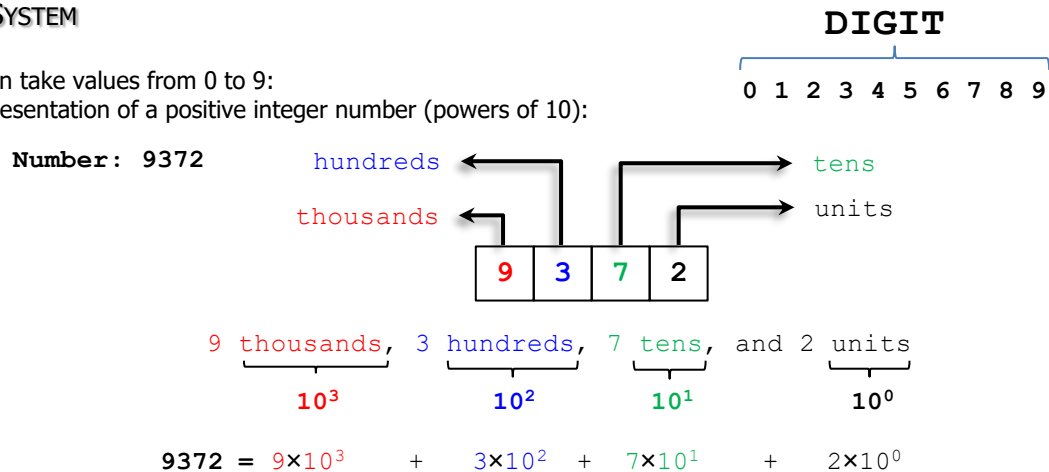


Unit 4 – Number Systems and Computer Arithmetic

UNSIGNED INTEGER NUMBERS

DECIMAL NUMBER SYSTEM

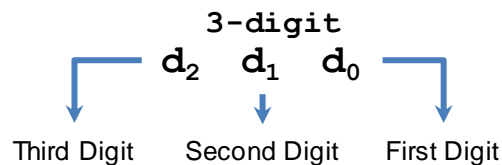
- A decimal digit can take values from 0 to 9:
- Digit-by-digit representation of a positive integer number (powers of 10):



POSITIONAL NUMBER REPRESENTATION

- Let's consider the numbers from 0 to 999. We represent these numbers with 3 digits (each *digit* being a number between 0 and 9). We show a 3-digit number using the positional number representation:

MATHEMATICAL REPRESENTATION



EXAMPLE



- The *positional number representation* allows us to express the decimal value using **powers of ten**: $d_2 \times 10^2 + d_1 \times 10^1 + d_0 \times 10^0$. Example:

Decimal Number	3-digit representation $d_2 d_1 d_0$	Powers of 10: $d_2 \times 10^2 + d_1 \times 10^1 + d_0 \times 10^0$
0	000	$0 \times 10^2 + 0 \times 10^1 + 0 \times 10^0$
9	009	$0 \times 10^2 + 0 \times 10^1 + 9 \times 10^0$
11	011	$0 \times 10^2 + 1 \times 10^1 + 1 \times 10^0$
25	025	$0 \times 10^2 + 2 \times 10^1 + 5 \times 10^0$
90	090	$0 \times 10^2 + 9 \times 10^1 + 0 \times 10^0$
128	128	$1 \times 10^2 + 2 \times 10^1 + 8 \times 10^0$
255	255	$2 \times 10^2 + 5 \times 10^1 + 5 \times 10^0$

Exercise: Write down the 3-digit and the powers of ten representations for the following numbers:

Decimal Number	3-digit representation	$d_2 \times 10^2 + d_1 \times 10^1 + d_0 \times 10^0$
5		
254		
100		
99		

General Case:

- Positional number representation for an integer positive number with n digits: $d_{n-1}d_{n-2} \dots d_1d_0$
Decimal Value:

$$D = \sum_{i=0}^{n-1} d_i \times 10^i = d_{n-1} \times 10^{n-1} + d_{n-2} \times 10^{n-2} + \dots + d_1 \times 10^1 + d_0 \times 10^0$$

- Example: 1098324 (7 digits). $1098324 = 1 \times 10^6 + 0 \times 10^5 + 9 \times 10^4 + 8 \times 10^3 + 3 \times 10^2 + 2 \times 10^1 + 4 \times 10^0$
203476 (6 digits). $203476 = 2 \times 10^5 + 0 \times 10^4 + 3 \times 10^3 + 4 \times 10^2 + 7 \times 10^1 + 6 \times 10^0$

Maximum value:

- The table presents the maximum attainable value for a given number of digits. What pattern do you find? Can you complete it for the highlighted cases (4 and 6)?

Number of digits	Maximum value	Range
1	$9 = 10^1 - 1$	$0 \rightarrow 9 \quad \equiv 0 \rightarrow 10^1 - 1$
2	$99 = 10^2 - 1$	$0 \rightarrow 99 \quad \equiv 0 \rightarrow 10^2 - 1$
3	$999 = 10^3 - 1$	$0 \rightarrow 999 \quad \equiv 0 \rightarrow 10^3 - 1$
4		
5	$99999 = 10^5 - 1$	$0 \rightarrow 99999 \quad \equiv 0 \rightarrow 10^5 - 1$
6		
...		
n	$999\dots999 = 10^n - 1$	$0 \rightarrow 999\dots999 \equiv 0 \rightarrow 10^n - 1$

- Maximum value for a number with 'n' digits:** Based on the table, the maximum decimal value for a number with 'n' digits is given by:

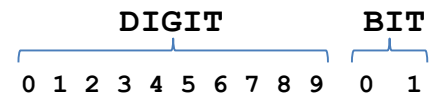
$$D = \underbrace{999\dots999}_{n \text{ digits}} = 9 \times 10^{n-1} + 9 \times 10^{n-2} + \dots + 9 \times 10^1 + 9 \times 10^0 = 10^n - 1$$

⇒ With 'n' digits, we can represent 10^n positive integer numbers from 0 to $10^n - 1$.

- With 7 digits, what is the range (starting from 0) of positive numbers that we can represent? How many different numbers can we represent?

BINARY NUMBER SYSTEM

- We are used to the decimal number system. However, there exist other number systems: octal, hexadecimal, vigesimal, binary, etc. In particular, binary numbers are very practical as they are used by digital computers. For binary numbers, the counterpart of the decimal digit (that can take values from 0 to 9) is the binary digit, or bit (that can take the value of 0 or 1).



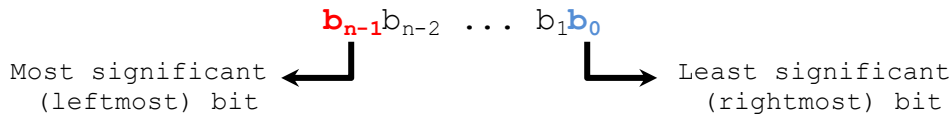
- Bit:** Unit of information that a computer uses to process and retrieve data. It can also be used as a Boolean variable (see Unit 1).
- Binary number:** This is represented by a string of bits using the positional number representation: $b_{n-1}b_{n-2} \dots b_1b_0$
- Converting a binary number into a decimal number:** The following figure depicts two cases: 2-bit numbers and 3-bit numbers. Note that the positional representation with powers of two let us obtain the decimal value (*integer positive*) of the binary number.

MATHEMATICAL REPRESENTATION	Binary number	Powers of 2:	Decimal Number
	b_1b_0	$b_1 \times 2^1 + b_0 \times 2^0$	
2-bit 	00	$0 \times 2^1 + 0 \times 2^0$	0
	01	$0 \times 2^1 + 1 \times 2^0$	1
	10	$1 \times 2^1 + 0 \times 2^0$	2
	11	$1 \times 2^1 + 1 \times 2^0$	3

MATHEMATICAL REPRESENTATION	Binary number	Powers of 2:	Decimal Number
	$b_2b_1b_0$	$b_2 \times 2^2 + b_1 \times 2^1 + b_0 \times 2^0$	
3-bit 	000	$0 \times 2^2 + 0 \times 2^1 + 0 \times 2^0$	0
	001	$0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$	1
	010	$0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0$	2
	011	$0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$	3
	100	$1 \times 2^2 + 0 \times 2^1 + 0 \times 2^0$	4
	101	$1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$	5
	110	$1 \times 2^2 + 1 \times 2^1 + 0 \times 2^0$	6
	111	$1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$	7

General case:

- Positional number representation for a binary number with 'n' bits:



The binary number can be converted to a positive decimal number by using the following formula:

$$D = \sum_{i=0}^{i=n-1} b_i \times 2^i = b_{n-1} \times 2^{n-1} + b_{n-2} \times 2^{n-2} + \dots + b_1 \times 2^1 + b_0 \times 2^0$$

- To avoid confusion, we usually write a binary number and attach a suffix '2': $(b_{n-1}b_{n-2} \dots b_1b_0)_2$
- Example: 6 bits: $(101011)_2 \equiv D = 1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 43$
4 bits: $(1011)_2 \equiv D = 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 11$
- Maximum value for a given number of bits. Complete the tables for the highlighted cases (4 and 6):

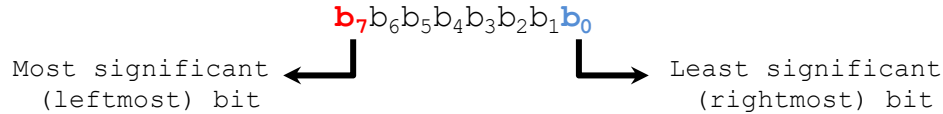
Number of bits	Maximum value	Range
1	$1_2 \equiv 2^1 - 1$	$0 \rightarrow 1_2 \equiv 0 \rightarrow 2^1 - 1$
2	$11_2 \equiv 2^2 - 1$	$0 \rightarrow 11_2 \equiv 0 \rightarrow 2^2 - 1$
3	$111_2 \equiv 2^3 - 1$	$0 \rightarrow 111_2 \equiv 0 \rightarrow 2^3 - 1$
4		
5	$11111_2 \equiv 2^5 - 1$	$0 \rightarrow 11111_2 \equiv 0 \rightarrow 2^5 - 1$
6		
...		
n	$111\dots111_2 \equiv 2^n - 1$	$0 \rightarrow 111\dots111_2 \equiv 0 \rightarrow 2^n - 1$

- Maximum value for 'n' bits:** The maximum binary number is given by an n-bit string of 1's: 111...111. Then, the maximum decimal number is given by:

$$D = \underbrace{111\dots111}_{n \text{ bits}} = 1 \times 2^{n-1} + 1 \times 2^{n-2} + \dots + 1 \times 2^1 + 1 \times 2^0 = 2^n - 1$$

\Rightarrow With 'n' bits, we can represent 2^n positive integer numbers, from 0 to $2^n - 1$.

- The case $n=8$ bits is of particular interest, as a string of 8 bits is called a byte. For 8-bit numbers, we have 256 numbers in the range 0 to $2^8-1 \equiv 0$ to 255.



- The table shows some examples:

Decimal Number	8-bit format $\text{b}_7 \text{b}_6 \text{b}_5 \text{b}_4 \text{b}_3 \text{b}_2 \text{b}_1 \text{b}_0$	$\text{b}_7 \times 2^7 + \text{b}_6 \times 2^6 + \text{b}_5 \times 2^5 + \text{b}_4 \times 2^4 + \text{b}_3 \times 2^3 + \text{b}_2 \times 2^2 + \text{b}_1 \times 2^1 + \text{b}_0 \times 2^0$
0	00000000	$0 \times 2^7 + 0 \times 2^6 + 0 \times 2^5 + 0 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 0 \times 2^0$
9	00001001	$0 \times 2^7 + 0 \times 2^6 + 0 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$
11	00001011	$0 \times 2^7 + 0 \times 2^6 + 0 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$
25	00011001	$0 \times 2^7 + 0 \times 2^6 + 0 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$
90	01011010	$0 \times 2^7 + 1 \times 2^6 + 0 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0$
128	10000000	$1 \times 2^7 + 0 \times 2^6 + 0 \times 2^5 + 0 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 0 \times 2^0$
255	11111111	$1 \times 2^7 + 1 \times 2^6 + 1 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$

Exercise: Convert the following binary numbers (positive integers) to their decimal values:

8-bit representation	$\text{b}_7 \times 2^7 + \text{b}_6 \times 2^6 + \text{b}_5 \times 2^5 + \text{b}_4 \times 2^4 + \text{b}_3 \times 2^3 + \text{b}_2 \times 2^2 + \text{b}_1 \times 2^1 + \text{b}_0 \times 2^0$	Decimal Number
00000001		
00001001		
10000101		
10000111		
11110011		

CONVERSION OF A NUMBER IN ANY BASE TO THE DECIMAL SYSTEM

- To convert a number of base 'r' ($r = 2, 3, 4, \dots$) to decimal, we use the following formula:

Number in base 'r': $(r_{n-1}r_{n-2} \dots r_1r_0)_r$

Conversion to decimal:

$$D = \sum_{i=0}^{n-1} r_i \times r^i = r_{n-1} \times r^{n-1} + r_{n-2} \times r^{n-2} + \dots + r_1 \times r^1 + r_0 \times r^0$$

Also, the maximum decimal value for a number in base 'r' with 'n' digits is:

$$D = rrr \dots rrr = r \times r^{n-1} + r_{n-2} \times r^{n-2} + \dots + r \times r^1 + r \times r^0 = r^n - 1$$

- Example: Base-8:

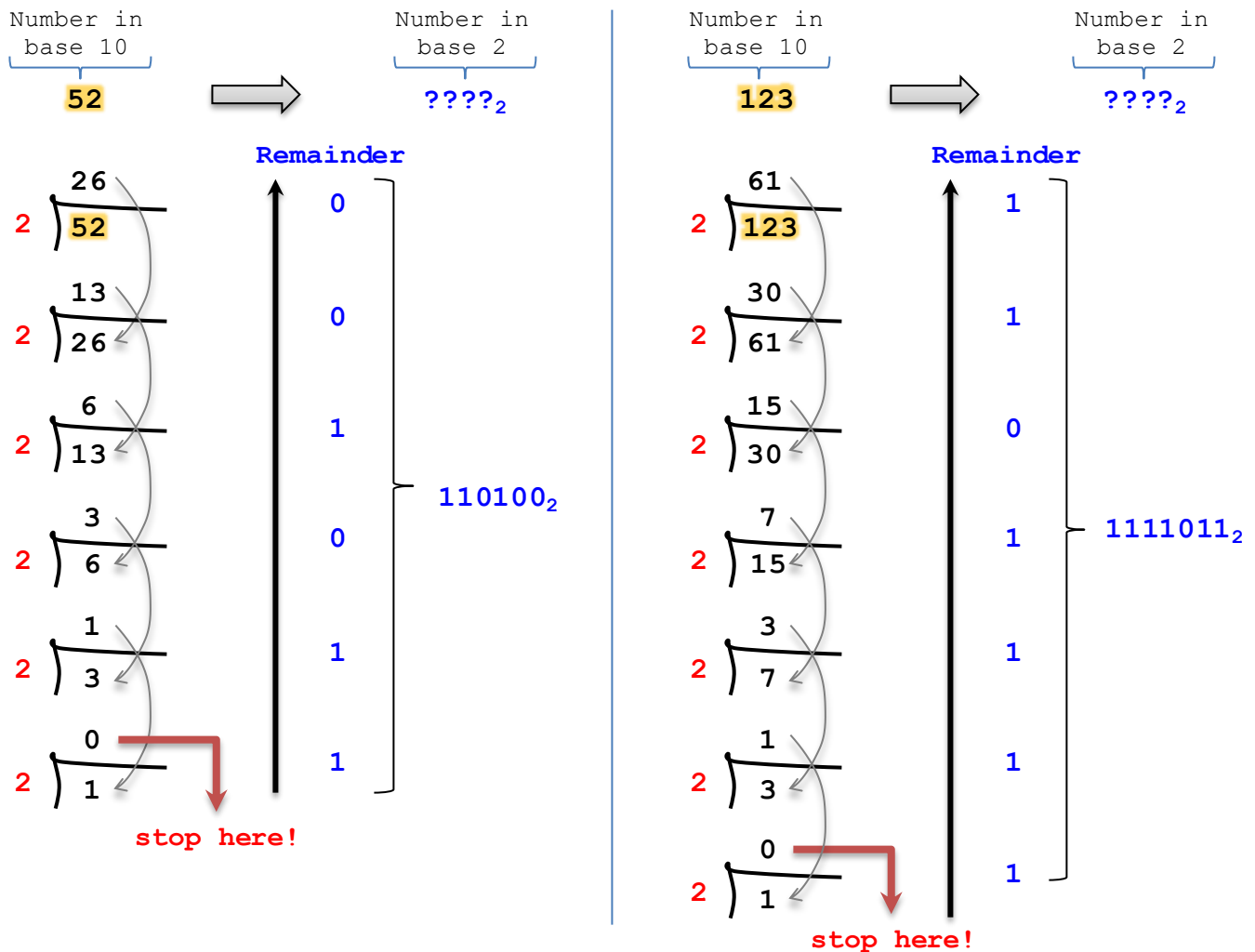
Number of digits	Maximum value	Range
1	$7_8 \equiv 8^1 - 1$	$0 \rightarrow 7_8 \equiv 0 \rightarrow 8^1 - 1$
2	$77_8 \equiv 8^2 - 1$	$0 \rightarrow 77_8 \equiv 0 \rightarrow 8^2 - 1$
3	$777_8 \equiv 8^3 - 1$	$0 \rightarrow 777_8 \equiv 0 \rightarrow 8^3 - 1$
...		
n	$777 \dots 777_8 \equiv 8^n - 1$	$0 \rightarrow 777 \dots 777_8 \equiv 0 \rightarrow 8^n - 1$

Examples:

- $(50632)_8$: Number in base 8 (octal system)
Number of digits: $n = 5$
Conversion to decimal: $D = 5 \times 8^4 + 0 \times 8^3 + 6 \times 8^2 + 3 \times 8^1 + 2 \times 8^0 = 20890$
- $(3102)_4$: Number in base 4 (quaternary system)
Number of digits: $n = 4$
Conversion to decimal: $D = 3 \times 4^3 + 1 \times 4^2 + 0 \times 4^1 + 2 \times 4^0 = 210$

CONVERSION OF DECIMAL (INTEGER POSITIVE) TO BINARY NUMBERS

- Examples:



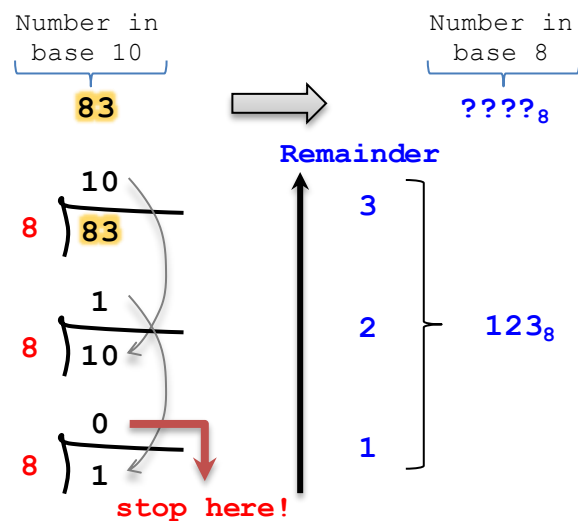
- Note that some numbers require fewer bits than others. If we want to use a specific bit representation, e.g., 8-bit, we just need to append zeros to the left until the 8 bits are completed. For example:

$$110100_2 = 00110100_2 \text{ (8-bit number)}$$

$$1111011_2 = 01111011_2 \text{ (8-bit number)}$$

- Actually, you can use this method to convert a decimal number into any other base. For example, if you want to convert it into a base-8 number, just divide by 8 and group the remainders.

- Example:** Converting a decimal number to base-8:



Exercise:

- Convert the following two decimal numbers to binary numbers. Fill in the blanks in the figure below.

Number in
base 10

63

Number in
base 2

????₂

Remainder

2

2

2

2

2

2

2

Number in
base 10

97

Number in
base 2

????₂

Remainder

2

2

2

2

2

2

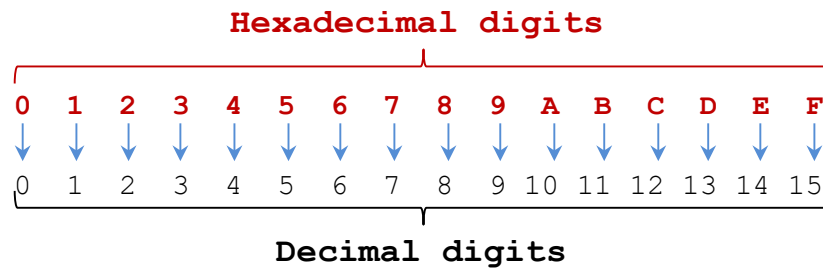
2

- Now, convert the following decimal numbers to binary numbers. The final binary number must have 8 bits (append zeros to the left to complete).

Decimal number	Binary number with 8 bits <small>b₇b₆b₅b₄b₃b₂b₁b₀</small>
40	
255	
111	
126	
9	

HEXADECIMAL NUMBER SYSTEM

- This is a very useful system as it is a short-hand notation for binary numbers
- In the decimal number system, a digit can take a value from 0 to 9.
- A hexadecimal digit is also called a *nibble*. A hexadecimal digit can take a value from 0 to 15. To avoid confusion, the numbers 10 to 15 are represented by letter (A-F):



- The following figure shows a 2-digit hexadecimal number. Note that the positional representation with powers of 16 let us obtain the decimal value (*integer positive*) of the hexadecimal number. This is the same as converting a hexadecimal number into a decimal number.

MATHEMATICAL REPRESENTATION 2-hexadecimal digits	Hex. number h_1h_0		Powers of 16: $h_1 \times 16^1 + h_0 \times 16^0$	Decimal Value
$\begin{array}{cc} \text{Second Digit} & \text{First Digit} \\ \downarrow & \downarrow \\ h_1 & h_0 \end{array}$	5A		$5 \times 16^1 + A \times 16^0$	90
	10		$1 \times 16^1 + 0 \times 16^0$	16
	08		$0 \times 16^1 + 8 \times 16^0$	8
	FB		$F \times 16^1 + B \times 16^0$	251
	3E		$3 \times 16^1 + E \times 16^0$	62
	A7		$A \times 16^1 + 7 \times 16^0$	167

- Note that when we use the letters A-F in the multiplications inside the powers of 16 representation (e.g., $A \times 16^1 + 7 \times 16^0$), we need to replace the hexadecimal symbol by its decimal value.

$$A = 10, B = 11, C = 12, D = 13, E = 14, F = 15$$

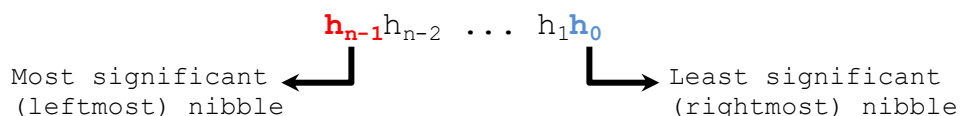
For example: $A \times 16^1 = (10) \times 16^1$.

EXERCISE: Convert the following hexadecimal numbers (positive integers) to their decimal values:

2-hex. digit representation	$h_1 \times 16^1 + h_0 \times 16^0$	Decimal Number
AB		
CE		
05		
70		
F0		
E9		

General case:

- Positional number representation for a hexadecimal number with 'n' nibbles (hexadecimal digits):



- To convert a hexadecimal number into a decimal, we apply the following formula:

Decimal Value (*integer positive*):

$$D = \sum_{i=0}^{n-1} h_i \times 16^i = h_{n-1} \times 16^{n-1} + h_{n-2} \times 16^{n-2} + \dots + h_1 \times 16^1 + h_0 \times 16^0$$

- To avoid confusion, it is sometimes customary to append the prefix '0x' to a hexadecimal number:

$$0xh_{n-1}h_{n-2}\dots h_1h_0$$

- Examples: $FD0A90: 0 \times FD0A90 \equiv F \times 16^5 + D \times 16^4 + 0 \times 16^3 + A \times 16^2 + 9 \times 16^1 + 0 \times 16^0$
 $0B871C: 0 \times 0B871C \equiv 0 \times 16^5 + B \times 16^4 + 8 \times 16^3 + 7 \times 16^2 + 1 \times 16^1 + C \times 16^0$
- The table presents the maximum attainable value for the given number of nibbles (hexadecimal digits). What pattern do you find? Can you complete it for the highlighted cases (4 and 6)?

Number of nibbles	Maximum value	Range
1	$F \equiv 16^1 - 1$	$0 \rightarrow F \equiv 0 \rightarrow 16^1 - 1$
2	$FF \equiv 16^2 - 1$	$0 \rightarrow FF \equiv 0 \rightarrow 16^2 - 1$
3	$FFF \equiv 16^3 - 1$	$0 \rightarrow FFF \equiv 0 \rightarrow 16^3 - 1$
4		
5	$FFFF \equiv 16^5 - 1$	$0 \rightarrow FFFFF \equiv 0 \rightarrow 16^5 - 1$
6		
...		
n	$FFF\dots FFF \equiv 16^n - 1$	$0 \rightarrow FFF\dots FFF \equiv 0 \rightarrow 16^n - 1$

- Maximum value for 'n' nibbles:** The maximum decimal value with 'n' nibbles is given by:

$$D = \underbrace{FFF\dots FFF}_{n \text{ nibbles}} = F \times 16^{n-1} + F \times 16^{n-2} + \dots + F \times 16^1 + F \times 16^0$$

$$= 15 \times 16^{n-1} + 15 \times 16^{n-2} + \dots + 15 \times 16^1 + 15 \times 16^0 = 16^n - 1$$

⇒ With 'n' nibbles, we can represent positive integer numbers from 0 to $16^n - 1$. (16^n numbers)

UNITS OF INFORMATION

Nibble	Byte	KB	MB	GB	TB
4 bits	8 bits	2^{10} bytes	2^{20} bytes	2^{30} bytes	2^{40} bytes

- Note that the nibble (4 bits) is one hexadecimal digit. Also, one byte (8 bits) is represented by two hexadecimal digits.
- While KB, MB, GB, TB (and so on) should be powers of 10 in the International System, it is customary in digital jargon to use powers of 2 to represent them as a matter of convenience. In microprocessor systems, memory size is determined by the number of addresses its address bus can handle (which is a power of 2). For example, for a memory with a 15-bit address line and 8 bits (1 byte) for each memory position, the memory size is 2^{15} bytes = $2^5 2^{10}$ bytes = 32 KB.
- To avoid conflicting definitions, it has been suggested to use KiB (kibibyte) for 2^{10} bytes, 1 MiB (mebibyte) for 2^{20} bytes, 1 GiB (gibibyte) for 2^{30} bytes, 1 TiB for 2^{40} bytes, and so on. However, this has not been fully adopted yet.
- Digital computers usually represent numbers utilizing a number of bits that is a multiple of 8. The simple hexadecimal to binary conversion allows us to quickly convert a string of bits that is a multiple of 8 into a string of hexadecimal digits.
- The size of the data bus in a processor represents the computing capacity of a processor, as the data bus size is the number of bits the processor can operate in one operation (e.g.: 8-bit, 16-bit, 32-bit processor). This is also usually expressed as a number of bits that is a multiple of 8.

RELATIONSHIP BETWEEN HEXADECIMAL AND BINARY NUMBERS

- Conversions between hexadecimal and binary systems are very common when dealing with digital computers. In this activity, we will learn how these 2 systems are related and how easy it is to convert between one and the other.
- Hexadecimal to binary:** We already know how to convert a hexadecimal number into a decimal number. We can then convert the decimal number into a binary number (using successive divisions).
- Binary to hexadecimal:** We can first convert the binary number to a decimal number. Then, using an algorithm similar to the one that converts decimals into binary, we can convert our decimal number into a hexadecimal number.

SIMPLE METHOD TO CONVERT HEXADECIMAL TO BINARY NUMBERS AND VICEVERSA

- The previous two conversion processes are too tedious. Fortunately, hexadecimal numbers have an interesting property that allows quick conversion of binary numbers to hexadecimal and viceversa.
- Binary to hexadecimal:** We group the binary numbers in groups of 4 (starting from the rightmost bit). If the last group of bits does not have four bits, we append zeros to the left. Then, we independently convert each group of 4 bits to its decimal value.
Notice that 4 bits can only take decimal values between 0 and $2^4-1 \equiv 0$ to 15, hence 4 bits represent only one hexadecimal digit. In other words, for each group of 4 bits, there are only 16 possible hexadecimal digits to pick from. The figure below shows an example.

Binary: $1011101_2 \rightarrow \underbrace{0101}_{5} \underbrace{1101}_{13}$
 decimal: 5 13
 hexadecimal: 5 D

Then: $01011101_2 = 0x5D$

Verification:

$$01011101_2 = 1 \times 2^6 + 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 1 \times 2^0 = 93$$

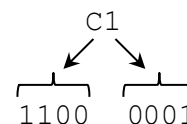
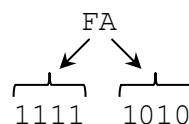
$$0x5D = 5 \times 16^1 + D \times 16^0 = 93$$

binary	dec	hex
0000	0	0
0001	1	1
0010	2	2
0011	3	3
0100	4	4
0101	5	5
0110	6	6
0111	7	7
1000	8	8
1001	9	9
1010	10	A
1011	11	B
1100	12	C
1101	13	D
1110	14	E
1111	15	F

- Exercise:** Group the following binary numbers in groups of 4 bits and obtain the hexadecimal representation. Use the table in the previous figure to pick the correspondent hexadecimal digit for each group of 4 bits.

Binary number	Hexadecimal number
1101111	
101	
1110	
110011	
11111110	
100001	

- Hexadecimal to binary:** It is basically the inverse of the process of converting a binary into a hexadecimal numbers. We pick each hexadecimal digit and convert it (always using 4 bits) to its 4-bit binary representation. The binary number is the concatenation of all resulting group of 4 bits.



$$0xFA = 11111010_2$$

$$0xC1 = 11000001_2$$

DO NOT discard these zeros
when concatenating!

Exercise:

- Convert these hexadecimal numbers to binary. Verify it by converting both the binary and hexadecimal number (they should match).

Hexadecimal number	Binary number	Decimal value
A10		
891		
43		
A2		
FACE		

- The reason hexadecimal numbers are popular is because hexadecimal numbers provide a short-hand notation for binary numbers.

OCTAL NUMBERS

- An octal digit can take between 0 and 7. This is another common number system in computers is base-8 (octal). The conversion between base-8 and base-2 resembles that of converting between base-16 and base-2. Here, we group binary numbers in 3-bit groups:

BINARY TO OCTAL

Binary: 1011101_2 \rightarrow $\underbrace{001}_{1}$ $\underbrace{011}_{3}$ $\underbrace{101}_{5}$

octal: 135_8

Then: $01011101_2 = 135_8$

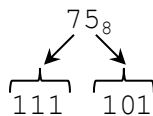
binary	dec	oct
000	0	0
001	1	1
010	2	2
011	3	3
100	4	4
101	5	5
110	6	6
111	7	7

Verification:

$$01011101_2 = 1 \times 2^6 + 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 1 \times 2^0 = 93$$

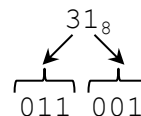
$$135_8 = 1 \times 8^2 + 3 \times 8^1 + 5 \times 8^0 = 93$$

OCTAL TO BINARY



$$75_8 = 111101_2$$

$$31_8 = 011001_2$$



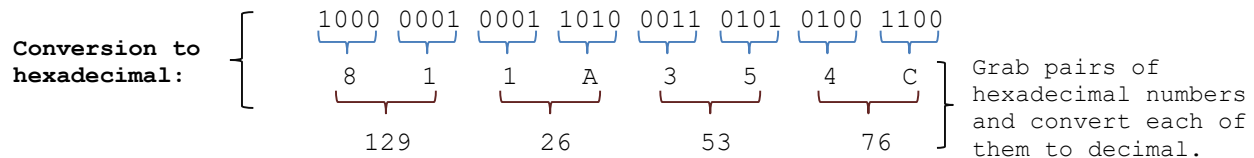
DO NOT discard these zeros when concatenating!

APPLICATIONS OF BINARY AND HEXADECIMAL REPRESENTATIONS

INTERNET PROTOCOL ADDRESS (IP ADDRESS):

- Hexadecimal numbers represent a compact way of representing binary numbers. The IP address is defined as a 32-bit number, but it is displayed as a concatenation of four decimal values separated by a dot (e.g., 129.26.53.76).
- The following figure shows how a 32-bit IP address expressed as a binary number is transformed into the standard IP address notation.

IP address (binary): 10000001000110100011010101001100



IP address (hex): 0x811A354C

IP address notation: 129.26.53.76

- The 32-bit IP address expressed as binary number is very difficult to read. So, we first convert the 32-bit binary number to a hexadecimal number.
 - The IP address expressed as a hexadecimal (0x811A354C) is a compact representation of a 32-bit IP address. This should suffice. However, it was decided to represent the IP address in a 'human-readable' notation. In this notation, we grab pairs of hexadecimal numbers and convert each of them individually to decimal numbers. Then we concatenate all the values and separate them by a dot.
 - Important:** Note that the IP address notation (decimal numbers) is NOT the decimal value of the binary number. It is rather a series of four decimal values, where each decimal value is obtained by independently converting each two hexadecimal digits to its decimal value.
- ✓ Given that each decimal number in the IP address can be represented by 2 hexadecimal digits (or 8 bits), what is the range (min. value, max. value) of each decimal number in the IP address?
With 8 bits, we can represent $2^8 = 256$ numbers from 0 to 255.
 - ✓ An IP address represents a unique device connected to the Internet. Given that the IP address has 32 bits (or 8 hexadecimal digits), how many numbers can be represented (i.e., how many devices can connect to the Internet)?
 $2^{32} = 4294967296$ devices.
 - ✓ The number of devices that can be connected to the Internet is huge, but considering the number of Internet-capable devices that exists in the entire world, it is becoming clear that 32 bits is not going to be enough. That is why the Internet Protocol is being currently extended to a new version (IPv6) that uses 128 bits for the addresses. With 128 bits, how many Internet-capable devices can be connected to the Internet?
 $2^{128} \approx 3.4 \times 10^{38}$ devices

REPRESENTING GRAYSCALE PIXELS

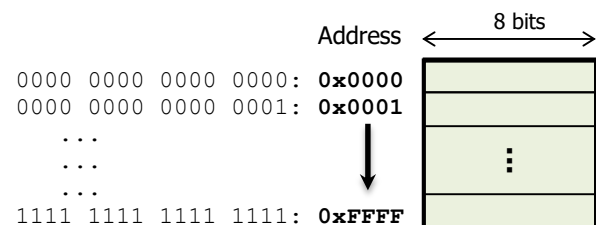
A grayscale pixel is commonly represented with 8 bits. So, a grayscale pixel value varies between 0 and 255, 0 being the darkest (black) and 255 being the brightest (white). Any value in between represents a shade of gray.



MEMORY ADDRESSES

The address bus size in processors is usually determined by the number of memory positions it can address. For example, if we have a microprocessor with an address bus of 16 bits, we can handle up to 2^{16} addresses. If the memory content is one byte wide, then the processor can handle up to 2^{16} bytes = 64KB.

Here, we use 16 bits per address, or 4 nibbles. The lowest address (in hex) is 0x0000 and highest address (in hex) is 0xFFFF.



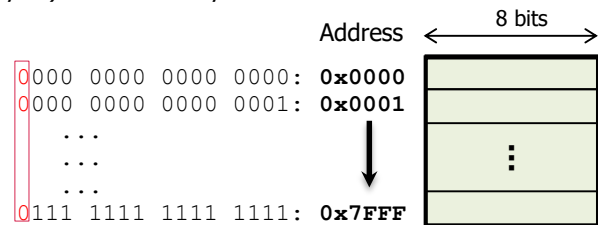
Examples:

- A microprocessor can only handle memory addresses from $0x0000$ to $0x7FFF$. What is the address bus size? If each memory position is one byte wide, what is the maximum size (in bytes) of the memory that we can connect?

We want to cover all the cases from $0x0000$ to $0x7FFF$:

The range from $0x0000$ to $0x7FFF$ is akin to all possible cases with 15 bits. Thus, the address bus size is **15 bits**.

We can handle $2^{15} \text{ bytes} = 32\text{KB}$ of memory.

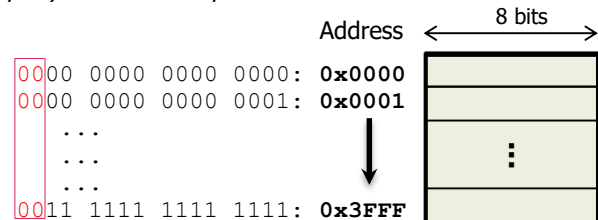


- A microprocessor can only handle memory addresses from $0x0000$ to $0x3FFF$. What is the address bus size? If each memory position is one byte wide, what is the maximum size (in bytes) of the memory that we can connect?

We want to cover all the cases from $0x0000$ to $0x3FFF$:

The range from $0x0000$ to $0x3FFF$ is akin to all possible cases with 14 bits. Thus, the address bus size is **14 bits**.

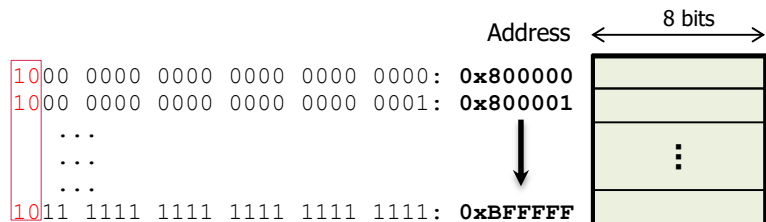
We can handle $2^{14} \text{ bytes} = 16\text{KB}$ of memory.



- A microprocessor has a 24-bit address line. We connect a memory chip to the microprocessor. The memory chip addresses are assigned the range $0x800000$ to $0xBFFFFFFF$. What is the minimum number of bits required to represent addresses in that individual memory chip? If each memory position is one byte wide, what is the memory size (in bytes)?

By looking at the binary numbers from $0x800000$ to $0xBFFFFFFF$, we notice that the addresses in that range require 24 bits. But all those addresses share the same first two MSBs: 10. Thus, if we were to use only that memory chip, we do not need those 2 bits, and we only need **22 bits**.

We can handle $2^{22} \text{ bytes} = 4\text{MB}$ of memory.



- A memory has a size of 512KB, where each memory content is 8-bits wide. How many bits do we need to address the contents of this memory?

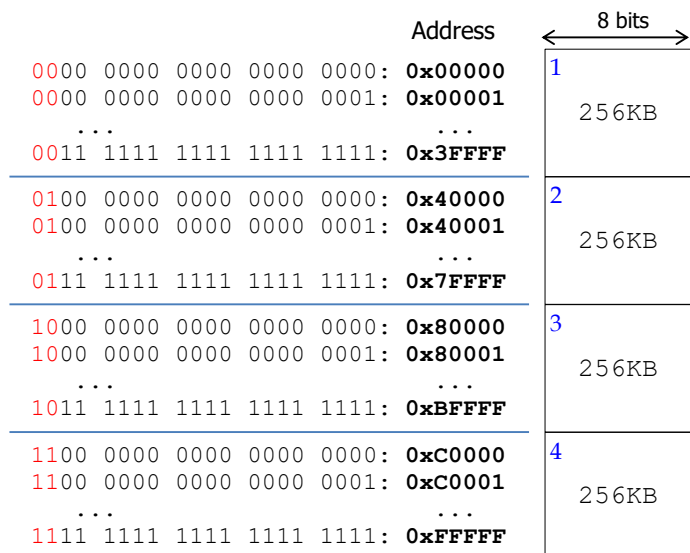
Recall that: $512\text{KB} = 2^{19} \text{ bytes}$. So, we need 19 bits to address the contents of this memory (address bus size = 19 bits)
In general, for a memory with N address positions, the number of bits to address those positions is given by: $\lceil \log_2 N \rceil$

- A 20-bit address line in a microprocessor with an 8-bit data bus handles 1 MB (2^{20} bytes) of data. We want to connect four 256 KB memory chips to the microprocessor. Provide the address ranges that each memory device will occupy.

Each memory chip can handle 256KB of memory. $256\text{KB} = 2^{18} \text{ bytes}$, requiring 18 bits for its address.

For a 20-bit address: we have 5 hexadecimal digits that go from $0x00000$ to $0xFFFFF$ (2^{20} memory positions).

We divide the 2^{20} memory positions into 4 contiguous groups, each with 2^{18} memory positions. The figure shows the optimal way of doing so: for each group, the 18 LSBs of the memory addresses correspond to the memory range of a 256 KB memory. And the 2 MSBs of the memory addresses are the same within a group. For a given memory address, we can quickly determine which group it belongs to by looking at its 2 MSBs.



UNSIGNED NUMBERS - ADDITION

- In the example, we add two 8-bit numbers using binary representation and hexadecimal representation (this is a short-hand notation). Note that every summation of two digits (binary or hexadecimal) generates a carry when the summation requires more than one digit. Also, note that c_0 is the *carry in* of the summation. c_0 is usually zero.
- The last carry (c_8 when $n=8$) is the *carry out* of the summation. If it is zero, it means that the summation can be represented with 8 bits. If it is one, it means that the summation requires more than 8 bits (in fact 9 bits); this is called an overflow. In the example, we add two numbers and overflow occurs: an extra bit (in red) is required to correctly represent the summation.

$$\begin{array}{r}
 \begin{array}{cccccccc}
 c_8 & c_7 & c_6 & c_5 & c_4 & c_3 & c_2 & c_1 & c_0 \\
 0 \times 3F & = & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & + \\
 0 \times B2 & = & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & \\
 \hline
 0 \times F1 & = & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 1 &
 \end{array}
 \quad
 \begin{array}{r}
 c_2 & c_1 & c_0 \\
 3 & F & + \\
 B & 2 & \\
 \hline
 F & 1 &
 \end{array}
 \\
 \\
 \begin{array}{r}
 \begin{array}{cccccccc}
 c_8 & c_7 & c_6 & c_5 & c_4 & c_3 & c_2 & c_1 & c_0 \\
 0 \times 3F & = & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & + \\
 0 \times C2 & = & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & \\
 \hline
 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 &
 \end{array}
 \quad
 \begin{array}{r}
 c_2 & c_1 & c_0 \\
 3 & F & + \\
 C & 2 & \\
 \hline
 1 & 0 & 1 &
 \end{array}
 \end{array}$$

DIGITAL CIRCUIT IMPLEMENTATION

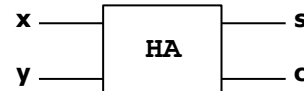
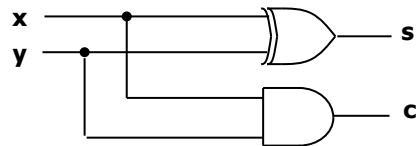
1-bit Addition:

- ✓ Addition of a bit without carry in: The circuit that implements this operation is called Half Adder (HA).

$$\begin{array}{r}
 x + \\
 y \\
 \hline
 \text{carry out} \leftarrow c \quad s \rightarrow \text{sum}
 \end{array}
 \quad
 \begin{array}{r}
 0 + \\
 0 \\
 \hline
 0 \ 0
 \end{array}
 \quad
 \begin{array}{r}
 0 + \\
 1 \\
 \hline
 0 \ 1
 \end{array}
 \quad
 \begin{array}{r}
 1 + \\
 0 \\
 \hline
 0 \ 1
 \end{array}
 \quad
 \begin{array}{r}
 1 + \\
 1 \\
 \hline
 1 \ 0
 \end{array}$$

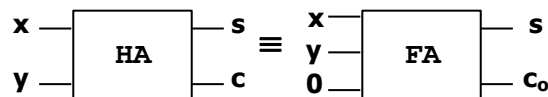
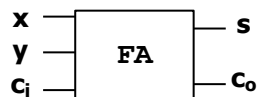
✓ Boolean Algebra is a very powerful tool for the implementation of digital circuits. Note how the 1-bit addition operation ' $x + y$ ' was mapped into Boolean equations for c and s . $c = xy$, $s = x \oplus y$, where x, y, c, s are Boolean variables.

x	y	c	s
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0



- ✓ Addition of a bit with carry in: The circuit that performs this operation is called Full Adder (FA).

$$\begin{array}{r}
 c_i \\
 x + \\
 y \\
 \hline
 \text{carry out} \leftarrow c_o \quad s \rightarrow \text{sum}
 \end{array}$$

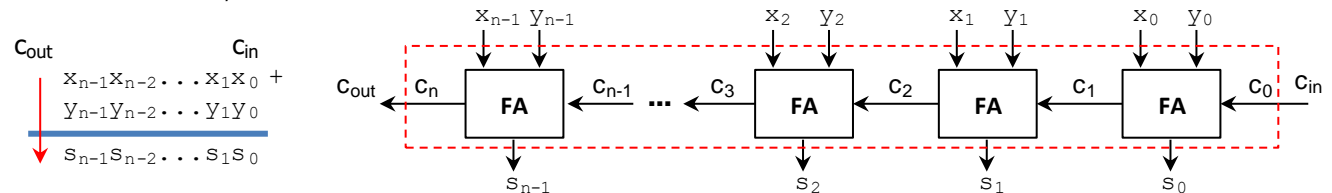


n-bit Addition (of two operands):

The figure on the right shows a 5-bit addition. To implement a circuit with the truth table method, we would need 11 inputs and 6 outputs (not practical!). A better method is to use a cascade of Full Adders.

$$\begin{array}{r}
 \begin{array}{cccccc}
 c_5 & c_4 & c_3 & c_2 & c_1 & c_0 \\
 15: & 0 & 1 & 1 & 1 & 1 & + \\
 10: & 0 & 1 & 0 & 1 & 0 & \\
 \hline
 25: & 1 & 1 & 0 & 0 & 1 &
 \end{array}
 \quad
 \begin{array}{r}
 C_{out} \quad C_{in} \\
 x_4x_3x_2x_1x_0 + \\
 y_4y_3y_2y_1y_0 \\
 \hline
 s_4s_3s_2s_1s_0
 \end{array}$$

For an n -bit addition, the circuit looks like:



Full Adder Design

x_i	y_i	c_i	c_{i+1}	s_i
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	0	1
1	0	1	1	0
1	1	0	1	0
1	1	1	1	1

c_i	$x_i y_i$	00	01	11	10
0		0	1	0	1
1		1	0	1	0

c_i	$x_i y_i$	00	01	11	10
0		0	0	1	0
1		0	1	1	1

$$s_i = \overline{x_i} \overline{y_i} c_i + x_i \overline{y_i} \overline{c_i} + \overline{x_i} y_i c_i + x_i y_i \overline{c_i}$$

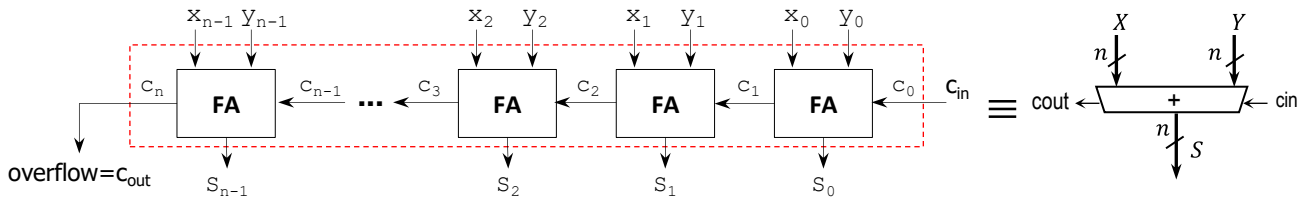
$$s_i = (x_i \oplus y_i) \overline{c_i} + \overline{(x_i \oplus y_i)} c_i$$

$$s_i = x_i \oplus y_i \oplus c_i$$

$$c_{i+1} = x_i y_i + x_i c_i + y_i c_i$$

ARITHMETIC OVERFLOW:

- In an n -bit addition operation, the operands and the result are represented with the same number of bits. Overflow occurs when the result requires more than n bits. For unsigned numbers, this is equivalent to the result being outside the n -bit range $[0, 2^n - 1]$, i.e., overflow occurs when the sum is greater than $2^n - 1$.
✓ Example: for $n = 4$ bits, the range is 0 to 15. Overflow occurs when the sum is greater than 15.
- In unsigned binary addition, a carry out of 1 means overflow. Thus, an overflow flag (or bit) can be specified as: $overflow = c_n = c_{out}$.
- To 'avoid' overflow in addition operation, a common technique is to sign-extend (this is zero-extension for unsigned numbers) the two summands. For example, if the two summands are 4-bits wide, then we add one more bit. So, we use 5 bits to represent the operands and the result.
- In general, if the two summands are n -bits wide, the result will have at most $n + 1$ bits ($2^n - 1 + 2^n - 1 = 2^{n+1} - 2$).

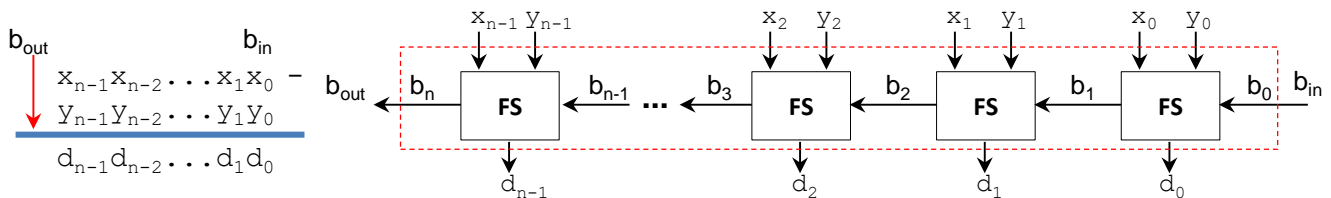


UNSIGNED NUMBERS - SUBTRACTION

- In the example, we subtract two 8-bit numbers using the binary and hexadecimal (short-hand notation) representations. A subtraction of two digits (binary or hexadecimal) generates a borrow when the difference is negative. So, we borrow 1 from the next digit so that the difference is positive. Recall that a borrow in a subtraction of two digits is an extra 1 that we need to subtract. Also, note that b_0 is the *borrow in* of the summation. This is usually zero.
- The last borrow (b_8 when $n=8$) is the *borrow out* of the subtraction. If it is zero, it means that the difference is positive and can be represented with 8 bits. If it is one, it means that the difference is negative, and we need to borrow 1 from the next digit. In the example, we subtract two 8-bit numbers, the result we have borrows 1 from the next digit.
- We can build an **n -bit subtractor** for unsigned numbers using Full Subtractor circuits. Subtraction for unsigned numbers only makes sense if the result is positive (or when doing multi-precision subtraction). In practice, subtraction is better performed in the 2's complement representation (for signed numbers).

$$\begin{array}{r}
 \begin{array}{cccccccc}
 b_8 & b_7 & b_6 & b_5 & b_4 & b_3 & b_2 & b_1 & b_0 \\
 0 \times 3A & = & 0 & 0 & 1 & 1 & 1 & 0 & 1 & 0 & - \\
 0 \times 2F & = & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 1 & - \\
 \hline
 0 \times 0B & = & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & \\
 \end{array}
 \quad
 \begin{array}{r}
 \begin{array}{ccc}
 b_2 & b_1 & b_0 \\
 3 & A & - \\
 2 & F & \\
 \hline
 0 & B &
 \end{array}
 \end{array}
 \end{array}$$

$$\begin{array}{r}
 \begin{array}{cccccccc}
 b_8 & b_7 & b_6 & b_5 & b_4 & b_3 & b_2 & b_1 & b_0 \\
 0 \times 3A & = & 0 & 0 & 1 & 1 & 1 & 0 & 1 & 0 & - \\
 0 \times 75 & = & 0 & 1 & 1 & 1 & 0 & 1 & 0 & 1 & - \\
 \hline
 0 \times C5 & = & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 1 &
 \end{array}
 \quad
 \begin{array}{r}
 \begin{array}{ccc}
 b_2 & b_1 & b_0 \\
 3 & A & - \\
 7 & 5 & \\
 \hline
 C & 5 &
 \end{array}
 \end{array}
 \end{array}$$



Full Subtractor Design

x_i	y_i	b_i	b_{i+1}	d_i
0	0	0	0	0
0	0	1	1	1
0	1	0	1	1
0	1	1	1	0
1	0	0	0	1
1	0	1	0	0
1	1	0	0	0
1	1	1	1	1

$x_i y_i$	00	01	11	10
b_i	0	1	0	1
1	1	0	1	0

$x_i y_i$	00	01	11	10
b_i	0	1	0	0
1	1	1	1	0

$$d_i = \overline{x_i} y_i \overline{b_i} + x_i \overline{y_i} \overline{b_i} + \overline{x_i} \overline{y_i} b_i + x_i y_i b_i$$

$$d_i = (x_i \oplus y_i) \overline{b_i} + (\overline{x_i \oplus y_i}) b_i$$

$$d_i = x_i \oplus y_i \oplus b_i$$

$$b_{i+1} = \overline{x_i} y_i + \overline{x_i} b_i + y_i b_i$$

SIGNED INTEGER NUMBERS

- For an n -bit number $b_{n-1}b_{n-2} \dots b_1b_0$, there exist three common signed representations: sign-and-magnitude, 1's complement, and 2's complement. In these three representations, the MSB always tells us whether the number is positive (MSB=0) or negative (MSB=1). These representations allow us to represent both positive and negative numbers.

SIGN-AND-MAGNITUDE (SM):

- Here, the sign and the magnitude (value) are represented separately. The MSB only represents the sign and the remaining $n - 1$ bits the magnitude.
- Example** ($n=4$): $0110 = +6$ $1110 = -6$

1'S COMPLEMENT (1C):

- Here, if the MSB=0, the number is positive and the remaining $n - 1$ bits represent the magnitude. If the MSB=1, the number is negative and the remaining $n - 1$ bits do not represent the magnitude. To invert the sign of a number in 1's complement representation, we apply the 1's complement operation to the number, which consists of inverting all the bits.
- Let $B = b_{n-1} \dots b_1b_0$ be a number represented in 1's complement. Let $K = k_{n-1} \dots k_1k_0$ represent $-B$. We get K by applying the 1's complement operation to B . K is also called the 1's complement of B (and viceversa).
- Definition:** The 1's complement of B is defined as $K = (2^n - 1) - B$, $n = \# \text{ of bits (including sign bit)}$, where $K = \sum_{i=0}^{n-1} k_i 2^i$ and $B = \sum_{i=0}^{n-1} b_i 2^i$. Note that K and B are treated as unsigned numbers in this formula. And $(2^n - 1)$ is the largest n -bit unsigned number. We can then show that the 1's complement operation amounts to inverting all the bits:

$$\sum_{i=0}^{n-1} k_i 2^i = (2^n - 1) - \sum_{i=0}^{n-1} b_i 2^i \rightarrow \sum_{i=0}^{n-1} (k_i + b_i) 2^i = 2^n - 1 \rightarrow k_i + b_i = 1, \forall i \Rightarrow k_i = \bar{b}_i$$

Example: Given $B = 01001_2 = 9$ in 1's complement, get the 1's complement representation of -9 using the formula:

$$\rightarrow K = (2^n - 1) - B = (2^5 - 1) - 9 = 22 = 10110_2.$$

Recall that the formula treats K and B as unsigned integers. So, $K = 22$ in unsigned representation, and $K = -9$ in 1's complement representation.

It is much simpler to just invert each bit!

- With n bits, we can represent $2^n - 1$ numbers from $-2^{n-1} + 1$ to $2^{n-1} - 1$. When using the 1's complement representation, it is mandatory to specify how many bits we are using.
- Examples:**
 - $+6=0110 \rightarrow -6=1001$, $+5=0101 \rightarrow -5=1010$, $+7=0111 \rightarrow -7=1000$.
 - If $-6=1001$, we get $+6$ by applying the 1's complement operation to $1001 \rightarrow +6 = 0110$
 - Get the 1's complement representation of 8: This is a positive number, thus the MSB=0. The remaining $n - 1$ bits represent the magnitude. The magnitude is represented with a minimum number of 4 bits as $8=1000_2$. Thus, using a minimum number of 5 bits, the number 8 in 1's complement representation is $8=01000_2$.
 - What is the decimal value of 1100 ? We first apply the 1's complement operation to 1100 , which results in 0011 (+3). Thus $1100=-3$.
 - What is the 1's complement representation of -4? We know that $+4=0100$. To get -4, we apply the 1's complement operation to 0100 , which results in 1011 . Thus $1011=-4$.

2'S COMPLEMENT (2C):

- Here, if the MSB=0, the remaining $n - 1$ bits represent the magnitude. If the MSB=1, the number is negative and the remaining $n - 1$ bits do not represent the magnitude. To invert the sign of a number in 2's complement representation, we apply the 2's complement operation to the number, which consists on inverting all the bits and add 1.
- Let $B = b_{n-1} \dots b_1b_0$ be a number represented in 2s complement. Let $K = k_{n-1} \dots k_1k_0$ represent $-B$. We get K by applying the 2's complement operation to B . K is also called the 2's complement of B (and viceversa).
- Definition:** The 2's complement of B is defined as $K = (2^n - 1) - B + 1$, $n = \# \text{ of bits (including sign bit)}$, where $K = \sum_{i=0}^{n-1} k_i 2^i$ and $B = \sum_{i=0}^{n-1} b_i 2^i$. Note that K and B are treated as unsigned numbers in this formula. We can see that we can first get the term $(2^n - 1) - B$ by inverting all the bits; then we add 1 to complete the equation.

Example: Given $B = 01001_2 = 9$ in 2's complement, get the 2's complement representation of -9 using the formula:

$$\rightarrow K = (2^n - 1) - B + 1 = (2^5 - 1) - 9 + 1 = 23 = 10111_2.$$

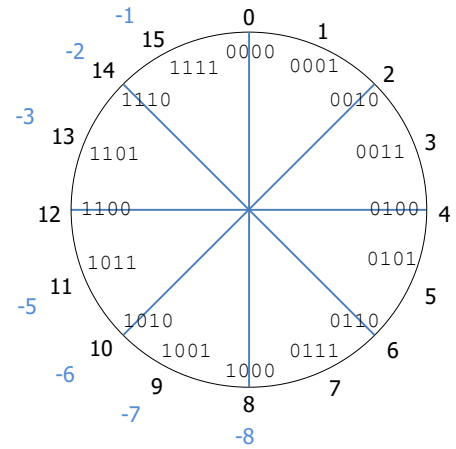
Recall that the formula treats K and B as unsigned integers. So, $K = 23$ in unsigned representation, and $K = -9$ in 2's complement representation.

It is much simpler to just invert each bit (i.e., apply 1's complement operation) and then add 1!

- With n bits, we can represent 2^n numbers from -2^{n-1} to $2^{n-1} - 1$. When using the 2's complement representation, it is mandatory to specify how many bits we are using.

▪ **Graphical interpretation of four-bit 2's complement numbers:**

If these 4-bit patterns represented unsigned integers, they would be numbers from 0 to 15. As they represent 2's complement integers, then the numbers range from -8 to 7.



▪ **Examples:**

- ✓ $+6=0110 \rightarrow -6=1010$, $+5=0101 \rightarrow -5=1011$, $+7=0111 \rightarrow -7=1001$.
- ✓ If $-6=1010$, we get $+6$ by applying the 2's complement operation to 1010 $\rightarrow +6 = 0110$
- ✓ Represent 12 in 2's complement: This is a positive number, \rightarrow MSB=0. The remaining $n-1$ bits represent the magnitude. We can get the magnitude with a minimum 4 bits: $12=1100_2$. Thus, using a minimum of 5 bits, the number 12 in 2's complement representation is $12=01100_2$.
- ✓ What is the decimal value of 1101? We first apply the 2's complement operation (or take the 2's complement) to 1101, which results in 0011 (+3). Thus $1101=-3$.
- ✓ What is the 2's complement representation of -4? We know that $+4=0100$. To get -4, we apply the 2's complement operation to 0100, which results in 1100. Thus $1100=-4$.

Getting the decimal value of a number in 2's complement representation:

- If the number B is positive, then MSB=0: $b_{n-1} = 0$.

$$\rightarrow B = \sum_{i=0}^{n-1} b_i 2^i = b_{n-1} 2^{n-1} + \sum_{i=0}^{n-2} b_i 2^i = \sum_{i=0}^{n-2} b_i 2^i \quad (a)$$

- If the number B is negative, $b_{n-1} = 1$ (MSB=1). If we take the 2's complement of B , we get K (which is a positive number). In 2's complement representation, K represents $-B$. Using $K = 2^n - B$ (K and B are treated as unsigned numbers):

$$\sum_{i=0}^{n-1} k_i 2^i = 2^n - \sum_{i=0}^{n-1} b_i 2^i$$

- We want to express $-K$ in terms of b_i , since the integer value $-K$ is the actual integer value of B .

$$-K = -\sum_{i=0}^{n-1} k_i 2^i = \sum_{i=0}^{n-1} b_i 2^i - 2^n = b_{n-1} 2^{n-1} + \sum_{i=0}^{n-2} b_i 2^i - 2^n = 2^{n-1}(b_{n-1} - 2) + \sum_{i=0}^{n-2} b_i 2^i$$

$$B = -K = 2^{n-1}(1 - 2) + \sum_{i=0}^{n-2} b_i 2^i = -2^{n-1} + \sum_{i=0}^{n-2} b_i 2^i \quad (b)$$

- Using (a) and (b), the formula for the decimal value of B (either positive or negative) is:

$$B = -b_{n-1} 2^{n-1} + \sum_{i=0}^{n-2} b_i 2^i$$

▪ **Examples:**

$$10110_2 = -2^4 + 2^2 + 2^1 = -10$$

$$11000_2 = -2^4 + 2^3 = -8$$

SUMMARY

- The following table summarizes the signed representations for a 4-bit number:

n=4: $b_3b_2b_1b_0$	SIGNED REPRESENTATION		
	Sign-and-magnitude	1's complement	2's complement
0 0 0 0	0	0	0
0 0 0 1	1	1	1
0 0 1 0	2	2	2
0 0 1 1	3	3	3
0 1 0 0	4	4	4
0 1 0 1	5	5	5
0 1 1 0	6	6	6
0 1 1 1	7	7	7
1 0 0 0	0	-7	-8
1 0 0 1	-1	-6	-7
1 0 1 0	-2	-5	-6
1 0 1 1	-3	-4	-5
1 1 0 0	-4	-3	-4
1 1 0 1	-5	-2	-3
1 1 1 0	-6	-1	-2
1 1 1 1	-7	0	-1
Range for n bits:	$[-(2^{n-1} - 1), 2^{n-1} - 1]$	$[-(2^{n-1} - 1), 2^{n-1} - 1]$	$[-2^{n-1}, 2^{n-1} - 1]$

- 1C and 2C are representations of signed numbers. 1C and 2C represent both negative and positive numbers. Do not confuse the 1C and 2C representations with the 1C and 2C operations.
- Note that the sign-and-magnitude and the 1's complement representations have a redundant representation for zero. This is not the case in 2's complement, which can represent an extra number.
- In 2C, the number -8 can be represented with 4 bits: $-8 = 1000$. To obtain +8, we apply the 2C operation to 1000, which results in 1000. But 1000 cannot be a positive number. This means that we require 5 bits to represent $+8 = 01000$.

SIGN EXTENSION

- Unsigned Numbers:** Here, if we want to use more bits, we just append zeros to the left.
Example: $12 = 1100_2$ with 4 bits. If we want to use 6 bits, then $12 = 001100_2$.
- Signed Numbers:**
 - ✓ **Sign-and-magnitude:** The MSB only represents the sign. If we want to use more bits, we append zeros to the left. The leftmost bit is always the sign.
Example: $-12 = 11100_2$ with 5 bits. If we want to use 7 bits, then $-12 = 1001100_2$.
 - ✓ **2's complement** (also applies to 1's complement): In many circumstances, we might want to represent numbers in 2's complement with a fixed number of bits. For example, the following two numbers require a minimum of 5 bits:
 $10111_2 = -2^4 + 2^2 + 2^1 + 2^0 = -9$ $01111_2 = 2^3 + 2^2 + 2^1 + 2^0 = +15$

What if we wanted to use 8 bits to represent them? In 2's complement, we need to sign-extend: If the number is positive, we append zeros to the left. If the number is negative, we attach ones to the left. In the example, note how we added three bits to the left in each case:

$$11110111_2 = -2^4 + 2^2 + 2^1 + 2^0 = -9 \quad 00001111_2 = 2^3 + 2^2 + 2^1 + 2^0 = +15$$

Demonstration of sign-extension in 2's complement:

- To increase the number of bits for representing a number, we append the MSB to the left as many times as needed:

$$\begin{aligned} \text{Examples: } & 00101_2 = 0000101_2 = 2^2 + 2^0 = 5 \\ & 10101_2 = 1110101_2 = -2^4 + 2^2 + 2^0 = -2^6 + 2^5 + 2^4 + 2^2 + 2^0 = -11 \end{aligned}$$

We can think of the sign-extended number as an m -bit number, where $m > n$:

$$b_{n-1} \dots b_{n-1} b_{n-1} b_{n-2} \dots b_0 = b_{m-1} \dots b_n b_{n-1} b_{n-2} \dots b_0, \text{ where: } b_i = b_{n-1}, i = n, n+1, \dots, m-1$$

- We need to demonstrate that $b_{n-1} b_{n-2} \dots b_0$ represents the same decimal number as $b_{n-1} \dots b_{n-1} b_{n-1} b_{n-2} \dots b_0$, i.e., that the sign-extension is correct for any $m > n$.

$$\text{We need that: } b_{m-1} \dots b_n b_{n-1} b_{n-2} \dots b_0 = b_{n-1} \dots b_{n-1} b_{n-1} b_{n-2} \dots b_0 = b_{n-1} b_{n-2} \dots b_0$$

Using the formula for 2's complement numbers:

$$\begin{aligned} -2^{m-1} b_{m-1} + \sum_{i=0}^{m-2} 2^i b_i &= -2^{n-1} b_{n-1} + \sum_{i=0}^{n-2} 2^i b_i \\ -2^{m-1} b_{m-1} + \sum_{i=n-1}^{m-2} 2^i b_i + \sum_{i=0}^{n-2} 2^i b_i &= -2^{n-1} b_{n-1} + \sum_{i=0}^{n-2} 2^i b_i \Rightarrow -2^{m-1} b_{m-1} + \sum_{i=n-1}^{m-2} 2^i b_i = -2^{n-1} b_{n-1} \\ -2^{m-1} b_{m-1} + b_{n-1} \sum_{i=n-1}^{m-2} 2^i &= -2^{n-1} b_{n-1}, \end{aligned}$$

$$\text{Recall: } \sum_{i=k}^l r^i = \frac{r^k - r^{l+1}}{1 - r}, r \neq 1 \rightarrow \sum_{i=k}^l 2^i = \frac{2^k - 2^{l+1}}{1 - 2} = 2^{l+1} - 2^k$$

Then:

$$\begin{aligned} -2^{m-1} b_{m-1} + b_{n-1} (2^{m-1} - 2^{n-1}) &= -2^{n-1} b_{n-1} \\ -2^{m-1} b_{m-1} + 2^{m-1} b_{n-1} - 2^{n-1} b_{n-1} &= -2^{n-1} b_{n-1} \therefore -2^{m-1} b_{m-1} = -2^{n-1} b_{n-1} \end{aligned}$$

ADDITION AND SUBTRACTION (SIGNED \equiv 2'S COMPLEMENT NUMBERS)

- We will use the 2's complement representation for signed numbers.
- The advantage of the 2's complement representation is that the summation can be carried out using the same circuitry as that of the unsigned summation. Here the operands can either be positive or negative.

ADDITION:

- We show addition examples of 4-bit signed numbers. Note that the *carry out* bit DOES NOT necessarily indicate overflow. In some cases, the carry out must be ignored, otherwise the result is incorrect.

$\begin{array}{r} +5 = 0101 \\ +2 = 0010 \\ \hline +7 = 0111 \\ \text{cout}=0 \end{array}$	$\begin{array}{r} -5 = 1011 \\ +2 = 0010 \\ \hline -3 = 1101 \\ \text{cout}=0 \end{array}$	$\begin{array}{r} +5 = 0101 \\ -2 = 1110 \\ \hline +3 = \text{X}0011 \\ \text{cout}=1 \end{array}$	$\begin{array}{r} -5 = 1011 \\ -2 = 1110 \\ \hline -7 = \text{X}1001 \\ \text{cout}=1 \end{array}$
--	--	--	--

- Now, we show addition examples of two 8-bit signed numbers. The *carry out* c_8 is not enough to determine overflow. Here, if $c_8 \neq c_7$ there is overflow. If $c_8 = c_7$, no overflow and we can ignore c_8 . Thus, the overflow bit is equal to $c_8 \text{ XOR } c_7$.
- Overflow:** For 8-bit operations, this occurs when the summation falls outside the 2's complement range for 8 bits: $[-2^7, 2^7 - 1]$. If there is no overflow, the carry out bit must not be part of the result.

$$\begin{array}{r} \text{c}_8 \text{ c}_7 \text{ c}_6 \text{ c}_5 \text{ c}_4 \text{ c}_3 \text{ c}_2 \text{ c}_1 \text{ c}_0 \\ +92 = 0 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ + \\ +78 = 0 \ 1 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \ + \\ \hline +170 = 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 0 \end{array}$$

overflow = $c_8 \oplus c_7 = 1 \rightarrow$ overflow!
 $+170 \notin [-2^7, 2^7 - 1] \rightarrow$ overflow!

$$\begin{array}{r} \text{c}_8 \text{ c}_7 \text{ c}_6 \text{ c}_5 \text{ c}_4 \text{ c}_3 \text{ c}_2 \text{ c}_1 \text{ c}_0 \\ -92 = 1 \ 0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ + \\ -78 = 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ + \\ \hline -170 = 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 1 \ 0 \end{array}$$

overflow = $c_8 \oplus c_7 = 1 \rightarrow$ overflow!
 $-170 \notin [-2^7, 2^7 - 1] \rightarrow$ overflow!

$$\begin{array}{r} \text{c}_8 \text{ c}_7 \text{ c}_6 \text{ c}_5 \text{ c}_4 \text{ c}_3 \text{ c}_2 \text{ c}_1 \text{ c}_0 \\ +92 = 0 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ + \\ -78 = 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ + \\ \hline +14 = \text{X} \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \end{array}$$

overflow = $c_8 \oplus c_7 = 0 \rightarrow$ no overflow
 $+14 \in [-2^7, 2^7 - 1] \rightarrow$ no overflow

$$\begin{array}{r} \text{c}_8 \text{ c}_7 \text{ c}_6 \text{ c}_5 \text{ c}_4 \text{ c}_3 \text{ c}_2 \text{ c}_1 \text{ c}_0 \\ -92 = 1 \ 0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ + \\ +78 = 0 \ 1 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \ + \\ \hline -14 = \text{X} \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 0 \end{array}$$

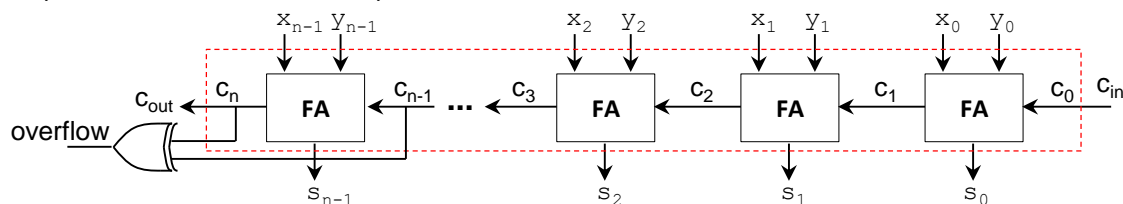
overflow = $c_8 \oplus c_7 = 0 \rightarrow$ no overflow
 $-14 \in [-2^7, 2^7 - 1] \rightarrow$ no overflow

- In general, for an n -bit number, overflow occurs when the summation falls outside the range $[-2^{n-1}, 2^{n-1} - 1]$. The overflow bit can quickly be computed as $\text{overflow} = c_n \oplus c_{n-1}$. Also, $c_{out} = c_n$.
- To avoid overflow, a common technique is to sign-extend the two summands. For example, for two 4-bits summands, we add an extra bit; so, we use 5 bits to represent our numbers.
- In general, if the two summands are n -bits wide, the result will have at most $n + 1$ bits.
- Recall that if there is no overflow in a summation result, the carry out bit must not be part of the result.

$\begin{array}{r} \text{c}_5 \text{ c}_4 \text{ c}_3 \text{ c}_2 \text{ c}_1 \text{ c}_0 \\ +7 = 0 \ 0 \ 1 \ 1 \ 1 \ + \\ +2 = 0 \ 0 \ 0 \ 1 \ 0 \ + \\ \hline +9 = 0 \ 1 \ 0 \ 0 \ 1 \end{array}$	$\begin{array}{r} \text{c}_5 \text{ c}_4 \text{ c}_3 \text{ c}_2 \text{ c}_1 \text{ c}_0 \\ -7 = 1 \ 1 \ 0 \ 0 \ 1 \ + \\ -2 = 1 \ 1 \ 1 \ 1 \ 0 \ + \\ \hline -9 = 1 \ 0 \ 1 \ 1 \ 1 \end{array}$
--	--

Digital Circuit - Addition

- The figure depicts an n -bit adder for 2's complement numbers:



SUBTRACTION

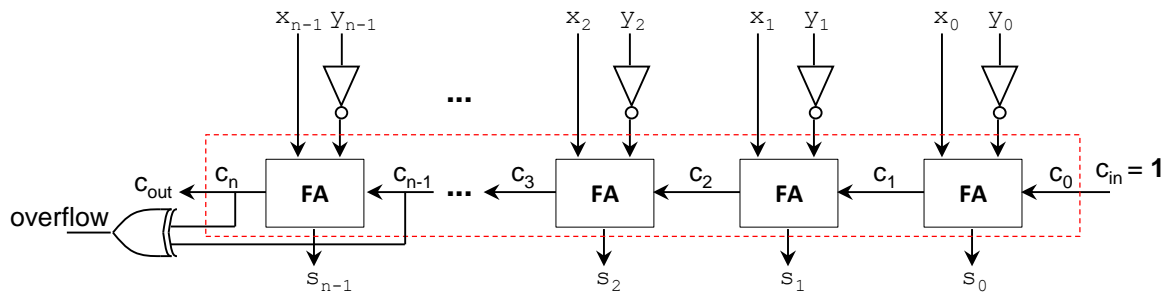
- Note that $A - B = A + 2C(B)$. To subtract two numbers represented in 2's complement arithmetic, we first apply the 2's complement operation to B (the subtrahend), and then add the numbers. So, in 2's complement arithmetic, subtraction is actually an addition of two numbers.
- The digital circuit for 2's complement subtraction is based on the adder. We account for the 2's complement operation for the subtrahend by inverting every bit in the subtrahend and by making the c_{in} bit equal to 1. Here, we give up the borrow in.

7 - 3 = 7 + (-3):

$+3 = 0011 \rightarrow -3 = 1101$

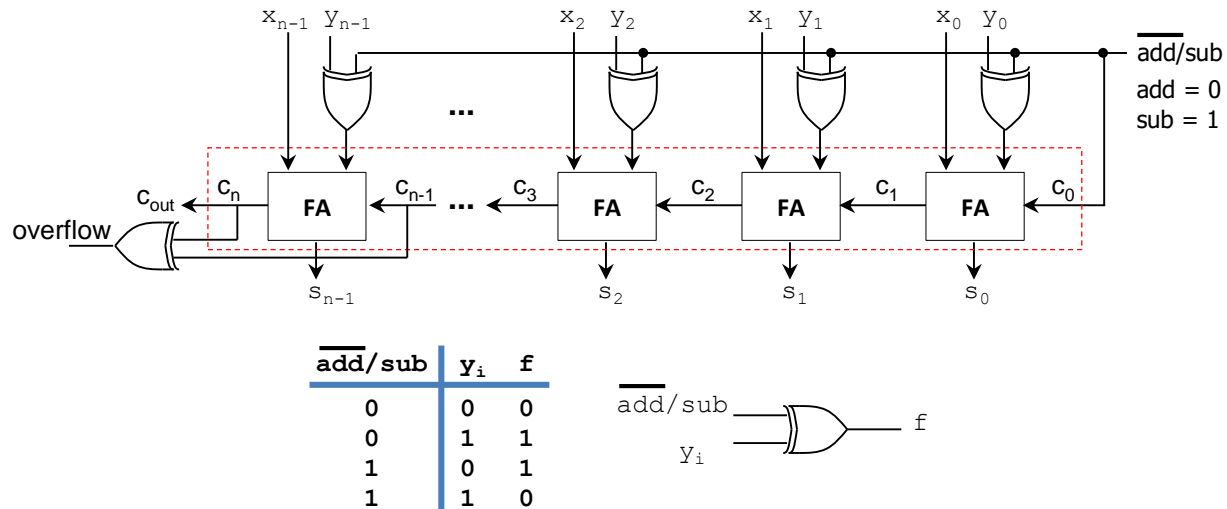
$$\begin{array}{r} \text{c}_4 \text{ c}_3 \text{ c}_2 \text{ c}_1 \text{ c}_0 \\ +7 = 0 \ 1 \ 1 \ 1 \ 1 \ + \\ -3 = 1 \ 1 \ 0 \ 1 \ 1 \ + \\ \hline +4 = 0 \ 1 \ 0 \ 0 \ 0 \end{array}$$

cout = 1
overflow = 0



Adder/Subtractor Unit for 2's complement numbers:

- We can combine the adder and subtractor in a single circuit if we are willing to give up the input c_{in} .



- If there is no carry in (or borrow in): The largest value (in magnitude) of addition of two n -bits operators is $-2^{n-1} + (-2^{n-1}) = -2^n$. In the case of subtraction, the largest value (in magnitude) is $-2^{n-1} - (2^{n-1} - 1) = -2^n + 1$, or $(2^{n-1} - 1) - (-2^{n-1}) = 2^n - 1$. Thus, the addition/subtraction of two n -bit operators needs at most $n + 1$ bits. $c_n = c_{out}$ is used in *multi-precision addition/subtraction*.

SUMMARY

- Addition/Subtraction of two n -bit numbers (no carry in or borrow in):

	UNSIGNED	SIGNED (2C)
Overflow bit	c_n	$c_n \oplus c_{n-1}$
Overflow occurs when:	$A + B \notin [0, 2^n - 1], c_n = 1$	$(A \pm B) \notin [-2^{n-1}, 2^{n-1} - 1], c_n \oplus c_{n-1} = 1$
Result range:	$[0, 2^{n+1} - 2]$	$A + B \in [-2^n, 2^n - 2], A - B \in [-2^n + 1, 2^n - 1]$
Result requires at most:	$n + 1$ bits	

- $c_n = c_{out}$ is used in *multi-precision addition/subtraction* for both signed and unsigned numbers.

MULTIPLICATION OF INTEGER NUMBERS

MULTIPLICATION OF UNSIGNED NUMBERS

- Simple operation: first, generate the products, then add up all the rows column by column (consider the carries).

$$\begin{array}{r}
 \begin{array}{cccc}
 a_3 & a_2 & a_1 & a_0 \\
 b_3 & b_2 & b_1 & b_0
 \end{array} \times \\
 \hline
 a_3b_0 & a_2b_0 & a_1b_0 & a_0b_0 \\
 a_3b_1 & a_2b_1 & a_1b_1 & a_0b_1 \\
 a_3b_2 & a_2b_2 & a_1b_2 & a_0b_2 \\
 a_3b_3 & a_2b_3 & a_1b_3 & a_0b_3
 \end{array}$$

p₇ p₆ p₅ p₄ p₃ p₂ p₁ p₀

$$\begin{array}{r}
 11 \times 13 \\
 \hline
 143
 \end{array}
 \rightarrow
 \begin{array}{r}
 \begin{array}{cccc}
 1 & 0 & 1 & 1 \\
 1 & 1 & 0 & 1
 \end{array} \times \\
 \hline
 \begin{array}{cccc}
 0 & 0 & 0 & 0 \\
 1 & 0 & 1 & 1 \\
 0 & 0 & 0 & 0 \\
 1 & 0 & 1 & 1
 \end{array}
 \end{array}$$

1 0 0 0 1 1 1 1

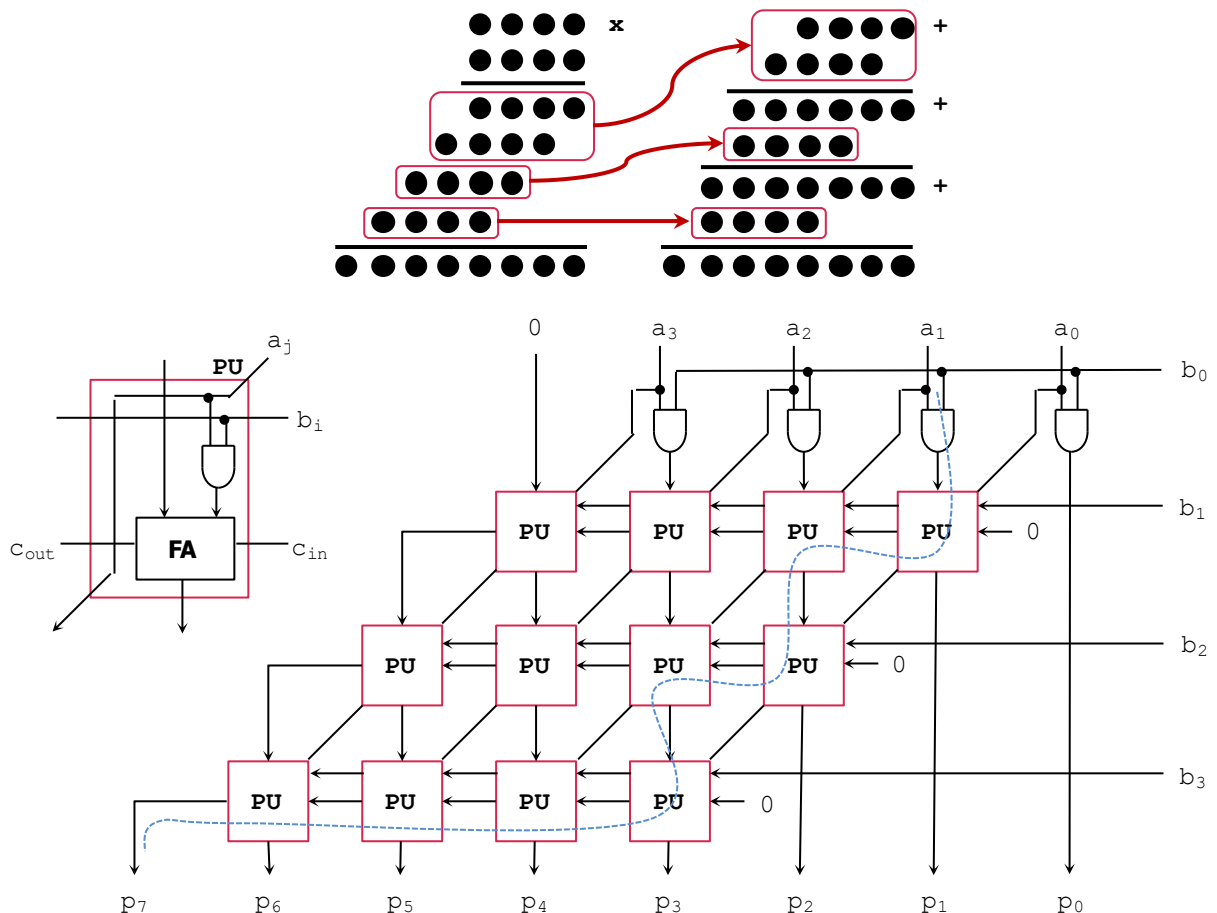
$$\begin{array}{r}
 13 \times 15 \\
 \hline
 195
 \end{array}
 \rightarrow
 \begin{array}{r}
 \begin{array}{cccc}
 1 & 1 & 0 & 1 \\
 1 & 1 & 1 & 1
 \end{array} \times \\
 \hline
 \begin{array}{cccc}
 1 & 0 & 0 & 1 \\
 1 & 1 & 0 & 1 \\
 1 & 1 & 0 & 1 \\
 1 & 1 & 0 & 1
 \end{array}
 \end{array}$$

1 1 0 0 0 0 1 1

- For two 4-bit unsigned numbers A and B, the multiplication $A \times B$ is given by: $A \times B = \left(\sum_{j=0}^3 a_j 2^j \right) \left(\sum_{j=0}^3 b_j 2^j \right)$
- This can be rewritten as: $A \times B = b_0 2^0 \left(\sum_{j=0}^3 a_j 2^j \right) + b_1 2^1 \left(\sum_{j=0}^3 a_j 2^j \right) + b_2 2^2 \left(\sum_{j=0}^3 a_j 2^j \right) + b_3 2^3 \left(\sum_{j=0}^3 a_j 2^j \right)$
- Note that if two operators are n -bits wide, the result (product) needs at most $2n$ bits. This is because the maximum possible value of the product is $(2^n - 1) \times (2^n - 1) = 2^{2n} - 2^{n+1} + 1$, which requires $2n$ bits to be represented.

COMBINATIONAL MULTIPLIER FOR UNSIGNED NUMBERS

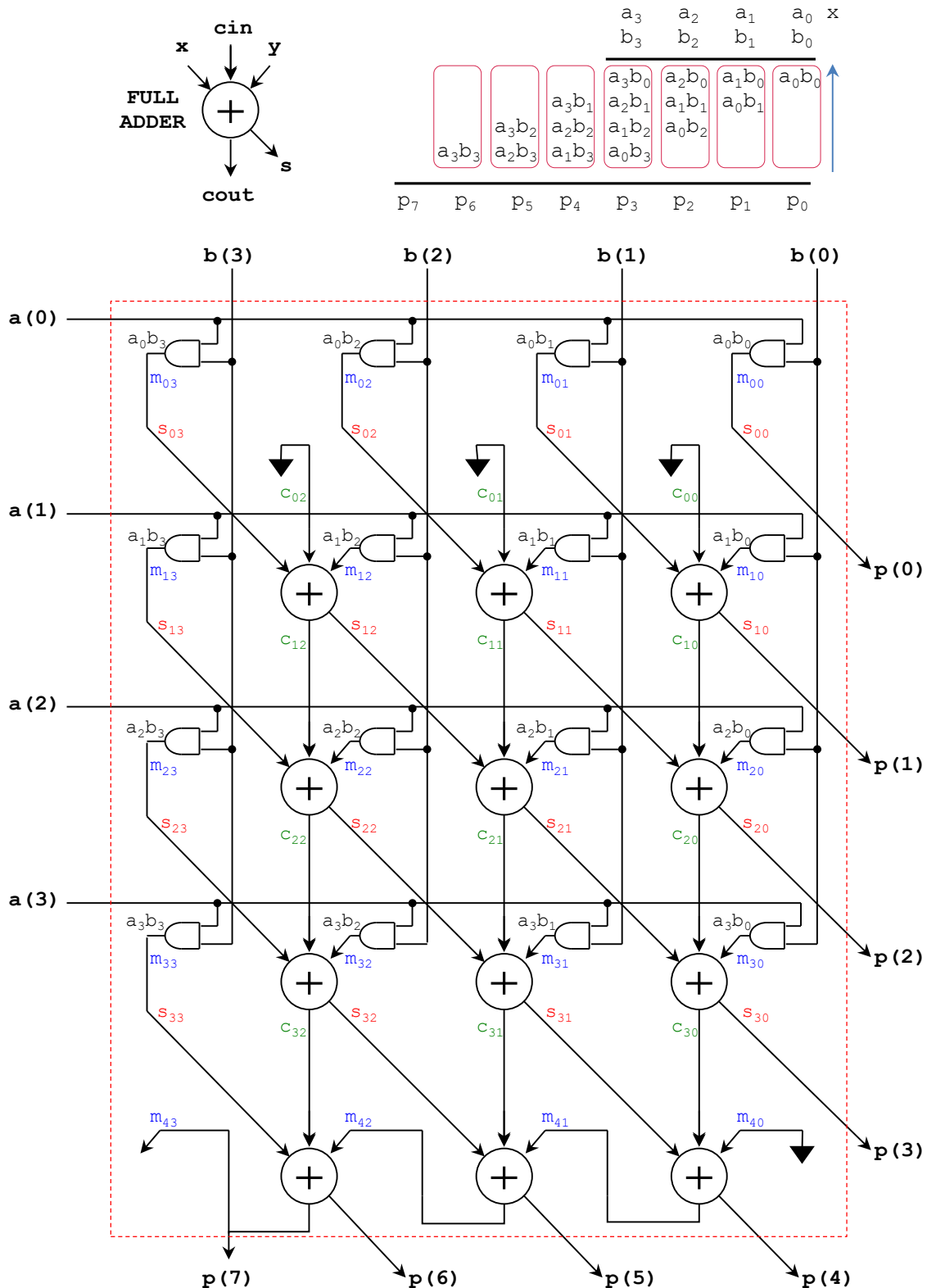
- The partial product computations are implemented by AND gates.
- A straightforward combinational implementation for the multiplication can be achieved by adding two partial products (rows) at each stage. This circuit, also called an Array Multiplier, is depicted in the figure:



- Though simple, this circuit has a large combinational delay from input to output. A signal must travel thru 8 Full Adders and an AND gate (longest path delay). Note that every stage (row) propagates the carries to the left.

ARRAY MULTIPLIER (OTHER) FOR UNSIGNED NUMBERS

- An alternative implementation of the multiplication operation of two 4-bit unsigned numbers is depicted in the figure: at every diagonal of the circuit, we add up all terms in a column of the multiplication.
- The propagation delay is reduced. A signal only needs to travel through 6 Full Adders and an AND gate (longest path delay). Note that every stage (row) does not propagate the carries to the left; instead, they are sent down to the next stage. Only the last stage propagates carries to the left.



MULTIPLICATION OF SIGNED NUMBERS

- A straightforward implementation consists of checking the sign of the multiplicand and multiplier. If one or both are negative, we change the sign by applying the 2's complement operation. This way, we are left with unsigned multiplication.
- As for the final output: if only one of the inputs was negative, then we modify the sign of the output. Otherwise, the result of the unsigned multiplication is the final output.

$\begin{array}{r} 101 \times 011 \\ 010 \times 010 \\ \hline 000 \\ 011 \\ 000 \\ \hline 000110 \\ \hline 111010 \end{array}$	$\begin{array}{r} 010 \times 010 \\ 110 \times 010 \\ \hline 000 \\ 010 \\ 000 \\ \hline 000100 \\ \hline 111100 \end{array}$	$\begin{array}{r} 111 \times 001 \\ 110 \times 010 \\ \hline 000 \\ 001 \\ 000 \\ \hline 000010 \\ \hline 000010 \end{array}$	$\begin{array}{r} 011 \times 011 \\ 010 \times 010 \\ \hline 000 \\ 011 \\ 000 \\ \hline 000110 \\ \hline 000110 \end{array}$
---	---	---	---

- Note:** If one of the inputs is -2^{n-1} , then the negative version is 2^{n-1} , which requires $n + 1$ bits. Here, we are allowed to use only n bits; in other words, we do not have to change its sign. This will not affect the final result since if we were to use $n + 1$ bits for 2^{n-1} , the MSB=0, which implies that i) the last row if full of zeros, or ii) there is an extra '0' to the MSB of every summand.

$\begin{array}{r} 100 \times 100 \\ 011 \times 011 \\ \hline 100 \\ 100 \\ 000 \\ \hline 001100 \\ \hline 110100 \end{array}$	$\begin{array}{r} 011 \times 011 \\ 100 \times 100 \\ \hline 000 \\ 000 \\ 011 \\ \hline 001100 \\ \hline 110100 \end{array}$	$\begin{array}{r} 100 \times 100 \\ 100 \times 100 \\ \hline 000 \\ 000 \\ 100 \\ \hline 010000 \\ \hline 010000 \end{array}$
---	---	---

- Final output: It requires $2n$ bits. Note that it is only because of the multiplication of -2^{n-1} by -2^{n-1} that we require those $2n$ bits (in 2's complement representation)

BINARY CODES

- We know that with n bits, we can represent 2^n numbers, from 0 to $2^n - 1$. This is a commonly used range. However, with 'n' bits, we can also represent 2^n numbers in any range.
- Moreover, with n bits we can represent 2^n different symbols. For example, in 24-bit color, each color is represented by 24 bits, providing 2^{24} distinct colors. Each color is said to have a binary code.
- $N = 5$ symbols. With 2 bits, only 4 symbols can be represented. With 3 bits, 8 symbols can be represented. Thus, the number of bits required is $n = 3 = \lceil \log_2 5 \rceil = \log_2 8$. Note that 8 is the power of 2 closest to $N=5$ that is greater than or equal to 5.
- In general, if we have N symbols to represent, the number of bits required is given by $\lceil \log_2 N \rceil$. For example:
 - ✓ Minimum number of bits to represent 70,000 colors: \rightarrow Number of bits: $\lceil \log_2 70000 \rceil = 17$ bits.
 - ✓ Minimum number of bits to represent numbers between 15,000 and 19,096: \rightarrow There are $19,096 - 15,000 + 1 = 4097$. Then, number of bits: $\lceil \log_2 4097 \rceil = 13$ bits.

7-bit US-ASCII character-encoding scheme: Each character is represented by 7 bits. Thus, the number of characters that can be represented is given by $2^7 = 128$. Each character is said to have a binary code.

Hex	Dec	Char	Hex	Dec	Char	Hex	Dec	Char	Hex	Dec	Char
0x00	0	NULL null	0x20	32	Space	0x40	64	@	0x60	96	`
0x01	1	SOH Start of heading	0x21	33	!	0x41	65	A	0x61	97	a
0x02	2	STX Start of text	0x22	34	"	0x42	66	B	0x62	98	b
0x03	3	ETX End of text	0x23	35	#	0x43	67	C	0x63	99	c
0x04	4	EOT End of transmission	0x24	36	\$	0x44	68	D	0x64	100	d
0x05	5	ENQ Enquiry	0x25	37	%	0x45	69	E	0x65	101	e
0x06	6	ACK Acknowledge	0x26	38	&	0x46	70	F	0x66	102	f
0x07	7	BELL Bell	0x27	39	'	0x47	71	G	0x67	103	g
0x08	8	BS Backspace	0x28	40	(0x48	72	H	0x68	104	h
0x09	9	TAB Horizontal tab	0x29	41)	0x49	73	I	0x69	105	i
0x0A	10	LF New line	0x2A	42	*	0x4A	74	J	0x6A	106	j
0x0B	11	VT Vertical tab	0x2B	43	+	0x4B	75	K	0x6B	107	k
0x0C	12	FF Form Feed	0x2C	44	,	0x4C	76	L	0x6C	108	l
0x0D	13	CR Carriage return	0x2D	45	-	0x4D	77	M	0x6D	109	m
0x0E	14	SO Shift out	0x2E	46	.	0x4E	78	N	0x6E	110	n
0x0F	15	SI Shift in	0x2F	47	/	0x4F	79	O	0x6F	111	o
0x10	16	DLE Data link escape	0x30	48	0	0x50	80	P	0x70	112	p
0x11	17	DC1 Device control 1	0x31	49	1	0x51	81	Q	0x71	113	q
0x12	18	DC2 Device control 2	0x32	50	2	0x52	82	R	0x72	114	r
0x13	19	DC3 Device control 3	0x33	51	3	0x53	83	S	0x73	115	s
0x14	20	DC4 Device control 4	0x34	52	4	0x54	84	T	0x74	116	t
0x15	21	NAK Negative ack	0x35	53	5	0x55	85	U	0x75	117	u
0x16	22	SYN Synchronous idle	0x36	54	6	0x56	86	V	0x76	118	v
0x17	23	ETB End transmission block	0x37	55	7	0x57	87	W	0x77	119	w
0x18	24	CAN Cancel	0x38	56	8	0x58	88	X	0x78	120	x
0x19	25	EM End of medium	0x39	57	9	0x59	89	Y	0x79	121	y
0x1A	26	SUB Substitute	0x3A	58	:	0x5A	90	Z	0x7A	122	z
0x1B	27	FSC Escape	0x3B	59	;	0x5B	91	[0x7B	123	{
0x1C	28	FS File separator	0x3C	60	<	0x5C	92	\	0x7C	124	
0x1D	29	GS Group separator	0x3D	61	=	0x5D	93]	0x7D	125	}
0x1E	30	RS Record separator	0x3E	62	>	0x5E	94	^	0x7E	126	~
0x1F	31	US Unit separator	0x3F	63	?	0x5F	95	_	0x7F	127	DEL

Unicode: This code can represent more than 110,000 characters and attempts to cover all world's scripts. A common character encoding is UTF-16, which uses 2 pair of 16-bit units: For most purposes, a 16 bit unit suffices ($2^{16} = 65536$ characters):

Θ (Greek theta symbol) = 03D1 Ω (Greek capital letter Omega): 03A9 Ж (Cyrillic capital letter zhe): 0416

BCD Code:

- In this coding scheme, decimal numbers are represented in binary form by independently encoding each decimal digit in binary form. Each digit requires 4 bits. Note that only values from 0 are 9 are represented here.
- This is a very useful code for input devices (e.g.: keypad). But it is not a coding scheme suitable for arithmetic operations. Also, notice that the binary numbers $1011_2(10)$ to $1111_2(15)$ are not used. Only 10 out of 16 values are used to encode each decimal digit.

BCD decimal #

0	0000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111
8	1000
9	1001

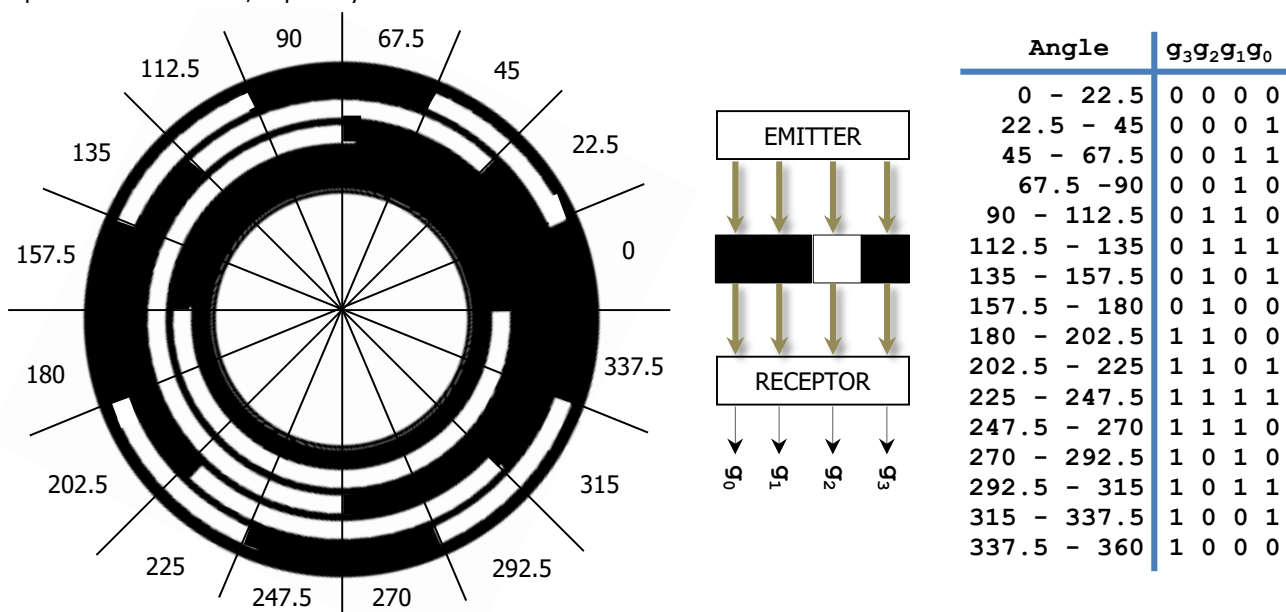
Examples:

- ✓ Decimal number **47**: This decimal number can be represented as a binary number: 101111_2 . In BCD format, this would be: **0100 0111**₂
 - ✓ Decimal number **58**: This decimal number can be represented as a binary number: 111010_2 . In BCD format, the binary representation would be: **0101 1000**₂
 - ✓ The BCD code is not the same as the binary number!
- There exist many other binary codes (e.g., reflective gray code, 6-3-1-1 code, 2-out-of-5 code) to represent decimal numbers. Usually, each of them is tailored to a specific application.

REFLECTIVE GRAY CODE:

g_1g_0	Decimal Number	$b_2b_1b_0$	$g_2g_1g_0$	$g_3g_2g_1g_0$	
0 0	0	0 0 0	0 0 0	0 0 0 0	$ \begin{array}{ccccccc} b_{n-1} & b_{n-2} & \dots & b_1 & b_0 \\ \downarrow & \downarrow & & \downarrow & \downarrow \\ g_{n-1} & g_{n-2} & & g_1 & g_0 \end{array} $
0 1	1	0 0 1	0 0 1	0 0 0 1	
1 1	2	0 1 0	0 1 1	0 0 1 1	$ \begin{array}{ccccccc} g_{n-1} & g_{n-2} & \dots & g_1 & g_0 \\ \downarrow & \downarrow & & \downarrow & \downarrow \\ b_{n-1} & b_{n-2} & & b_1 & b_0 \end{array} $
1 0	3	0 1 1	0 1 0	0 0 1 0	
	4	1 0 0	1 1 0	0 1 1 0	
	5	1 0 1	1 1 1	0 1 1 1	
	6	1 1 0	1 0 1	0 1 0 1	
	7	1 1 1	1 0 0	0 1 0 0	
				1 1 0 0	
				1 1 0 1	
				1 1 1 1	
				1 1 1 0	
				1 0 1 0	
				1 0 1 1	
				1 0 0 1	
				1 0 0 0	

- Application:** Measuring angular position with 4-bit RGC. 4 beams are emitted along an axis. When a light beam passes (transparent spots, represented as whites), we get a logical 1, 0 otherwise. The RGC encoding makes that between areas only one bit changes, thereby reducing the possibility of an incorrect reading (especially when the beam between adjacent areas). For example: from 0001 to 0011 only one bit flips. If we used 0001 to 0010, two bits would flip: that would be prone to more errors, especially when the beams are close to the line where the two areas meet.



INTRODUCTION TO FIXED-POINT ARITHMETIC

- We have been representing positive integer numbers. But what if we wanted to represent numbers with fractional parts?
- Fixed-point arithmetic: Binary representation of positive decimal numbers with fractional parts.

Given the following binary number:

$$(b_{n-1}b_{n-2} \dots b_1b_0.b_{-1}b_{-2} \dots b_{-k})_2$$

Formula to convert it to decimal:

$$D = \sum_{i=-k}^{n-1} b_i \times 2^i = b_{n-1} \times 2^{n-1} + b_{n-2} \times 2^{n-2} + \dots + b_1 \times 2^1 + b_0 \times 2^0 + b_{-1} \times 2^{-1} + b_{-2} \times 2^{-2} + \dots + b_{-k} \times 2^{-k}$$

Conversion from binary to hexadecimal (or octal): (unsigned numbers)

Binary: $1101.11_2 \rightarrow 001 \ 101.110$

octal: $\downarrow \quad \downarrow \quad \downarrow$
 $1 \quad 5 \ . \ 6$

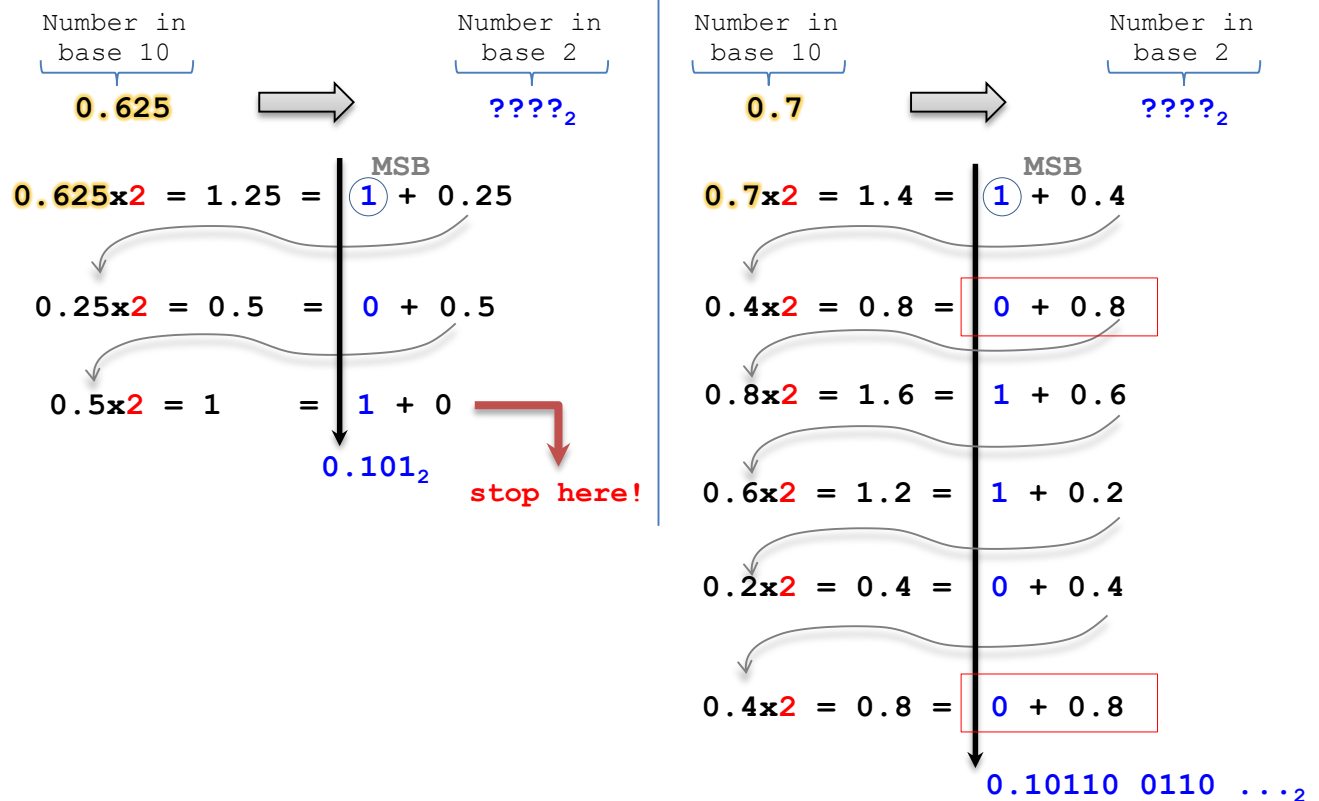
Binary: $10101.10101_2 \rightarrow 0001 \ 0101.1010 \ 1000$

hexadecimal: $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $1 \quad 5 \ . \ A \quad 8$

- Example:** (unsigned number)

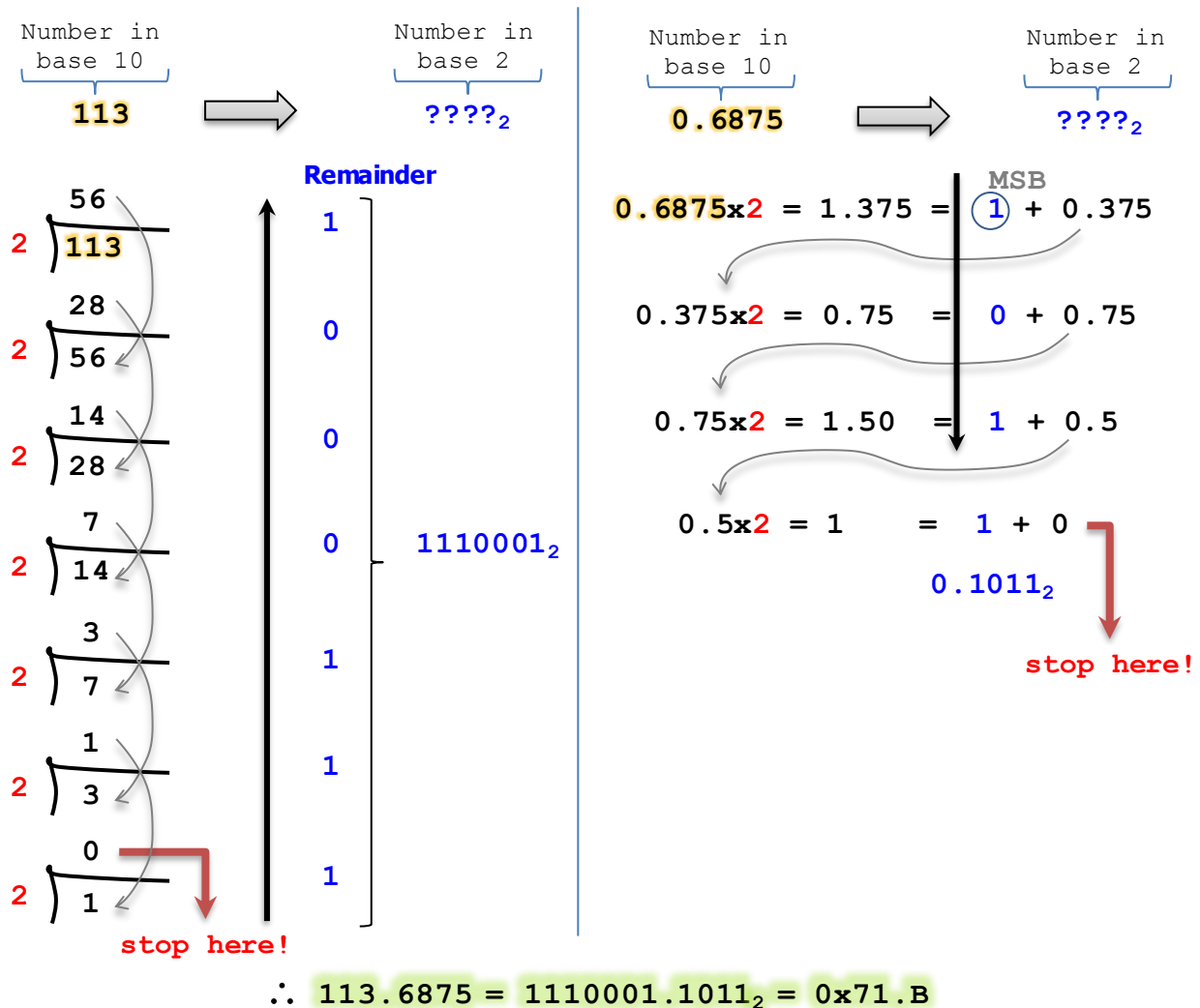
$$1011.101_2 = 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} = 11.625$$

- Example:** Now, what if we have a decimal number with fractional part? What we do is we divide the integer part and the fractional part. We obtain the binary representation of the integer part using what we know. As for the fractional part, what we do is successive multiplications by 2, the resulting integer parts resulting is the result.



- Example (signed number):** Convert -379.21875 to the 2's complement representation.
First, we get the binary representation of +379.21875, and then apply the 2's complement operation to that result.
 - ✓ $379 = 101111011_2$. In 2's complement: $379 = 0101111011_2$, $0.21875 = 0.00111_2$.
 - ✓ Then: $379.21875 = 0101111011.00111_2$. This is the 2's complement representation of 379.21875.
 - ✓ Finally, we get -379.2185 by getting the 2's complement of the previous result: $-379.21875 = 1010000100.11001_2 = 0xE84.C8$ (to convert to hexadecimal, we append zeros to the LSB and sign-extend the MSB)

- **Example (unsigned number):** Convert 113.6875 to the binary representation (unsigned integer).



PRACTICE EXERCISES

- Determine the radix r such that:
 - ✓ $645_r = 327_{10}$
 - ✓ $150_r + 256_r = 426_r$
 - ✓ $45_r \times 6_r = 450_r$
- Represent the unsigned decimal number 1237 in the following formats. Determine the number of bits required in each case.
 - ✓ BCD
 - ✓ Binary
 - ✓ 7-bit ASCII
- Complete the following table:

REPRESENTATION			
Decimal	Sign-and-magnitude	1's complement	2's complement
		010011	
			1
	11000		
			10000
-15			
		111	

- Convert the following decimal numbers to their 2's complement representations: binary and hexadecimal.
-93.65625, -256.6875, 31.6875, -138.625.