

1	2	4	8	16	32	64	128	256	512	1024
2 ⁰	2 ¹	2 ²	2 ³	2 ⁴	2 ⁵	2 ⁶	2 ⁷	2 ⁸	2 ⁹	2 ¹⁰

CS 61C

Number Representation

Fall 2018

Discussion 1: August 27, 2018

Notes

1 Unsigned Integers

lets represent

$$1028 = 1 \cdot 10^3 + 0 \cdot 10^2 + 2 \cdot 10^1 + 8 \cdot 10^0$$

[1.1] If we have an n -digit unsigned numeral $d_{n-1}d_{n-2}\dots d_0$ in radix (or base) r , then the value of that numeral is $\sum_{i=0}^{n-1} r^i d_i$, which is just fancy notation to say that instead of a 10's or 100's place we have an r 's or r^2 's place. For the three radices, binary, decimal, and hex, we just let r be 2, 10, and 16, respectively.

We don't have calculators during exams, so let's try this by hand. Recall that our preferred tool for writing large numbers is the IEC prefixing system:

- Ki (Kibi) = 2^{10} · Gi (Gibi) = 2^{30} · Pi (Pebi) = 2^{50} · Zi (Zebi) = 2^{70}
- Mi (Mebi) = 2^{20} · Ti (Tebi) = 2^{40} · Ei (Exbi) = 2^{60} · Yi (Yobi) = 2^{80}

(a) Convert the following numbers from their initial radix into the other two common radices:

mon radices:

1) 0b10010011 = 147 = 0x93

$$1 + 2 + 2^4 + 2^7 = 1 + 2 + 16 + 128 = 147$$

2) 63 = 0b00111111 = 0x3F

Trick: know $64 = 2^6 = 0100\ 0000$ so $64-1 = 0011\ 1111$

3) 0b00100100 = 36 = 0x24

$$2^2 + 2^5 = 32 + 4 = 36$$

4) 0 = 0b0 = 0x0

$$32 + 16 = 48; 32 + 8 = 40; 32 + 4 = 36$$

5) 39 = 0b0010 0111 = 0x27

$$36 + 2 = 38; 38 + 1 = 39 \text{ so } 2^0 + 2^1 + 2^2 + 2^5$$

6) 437 = 0b0001 1011 0101 = 0x1B5

Do divex

7) 0x0123 = 0b0000 0001 0010 0011 = 291

$$2^0 + 2^1 + 2^5 + 2^8 = 1 + 2 + 32 + 256 = 291$$

(b) Convert the following numbers from hex to binary:

1) 0xD3AD = 0b1101 0011 1010 1101 = 54189

2) 0xB33F = 0b1011 0011 0011 1111 = 45887

3) 0x7EC4 = 0b0111 1110 1100 0100 = 32452

(c) Write the following numbers using IEC prefixes:

1. $2^{16} = 64 \text{ Ki}$

2. $2^{27} = 128 \text{ Mi}$

3. $2^{43} = 8 \text{ Ti}$

4. $2^{36} = 64 \text{ Gi}$

5. $2^{34} = 16 \text{ Gi}$

6. $2^{61} = 2 \text{ Ei}$

7. $2^{47} = 128 \text{ Ti}$

8. $2^{58} = 256 \text{ Pi}$

(d) Write the following numbers as powers of 2:

1. $2 \text{ Ki} = 2^{11}$

2. $512 \text{ Ki} = 2^{19}$

3. $16 \text{ Mi} = 2^{24}$

4. $256 \text{ Pi} = 2^{58}$

5. $64 \text{ Gi} = 2^{36}$

6. $128 \text{ Ei} = 2^{67}$

$$\begin{array}{r} 218 \\ 2 \overline{) 437} \\ \underline{86} \\ 17 \\ \underline{16} \\ 1 \end{array}$$

LSB

$$\begin{array}{r} 109 \\ 2 \overline{) 218} \\ \underline{218} \\ 0 \end{array}$$

$$\begin{array}{r} 54 \\ 2 \overline{) 109} \\ \underline{108} \\ 1 \end{array}$$

$$\begin{array}{r} 27 \\ 2 \overline{) 54} \\ \underline{54} \\ 0 \end{array}$$

$$\begin{array}{r} 13 \\ 2 \overline{) 27} \\ \underline{26} \\ 1 \end{array}$$

$$\begin{array}{r} 6 \\ 2 \overline{) 13} \\ \underline{12} \\ 1 \end{array}$$

$$\begin{array}{r} 3 \\ 2 \overline{) 6} \\ \underline{6} \\ 0 \end{array}$$

$$\begin{array}{r} 1 \\ 2 \overline{) 3} \\ \underline{2} \\ 1 \end{array}$$

$$\begin{array}{r} 0 \\ 2 \overline{) 1} \\ \underline{0} \\ 1 \end{array}$$

MSB

All hold zero

Use floor division + get remainder. Read from

So 00110110101 This is general idea from going from higher radix to lower radix.

Hex \rightarrow Binary: group into 4 ones + convert directly.
 Binary to hex:

A=10	D=13	1010 = 8+2=10=A
B=11	E=14	B=11 = 1011
C=12	F=15	

Alternative for 437:

$$\begin{array}{r} 512 - 437 = 75 \\ \downarrow \\ 100000000 \end{array}$$

$$\begin{array}{r} 75 \\ \downarrow \\ 1001011 \end{array}$$

$$\begin{array}{r} 512 - 75 = 437 \\ \begin{array}{r} 100000000 \\ + 1001011 \\ \hline 1001011011 \end{array} \end{array}$$

2 Signed Integers

2.1 Unsigned binary numbers work for natural numbers, but many calculations use negative numbers as well. To deal with this, a number of different schemes have been used to represent signed numbers, but we will focus on two's complement, as it is the standard solution for representing signed integers.

- Most significant bit has a negative value, all others are positive. So the value of an n -digit two's complement number can be written as $\sum_{i=0}^{n-2} 2^i d_i - 2^{n-1} d_{n-1}$.
- Otherwise exactly the same as unsigned integers.
- A neat trick for flipping the sign of a two's complement number: flip all the bits and add 1.
- Addition is exactly the same as with an unsigned number.
- Only one 0, and it's located at 0b0.

For questions (a) through (c), assume an 8-bit integer and answer each one for the case of an unsigned number, biased number with a bias of -127, and two's complement number. Indicate if it cannot be answered with a specific representation.

(a) What is the largest integer? The largest integer's representation + 1?

1. Unsigned? 255, 0

2. Biased? 128, -127

3. Two's Complement? 127, -128

(b) How would you represent the numbers 0, 1, and -1?

1. Unsigned? 0b0000 0000, 0b0000 0001, N/A

2. Biased? 0b0111 1111, 0b1000 0000, 0b0111 1110

3. Two's Complement? 0b0000 0000, 0b0000 0001, 0b1111 1111

(c) How would you represent 17 and -17?

1. Unsigned? 0b0001 0001, N/A

2. Biased? 0b1001 0000, 0b0110 1110

3. Two's Complement? 0b0001 0001, 0b1110 1111

(d) What is the largest integer that can be represented by any encoding scheme that only uses 8 bits?

There is no such integer. For example, an arbitrary 8-bit mapping could choose to represent the numbers from 1 to 256 instead of 0 to 255.

(e) Prove that the two's complement inversion trick is valid (i.e. that x and $\bar{x} + 1$ sum to 0).

Note that for any x we have $x + \bar{x} = 0b1 \dots 1$. A straightforward hand calculation shows that $0b1 \dots 1 + 0b1 = 0$.

$$\text{Ex: } -17: x = 11101111 \quad \bar{x} = 00010000$$

$$x + \bar{x} = 11111111 \quad + \quad x + \bar{x} + 1 = 11111111 + 00000001 = 100000000 \leftarrow \text{all zeros}$$

8 bits means

2^8 possible numbers
= 256 numbers

$$\text{largest number} + \text{bias} = 255 - 127 = 128$$

Two's complement has one positive zero, zero takes up the extra positive number.

cannot represent negative number in unsigned.

$$0: 0 - \text{bias} = -(-127) = 127 \quad 1 - (-127) = 128 \quad -1 - (-127) = 126$$

$$2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0 = 128 + 127 = -1$$

$$\text{or flip bits of } 00000001 \text{ and add } 1$$

$$00000001 \rightarrow 11111110 + 1 = 11111111$$

$$17 =$$

$$00010001$$

$$\bar{x} = 11101110 + 1$$

$$-17 = 11101111$$

- (f) Explain where each of the three radices shines and why it is preferred over other bases in a given context.

Decimal is the preferred radix for human hand calculations, likely related to the fact that humans have 10 fingers.

Binary numerals are particularly useful for computers. Binary signals are less likely to be garbled than higher radix signals, as there is more distance (voltage or current) between valid signals. Additionally, binary signals are quite convenient to design circuits with, as well see later in the course.

Hexadecimal numbers are a convenient shorthand for displaying binary numbers, owing to the fact that one hex digit corresponds exactly to four binary digits.

3 Counting

Do: 60 over 2
3 min + 2 min

3.1 Bitstrings can be used to represent more than just numbers. In fact, we use bitstrings to represent *everything* inside a computer. And, because we don't want to be wasteful with bits it is important that to remember that n bits can be used to represent 2^n distinct things. For each of the following questions, answer with the minimum number of bits possible.

- (a) How many bits do we need to represent a variable that can only take on the values 0, π or e ? $1 \text{ bit} = 2 \text{ options}$

2

$2 \text{ bits} = 4 \text{ options}$

- (b) If we need to address 3 TiB of memory and we want to address every byte of memory, how long does an address need to be?

42 bits

$2^0 < 3 \text{ so } 2^1 < 3 \text{ so } 2^2 > 3 \text{ so } 2^2 \cdot 2^{40} = 2^{42}$
we need to be able to have at least 3 TiB so cannot use 2^0 or 2^1 that's why we use 2^2

- (c) If the only value a variable can take on is e , how many bits are needed to represent it?

0 This may sound weird since you would need a bit to represent this. The reason it is zero is because you always know the value so you do not need a bit to represent this. & since the variable can take only one value, we do not need any bits to represent the single value