# PURDUE UNIVERSITY®

**CS 50011: Introduction to Systems II**

**Lecture 4: Introduction to Assembly**

Prof. Jeff Turkstra

# Lecture 03

- History
- Background
- x86
  - Syntax
  - Operands
  - Addressing modes
  - Data types
  - Instructions

# A software hierarchy

Java, Pascal, FORTRAN, etc.

Portable among different computers with some effort; some machine-dependent features are visible; performance

C

Assembly languages: a distinct one for each computer, hardware dependence, just a few abstractions

| VAX | X86 | MIPS | ARM | ... |

Machine languages: one per computer; NO abstractions

| Variable length bit strings | Variable length bit strings | 32- & 64-bit strings | 32- & 64-bit strings | ... |

# Assembly language

- All (somewhat) different
- Many assembly languages share the same fundamental structure
  - Why?
- Typical assembly language statement syntax and corresponding machine code in hex...
  ```
          label: op result, operand1, operand2
  0x004005F9 0x23CC803C
  ```
- Label is symbolic (an abstraction) for a memory address
- "op" is a mnemonic for the operation

# **Assembly is two-pass**

- Initial pass of assembler resolves memory addresses for all labels
  - Even (especially) forward references
  - Symbol table
- Second pass emits machine code bitstrings
  - Translates mnemonics, register names, etc
  - Uses symbol table to fill in offset bit field
    - Offset = branch_target - current_addr

# Why?

- Many languages are one-pass
  - C, for example
    - Have to prototype functions, declare/define
- Would have to manually determine instruction addresses and branch targets
- Changing the code often changes all of the offsets and addresses
- Impractical

# Opcodes

- Set of opcode-field bit strings defines what the processor circuit can do

- Different processors have different sets of opcodes

- Assembly language defines a memorable symbolic name of a few characters for each opcode, a mnemonic

- No agreement on opcode mnemonics across assembly languages

# Readability

- Assembly is easy to write but hard to follow
- Comments are essential
  - Block comment – explain the purpose of a section of code, detail the use of registers and memory
  - Line comment – explains each instruction
- Comment usually starts with a delimiter, runs to end of line
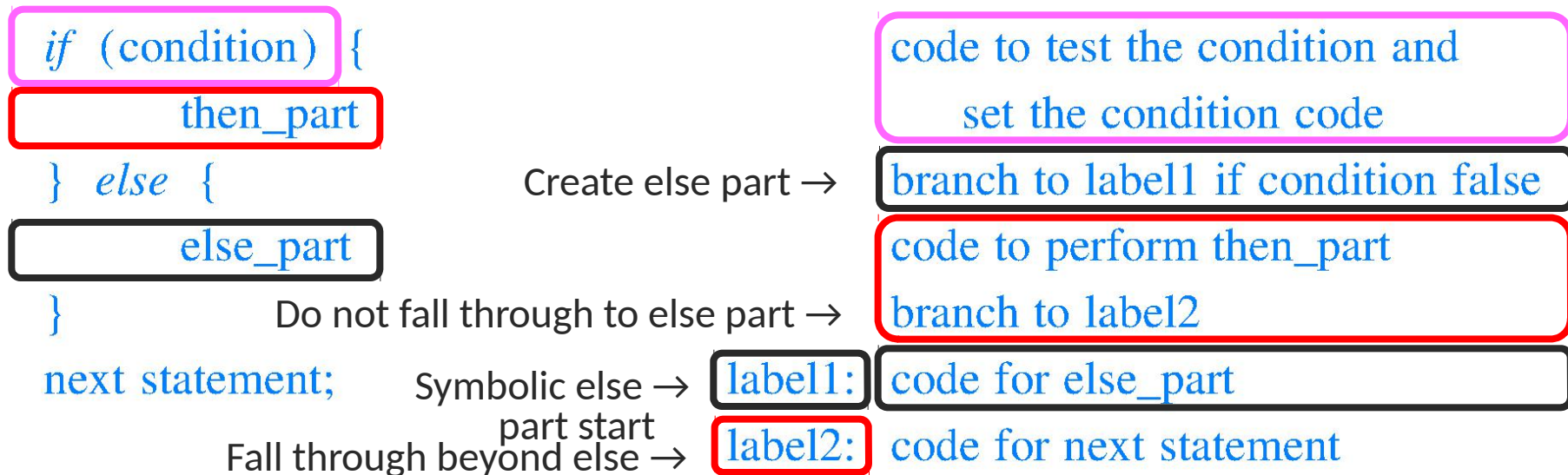- Best strategy: comment every line

# Example

```
###############################################
# Search linked list of free memory blocks to find  #
# a block of size N bytes or greater. Pointer must   #
# be in r3 and N in r4. Code destroys contents of    #
# r5, which is used to walk the list.                #
###############################################
        ld    r5,r3    # load address of list into r5
loop_1: cmp   r5,0     # test to see if at list end
        bz    notfnd   # if reached end go to notfnd
```

# Coding IF-THEN-ELSE in assembly

*if* (condition) {

       then_part

} *else* {　　　　　　　　　Create else part →　　branch to label1 if condition false

       else_part

}　　　　　Do not fall through to else part →

next statement;　　Symbolic else →　label1:　code for else_part
　　　　　　　　　part start
　　　　Fall through beyond else →　label2:　code for next statement

code to test the condition and
　　　　　set the condition code

code to perform then_part

branch to label2

"Fall through" means to fetch at the default next instruction location; must code two exceptions for if-then-else

**Figure 9.2** (a) An *if-then-else* statement used in a high-level language, and (b) the equivalent assembly language code.
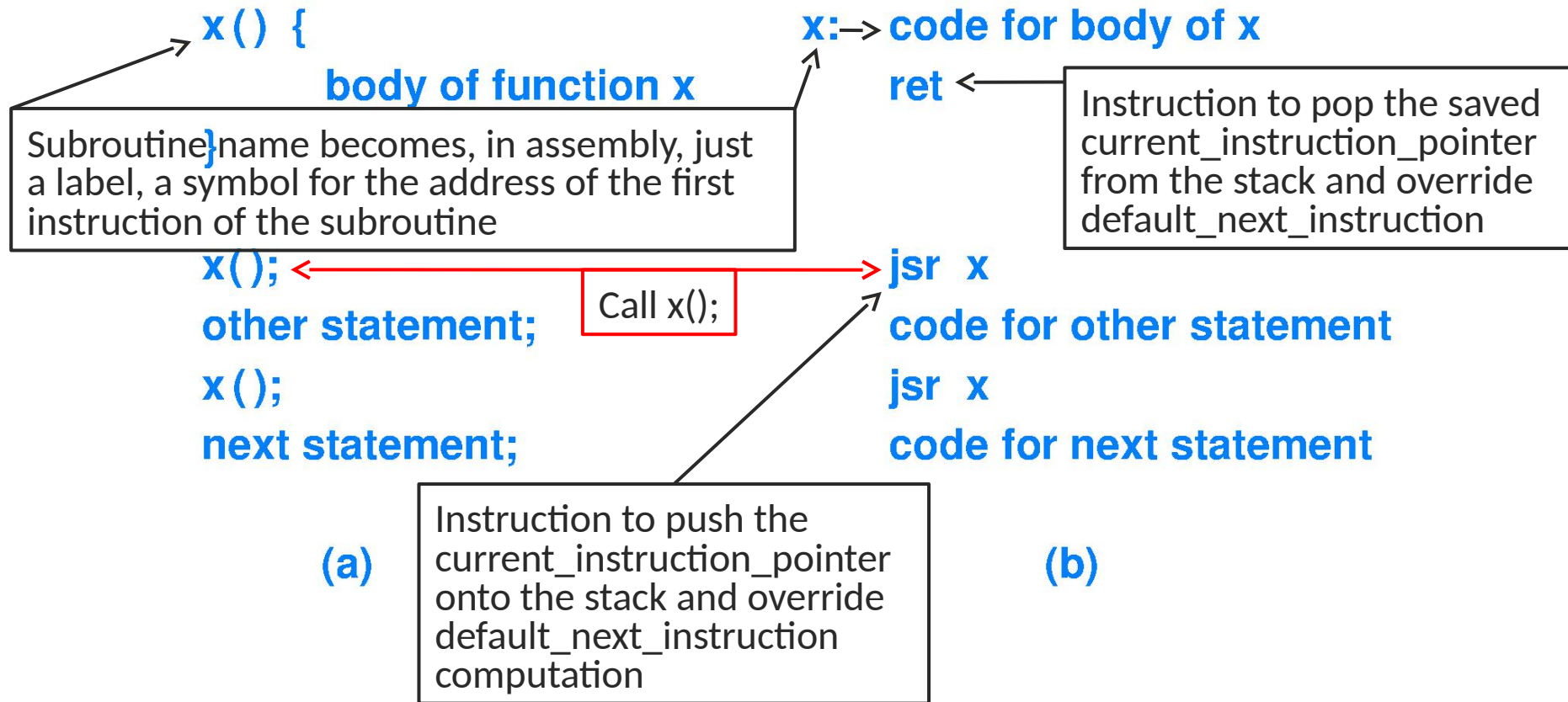
# Subroutine call in assembly

x() {

    **body of function x**

}

Subroutine name becomes, in assembly, just a label, a symbol for the address of the first instruction of the subroutine

x();

other statement;

x();

next statement;

**(a)**

Call x();

Instruction to push the current_instruction_pointer onto the stack and override default_next_instruction computation

x:→ **code for body of x**

    **ret**

Instruction to pop the saved current_instruction_pointer from the stack and override default_next_instruction

**jsr  x**

**code for other statement**

**jsr  x**

**code for next statement**

**(b)**

**Figure 9.5** (a) A declaration for procedure *x* and two invocations in a high-level language, and (b) the assembly language equivalent.

# Language specifics

- Documentation
  - Operand order
  - Register naming
  - Syntax
    - Immediate values, register values, memory, etc
- Assembly language does not provide any program control structures, nor enforce any coding style
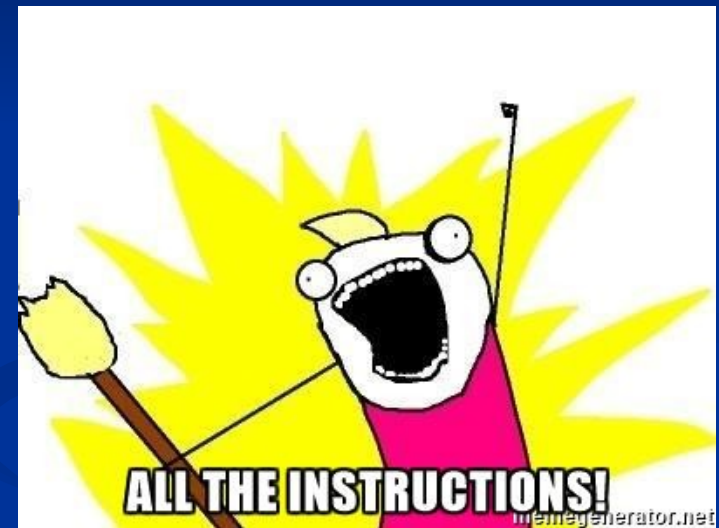
# Intel documentation

- Volume 1: Basic Architecture
  - 482 pages
  - 19 Chapters
  - Includes basic execution environment as well as summary of instructions
    - Groups instructions for programming
      - MMX, SIMD, SSE, etc

- Volume 2: Instruction Set Reference A-Z
  - 2234 pages
  - "Only" 6 chapters
  - Instruction format
  - All of the instructions
  - Safer Mode Extensions

- Volume 3: System Programming Guide
  - 1660 pages
  - 43 Chapters
  - Everything the hardware does to support an OS and how to use it

# CPUs have errata

- Ever hear of the original Pentium floating point bug?
    - Could have been errata, but the press picked it up
- Ever find a compiler error?
- Imagine finding a hardware error
    - Probably involves premature baldness
        - Possibly temporary

# x86 Assembly

- Unfortunately, x86 is arguably the most complex assembly language around
  - MOV is even Turing complete
- Exposure to most common instructions
  - Focus on ability to read assembled C programs
  - Maybe a little writing

Differences with x86_64

# **The Intel Legacy**

- Started with 4004
  - 4-bit processor
- 8086, first x86 CPU
  - 16-bits
  - June 8, 1978
  - 5MHz, 8MHz, and 10MHz
- 80186, 80286
- 80386 (SX/DX), 80486 (SX/DX/DX2/etc)

# Pentium

- MMX
- SSE, SSE2, SSE3
- X86-64
- AMD-V
- Intel VT-x
- etc
- …and it's all backwards compatible

# **Fortunately**

- Some analyses claim only 14 instructions account for 90% of compiled code

# Assembly is symbolic

- label: mnemonic arg1, arg2, arg3
  - Zero to three args
  - Right is source, left is destination
- Mnemonic may represent different (multiple) opcodes

# Remember



**Figure 5.1** The general instruction format that many processors use. The opcode at the beginning of an instruction determines exactly which operands follow.

# 64-bit prefix ordering

| Legacy Prefixes | REX Prefix | Opcode | ModR/M | SIB | Displacement | Immediate |
|---|---|---|---|---|---|---|
| Grp 1, Grp 2, Grp 3, Grp 4 (optional) | (optional) | 1-, 2-, or 3-byte opcode | 1 byte (if required) | 1 byte (if required) | Address displacement of 1, 2, or 4 bytes | Immediate data of 1, 2, or 4 bytes or none |

mov    rcx,0x4004e0
48 c7 c1 e0 04 40 00

48: REX.W prefix: 64-bit operand
c7: MOV
c1: ecx (but really rcx)
e0044000: 004004e0

# REX prefix

| Field Name | Bit Position | Definition |
|---|---|---|
| - | 7:4 | 0100 |
| W | 3 | 0 = Operand size determined by CS.D |
| | | 1 = 64 Bit Operand Size |
| R | 2 | Extension of the ModR/M reg field |
| X | 1 | Extension of the SIB index field |
| B | 0 | Extension of the ModR/M r/m field, SIB base field, or Opcode reg field |

# Wat?

| r8(/r)<br>r16(/r)<br>r32(/r)<br>mm(/r)<br>xmm(/r)<br>(In decimal) /digit (Opcode)<br>(In binary) REG = | | | AL<br>AX<br>EAX<br>MM0<br>XMM0<br>0<br>000 | CL<br>CX<br>ECX<br>MM1<br>XMM1<br>1<br>001 | DL<br>DX<br>EDX<br>MM2<br>XMM2<br>2<br>010 | BL<br>BX<br>EBX<br>MM3<br>XMM3<br>3<br>011 | AH<br>SP<br>ESP<br>MM4<br>XMM4<br>4<br>100 | CH<br>BP<br>EBP<br>MM5<br>XMM5<br>5<br>101 | DH<br>SI<br>ESI<br>MM6<br>XMM6<br>6<br>110 | BH<br>DI<br>EDI<br>MM7<br>XMM7<br>7<br>111 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Effective Address** | **Mod** | **R/M** | **Value of ModR/M Byte (in Hexadecimal)** | | | | | | | |
| [EAX]<br>[ECX]<br>[EDX]<br>[EBX]<br>[--][--]$^1$<br>disp32$^2$<br>[ESI]<br>[EDI] | 00 | 000<br>001<br>010<br>011<br>100<br>101<br>110<br>111 | 00<br>01<br>02<br>03<br>04<br>05<br>06<br>07 | 08<br>09<br>0A<br>0B<br>0C<br>0D<br>0E<br>0F | 10<br>11<br>12<br>13<br>14<br>15<br>16<br>17 | 18<br>19<br>1A<br>1B<br>1C<br>1D<br>1E<br>1F | 20<br>21<br>22<br>23<br>24<br>25<br>26<br>27 | 28<br>29<br>2A<br>2B<br>2C<br>2D<br>2E<br>2F | 30<br>31<br>32<br>33<br>34<br>35<br>36<br>37 | 38<br>39<br>3A<br>3B<br>3C<br>3D<br>3E<br>3F |
| [EAX]+disp8$^3$<br>[ECX]+disp8<br>[EDX]+disp8<br>[EBX]+disp8<br>[--][--]+disp8<br>[EBP]+disp8<br>[ESI]+disp8<br>[EDI]+disp8 | 01 | 000<br>001<br>010<br>011<br>100<br>101<br>110<br>111 | 40<br>41<br>42<br>43<br>44<br>45<br>46<br>47 | 48<br>49<br>4A<br>4B<br>4C<br>4D<br>4E<br>4F | 50<br>51<br>52<br>53<br>54<br>55<br>56<br>57 | 58<br>59<br>5A<br>5B<br>5C<br>5D<br>5E<br>5F | 60<br>61<br>62<br>63<br>64<br>65<br>66<br>67 | 68<br>69<br>6A<br>6B<br>6C<br>6D<br>6E<br>6F | 70<br>71<br>72<br>73<br>74<br>75<br>76<br>77 | 78<br>79<br>7A<br>7B<br>7C<br>7D<br>7E<br>7F |
| [EAX]+disp32<br>[ECX]+disp32<br>[EDX]+disp32<br>[EBX]+disp32<br>[--][--]+disp32<br>[EBP]+disp32<br>[ESI]+disp32<br>[EDI]+disp32 | 10 | 000<br>001<br>010<br>011<br>100<br>101<br>110<br>111 | 80<br>81<br>82<br>83<br>84<br>85<br>86<br>87 | 88<br>89<br>8A<br>8B<br>8C<br>8D<br>8E<br>8F | 90<br>91<br>92<br>93<br>94<br>95<br>96<br>97 | 98<br>99<br>9A<br>9B<br>9C<br>9D<br>9E<br>9F | A0<br>A1<br>A2<br>A3<br>A4<br>A5<br>A6<br>A7 | A8<br>A9<br>AA<br>AB<br>AC<br>AD<br>AE<br>AF | B0<br>B1<br>B2<br>B3<br>B4<br>B5<br>B6<br>B7 | B8<br>B9<br>BA<br>BB<br>BC<br>BD<br>BE<br>BF |
| EAX/AX/AL/MM0/XMM0<br>ECX/CX/CL/MM/XMM1<br>EDX/DX/DL/MM2/XMM2<br>EBX/BX/BL/MM3/XMM3<br>ESP/SP/AH/MM4/XMM4<br>EBP/BP/CH/MM5/XMM5<br>ESI/SI/DH/MM6/XMM6<br>EDI/DI/BH/MM7/XMM7 | 11 | 000<br>001<br>010<br>011<br>100<br>101<br>110<br>111 | C0<br>C1<br>C2<br>C3<br>C4<br>C5<br>C6<br>C7 | C8<br>C9<br>CA<br>CB<br>CC<br>CD<br>CE<br>CF | D0<br>D1<br>D2<br>D3<br>D4<br>D5<br>D6<br>D7 | D8<br>D9<br>DA<br>DB<br>DC<br>DD<br>DE<br>DF | E0<br>E1<br>E2<br>E3<br>E4<br>E5<br>E6<br>E7 | E8<br>E9<br>EA<br>EB<br>EC<br>ED<br>EE<br>EF | F0<br>F1<br>F2<br>F3<br>F4<br>F5<br>F6<br>F7 | F8<br>F9<br>FA<br>FB<br>FC<br>FD<br>FE<br>FF |

# Syntax

- Intel
  - [base + index*scale + disp]
    call DWORD PTR [rbx+rsi*4-0xe8]
    mov rax, DWORD PTR [rbp+0x8]
    lea rax, [rbx-0xe8]
- AT&T
  - disp(base, index, scale)
    call *-0xe8(%rbx,%rsi,4)
    mov 0x8(%rbp), %rax
    lea -0xe8(%rbx), %rax

# Intel vs. AT&T syntax

- Intel
  - Destination comes first
    ```
    mov rbp, rsp
    add rax, 0x14
    ```
- AT&T
  - Reverse
    ```
    mov %rsp, %rbp
    add $0x14, %rsp
    ```
  - Registers prefixed with %, immediate $

# Registers

- EIP/RIP
- (E|R)[ABCD]X
  - A: Accumulator
  - B: Base
  - C: Counter
  - D: Data
- ESI, EDI: source and destination pointers for string operations
  - Based off DS in compatibility mode
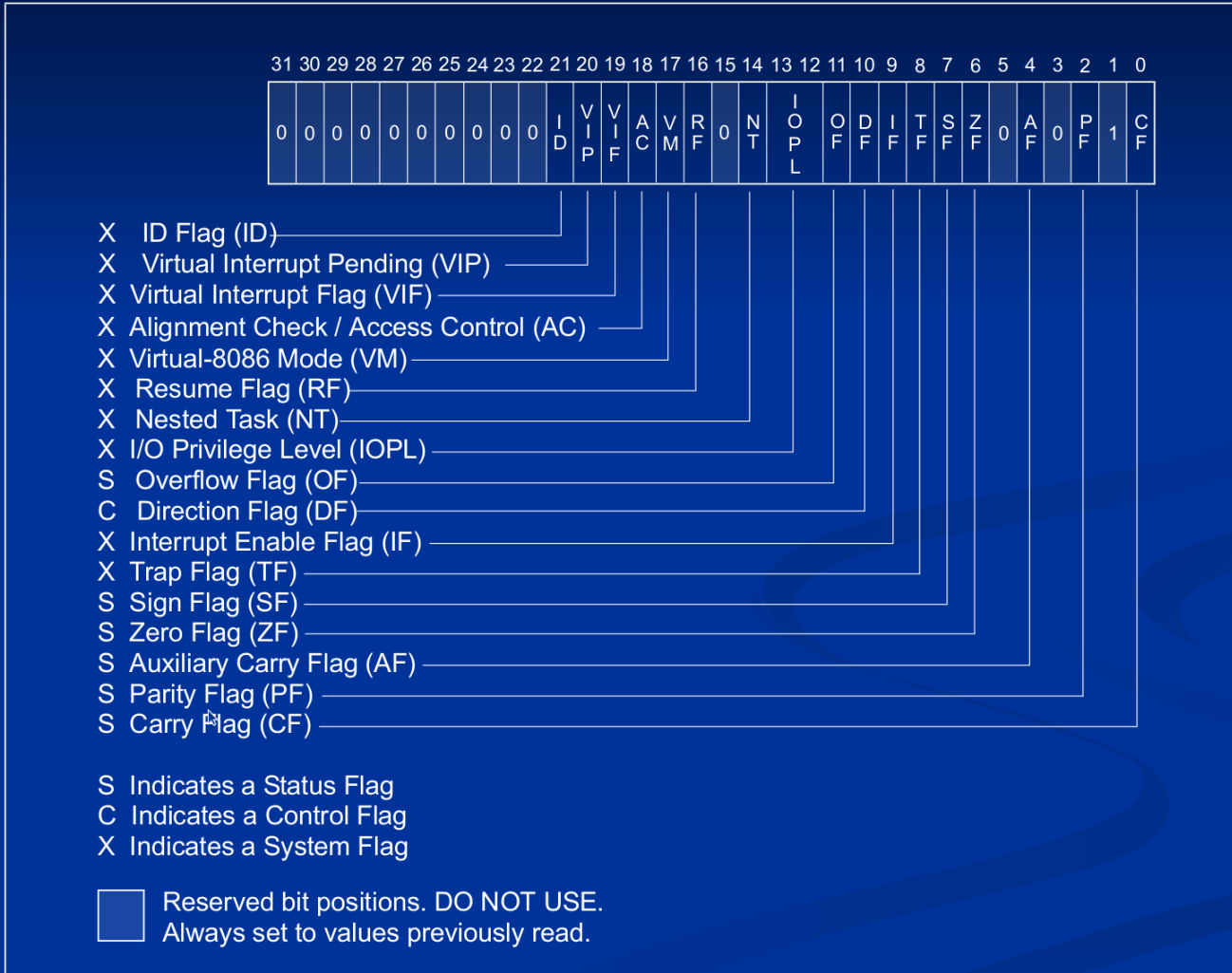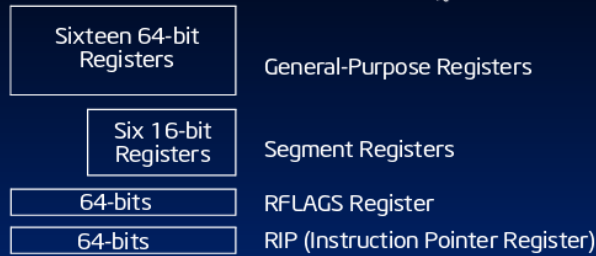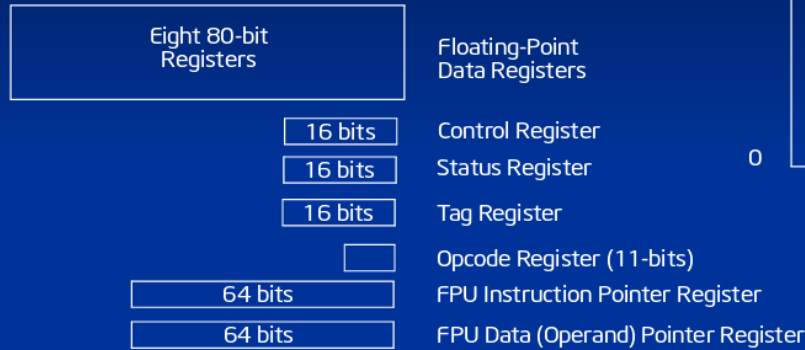- ESP, EBP
  - SS segment

# EFLAGS/RFLAGS

| 31 | 30 | 29 | 28 | 27 | 26 | 25 | 24 | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ID | VIP | VIF | AC | VM | RF | 0 | NT | IOPL | | OF | DF | IF | TF | SF | ZF | 0 | AF | 0 | PF | 1 | CF |

X   ID Flag (ID)
X   Virtual Interrupt Pending (VIP)
X   Virtual Interrupt Flag (VIF)
X   Alignment Check / Access Control (AC)
X   Virtual-8086 Mode (VM)
X   Resume Flag (RF)
X   Nested Task (NT)
X   I/O Privilege Level (IOPL)
S   Overflow Flag (OF)
C   Direction Flag (DF)
X   Interrupt Enable Flag (IF)
X   Trap Flag (TF)
S   Sign Flag (SF)
S   Zero Flag (ZF)
S   Auxiliary Carry Flag (AF)
S   Parity Flag (PF)
S   Carry Flag (CF)

S   Indicates a Status Flag
C   Indicates a Control Flag
X   Indicates a System Flag

Reserved bit positions. DO NOT USE.
Always set to values previously read.

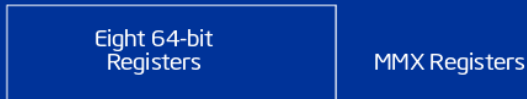Figure 3-8.  EFLAGS Register

**Basic Program Execution Registers**

**Address Space**

| Sixteen 64-bit Registers | General-Purpose Registers |

$2^{64} -1$

| Six 16-bit Registers | Segment Registers |

| 64-bits | RFLAGS Register |
| 64-bits | RIP (Instruction Pointer Register) |

**FPU Registers**

| Eight 80-bit Registers | Floating-Point Data Registers |

0

| 16 bits | Control Register |
| 16 bits | Status Register |
| 16 bits | Tag Register |
| | Opcode Register (11-bits) |
| 64 bits | FPU Instruction Pointer Register |
| 64 bits | FPU Data (Operand) Pointer Register |

**MMX Registers**

**Bounds Registers**

| Eight 64-bit Registers | MMX Registers |

| Four 128-bit Registers |

| BNDCFGU | | BNDSTATUS |

**XMM Registers**

| Sixteen 128-bit Registers | XMM Registers |

| 32-bits | MXCSR Register |

**YMM Registers**

| Sixteen 256-bit Registers | YMM Registers |

© 2017 Dr. Jeffrey A. Turkstra

# Operand Addressing

- Data for a source operand can be found in...
  - The instruction itself (immediate)
  - A register
  - A memory location
  - An I/O port
- A destination operand can be:
  - A register
  - A memory location
  - An I/O port

# Immediate operands

- Example: ADD EAX, 14

- All arithmetic instructions permit an immediate source operand.
- Max value varies, never larger than an unsigned doubleword integer ($2^{32}$)

# Register operands

- 64-bit general-purpose registers:
  - RAX, RBX, RCX, RDX, RSI, RDI, RSP, RBP, R8-R15
- 32-bit general-purpose registers:
  - EAX, EBX, ECX, EDX, ESI, EDI, ESP, EBP, R8D-R15D
- 16-bit general-purpose registers
- 8-bit general-purpose registers
- Segment registers

- RFLAGS
- FPU registers
- MMX, XMM, Control, Debug, and MSR registers
- RDX:RAX register pair (128-bit operand)

# Memory operands

- Segment selector and offset

```
15          0  63                        0
Segment        Offset (or Linear Address)
Selector
```

- 64-bit mode segmentation is generally disabled (flat 64-bit linear address space)
  - CS, DS, ES, SS are 0
  - FS and GS can be used as additional base registers

# Memory offset

- Displacement: 8, 16, or 32-bits
  - Direct, static value
- Base and Index
  - Values from general-purpose registers
- Scale factor
  - 2, 4, or 8
  - Multiplies Index
- RIP + Displacement
- Result is called an effective address

# 64-bit prefix ordering

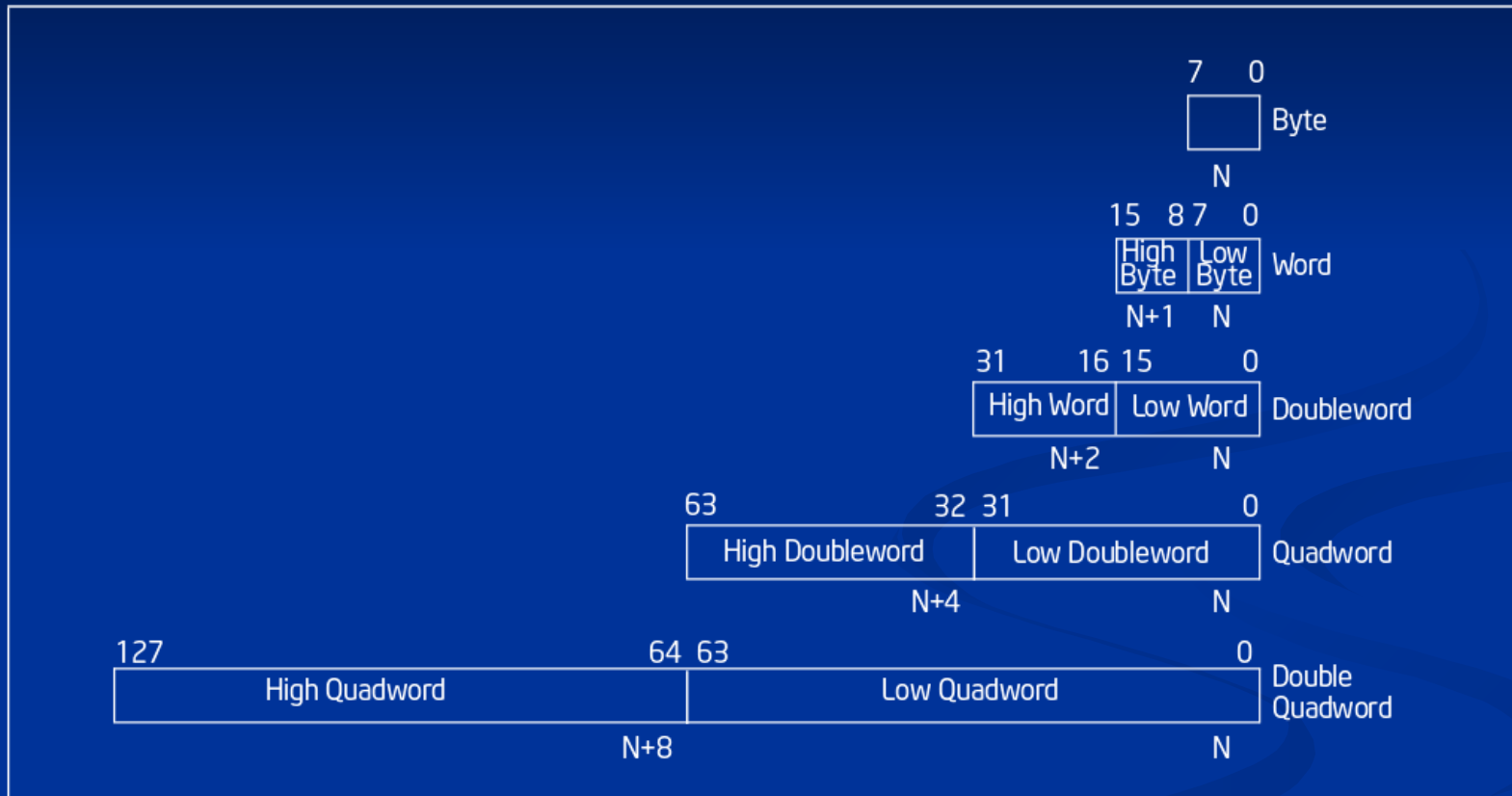| Legacy Prefixes | REX Prefix | Opcode | ModR/M | SIB | Displacement | Immediate |
|---|---|---|---|---|---|---|
| Grp 1, Grp 2, Grp 3, Grp 4 (optional) | (optional) | 1-, 2-, or 3-byte opcode | 1 byte (if required) | 1 byte (if required) | Address displacement of 1, 2, or 4 bytes | Immediate data of 1, 2, or 4 bytes or none |

# SIB?

- Scale
- Index
- Base

# Effective address computation



| Base | | Index | Scale | Displacement |
|------|---|-------|-------|--------------|
| EAX | | EAX | | None |
| EBX | | EBX | 1 | |
| ECX | | ECX | | 8-bit |
| EDX | + | EDX | 2 | + |
| ESP | | EBP | 4 | 16-bit |
| EBP | | ESI | | |
| ESI | | EDI | 8 | 32-bit |
| EDI | | | | |

Offset = Base + (Index * Scale) + Displacement

# Data types

# LEA

- LEA, the only instruction that performs memory addressing calculations but doesn't actually address memory. LEA accepts a standard memory addressing operand, but does nothing more than store the calculated memory offset in the specified register, which may be any general purpose register.

- What does that give us? Two things that ADD doesn't provide:

  - the ability to perform addition with either two or three operands, and

  - the ability to store the result in any register; not just one of the source operands.

# What about 32-bits

- Many systems now are x86_64
- BUT, they can run a lot of 32-bit software
  - "Compatibility mode"
  - Segment registers actually matter
  - Relies on 32-bit registers/addresses/etc
- x86_64 CPUs can switch in and out of compatibility mode with ease
  - Consider system calls for a 64-bit kernel running a 32-bit program

# Instruction set

- Data transfer instructions
- Binary arithmetic
- Decimal arithmetic
- Logical
- Shift and rotate
- Bit and byte
- Control
- String

- Flag control ([ER]FLAG)
- Segment registers
- Miscellaneous

# Data transfer instructions

- Move data between memory and registers
  - Can be conditional
  - Includes stack access
- CMOV and friends
- XCHG
- BSWAP
- PUSH, PUSHA
- POP, POPA

# MOV

- Register to register
- Memory to register
- Register to Memory
- Never memory to memory
  - Remember DMA?

# Binary arithmetic instructions

- Basic binary integer computations
- ADD
- SUB
- IMUL, IDIV
- MUL, DIV
- INC, DEC, NEG
- CMP

# Decimal arithmetic instructions

- Manipulate BCD data
- Invalid in 64-bit mode

# Logical, shift and rotate instructions

- AND, OR, XOR, NOT
- SAR, SHR, SAL, SHL
- ROR, ROL, RCR, RCL

# Bit and byte instructions

- BT
- BTS, BTR
  - Semaphores
- SETE, SETZ and friends
- TEST
- CRC32, POPCNT

# Control transfer instructions

- JMP
- JE, JZ, JNE, JNZ
- CALL, RET
- INT, IRET
- ENTER, LEAVE

# String instructions

- MOVS, MOVSB
  - B/W/D: byte, word, doubleword
- CMPS, CMPSB

# Flag control instructions

- STC, CLC
- STD, CLD
- LAHF, SAHF
- PUSHF, PUSHFD
- POPF, POPFD
- STI
- CLI

# Questions?