

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Floating Point II, x86-64 Intro

CSE 351 Autumn 2017

**Instructor:**  
Justin Hsia

**Teaching Assistants:**  
 Lucas Wotton   Michael Zhang   Parker DeWilde   Ryan Wong  
 Sam Gehman   Sam Wolfson   Savanna Yee   Vinny Palaniappan

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017


## Administrivia

- ❖ Lab 1 due on Friday (10/13)
  - Submit `bits.c`, `pointer.c`, `lab1reflect.txt`
- ❖ Homework 2 due next Friday (10/20)
  - On Integers, Floating Point, and x86-64
- ❖ Section tomorrow on Integers and Floating Point
- ❖ Peer Instruction Questions are for your benefit!
  - TAs are scattered about as well to help

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Floating point topics

- ❖ Fractional binary numbers
- ❖ IEEE floating-point standard
- ❖ **Floating-point operations and rounding**
- ❖ Floating-point in C
- ❖ There are many more details that we won't cover
  - It's a 58-page standard...



UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

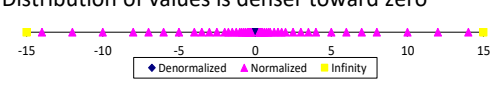
## Floating Point Encoding Summary

Exponent	Mantissa	Meaning
0x00	0	$\pm 0$
0x00	non-zero	$\pm$ denorm num
0x01 – 0xFE	anything	$\pm$ norm num
0xFF	0	$\pm \infty$
0xFF	non-zero	NaN

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017


## Distribution of Values

- ❖ What ranges are NOT representable?
  - Between largest norm and infinity **Overflow** (Exp too large)
  - Between zero and smallest denorm **Underflow** (Exp too small)
  - Between norm numbers? **Rounding**
- ❖ Given a FP number, what's the bit pattern of the next largest representable number?
  - What is this "step" when Exp = 0?
  - What is this "step" when Exp = 100?
- ❖ Distribution of values is denser toward zero



UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Floating Point Operations: Basic Idea

$$\text{Value} = (-1)^S \times 2^{\text{Exponent}} \times \text{Mantissa}$$


- ❖  $x \oplus_{\text{E}} y = \text{Round}(x + y)$
- ❖  $x \otimes_{\text{E}} y = \text{Round}(x * y)$
- ❖ Basic idea for floating point operations:
  - First, **compute the exact result**
  - Then **round** the result to make it fit into desired precision:
    - Possibly over/underflow if exponent outside of range
    - Possibly drop least-significant bits of mantissa to fit into M bit vector



UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Floating Point and the Programmer

```
#include <stdio.h>

int main(int argc, char* argv[]) {
    float f1 = 1.0;
    float f2 = 0.0;
    int i;
    for (i = 0; i < 10; i++)
        f2 += 1.0/10.0;

    printf("0x%08x 0x%08x\n", *(int*)&f1, *(int*)&f2);
    printf("f1 = %10.9f\n", f1);
    printf("f2 = %10.9f\n", f2);

    f1 = 1E30;
    f2 = 1E-30;
    float f3 = f1 + f2;
    printf("f1 == f3? %s\n", f1 == f3 ? "yes" : "no" );

    return 0;
}
```

```
$ ./a.out
0x3f800000 0x3f800001
f1 = 1.000000000
f2 = 1.000000119
f1 == f3? yes
```

13

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Floating Point Summary

- Floats also suffer from the fixed number of bits available to represent them
  - Can get overflow/underflow
  - "Gaps" produced in representable numbers means we can lose precision, unlike ints
    - Some "simple fractions" have no exact representation (e.g. 0.2)
    - "Every operation gets a slightly wrong result"
- Floating point arithmetic not associative or distributive
  - Mathematically equivalent ways of writing an expression may compute different results
- Never** test floating point values for equality!
- Careful** when converting between ints and floats!

14

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Number Representation Really Matters

- 1991:** Patriot missile targeting error
  - clock skew due to conversion from integer to floating point
- 1996:** Ariane 5 rocket exploded (\$1 billion)
  - overflow converting 64-bit floating point to 16-bit integer
- 2000:** Y2K problem
  - limited (decimal) representation: overflow, wrap-around
- 2038:** Unix epoch rollover
  - Unix epoch = seconds since 12am, January 1, 1970
  - signed 32-bit integer representation rolls over to TMin in 2038
- Other related bugs:**
  - 1982: Vancouver Stock Exchange 10% error in less than 2 years
  - 1994: Intel Pentium FDIV (floating point division) HW bug (\$475 million)
  - 1997: USS Yorktown "smart" warship stranded: divide by zero
  - 1998: Mars Climate Orbiter crashed: unit mismatch (\$193 million)

15

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Roadmap

C:

```
car *c = malloc(sizeof(car));
c->miles = 100;
c->gals = 17;
float mpg = get_mpg(c);
free(c);
```

Java:

```
Car c = new Car();
c.setMiles(100);
c.setGals(17);
float mpg = c.getMPG();
```

Assembly language:

```
get_mpg:
    pushq %rbp
    movq %rsp, %rbp
    ...
    popq %rbp
    ret
```

Machine code:

```
01110100000011000
100011010000010000000010
1000100111000010
110000011111101000011111
```

OS:

Windows 10, OS X Yosemite, Linux

Computer system:

Memory & data  
Integers & floats  
**x86 assembly**  
Procedures & stacks  
Executables  
Arrays & structs  
Memory & caches  
Processes  
Virtual memory  
Memory allocation  
Java vs. C

16

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Translation

Code Time      Compile Time      Run Time

```
graph LR
    subgraph Code_Time [Code Time]
        User((User)) --> C[C: User program in C]
    end
    subgraph Compile_Time [Compile Time]
        C --> CC[C compiler]
        CC --> AS[Assembler]
    end
    subgraph Run_Time [Run Time]
        AS --> HW[Hardware]
    end
    C -- ".c file" --> CC
    AS -- ".exe file" --> HW
```

What makes programs run fast(er)?

17

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## HW Interface Affects Performance

Source code: Different applications or algorithms

Compiler: Perform optimizations, generate instructions

Architecture: Instruction set

Hardware: Different implementations

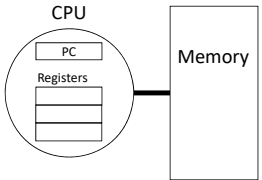
```
graph LR
    subgraph Source_Code [Source code]
        C[C Language]
        PA[Program A]
        PB[Program B]
        YP[Your program]
    end
    subgraph Compiler [Compiler]
        GCC[GCC]
        Clang[Clang]
    end
    subgraph Architecture [Architecture]
        x86[x86-64]
        ARMv8[ARMv8 AArch64/A64]
    end
    subgraph Hardware [Hardware]
        IP4[Intel Pentium 4]
        IC2[Intel Core 2]
        IC7[Intel Core i7]
        AO[AMD Opteron]
        AA[AMD Athlon]
        ACA53[ARM Cortex-A53]
        AA7[Apple A7]
    end
    C --> GCC
    C --> Clang
    GCC --> x86
    GCC --> ARMv8
    Clang --> x86
    Clang --> ARMv8
    x86 --> IP4
    x86 --> IC2
    x86 --> IC7
    x86 --> AO
    x86 --> AA
    ARMv8 --> ACA53
    ARMv8 --> AA7
```

18

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Instruction Set Architectures

- ❖ The ISA defines:
  - The system's state (e.g. registers, memory, program counter)
  - The instructions the CPU can execute
  - The effect that each of these instructions will have on the system state



19

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Instruction Set Philosophies

- ❖ **Complex Instruction Set Computing (CISC):** Add more and more elaborate and specialized instructions as needed
  - Lots of tools for programmers to use, but hardware must be able to handle all instructions
  - x86-64 is CISC, but only a small subset of instructions encountered with Linux programs
- ❖ **Reduced Instruction Set Computing (RISC):** Keep instruction set small and regular
  - Easier to build fast hardware
  - Let software do the complicated operations by composing simpler ones

20

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## General ISA Design Decisions

- ❖ Instructions
  - What instructions are available? What do they do?
  - How are they encoded?
- ❖ Registers
  - How many registers are there?
  - How wide are they?
- ❖ Memory
  - How do you specify a memory location?

21

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Mainstream ISAs

x86	ARM	MIPS
<b>Designer</b> Intel, AMD <b>Bits</b> 16-bit, 32-bit and 64-bit <b>Introduced</b> 1978 (16-bit), 1985 (32-bit), 2003 (64-bit) <b>Design</b> CISC <b>Type</b> Register-memory <b>Encoding</b> Variable (1 to 15 bytes) <b>Endianness</b> Little	<b>Designer</b> ARM Holdings <b>Bits</b> 32-bit, 64-bit <b>Introduced</b> 1985; 31 years ago <b>Design</b> RISC <b>Type</b> Register-Register <b>Encoding</b> AArch64/AArch32 use 32-bit instructions, T32 (Thumb-2) uses mixed 16- and 32-bit instructions. ARMv7 user-space compatibility. <sup>[1]</sup> <b>Endianness</b> Bi (little as default)	<b>Designer</b> MIPS Technologies, Inc. <b>Bits</b> 64-bit (32-64) <b>Introduced</b> 1981; 35 years ago <b>Design</b> RISC <b>Type</b> Register-Register <b>Encoding</b> Fixed <b>Endianness</b> Bi

Macbooks & PCs (Core i3, i5, i7, M)  
[x86-64 Instruction Set](#)

Smartphone-like devices (iPhone, iPad, Raspberry Pi)  
[ARM Instruction Set](#)

Digital home & networking equipment (Blu-ray, PlayStation 2)  
[MIPS Instruction Set](#)

22

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

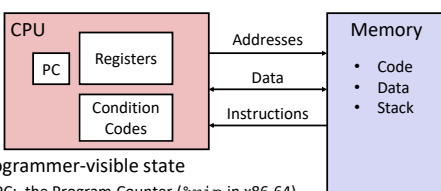
## Definitions

- ❖ **Architecture (ISA):** The parts of a processor design that one needs to understand to write assembly code
  - "What is directly visible to software"
- ❖ **Microarchitecture:** Implementation of the architecture
  - CSE/EE 469, 470
- ❖ Are the following part of the architecture?
  - Number of registers?
  - How about CPU frequency?
  - Cache size? Memory size?

23

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Assembly Programmer's View



- ❖ **Programmer-visible state**
  - PC: the Program Counter (`%rip` in x86-64)
    - Address of next instruction
  - Named registers
    - Together in "register file"
    - Heavily used program data
  - Condition codes
    - Store status information about most recent arithmetic operation
    - Used for conditional branching
- ❖ **Memory**
  - Byte-addressable array
  - Code and user data
  - Includes *the Stack* (for supporting procedures)

24

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## x86-64 Assembly “Data Types”

- ❖ Integral data of 1, 2, 4, or 8 bytes
  - Data values
  - Addresses (untyped pointers)
- ❖ Floating point data of 4, 8, 10 or 2x8 or 4x4 or 8x2
  - Different registers for those (e.g. %xmm1, %ymm2)
  - Come from *extensions to x86* (SSE, AVX, ...)
- ❖ No aggregate types such as arrays or structures
  - Just contiguously allocated bytes in memory
- ❖ Two common syntaxes
  - “AT&T”: used by our course, slides, textbook, gnu tools, ...
  - “Intel”: used by Intel documentation, Intel tools, ...
  - Must know which you’re reading

} Not covered in 351

25

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## What is a Register?

- ❖ A location in the CPU that stores a small amount of data, which can be accessed very quickly (once every clock cycle)
- ❖ Registers have *names*, not *addresses*
  - In assembly, they start with % (e.g. %rsi)
- ❖ Registers are at the heart of assembly programming
  - They are a precious commodity in all architectures, but *especially* x86

26

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## x86-64 Integer Registers – 64 bits wide

%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 2 & 1 bytes)

27

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Some History: IA32 Registers – 32 bits wide

%eax	%ax	%ah	%al	accumulate
%ecx	%cx	%ch	%cl	counter
%edx	%dx	%dh	%dl	data
%ebx	%bx	%bh	%bl	base
%esi	%si			source index
%edi	%di			destination index
%esp	%sp			stack pointer
%ebp	%bp			base pointer

16-bit virtual registers (backwards compatibility)      Name Origin (mostly obsolete)

28

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Memory vs. Registers

❖ Addresses	vs.	Names
▪ 0x7FFFD024C3DC		%rdi
❖ Big	vs.	Small
▪ ~ 8 GiB		(16 x 8 B) = 128 B
❖ Slow	vs.	Fast
▪ ~50-100 ns		sub-nanosecond timescale
❖ Dynamic	vs.	Static
▪ Can “grow” as needed while program runs		fixed number in hardware

29

UNIVERSITY of WASHINGTON L07: Floating Point II, x86-64 Intro CSE351, Autumn 2017

## Three Basic Kinds of Instructions

- 1) Transfer data between memory and register
  - **Load** data from memory into register
    - %reg = Mem[address]
  - **Store** register data into memory
    - Mem[address] = %reg
- 2) Perform arithmetic operation on register or memory data
  - $c = a + b;$        $z = x \ll y;$        $i = h \& g;$
- 3) Control flow: what instruction to execute next
  - Unconditional jumps to/from procedures
  - Conditional branches

**Remember:** Memory is indexed just like an array of bytes!

30

UNIVERSITY of WASHINGTON
L07: Floating Point II, x86-64 Intro
CSE351, Autumn 2017

## Operand types

- ❖ **Immediate:** Constant integer data
  - Examples: `$0x400`, `$-533`
  - Like C literal, but prefixed with ``$'`
  - Encoded with 1, 2, 4, or 8 bytes depending on the instruction
- ❖ **Register:** 1 of 16 integer registers
  - Examples: `%rax`, `%r13`
  - But `%rsp` reserved for special use
  - Others have special uses for particular instructions
- ❖ **Memory:** Consecutive bytes of memory at a computed address
  - Simplest example: `(%rax)`
  - Various other "address modes"

`%rax`

`%rcx`

`%rdx`

`%rbx`

`%rsi`

`%rdi`

`%rsp`

`%rbp`

`%rN`

31

UNIVERSITY of WASHINGTON
L07: Floating Point II, x86-64 Intro
CSE351, Autumn 2017

## Summary

- ❖ Converting between integral and floating point data types *does* change the bits
  - Floating point rounding is a HUGE issue!
    - Limited mantissa bits cause inaccurate representations
    - Floating point arithmetic is NOT associative or distributive
- ❖ x86-64 is a complex instruction set computing (CISC) architecture
- ❖ **Registers** are named locations in the CPU for holding and manipulating data
  - x86-64 uses 16 64-bit wide registers
- ❖ Assembly operands include immediates, registers, and data at specified memory locations

32