

Princeton University

Computer Science 217: Introduction to Programming Systems



Goals of this Lecture

Number Systems and Number Representation

Q: Why do computer programmers confuse Christmas and Halloween?

A: Because $25_{10} = 31_{10}$



1

- Help you learn (or refresh your memory) about:
- The binary, hexadecimal, and octal number systems
 - Finite representation of unsigned integers
 - Finite representation of signed integers
 - Finite representation of rational numbers (if time)

Why?

- A power programmer must know number systems and data representation to fully understand C's primitive data types

Primitive values and the operations on them

2

Agenda

Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational numbers (if time)

3

The Decimal Number System

Name

- "decem" (Latin) \Rightarrow ten

Characteristics

- Ten symbols
 - 0 1 2 3 4 5 6 7 8 9
- Positional
 - $2945 \neq 2495$
 - $2945 = (2*10^3) + (9*10^2) + (4*10^1) + (5*10^0)$

(Most) people use the decimal number system

Why?

4

The Binary Number System

binary

adjective: being in a state of one of two mutually exclusive conditions such as on or off, true or false, molten or frozen, presence or absence of a signal. From Late Latin *bimanus* ("consisting of two").

Characteristics

- Two symbols
 - 0 1
- Positional
 - $1010_B \neq 1100_B$

Most (digital) computers use the binary number system

Terminology

- Bit: a binary digit
- Byte: (typically) 8 bits



5



Decimal Binary

0	0
1	1
2	10
3	11
4	100
5	101
6	110
7	111
8	1000
9	1001
10	1010
11	1011
12	1100
13	1101
14	1110
15	1111
	...

6

Decimal-Binary Conversion

Binary to decimal: expand using positional notation

$$\begin{aligned}100101_8 &= (1 * 2^5) + (0 * 2^4) + (0 * 2^3) + (1 * 2^2) + (0 * 2^1) + (1 * 2^0) \\&= 32 + 0 + 0 + 4 + 0 + 1 \\&= 37\end{aligned}$$

7

Integer Decimal-Binary Conversion

Binary to decimal: expand using positional notation

$$\begin{aligned}100101_8 &= (1 * 2^5) + (0 * 2^4) + (0 * 2^3) + (1 * 2^2) + (0 * 2^1) + (1 * 2^0) \\&= 32 + 0 + 0 + 4 + 0 + 1 \\&= 37\end{aligned}$$

8

Integer
Binary to decimal: expand using positional notation

$$\begin{aligned}100101_8 &= (1 * 2^5) + (0 * 2^4) + (0 * 2^3) + (1 * 2^2) + (0 * 2^1) + (1 * 2^0) \\&= 32 + 0 + 0 + 4 + 0 + 1 \\&= 37\end{aligned}$$

These are integers

They exist as their pure selves
no matter how we might choose
to represent them with our
fingers or toes

8

Integer-Binary Conversion

Integer to binary: do the reverse

- Determine largest power of 2 ≤ number; write template

$$37 = (\underline{2} * 2^5) + (\underline{2} * 2^4) + (\underline{2} * 2^3) + (\underline{2} * 2^2) + (\underline{2} * 2^1) + (\underline{2} * 2^0)$$

- Fill in template

$$\begin{array}{r}37 = (\underline{1} * 2^5) + (0 * 2^4) + (\underline{0} * 2^3) + (\underline{1} * 2^2) + (0 * 2^1) + (\underline{1} * 2^0) \\ \underline{-32} \\ \hline 5 \\ \underline{-4} \\ \hline 1 \\ \underline{-1} \\ \hline 0\end{array}$$

100101₈

9

Integer-Binary Conversion

Integer to binary shortcut

- Repeatedly divide by 2, consider remainder

$$\begin{array}{r}37 / 2 = 18 \text{ R } 1 \\ 18 / 2 = 9 \text{ R } 0 \\ 9 / 2 = 4 \text{ R } 1 \\ 4 / 2 = 2 \text{ R } 0 \\ 2 / 2 = 1 \text{ R } 0 \\ 1 / 2 = 0 \text{ R } 1\end{array}$$

Read from bottom
to top: 100101₈

10

Integer-Binary Conversion

- Repeatedly divide by 2, consider remainder

$$\begin{array}{r}37 / 2 = 18 \text{ R } 1 \\ 18 / 2 = 9 \text{ R } 0 \\ 9 / 2 = 4 \text{ R } 1 \\ 4 / 2 = 2 \text{ R } 0 \\ 2 / 2 = 1 \text{ R } 0 \\ 1 / 2 = 0 \text{ R } 1\end{array}$$

Read from bottom
to top: 100101₈

10

The Hexadecimal Number System

Name

- “hexa” (Greek) ⇒ six
- “dece[m]” (Latin) ⇒ ten

Characteristics

- Sixteen symbols
 - 0 1 2 3 4 5 6 7 8 9 A B C D E F
- Positional
 - A13D_H ≠ 3DA1_H

Computer programmers often use the hexadecimal number system

Why?

Decimal-Hexadecimal Equivalence

Decimal	Hex	Decimal	Hex	Decimal	Hex
0	0	16	10	32	20
1	1	17	11	33	21
2	2	18	12	34	22
3	3	19	13	35	23
4	4	20	14	36	24
5	5	21	15	37	25
6	6	22	16	38	26
7	7	23	17	39	27
8	8	24	18	40	28
9	9	25	19	41	29
10	A	26	1A	42	2A
11	B	27	1B	43	2B
12	C	28	1C	44	2C
13	D	29	1D	45	2D
14	E	30	1E	46	2E
15	F	31	1F	47	2F
					...

11

12

Integer-Hexadecimal Conversion

Hexadecimal to integer: expand using positional notation

$$\begin{aligned}25_{16} &= (2 * 16^1) + (5 * 16^0) \\&= 32 + 5 \\&= 37\end{aligned}$$

Integer to hexadecimal: use the shortcut

$$\begin{array}{r}37 \quad / \quad 16 = 2 \text{ R } 5 \\2 \quad / \quad 16 = 0 \text{ R } 2\end{array}$$

↑ Read from bottom
to top: 25_{16}

Binary-Hexadecimal Conversion

Observation: $16^1 = 2^4$

- Every 1 hexadecimal digit corresponds to 4 binary digits

Binary to hexadecimal

$$\begin{array}{cccccc}101 & 0000 & 1001 & 1110 & 10_8 \\A & 1 & 3 & D_H\end{array}$$

Hexadecimal to binary

$$\begin{array}{cccccc}A & 1 & 3 & D_H \\101 & 0000 & 1001 & 1110 & 10_8\end{array}$$

Is it clear why programmers often use hexadecimal?

14

The Octal Number System

Name

- “octo” (Latin) ⇒ eight

Characteristics

- Eight symbols
 - 0 1 2 3 4 5 6 7
- Positional
 - $1743_8 \neq 7314_8$



Computer programmers often use the octal number system

(So does Mickey Mouse!)

15

Agenda

Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational numbers (if time)

16

Unsigned Data Types: Java vs. C

Java has type:

- int
 - Can represent signed integers
- signed int
 - Can represent signed integers
- int
 - Same as signed int
- unsigned int
 - Can represent only unsigned integers

Mathematics

- Range is 0 to ∞

Computer programming

- Range limited by computer's word size
 - Word size is n bits ⇒ range is 0 to $2^n - 1$
 - Exceed range ⇒ overflow

CourseLab computers

- $n = 64$, so range is 0 to $2^{64} - 1$ (huge!)

Pretend computer

- $n = 4$, so range is 0 to $2^4 - 1$ (15)

Hereafter, assume word size = 4

- All points generalize to word size = 64, word size = n

17

Representing Unsigned Integers

18

Representing Unsigned Integers

On pretend computer

Unsigned Integer	Rep.
0	0000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111
8	1000
9	1001
10	1010
11	1011
12	1100
13	1101
14	1110
15	1111

19

Adding/subtracting binary numbers

Addition

$$\begin{array}{r} 0011 \\ + 1010 \\ \hline \end{array}$$

Subtraction

$$\begin{array}{r} 0011 \\ - 1010 \\ \hline 0111 \end{array}$$

20

Adding Unsigned Integers

Addition

$$\begin{array}{r} 1 \\ 0011_2 \\ + 1010_2 \\ \hline 1101_2 \end{array}$$

Start at right column
Proceed leftward
Carry 1 when necessary

$$\begin{array}{r} 1 \\ 0111_2 \\ + 1010_2 \\ \hline 0001_2 \end{array}$$

Beware of overflow

Results are mod 2^4

21

Subtracting Unsigned Integers

Subtraction

$$\begin{array}{r} 111 \\ 1010_2 \\ - 0111_2 \\ \hline 0011_2 \end{array}$$

Start at right column
Proceed leftward
Borrow when necessary

$$\begin{array}{r} 1 \\ 0011_2 \\ - 1010_2 \\ \hline 1001_2 \end{array}$$

Beware of overflow

Results are mod 2^4

22

Shifting Unsigned Integers

Bitwise right shift ($>>$ in C): fill on left with zeros

$$\begin{array}{r} 10 >> 1 \Rightarrow 5 \\ 1010_2 \quad 0101_2 \\ \hline \end{array}$$

$$\begin{array}{r} 10 >> 2 \Rightarrow 2 \\ 1010_2 \quad 0010_2 \\ \hline \end{array}$$

Bitwise NOT (\sim in C)

$$\begin{array}{r} \sim 10 \Rightarrow 5 \\ 1010_2 \quad 0101_2 \\ \hline \end{array}$$

Bitwise AND ($\&$ in C)

$$\begin{array}{r} \bullet \text{ Logical AND corresponding bits} \\ 10 \quad 1010_2 \\ \& 7 \quad 0111_2 \\ \hline 2 \quad 0010_2 \end{array}$$

23

Other Operations on Unsigned Ints

• Flip each bit

$$\begin{array}{r} \sim 10 \Rightarrow 5 \\ 1010_2 \quad 0101_2 \\ \hline \end{array}$$

• Logical AND corresponding bits

$$\begin{array}{r} 5 << 1 \Rightarrow 10 \\ 0101_2 \quad 1010_2 \\ \hline \end{array}$$

$$\begin{array}{r} 3 << 2 \Rightarrow 12 \\ 0011_2 \quad 1100_2 \\ \hline \end{array}$$

Useful for setting selected bits to 0

$$\begin{array}{r} \bullet \text{ What is the effect arithmetically? (No fair looking ahead)} \\ 5 << 1 \Rightarrow 10 \\ 0101_2 \quad 1010_2 \\ \hline \end{array}$$

$$\begin{array}{r} 3 << 2 \Rightarrow 12 \\ 0011_2 \quad 1100_2 \\ \hline \end{array}$$

Results are mod 2^4

24

Other Operations on Unsigned Ints

Aside: Using Bitwise Ops for Arith

- Bitwise OR: (`|` in C)
- Logical OR corresponding bits

10	1010 _b
	0001 _b
--	---
11	1011 _b

Useful for setting selected bits to 1

Bitwise exclusive OR (`^` in C)

- Logical exclusive OR corresponding bits

10	1010 _b
^	^ 1010 _b
--	---
0	0000 _b

$x \wedge x$ sets all bits to 0

25

26

Can use `<<`, `>>`, and `&` to do some arithmetic efficiently

- $x * 2^y == x << y$ Fast way to multiply by a power of 2
- $x / 2^y == x >> y$ Fast way to divide `unsigned` by power of 2
- $x \% 2^y == x \& (2^y - 1)$ Fast way to mod by a power of 2
- $x \% 4 = 13 \% 2^2 = 13 \& (2^2 - 1) = 13 \& 3 \Rightarrow 1$

13	1101 _b
\& 3	\& 0011 _b
--	---
1	0001 _b

Aside: Example C Program

```
#include <stdio.h>
#include <stdlib.h>
int main(void)
{
    unsigned int n;
    unsigned int count;
    printf("Enter an unsigned integer: ");
    if (scanf ("%u", &n) != 1)
        fprintf(stderr, "Error: Expect unsigned int.\n");
    exit(EXIT_FAILURE);
}
```

```
for (count = 0; n > 0; n = n >> 1)
```

```
    count += (n & 1);
```

```
printf("%u\n", count);
```

```
return 0;
```

```
}
```

```
What does it write?
```

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

Signed Magnitude (cont.)



Ones' Complement



Integer	Rep
-7	1111
-6	1110
-5	1101
-4	1100
-3	1011
-2	1010
-1	1001
0	1000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111

31

- Computing negative**
 $\text{neg}(x) = \text{flip high order bit of } x$
 $\text{neg}(\text{0101}_B) = \text{1101}_B$
 $\text{neg}(\text{1101}_B) = \text{0101}_B$
- Pros and cons**
- + easy for people to understand
 - + symmetric
 - two representations of zero
 - can't use the same "add" algorithm for both signed and unsigned numbers

32

Integer	Rep
-7	1000
-6	1001
-5	1010
-4	1011
-3	1100
-2	1101
-1	1110
0	1111
1	0000
2	0001
3	0010
4	0011
5	0100
6	0101
7	0110

32

Definition

$$\begin{aligned}\text{High-order bit has weight -7} \\ \text{1010}_B &= (1 * -7) + (0 * 4) + (1 * 2) + (0 * 1) \\ &= -5\end{aligned}$$
$$\begin{aligned}\text{0010}_B &= (0 * -7) + (0 * 4) + (1 * 2) + (0 * 1) \\ &= 2\end{aligned}$$

Ones' Complement (cont.)



Computing negative

$$\begin{aligned}\text{neg}(x) &= \sim X \\ \text{neg}(\text{0101}_B) &= \text{1010}_B \\ \text{neg}(\text{1010}_B) &= \text{0101}_B\end{aligned}$$

Computing negative (alternative)

$$\begin{aligned}\text{neg}(x) &= 1111_B - X \\ \text{neg}(\text{0101}_B) &= 1111_B - \text{0101}_B \\ &= \text{1010}_B\end{aligned}$$
$$\begin{aligned}\text{neg}(\text{1010}_B) &= 1111_B - \text{1010}_B \\ &= \text{0101}_B\end{aligned}$$

Pros and cons

- + symmetric
- two reps of zero
- can't use the same "add" algorithm for both signed and unsigned numbers

33

Integer	Rep
-8	1000
-7	1001
-6	1010
-5	1011
-4	1100
-3	1101
-2	1110
-1	1111
0	0000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111

33



Definition

$$\begin{aligned}\text{High-order bit has weight -8} \\ \text{1010}_B &= (1 * -8) + (0 * 4) + (1 * 2) + (0 * 1) \\ &= -6\end{aligned}$$
$$\begin{aligned}\text{0010}_B &= (0 * -8) + (0 * 4) + (1 * 2) + (0 * 1) \\ &= 2\end{aligned}$$

34

Two's Complement



Definition

$$\begin{aligned}\text{High-order bit has weight -8} \\ \text{1010}_B &= (1 * -8) + (0 * 4) + (1 * 2) + (0 * 1) \\ &= -6\end{aligned}$$
$$\begin{aligned}\text{0010}_B &= (0 * -8) + (0 * 4) + (1 * 2) + (0 * 1) \\ &= 2\end{aligned}$$

34

Two's Complement (cont.)



- Almost all computers use two's complement to represent signed integers
- Why?
- Arithmetic is easy
 - Will become clear soon
- Hereafter, assume two's complement representation of signed integers

35

Integer	Rep
-8	1000
-7	1001
-6	1010
-5	1011
-4	1100
-3	1101
-2	1110
-1	1111
0	0000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111

35

Computing negative

$$\begin{aligned}\text{neg}(x) &= \sim x + 1 \\ \text{neg}(\text{0101}_B) &= \text{onescomp}(x) + 1 \\ \text{neg}(\text{1010}_B) &= 1010_B + 1 = \text{1011}_B \\ \text{neg}(\text{1011}_B) &= 0100_B + 1 = \text{0101}_B\end{aligned}$$

Pros and cons

- not symmetric
- + one representation of zero
- + same algorithm adds unsigned numbers or signed numbers

36

Adding Signed Integers



Subtracting Signed Integers

$$\begin{array}{r} \text{pos + pos} \\ \hline \begin{array}{r} 11 \\ 3 + 3 \\ \hline - \\ 6 \end{array} \end{array}$$

$$\begin{array}{r} \text{pos + neg} \\ \hline \begin{array}{r} 1111 \\ 0011_0 \\ + 1111_0 \\ \hline - \\ 2 \end{array} \end{array}$$

$$\begin{array}{r} \text{neg + neg} \\ \hline \begin{array}{r} 11 \\ -3 + -2 \\ \hline -5 \end{array} \end{array}$$

$$\begin{array}{r} \text{pos + pos (overflow)} \\ \hline \begin{array}{r} 111 \\ 7 \\ + 1 \\ \hline -8 \end{array} \end{array}$$

$$\begin{array}{r} \text{pos + neg (overflow)} \\ \hline \begin{array}{r} 11 \\ -6 + -5 \\ \hline 5 \end{array} \end{array}$$

$$\begin{array}{r} \text{neg + neg (overflow)} \\ \hline \begin{array}{r} 11 \\ -6 + -5 \\ \hline 5 \end{array} \end{array}$$

Perform subtraction with borrows
or
Compute two's comp and add

$$\begin{array}{r} \text{pos + pos} \\ \hline \begin{array}{r} 11 \\ 3 + 3 \\ \hline - \\ 6 \end{array} \end{array} \quad \begin{array}{r} \text{pos + neg} \\ \hline \begin{array}{r} 1111 \\ 0011_0 \\ + 1111_0 \\ \hline - \\ 2 \end{array} \end{array} \quad \begin{array}{r} \text{neg + neg} \\ \hline \begin{array}{r} 11 \\ -3 + -2 \\ \hline -5 \end{array} \end{array}$$

$$\begin{array}{r} \text{pos + pos} \\ \hline \begin{array}{r} 111 \\ 7 \\ + 1 \\ \hline -8 \end{array} \end{array} \quad \begin{array}{r} \text{pos + neg} \\ \hline \begin{array}{r} 11 \\ -6 + -5 \\ \hline 5 \end{array} \end{array} \quad \begin{array}{r} \text{neg + neg} \\ \hline \begin{array}{r} 11 \\ -6 + -5 \\ \hline 5 \end{array} \end{array}$$

How would you detect overflow programmatically?

See Bryant & O'Hallaron book for much more info

Negating Signed Integers



Question: Why does two's comp arithmetic work?

Answer: $[-b] \bmod 2^4 = [\text{twoscomp}(b)] \bmod 2^4$

$$\begin{aligned} [-b] \bmod 2^4 &= [2^4 - b] \bmod 2^4 \\ &= [2^4 - 1 - b + 1] \bmod 2^4 \\ &= [(2^4 - 1 - b) + 1] \bmod 2^4 \\ &= [\text{onescomp}(b) + 1] \bmod 2^4 \\ &= [\text{twoscomp}(b)] \bmod 2^4 \end{aligned}$$

See Bryant & O'Hallaron book for much more info

39

Subtracting Signed Integers: Math



And so:
 $[a - b] \bmod 2^4 = [a + \text{twoscomp}(b)] \bmod 2^4$

$$\begin{aligned} [a - b] \bmod 2^4 &= [a + 2^4 - b] \bmod 2^4 \\ &= [a + 2^4 - 1 - b + 1] \bmod 2^4 \\ &= [a + (2^4 - 1 - b) + 1] \bmod 2^4 \\ &= [a + \text{onescomp}(b) + 1] \bmod 2^4 \\ &= [a + \text{twoscomp}(b)] \bmod 2^4 \end{aligned}$$

See Bryant & O'Hallaron book for much more info

40

Shifting Signed Integers



Bitwise left shift (<< in C): fill on right with zeros

$$\begin{array}{r} 3 << 1 \Rightarrow 6 \\ 0011_0 \\ 0110_0 \end{array}$$

$$\begin{array}{r} -3 << 1 \Rightarrow -6 \\ 1101_0 \\ -1010_0 \end{array}$$

Bitwise arithmetic right shift: fill on left with sign bit

$$\begin{array}{r} 6 >> 1 \Rightarrow 3 \\ 0110_0 \\ 0011_0 \end{array}$$

$$\begin{array}{r} -6 >> 1 \Rightarrow -3 \\ 1010_0 \\ -1101_0 \end{array}$$

Results are mod 2^4

Shifting Signed Integers (cont.)



Bitwise logical right shift: fill on left with zeros

$$\begin{array}{r} 6 >> 1 \Rightarrow 3 \\ 0110_0 \\ 0011_0 \end{array}$$

$$\begin{array}{r} -6 >> 1 \Rightarrow 5 \\ 1010_0 \\ 0101_0 \end{array}$$

In C, right shift (>) could be logical or arithmetic

- Not specified by C90 standard
 - Compiler designer decides

Best to avoid shifting signed integers

(if you must shift signed integers, make sure you're on a 2's complement machine, and do a bitwise-and after shifting)
(Java does this better, with two operators: >> >>>)



What is the effect arithmetically?

What is the effect arithmetically?

What is the effect arithmetically?

41

42

Shifting Signed Integers (cont.)



Other Operations on Signed Ints



Bitwise NOT (`~` in C)

- Same as with unsigned ints

Bitwise AND (`&` in C)

- Same as with unsigned ints

Bitwise OR: (`|` in C)

- Same as with unsigned ints

Bitwise exclusive OR (`^` in C)

- Same as with unsigned ints

(if you must shift signed integers, make sure you're on a 2's complement machine, and do a bitwise-and after shifting)

43

Is it after 1980?
OK, then we're surely
two's complement



Agenda

Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational numbers (if time)

45



Rational Numbers

Mathematics

- A **rational** number is one that can be expressed as the ratio of two integers
 - Unbounded range and precision
- Finite range and precision
 - Approximate using **floating point** number
 - Binary point "floats" across bits

46



Rational Numbers

IEEE Floating Point Representation



Common finite representation: IEEE floating point

- More precisely: ISO/IEEE 754 standard

Using 32 bits (type float in C):

- 1 bit: sign (0⇒positive, 1⇒negative)
- 8 bits: exponent + 127
- 23 bits: binary fraction of the form $1.ddd\ldots dddddd$
- Using 64 bits (type double in C):
 - 1 bit: sign (0⇒positive, 1⇒negative)
 - 11 bits: exponent + 1023
 - 52 bits: binary fraction of the form $1.ddd\ldots dddddd$

- 1.10110110000000000000000000000000_B
- $1 + (1 * 2^{-1}) + (0 * 2^{-2}) + (1 * 2^{-3}) + (1 * 2^{-4}) + (0 * 2^{-5}) + (1 * 2^{-6}) + (1 * 2^{-7}) = 1.7109375$

Number:

$$\bullet -1.7109375 * 2^4 = -27.375$$

Sign (1 bit):

- 1 ⇒ negative

Exponent (8 bits):

$$\bullet 10000011_8 = 131$$
$$\bullet 131 - 127 = 4$$

Fraction (23 bits):

$$\bullet \text{also called "mantissa"} \\ \bullet 1.10110110000000000000000000000000_8 \\ \bullet 1 + (1 * 2^{-1}) + (0 * 2^{-2}) + (1 * 2^{-3}) + (1 * 2^{-4}) + (0 * 2^{-5}) + (1 * 2^{-6}) + (1 * 2^{-7}) = 1.7109375$$

47

11000001110110000000000000000000

32-bit representation

48

When was floating-point invented?

Answer: long before computers!

mantissa

noun
decimal part of a logarithm, 1865, from Latin *mantissa* "a worthless addition, make-weight," perhaps a Gaulish word introduced into Latin via Etruscan (cf. Old Irish *meit*, Welsh *maint* "size").

COMMON LOGARITHMS Log ₁₀ x										
x	0	1	2	3	4	5	6	7	8	
.50	-0.6990	0.6998	1.097	1.7016	2.024	2.033	2.042	2.050	2.059	2.067
.51	-0.6976	0.6984	1.093	1.701	2.010	2.018	2.026	2.035	2.043	2.052
.52	-0.6960	0.6968	1.088	1.717	2.015	2.019	2.023	2.030	2.038	2.045
.53	-0.6943	0.6943	1.081	1.729	2.016	2.025	2.030	2.039	2.047	2.054
.54	-0.6924	0.6924	1.074	1.740	2.048	2.036	2.064	2.072	2.080	2.088
.55	-0.6904	0.6904	1.067	1.749	2.047	2.045	2.063	2.071	2.079	2.096
.56	-0.6882	0.6882	1.060	1.749	2.047	2.045	2.062	2.070	2.078	2.094
.57	-0.6859	0.6859	1.053	1.754	2.052	2.049	2.067	2.075	2.083	2.102
.58	-0.6834	0.6834	1.046	1.764	2.049	2.052	2.064	2.072	2.079	2.098
.59	-0.6809	0.6809	1.039	1.773	2.051	2.057	2.067	2.075	2.083	2.107

Floating Point Warning

- Decimal number system can represent only some rational numbers with finite digit count
 - Example: 1/3

- Binary number system can represent only some rational numbers with finite digit count
 - Example: 1/5

Beware of roundoff error

- Error resulting from inexact representation
 - Can accumulate

CS

Stanford

CS

Stanford</