
%rsi	%esi	%r12	%r12d
%rdi	%edi	%r13	%r13d
%rsp	%esp	%r14	%r14d
%rbp	%ebp	%r15	%r15d

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

Some History: IA32 Registers



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

26

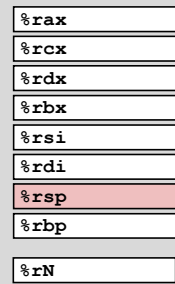
Moving Data

■ 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
 - Some 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"



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

27

movq Operand Combinations

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

Cannot do memory-memory transfer with a single instruction

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

28

Simple Memory Addressing Modes

- **Normal (R) Mem[Reg[R]]**
 - Register R specifies memory address
 - Like 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`

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

29

Example of Simple Addressing Modes

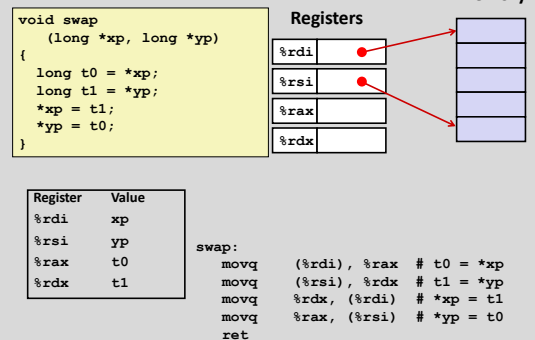
```
void swap
(long *xp, long *yp)
{
    long t0 = *xp;
    long t1 = *yp;
    *xp = t1;
    *yp = t0;
}
```

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

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

30

Understanding Swap()



Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

31

Understanding Swap()

Registers		Memory	
%rdi	0x120	123	0x120
%rsi	0x100		0x118
%rax			0x110
%rdx		456	0x108
			0x100

```

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		Memory	
%rdi	0x120	123	0x120
%rsi	0x100		0x118
%rax	123		0x110
%rdx		456	0x108
			0x100

```

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		Memory	
%rdi	0x120	123	0x120
%rsi	0x100		0x118
%rax	123		0x110
%rdx	456	456	0x108
			0x100

```

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		Memory	
%rdi	0x120	456	0x120
%rsi	0x100		0x118
%rax	123		0x110
%rdx	456	456	0x108
			0x100

```

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		Memory	
%rdi	0x120	456	0x120
%rsi	0x100		0x118
%rax	123		0x110
%rdx	456	123	0x108
			0x100

```

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

```

Complete Memory Addressing Modes

Most General Form

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

- 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

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

Address Computation Examples

<code>%rdx</code>	<code>0xf000</code>
<code>%rcx</code>	<code>0x0100</code>

Expression	Address Computation	Address
<code>0x8(%rdx)</code>		
<code>(%rdx,%rcx)</code>		
<code>(%rdx,%rcx,4)</code>		
<code>0x80(,%rdx,2)</code>		

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

38

Address Computation Examples

<code>%rdx</code>	<code>0xf000</code>
<code>%rcx</code>	<code>0x0100</code>

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

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

39

Today: Machine Programming I: Basics

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

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

41

Address Computation Instruction

- leaq Src, Dest**
 - Src is address mode expression
 - Set Dest 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, or 8
- Example**

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

Converted to ASM by compiler:

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

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

42

Some Arithmetic Operations

Two Operand Instructions:

Format **Computation**

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

*Also called salq
Arithmetic
Logical*

- Watch out for argument order!
- No distinction between signed and unsigned int (why?)

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

43

Some Arithmetic Operations

One Operand Instructions

```
incq Dest    Dest = Dest + 1
decq Dest    Dest = Dest - 1
negq Dest    Dest = - Dest
notq Dest    Dest = ~Dest
```

- See book for more instructions

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

44

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
    addq    %rdx, %rax
    leaq    (%rsi,%rsi,2), %rdx
    salq    $4, %rdx
    leaq    4(%rdi,%rdx), %rcx
    imulq    %rcx, %rax
    ret
```

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)
%rdi	Argument x
%rsi	Argument y
%rdx	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