

# Programming Assignment 1

695.744

T. McGuire

Johns Hopkins University

## Introduction

Throughout the semester, you will be utilizing tools such as **Ghidra**, **IDA Pro**, **objdump**, **Windbg**, **gdb** and other debugging/disassembling utilities. One feature these tools share is the ability to turn *machine code* into human readable *assembly language*. It is important to have an understanding of how these tools work. In doing so, you will have a better idea of how to make the tool work better for you. The goal of this project is to create a program that can turn *machine code* into *assembly language*.

## Deliverables

1. Brief discussion (about 1 page) on the strengths and weaknesses of the recursive descent and linear sweep algorithms. What makes tools like IDA and Ghidra powerful disassemblers?
2. Your source code for a disassembler for a small subset of the Intel Instruction Set, as described in the remainder of this assignment.

# 1 Requirements

Your disassembler must:

- Be written in any of the following programming languages: C, C++, Go, Rust, Java, Python. Please ask the instructor if you have issues with this requirement.
- Not crash on any (in)valid inputs.
- Use either the linear sweep or recursive descent algorithm. Most students choose linear sweep.
- Print disassembled instructions to standard output.
- Handle jumping/calling forwards and backwards, adding labels where appropriate with the following form (see Example 2 below).

**offset\_XXXXXXXXh:**

- Handle unknown opcodes by printing the address, the byte, and the assembly as follows (see skeleton code for an example).

**db <byte>**

- Work on the supplied examples in addition to other test files that are not supplied.
- Implement only the given opcodes detailed in the Supported Mnemonics section.
- Implements both **SIB** and **MODRM** bytes.
- Have the input file specified using the “**-i**” command-line option.
- Display only addresses, instruction machine code (i.e. the bytes that make up the instruction), disassembled instructions/data, and labels.

# 2 Assumptions

The following assumptions may be made when writing your disassembler:

- The input file is a binary that contains some x86 machine code.
- Code starts at offset 0 in the given file. Do not worry about headers that are generally added by linkers.
- You will start the assignment the week it is assigned.
- If you have any questions, feel free to reach out to me! I’m always happy to help clarify and point you in the right direction.

### 3 Supported Mnemonics

For all instructions below, remember that addressing modes **0b00**, **0b01**, and **0b10** with an **r/m32** operand value of **0b100** indicates that a **SIB** follows the **MODRM** byte. In the Intel SDM this is noted by an Effective Address of **[--][--]**. Examples are given in the next section and you may reference the Intel Instruction Tutorial for more details.

All register references will be 32-bit references. For example, you do not need to handle “**mov dl, byte [ ebx ]**”, you only need to handle “**mov edx, dword [ ebx ]**”. An immediate will be a 32-bit value while the displacement *may* be 8-bits or 32-bits in size. The only exceptions to this are the “**retn imm16**” and “**retf imm16**” instructions.

<b>add</b>	<b>jmp</b>	<b>pop</b>
<b>and</b>	<b>jz/jnz</b>	<b>push</b>
<b>call</b>	<b>lea</b>	<b>repne cmpsd</b>
<b>clflush</b>	<b>mov</b>	<b>retf</b>
<b>cmp</b>	<b>movsd</b>	<b>retn</b>
<b>dec</b>	<b>nop</b>	<b>sub</b>
<b>idiv</b>	<b>not</b>	<b>test</b>
<b>inc</b>	<b>or</b>	<b>xor</b>

## 4 Instruction Details

### 4.1 add/and/mov/not/or/pop/push/sub/etc.

For these and similar instructions, you must implement (where applicable):

```

add r/m32, imm32
add r32, imm32
add r32, [ disp32 ]
add r32, r/m32
add r32, [ r/m32 + disp8 ]
add r32, [ r/m32 + disp32 ]
add r32, [ r/m32*1 + disp32 ]
add r32, [ r/m32*2 + disp32 ]
add r32, [ r/m32*4 + disp32 ]
add r32, [ r/m32*8 + disp32 ]
add r32, [ r/m32*1 + r32 + disp32 ]
add r32, [ r/m32*2 + r32 + disp32 ]
add r32, [ r/m32*4 + r32 + disp32 ]
add r32, [ r/m32*8 + r32 + disp32 ]
add r/m32, r32
add [ disp32 ], imm32
add [ disp32 ], r32
add [ r/m32 ], imm32
add [ r/m32 + disp8 ], imm32
add [ r/m32 + disp32 ], imm32
add [ r/m32*1 + disp32 ], imm32
add [ r/m32*2 + disp32 ], imm32
add [ r/m32*4 + disp32 ], imm32
add [ r/m32*8 + disp32 ], imm32
add [ r/m32*1 + r32 + disp32 ], imm32
add [ r/m32*2 + r32 + disp32 ], imm32
add [ r/m32*4 + r32 + disp32 ], imm32
add [ r/m32*8 + r32 + disp32 ], imm32
add [ r/m32 + disp8 ], r32
add [ r/m32 + disp32 ], r32
add [ r/m32*1 + disp32 ], r32
add [ r/m32*2 + disp32 ], r32
add [ r/m32*4 + disp32 ], r32
add [ r/m32*8 + disp32 ], r32
add [ r/m32*1 + r32 + disp32 ], r32
add [ r/m32*2 + r32 + disp32 ], r32
add [ r/m32*4 + r32 + disp32 ], r32
add [ r/m32*8 + r32 + disp32 ], r32

```

## 4.2 jz/jnz/jmp

You must implement the following:

```
jz rel8
jz rel32
jnz rel8
jnz rel32
jmp rel8
jmp rel32
jmp r/m32
jmp [ disp32 ]
jmp [ r/m32 + disp8 ]
jmp [ r/m32 + disp32 ]
jmp [ r/m32*1 + disp32 ]
jmp [ r/m32*2 + disp32 ]
jmp [ r/m32*4 + disp32 ]
jmp [ r/m32*8 + disp32 ]
jmp [ r/m32*1 + r32 + disp32 ]
jmp [ r/m32*2 + r32 + disp32 ]
jmp [ r/m32*4 + r32 + disp32 ]
jmp [ r/m32*8 + r32 + disp32 ]
```

## 4.3 movsd

Recall that the “d” in “**movsd**” refers to the data size. In this case, it is a DWORD or 32-bit value. Thus, in the Intel Manual we are looking for “**movs m32, m32**”.

*Note: This is a string operation and NOT the “move scalar double-precision” operation.*

## 4.4 repne cmpsd

Recall that the “d” in “**cmpsd**” refers to the data size. In this case, it is a DWORD or 32-bit value. Thus, in the Intel Manual we are looking for “**repne cmps m32, m32**”.

## 4.5 retf/retn

For this instruction family (listed as just “**ret**” in the Intel Instruction Manual), you must implement the following:

```
retf
retf imm16
retn
retn imm16
```

*Note: “**retf**” refers to “return far” and “**retn**” refers to “return near.”*

## 5 Sample Output

The following samples show how the input file is passed to your program and how output is to be formatted. The name of your program and the method by which it is invoked may be different depending on the language you use to implement it.

### Example 1: No jumps

```
$ disasm -i nojump.o
00000000: 31C0          xor eax,eax
00000002: 01C8          add eax,ecx
00000004: 01D0          add eax,edx
00000006: 55           push ebp
00000007: 89E5          mov ebp,esp
00000009: 52           push edx
0000000A: 51           push ecx
0000000B: B844434241    mov eax,0x41424344
00000010: 8B9508000000  mov edx,[ebp+0x00000008]
00000016: 8B8D0C000000  mov ecx,[ebp+0x0000000c]
0000001C: 01D1          add ecx,edx
0000001E: 89C8          mov eax,ecx
00000020: 5A           pop edx
00000021: 59           pop ecx
00000022: 5D           pop ebp
00000023: C20800       retn 0x0008
```

### Example 2: With conditional jump

```
$ disasm -i condjump.o
00000000: 55           push ebp
00000001: 89E5          mov ebp,esp
00000003: 52           push edx
00000004: 51           push ecx
00000005: 39D1          cmp ecx,edx
00000007: 740F          jz offset_00000018h
00000009: B844434241    mov eax,0x41424344
0000000E: 8B5508        mov edx,[ebp+0x00000008]
00000011: 8B4D0C        mov ecx,[ebp+0x0000000c]
00000014: 01D1          add ecx,edx
00000016: 89C8          mov eax,ecx
offset_00000018h:
00000018: 5A           pop edx
00000019: 59           pop ecx
0000001A: 5D           pop ebp
0000001B: C20800       retn 0x0008
```

## 6 Complete Opcode and Addressing Mode Requirements

Mnemonic/Syntax	Opcode	Addressing Modes
<b>add</b> <b>eax</b> , <b>imm32</b>	<b>0x05 id</b>	MODR/M Not Required
<b>add</b> <b>r/m32</b> , <b>imm32</b>	<b>0x81 /0 id</b>	00/01/10/11
<b>add</b> <b>r/m32</b> , <b>r32</b>	<b>0x01 /r</b>	00/01/10/11
<b>add</b> <b>r32</b> , <b>r/m32</b>	<b>0x03 /r</b>	00/01/10/11
<b>and</b> <b>eax</b> , <b>imm32</b>	<b>0x25 id</b>	MODR/M Not Required
<b>and</b> <b>r/m32</b> , <b>imm32</b>	<b>0x81 /4 id</b>	00/01/10/11
<b>and</b> <b>r/m32</b> , <b>r32</b>	<b>0x21 /r</b>	00/01/10/11
<b>and</b> <b>r32</b> , <b>r/m32</b>	<b>0x23 /r</b>	00/01/10/11
<b>call</b> <b>rel32</b>	<b>0xE8 cd</b>	Note: treat <b>cd</b> as <b>id</b>
<b>call</b> <b>r/m32</b>	<b>0xFF /2</b>	00/01/10/11
<b>clflush</b> <b>m8</b>	<b>0x0F 0xAE /7</b>	00/01/10 Note: <b>m8</b> can be a [ <b>disp32</b> ] only, a [ <b>reg</b> ], a [ <b>reg + disp8</b> ], or a [ <b>reg + disp32</b> ]. Addressing mode 11 is illegal.
<b>cmp</b> <b>eax</b> , <b>imm32</b>	<b>0x3D id</b>	MODR/M Not Required
<b>cmp</b> <b>r/m32</b> , <b>imm32</b>	<b>0x81 /7 id</b>	00/01/10/11
<b>cmp</b> <b>r/m32</b> , <b>r32</b>	<b>0x39 /r</b>	00/01/10/11
<b>cmp</b> <b>r32</b> , <b>r/m32</b>	<b>0x3B /r</b>	00/01/10/11
<b>dec</b> <b>r/m32</b>	<b>0xFF /1</b>	00/01/10/11
<b>dec</b> <b>r32</b>	<b>0x48 + rd</b>	MODR/M Not Required
<b>idiv</b> <b>r/m32</b>	<b>0xF7 /7</b>	00/01/10/11
<b>inc</b> <b>r/m32</b>	<b>0xFF /0</b>	00/01/10/11
<b>inc</b> <b>r32</b>	<b>0x40 + rd</b>	MODR/M Not Required

<b>jmp rel8</b>	<b>0xEB cb</b>	Note: treat <b>cb</b> as <b>ib</b>
<b>jmp rel32</b>	<b>0xE9 cd</b>	Note: treat <b>cd</b> as <b>id</b>
<b>jmp r/m32</b>	<b>0xFF /4</b>	00/01/10/11
<b>jz rel8</b>	<b>0x74 cb</b>	Note: treat <b>cb</b> as <b>ib</b>
<b>jz rel32</b>	<b>0x0f 0x84 cd</b>	Note: treat <b>cd</b> as <b>id</b>
<b>jnz rel8</b>	<b>0x75 cb</b>	Note: treat <b>cb</b> as <b>ib</b>
<b>jnz rel32</b>	<b>0x0f 0x85 cd</b>	Note: treat <b>cd</b> as <b>id</b>
<b>lea r32, m</b>	<b>0x8D /r</b>	00/01/10 Note: <b>m</b> can be a <b>[disp32]</b> only, a <b>[reg]</b> , a <b>[reg + disp8]</b> , or a <b>[reg + disp32]</b> . Addressing mode 11 is illegal.
<b>mov r32, imm32</b>	<b>0xB8+rd id</b>	MODR/M Not Required
<b>mov r/m32, imm32</b>	<b>0xC7 /0 id</b>	00/01/10/11
<b>mov r/m32, r32</b>	<b>0x89 /r</b>	00/01/10/11
<b>mov r32, r/m32</b>	<b>0x8B /r</b>	00/01/10/11
<b>movsd</b>	<b>0xA5</b>	MODR/M Not Required
<b>nop</b>	<b>0x90</b>	MODR/M Not Required Note: this is really <b>xchg eax, eax</b>
<b>not r/m32</b>	<b>0xF7 /2</b>	00/01/10/11
<b>or eax, imm32</b>	<b>0x0D id</b>	MODR/M Not Required
<b>or r/m32, imm32</b>	<b>0x81 /1 id</b>	00/01/10/11
<b>or r/m32, r32</b>	<b>0x09 /r</b>	00/01/10/11
<b>or r32, r/m32</b>	<b>0x0B /r</b>	00/01/10/11
<b>pop r/m32</b>	<b>0x8F /0</b>	00/01/10/11
<b>pop r32</b>	<b>0x58 + rd</b>	MODR/M Not Required



<b>push r/m32</b>	<b>0xFF /6</b>	00/01/10/11
<b>push r32</b>	<b>0x50 + rd</b>	MODR/M Not Required
<b>push imm32</b>	<b>0x68 id</b>	MODR/M Not Required
<b>repne cmpsd</b>	<b>0xF2 0xA7</b>	MODR/M Not Required Note: <b>0xF2</b> is the <b>repne</b> prefix
<b>retf</b>	<b>0xCB</b>	MODR/M Not Required
<b>retf imm16</b>	<b>0xCA iw</b>	MODR/M Not Required Note: <b>iw</b> is a 16-bit immediate
<b>retn</b>	<b>0xC3</b>	MODR/M Not Required
<b>retn imm16</b>	<b>0xC2 iw</b>	MODR/M Not Required Note: <b>iw</b> is a 16-bit immediate
<b>sub eax, imm32</b>	<b>0x2D id</b>	MODR/M Not Required
<b>sub r/m32, imm32</b>	<b>0x81 /5 id</b>	00/01/10/11
<b>sub r/m32, r32</b>	<b>0x29 /r</b>	00/01/10/11
<b>sub r32, r/m32</b>	<b>0x2B /r</b>	00/01/10/11
<b>test eax, imm32</b>	<b>0xA9 id</b>	MODR/M Not Required
<b>test r/m32, imm32</b>	<b>0xF7 /0 id</b>	00/01/10/11
<b>test r/m32, r32</b>	<b>0x85 /r</b>	00/01/10/11
<b>xor eax, imm32</b>	<b>0x35 id</b>	MODR/M Not Required
<b>xor r/m32, imm32</b>	<b>0x81 /6 id</b>	00/01/10/11
<b>xor r/m32, r32</b>	<b>0x31 /r</b>	00/01/10/11
<b>xor r32, r/m32</b>	<b>0x33 /r</b>	00/01/10/11