

Let us give Mr C a Complement

ESC101: Fundamentals of Computing

Purushottam Kar

The Binary Number System

2



The Binary Number System

Just two digits – 0 and 1

2



The Binary Number System

2

Just two digits – 0 and 1

--	--	--	--	--	--



The Binary Number System

2

Just two digits – 0 and 1

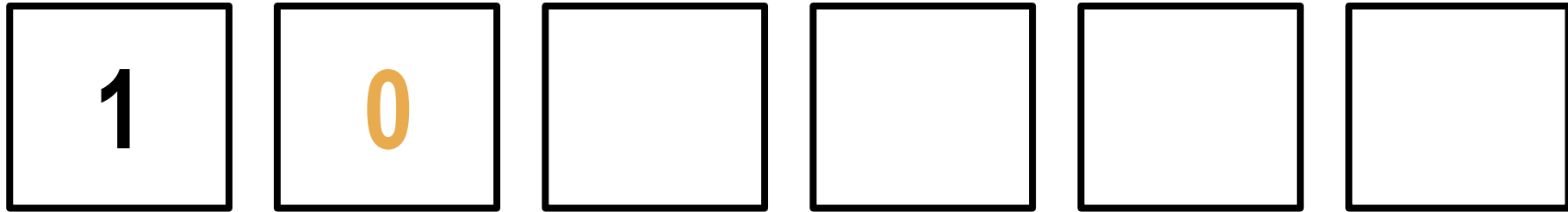
1					
---	--	--	--	--	--



The Binary Number System

2

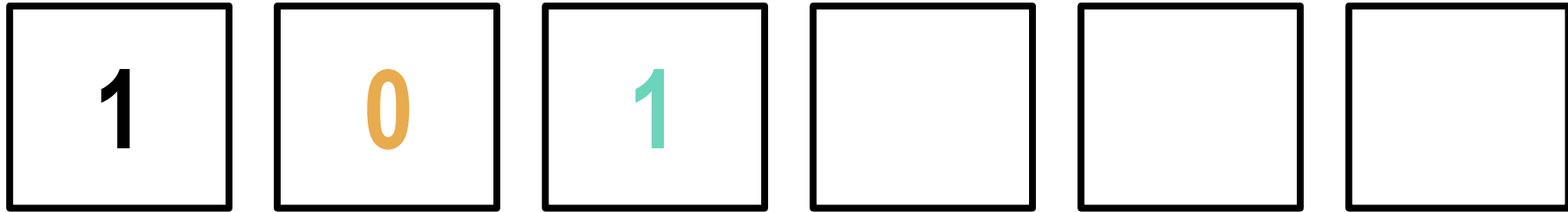
Just two digits – 0 and 1



The Binary Number System

2

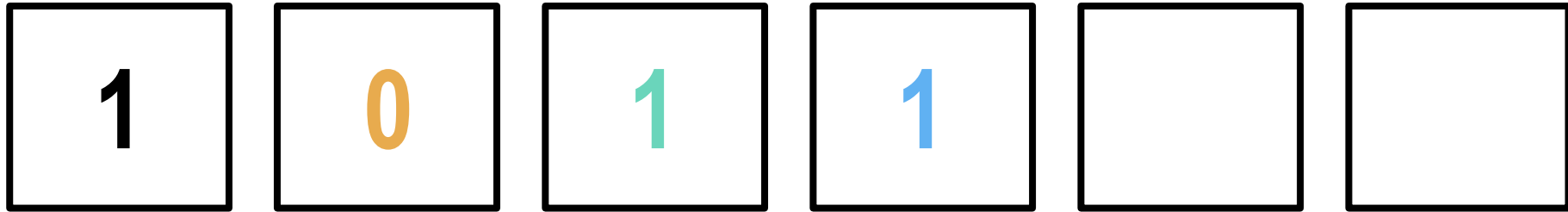
Just two digits – 0 and 1



The Binary Number System

2

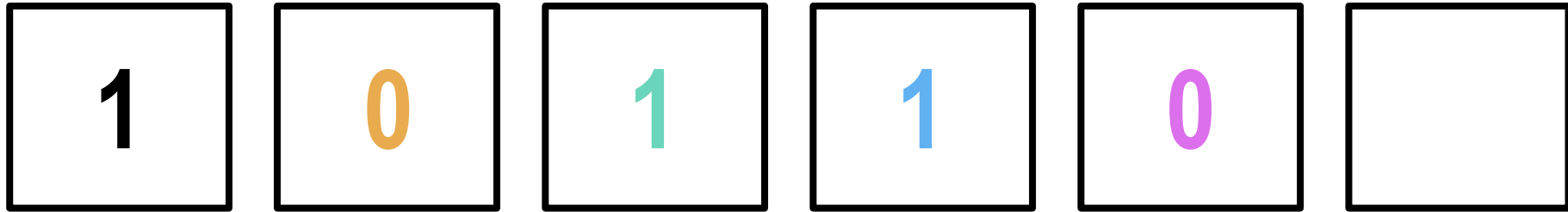
Just two digits – 0 and 1



The Binary Number System

2

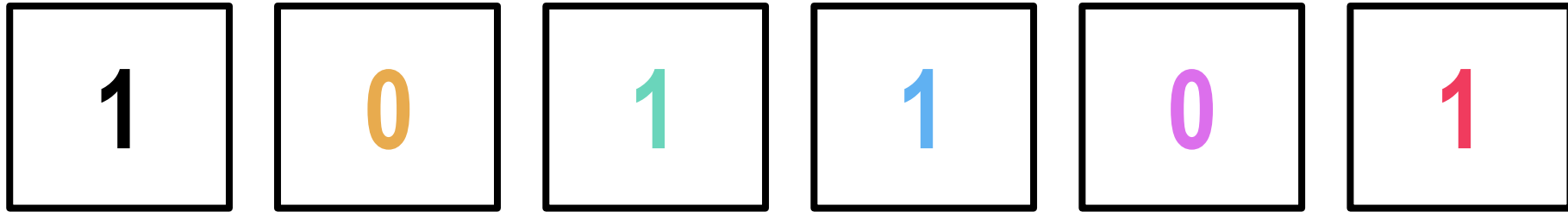
Just two digits – 0 and 1



The Binary Number System

2

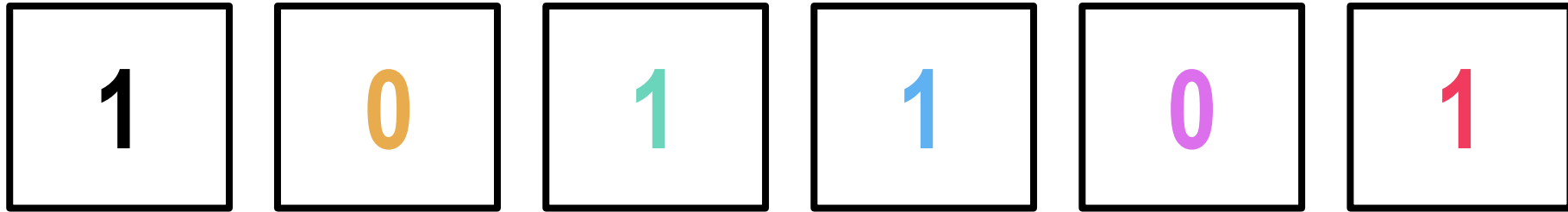
Just two digits – 0 and 1



The Binary Number System

2

Just two digits – 0 and 1



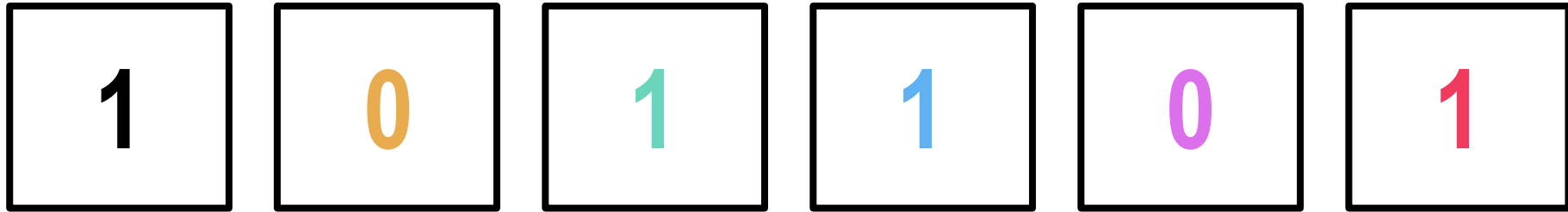
$$= 1 \times 2^0 + 0 \times 2^1 + 1 \times 2^2 + 1 \times 2^3 + 0 \times 2^4 + 1 \times 2^5$$



The Binary Number System

2

Just two digits – 0 and 1



$$= 1 \times 2^0 + 0 \times 2^1 + 1 \times 2^2 + 1 \times 2^3 + 0 \times 2^4 + 1 \times 2^5$$

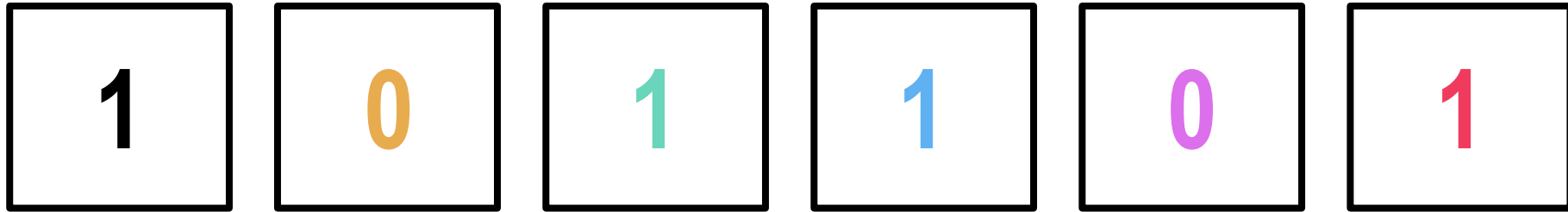
$$= 1 + 0 + 4 + 8 + 0 + 32$$



The Binary Number System

2

Just two digits – 0 and 1



$$= 1 \times 2^0 + 0 \times 2^1 + 1 \times 2^2 + 1 \times 2^3 + 0 \times 2^4 + 1 \times 2^5$$

$$= 1 + 0 + 4 + 8 + 0 + 32$$

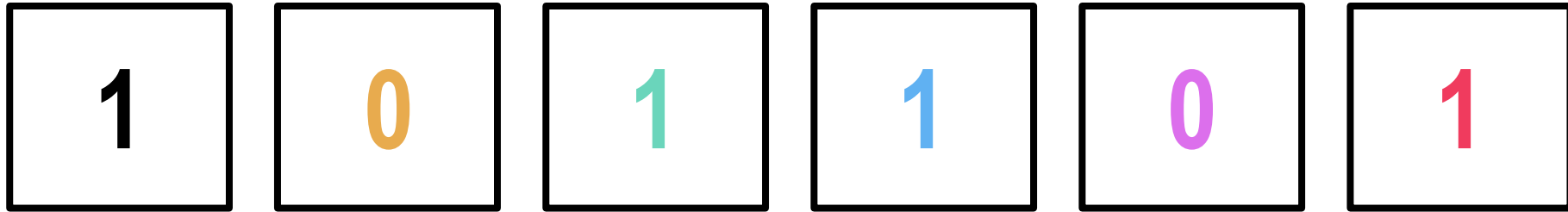
$$= 45$$



The Binary Number System

2

Just two digits – 0 and 1



$$= 1 \times 2^0 + 0 \times 2^1 + 1 \times 2^2 + 1 \times 2^3 + 0 \times 2^4 + 1 \times 2^5$$

$$= 1 + 0 + 4 + 8 + 0 + 32$$

$$= 45$$

We will today see how the binary system is used to represent negative and fractional numbers



Bitwise Operators

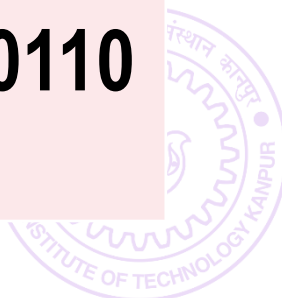
3



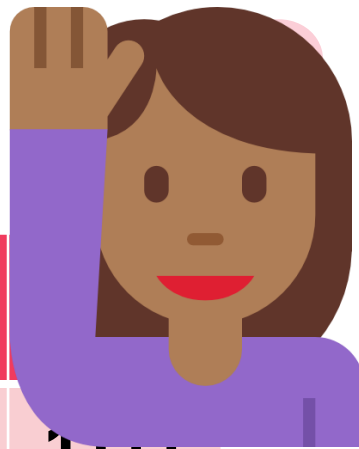
Bitwise Operators

3

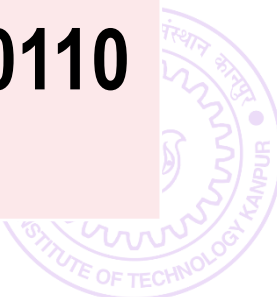
Operation	C Code	a	b	c	d	e	f
BITWISE AND	c = a & b	0000	1111	0000	1111	1111	1111
BITWISE OR	d = a b	0101	1100	0100	1101	1001	1010
BITWISE XOR	e = a ^ b	1010	1110	1010	1110	0100	0101
BITWISE COMPLEMENT	f = ~a	1001	0111	0001	1111	1110	0110



Bitwise Operators



Operation	C Code	a	b	c	d	e	f
BITWISE AND	c = a & b	0000	1111	0000	1111	1111	1111
BITWISE OR	d = a b	0101	1100	0100	1101	1001	1010
BITWISE XOR	e = a ^ b	1010	1110	1010	1110	0100	0101
BITWISE COMPLEMENT	f = ~a	1001	0111	0001	1111	1110	0110

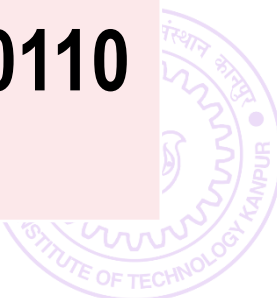


Bitwise Operators

Note: the word is complement and not compliment 😊



Operation	C Code	a	b	c	d	e	f
BITWISE AND	c = a & b	0000	1111	0000	1111	1111	1111
BITWISE OR	d = a b	0101	1100	0100	1101	1001	1010
BITWISE XOR	e = a ^ b	1010	1110	1010	1110	0100	0101
BITWISE COMPLEMENT	f = ~a	1001	0111	0001	1111	1110	0110

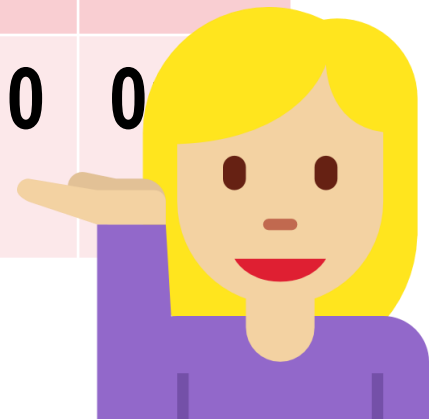


Bitwise Operators

Note: the word is complement and not compliment 😊



Operation	C Code	a	b	c	d	e	f
BITWISE AND	c = a & b	0000	1111	0000	1111	1111	1111
BITWISE OR	d = a b	0101	1100	0100	1101	1001	1010
BITWISE XOR	e = a ^ b	1010	1110	1010	1110	0100	0101
BITWISE COMPLEMENT	f = ~a	1001	0111	0001	1111	1110	0100



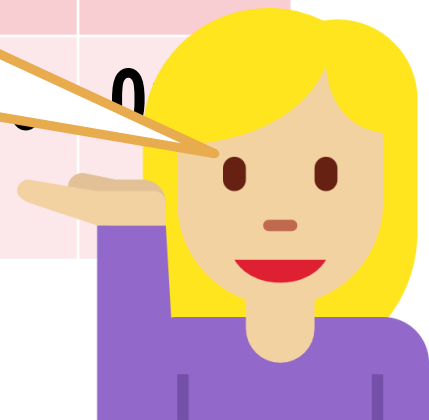
Bitwise Operators

Note: the word is complement and not compliment 😊



Operation	C Code	a	b	c	d	e	f
BITWISE AND	c = a & b	0000	1111	0000	1111	1111	1111
BITWISE OR	d = a b	0101	1100	0100	1101	1001	1010
BITWISE XOR	e = a ^ b	1010	1110	1010	1110	0100	0101
BITWISE COMPLEMENT	f = ~a	1010	0000	0101	0000	1111	0101

The word compliment means to praise someone, complement means to enhance something



Bitwise And Operator &

4



Bitwise And Operator &

4

The output of bitwise AND is 1 if the corresponding bits of two operands is 1. If either bit of an operand is 0, the result of corresponding bit is evaluated to 0



Bitwise And Operator &

4

The output of bitwise AND is 1 if the corresponding bits of two operands is 1. If either bit of an operand is 0, the result of corresponding bit is evaluated to 0

In C Programming, bitwise AND operator is denoted by &



Bitwise And Operator &

4

The output of bitwise AND is 1 if the corresponding bits of two operands is 1. If either bit of an operand is 0, the result of corresponding bit is evaluated to 0

In C Programming, bitwise AND operator is denoted by &

12 = 00001100 (In Binary)

25 = 00011001 (In Binary)

Bitwise AND of 12 and 25

0000 1100

& 0001 1001

0000 1000 = 8 (In decimal)



Bitwise And Operator &

4

The output of bitwise AND is 1 if the corresponding bits of two operands is 1. If either bit of an operand is 0, the result of corresponding bit is evaluated to 0

In C Programming, bitwise AND operator is denoted by &

```
12 = 00001100 (In Binary)
25 = 00011001 (In Binary)
Bitwise AND of 12 and 25
  0000 1100
& 0001 1001
  _____
  0000 1000 = 8 (In decimal)
```

```
#include <stdio.h>
int main(){
    int a = 12, b = 25;
    printf("Output = %d", a & b);
    return 0;
}
```



Bitwise OR Operator |

5



Bitwise OR Operator |

5

The output of bitwise OR is 1 if at least one corresponding bit of two operands is 1



Bitwise OR Operator |

5

The output of bitwise OR is 1 if at least one corresponding bit of two operands is 1

In C Programming, bitwise OR operator is denoted by |



Bitwise OR Operator |

5

The output of bitwise OR is 1 if at least one corresponding bit of two operands is 1

In C Programming, bitwise OR operator is denoted by |

12 = 00001100 (In Binary)

25 = 00011001 (In Binary)

Bitwise OR of 12 and 25

0000 1100

| 0001 1001

00011101 = 29 (In decimal)



Bitwise OR Operator |

5

The output of bitwise OR is 1 if at least one corresponding bit of two operands is 1

In C Programming, bitwise OR operator is denoted by |

```
12 = 00001100 (In Binary)
25 = 00011001 (In Binary)
Bitwise OR of 12 and 25
  0000 1100
| 0001 1001
  _____
  00011101 = 29 (In decimal)
```

```
#include <stdio.h>
int main(){
    int a = 12, b = 25;
    printf("Output = %d", a | b);
    return 0;
}
```



Bitwise XOR Operator \wedge

6



Bitwise XOR Operator \wedge

6

The result of bitwise XOR operator is 1 if the corresponding bits of two operands are opposite i.e. one is 1 and the other is 0



Bitwise XOR Operator \wedge

6

The result of bitwise XOR operator is 1 if the corresponding bits of two operands are opposite i.e. one is 1 and the other is 0

In C Programming, bitwise XOR operator is denoted by \wedge



Bitwise XOR Operator ^

6

The result of bitwise XOR operator is 1 if the corresponding bits of two operands are opposite i.e. one is 1 and the other is 0

In C Programming, bitwise XOR operator is denoted by ^

12 = 00001100 (In Binary)

25 = 00011001 (In Binary)

Bitwise XOR of 12 and 25

00001100

^ 00011001

00010101 = 21 (In decimal)



Bitwise XOR Operator ^

6

The result of bitwise XOR operator is 1 if the corresponding bits of two operands are opposite i.e. one is 1 and the other is 0

In C Programming, bitwise XOR operator is denoted by ^

```
12 = 00001100 (In Binary)
25 = 00011001 (In Binary)
Bitwise XOR of 12 and 25
  00001100
^ 00011001
  _____
  00010101 = 21 (In decimal)
```

```
#include <stdio.h>
int main(){
    int a = 12, b = 25;
    printf("Output = %d", a^b);
    return 0;
}
```



Bitwise Complement Operator ~

7



Bitwise Complement Operator ~

7

A unary operator that simply flips each bit of the input



Bitwise Complement Operator ~

7

A unary operator that simply flips each bit of the input

In C Programming, bitwise complement operator is denoted by ~



Bitwise Complement Operator ~

7

A unary operator that simply flips each bit of the input

In C Programming, bitwise complement operator is denoted by ~

12 = 0000 0000 0000 0000 0000 0000 0000 1100

Bitwise complement of 12 00001100

~ 0000 0000 0000 0000 0000 0000 0000 1100

1111 1111 1111 1111 1111 1111 1111 0011

= -13 (decimal)



Bitwise Complement Operator ~

7

A unary operator that simply flips each bit of the input

In C Programming, bitwise complement operator is denoted by ~

12 = 0000 0000 0000 0000 0000 0000 0000 1100

Bitwise complement of 12 00001100

~ 0000 0000 0000 0000 0000 0000 0000 1100

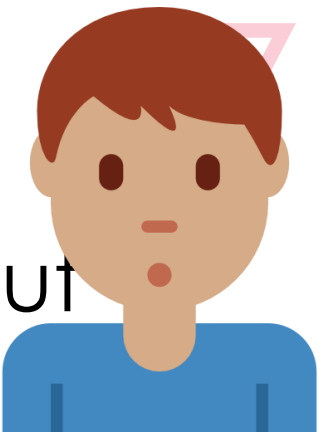
1111 1111 1111 1111 1111 1111 1111 0011

= -13 (decimal)

```
#include <stdio.h>
int main(){
    int a = 12;
    printf("Output = %d", ~a);
    return 0;
}
```



Bitwise Complement Operator ~



A unary operator that simply flips each bit of the input

In C Programming, bitwise complement operator is denoted by ~

12 = 0000 0000 0000 0000 0000 0000 0000 1100

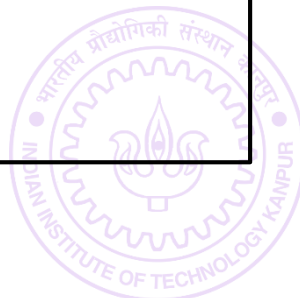
Bitwise complement of 12 00001100

~ 0000 0000 0000 0000 0000 0000 0000 1100

1111 1111 1111 1111 1111 1111 1111 0011

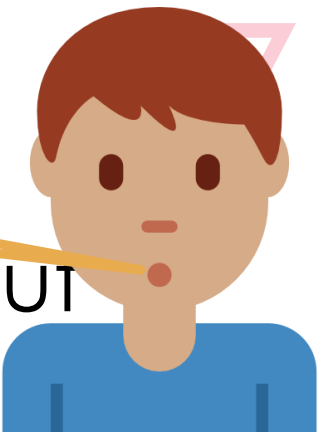
= -13 (decimal)

```
#include <stdio.h>
int main(){
    int a = 12;
    printf("Output = %d", ~a);
    return 0;
}
```



Bitwise Complement

Why does flipping bits of 12 generate -13?



A unary operator that simply flips each bit of the input

In C Programming, bitwise complement operator is denoted by ~

12 = 0000 0000 0000 0000 0000 0000 0000 1100

Bitwise complement of 12 00001100

~ 0000 0000 0000 0000 0000 0000 0000 1100

1111 1111 1111 1111 1111 1111 1111 0011

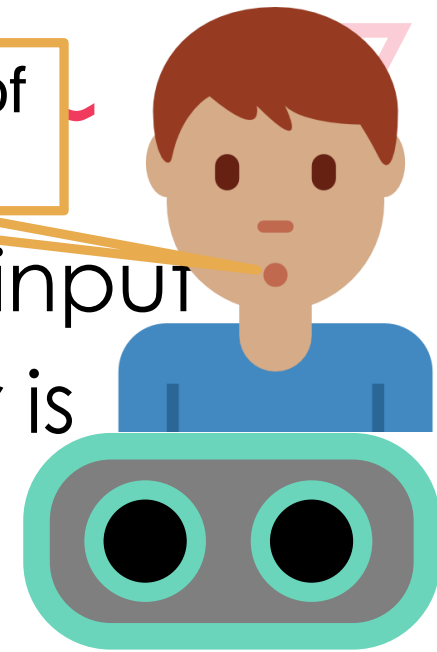
= -13 (decimal)

```
#include <stdio.h>
int main(){
    int a = 12;
    printf("Output = %d", ~a);
    return 0;
}
```



Bitwise Complement

Why does flipping bits of 12 generate -13?



A unary operator that simply flips each bit of the input

In C Programming, bitwise complement operator is denoted by ~

12 = 0000 0000 0000 0000 0000 0000 0000 1100

Bitwise complement of 12 00001100

~ 0000 0000 0000 0000 0000 0000 0000 1100

1111 1111 1111 1111 1111 1111 1111 0011

= -13 (decimal)

```
#include <stdio.h>
int main(){
    int a = 12;
    printf("Output = %d", ~a);
    return 0;
}
```



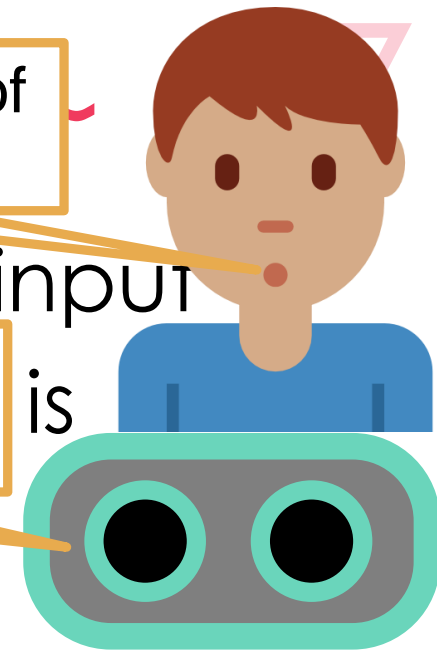
Bitwise Complement

Why does flipping bits of 12 generate -13?

A unary operator that simply flips each bit of the input

In C Programming, bitwise complement is denoted by ~

Wait just a sec – I will teach you how negative numbers are represented



12 = 0000 0000 0000 0000 0000 0000 0000 1100

Bitwise complement of 12 00001100

~ 0000 0000 0000 0000 0000 0000 0000 1100

1111 1111 1111 1111 1111 1111 1111 0011

= -13 (decimal)

```
#include <stdio.h>
int main(){
    int a = 12;
    printf("Output = %d", ~a);
    return 0;
}
```



Right Shift Operator >>

8



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations

Bits that “fall off” from the right most end are lost



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations

Bits that “fall off” from the right most end are lost

Blank spaces in the leftmost positions are filled with 0s



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations

Bits that “fall off” from the right most end are lost

Blank spaces in the leftmost positions are filled with 0s

212 = 0000 0000 0000 0000 0000 0000 1101 0100



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations

Bits that “fall off” from the right most end are lost

Blank spaces in the leftmost positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations

Bits that “fall off” from the right most end are lost

Blank spaces in the leftmost positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101$



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations

Bits that “fall off” from the right most end are lost

Blank spaces in the leftmost positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101$

$212 \gg 6 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0011$



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations

Bits that “fall off” from the right most end are lost

Blank spaces in the leftmost positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101$

$212 \gg 6 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0011$

$212 \gg 3 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0001\ 1010$



Right Shift Operator >>

8

Right shift operator shifts all bits towards right by a certain number of locations

Bits that “fall off” from the right most end are lost

Blank spaces in the leftmost positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \gg 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101$

$212 \gg 6 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0011$

$212 \gg 3 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0001\ 1010$

Right shift by k is equivalent to integer division with 2^k 😊

Left Shift Operator <<

9



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost

Blank spaces in the right positions are filled with 0s



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost

Blank spaces in the right positions are filled with 0s

212 = 0000 0000 0000 0000 0000 0000 1101 0100



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost

Blank spaces in the right positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost

Blank spaces in the right positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100\ 0000$



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost

Blank spaces in the right positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100\ 0000$

$212 \ll 6 = 0000\ 0000\ 0000\ 0000\ 0011\ 0101\ 0000\ 0000$



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost

Blank spaces in the right positions are filled with 0s

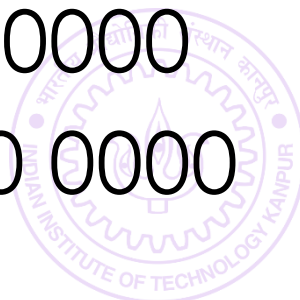
$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100\ 0000$

$212 \ll 6 = 0000\ 0000\ 0000\ 0000\ 0011\ 0101\ 0000\ 0000$

$212 \ll 28 = 0100\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000$



Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost

Blank spaces in the right positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100\ 0000$

$212 \ll 6 = 0000\ 0000\ 0000\ 0000\ 0011\ 0101\ 0000\ 0000$

$212 \ll 28 = 0100\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000$

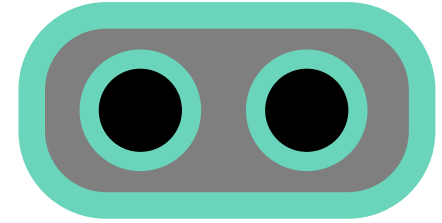
Left shift by k is equivalent to integer multiplication with 2^k

Left Shift Operator <<

9

Left shift operator shifts all bits towards left by a certain number of locations

Bits that “fall off” from the left most end are lost



Blank spaces in the right positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100\ 0000$

$212 \ll 6 = 0000\ 0000\ 0000\ 0000\ 0011\ 0101\ 0000\ 0000$

$212 \ll 28 = 0100\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000$

Left shift by k is equivalent to integer multiplication with 2^k

Left Shift Op

However, be warned. If you left shift too much, multiplication with 2^k may give you a number too large to be represented by an int. You may get garbage results.

Left shift operator
number of locations

Bits that “fall off” from the left most end are lost

Blank spaces in the right positions are filled with 0s

$212 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

$212 \ll 0 = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100$

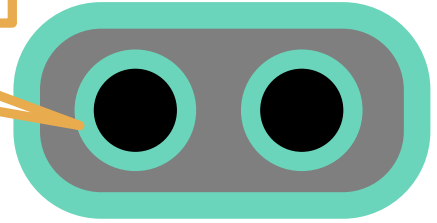
$212 \ll 4 = 0000\ 0000\ 0000\ 0000\ 0000\ 1101\ 0100\ 0000$

$212 \ll 6 = 0000\ 0000\ 0000\ 0000\ 0011\ 0101\ 0000\ 0000$

$212 \ll 28 = 0100\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000$

Left shift by k is equivalent to integer multiplication with 2^k

certain



Example use of bitwise operators 10



Example use of bitwise operators

10

Can use “masks” to extract certain bits of a number



Example use of bitwise operators 10

Can use “masks” to extract certain bits of a number

Suppose I want to look at the last 6 bits of a number a



Example use of bitwise operators 10

Can use “masks” to extract certain bits of a number

Suppose I want to look at the last 6 bits of a number a

Create a mask with only last bits set to 1 and take & with a



Example use of bitwise operators 10

Can use “masks” to extract certain bits of a number

Suppose I want to look at the last 6 bits of a number a

Create a mask with only last bits set to 1 and take & with a

```
int a = 427;  
int p = 1;  
int q = p << 6;  
int m = q - 1;  
int r = a & m;  
printf("%d", r); // 43
```



Example use of bitwise operators 10

Can use “masks” to extract certain bits of a number

Suppose I want to look at the last 6 bits of a number a

Create a mask with only last bits set to 1 and take & with a

```
a = 0000 0000 0000 0000 0000 0001 1010 1011
p = 0000 0000 0000 0000 0000 0000 0000 0001
q = 0000 0000 0000 0000 0000 0000 0100 0000
m = 0000 0000 0000 0000 0000 0000 0011 1111
r = 0000 0000 0000 0000 0000 0000 0010 1011
```

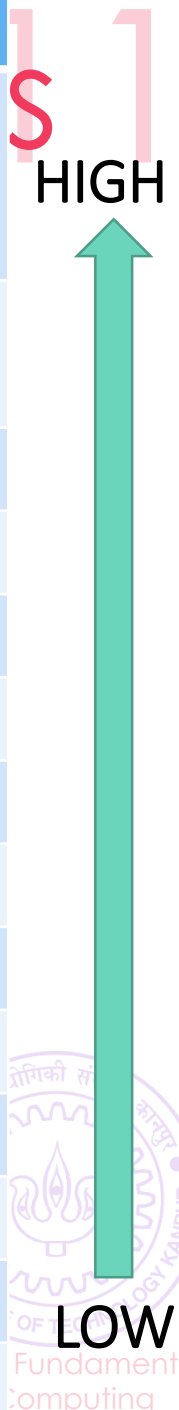
```
int a = 427;
int p = 1;
int q = p << 6;
int m = q - 1;
int r = a & m;
printf("%d", r); // 43
```



BODMAS table has more members



Operator Name	Symbol/Sign	Associativity
Brackets (array subscript), Post increment/decrement, structure field selection (dot and arrow)	(), [] ++, --, ., ->	Left
Unary negation, Pre-increment/decrement, NOT, (de)reference, sizeof, bitwise complement	-, ++, --, !, *, &, sizeof, ~	Right
Multiplication/division/remainder	*, /, %	Left
Addition/subtraction	+, -	Left
Bitwise left shift and right shift	<<, >>	
Relational	<, <=, >, >=	Left
Relational	==, !=	Left
Bitwise AND	&	Left
Bitwise XOR	^	Left
Bitwise OR		Left
Logical AND	&&	Left
Logical OR		Left
Ternary Conditional	? :	Right
Assignment, Compound assignment	=, +=, -=, *=, /=, % =, &=, ^=, =, <<=, >>=	Right

Operator Name	Symbol/Sign	Associativity	
Brackets (array subscript), Post increment/decrement, structure field selection (dot and arrow)	(), [] ++, --, ., ->	Left	
Unary negation, Pre-increment/decrement, NOT, (de)reference, sizeof, bitwise complement	~, ++, --, !, *, &, sizeof, ~	Right	
Multiplication/division/remainder	*, /, %	Left	
Addition/subtraction	+, -	Left	
Bitwise left shift and right shift	<<, >>		
Relational	<, <=, >, >=	Left	
Relational	==, !=	Left	
Bitwise AND	&	Left	
Bitwise XOR	^	Left	
Bitwise OR		Left	
Logical AND	&&	Left	
Logical OR		Left	
Ternary Conditional	? :	Right	
Assignment, Compound assignment	=, +=, -=, *=, /=, % =, &=, ^=, =, <<=, >>=	Right	

One's Compliment

12



One's Complement

12

The one's complement of a binary number) is simply the bitwise complement of that binary number



One's Complement

12

The one's complement of a binary number) is simply the bitwise complement of that binary number

A long time ago (25-30 years ago) one's complement used to be used to represent negative numbers



One's Complement

12

The one's complement of a binary number) is simply the bitwise complement of that binary number

A long time ago (25-30 years ago) one's complement used to be used to represent negative numbers

35 as a 4 byte int is represented as

0000 0000 0000 0000 0000 0000 0010 0011



One's Complement

12

The one's complement of a binary number) is simply the bitwise complement of that binary number

A long time ago (25-30 years ago) one's complement used to be used to represent negative numbers

35 as a 4 byte int is represented as

0000 0000 0000 0000 0000 0000 0010 0011

So, in those old computers, -35 used to be represented as

1111 1111 1111 1111 1111 1111 1101 1100



One's Complement

12

The one's complement of a binary number is simply the bitwise complement of that binary number

A long time ago (25-30 years ago) one's complement used to be used to represent negative numbers

35 as a 4 byte int is represented as

0000 0000 0000 0000 0000 0000 0010 0011

So, in those old computers, -35 used to be represented as

1111 1111 1111 1111 1111 1111 1101 1100

Note that $b + \sim b = 11111111 \ 11111111 \ 11111111 \ 11111111$



One's Complement

12

The bit If we have n bits, then using one's complement, we can represent numbers between $-(2^{n-1} - 1)$ and $+(2^{n-1} - 1)$

A long time ago (25-30 years ago) one's complement used to be used to represent negative numbers

35 as a 4 byte int is represented as

0000 0000 0000 0000 0000 0000 0010 0011

So, in those old computers, -35 used to be represented as

1111 1111 1111 1111 1111 1111 1101 1100

Note that $b + \sim b = 11111111 \ 11111111 \ 11111111 \ 11111111$



One's Complement

12

The bit

If we have n bits, then using one's complement, we can represent numbers between $-(2^{n-1} - 1)$ and $+(2^{n-1} - 1)$

The first bit acts as a sign bit – if the first bit is 1, it is treated as a negative number, if the first bit is 0, it is treated as a positive number

0000 0000 0000 0000 0000 0000 0010 0011

So, in those old computers, -35 used to be represented as
1111 1111 1111 1111 1111 1111 1101 1100

Note that $b + \sim b = 11111111 11111111 11111111 11111111$



One's Complement

12

The bit

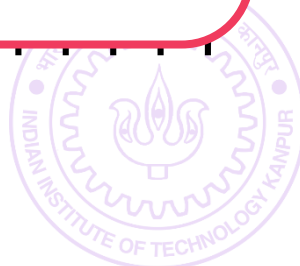
If we have n bits, then using one's complement, we can represent numbers between $-(2^{n-1} - 1)$ and $+(2^{n-1} - 1)$

The first bit acts as a sign bit – if the first bit is 1, it is treated as a negative number, if the first bit is 0, it is treated as a positive number

Largest positive number is 01111111 11111111 11111111 11111111
Smallest negative number is 10000000 00000000 00000000 00000000

Weird thing – negative 0 😊

11111111 11111111 11111111 11111111



One's Complement

12

The first bit
If we have n bits, then using one's complement, we can represent numbers between $-(2^{n-1} - 1)$ and $+(2^{n-1} - 1)$

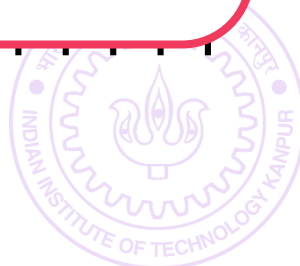
The first bit acts as a sign bit – if the first bit is 1, it is treated as a negative number, if the first bit is 0, it is treated as a positive number

Largest positive number is 01111111 11111111 11111111 11111111
Smallest negative number is 10000000 00000000 00000000 00000000

Weird thing – negative 0 😊

11111111 11111111 11111111 11111111

Used no more. These days, computers use two's complement to represent negative numbers



Two's Complement

86



Two's Complement

86

Advantages over one's complement



Two's Complement

86

Advantages over one's complement

Arithmetic circuits inside microprocessors do not need to look at sign bits



Two's Complement

86

Advantages over one's complement

Arithmetic circuits inside microprocessors do not need to look at sign bits
There is no weirdness about positive zero and negative zero



Two's Complement

86

Advantages over one's complement

Arithmetic circuits inside microprocessors do not need to look at sign bits

There is no weirdness about positive zero and negative zero

Two's complement of an n -bit binary number is the number which when added to this number, gives 2^n



Two's Complement

86

Advantages over one's complement

Arithmetic circuits inside microprocessors do not need to look at sign bits

There is no weirdness about positive zero and negative zero

Two's complement of an n -bit binary number is the number which when added to this number, gives 2^n

This means two's complement of b is $2^n - b$



Two's Complement

86

Advantages over one's complement

Arithmetic circuits inside microprocessors do not need to look at sign bits

There is no weirdness about positive zero and negative zero

Two's complement of an n -bit binary number is the number which when added to this number, gives 2^n

This means two's complement of b is $2^n - b$

Recall that $b + \sim b = \text{all ones} = 2^n - 1$ i.e. two's complement of b is $2^n - b = \sim b + 1$



Two's Complement

86

Advantages over one's complement

Arithmetic circuits inside microprocessors do not need to look at sign bits

There is no weirdness about positive zero and negative zero

Two's complement of an n -bit binary number is the number which when added to this number, gives 2^n

This means two's complement of b is $2^n - b$

Recall that $b + \sim b = \text{all ones} = 2^n - 1$ i.e. two's complement of b is $2^n - b = \sim b + 1$

Easier way of calculating two's complement – take the one's complement and add 1 to the binary string 😊



Two's Complement

86

Advantages over one's complement

Arithmetic circuits inside microprocessors do not need to look at sign bits

There is no weirdness about positive zero and negative zero

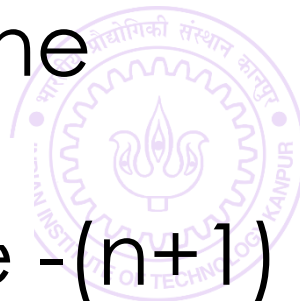
Two's complement of an n -bit binary number is the number which when added to this number, gives 2^n

This means two's complement of b is $2^n - b$

Recall that $b + \sim b = \text{all ones} = 2^n - 1$ i.e. two's complement of b is $2^n - b = \sim b + 1$

Easier way of calculating two's complement – take the one's complement and add 1 to the binary string ☺

So for any integer n , bitwise complement of n will be $-(n+1)$



Two's Complement

86

These days, the two's complement of a number represents the negative of that number e.g. -9 is represented as two's complement of 9

there is no weirdness about positive zero and negative zero

Two's complement of an n -bit binary number is the number which when added to this number, gives 2^n

This means two's complement of b is $2^n - b$

Recall that $b + \sim b = \text{all ones} = 2^n - 1$ i.e. two's complement of b is $2^n - b = \sim b + 1$

Easier way of calculating two's complement – take the one's complement and add 1 to the binary string ☺

So for any integer n , bitwise complement of n will be $-(n+1)$



Two's Complement

86

These days, the two's complement of a number represents the negative of that number e.g. -9 is represented as two's complement of 9

Tw
wh If we have n bits, then using one's complement, we can represent numbers between -2^{n-1} and $+(2^{n-1} - 1)$

This means two's complement of b is $2^n - b$

Recall that $b + \sim b = \text{all ones} = 2^n - 1$ i.e. two's complement of b is $2^n - b = \sim b + 1$

Easier way of calculating two's complement – take the one's complement and add 1 to the binary string ☺

So for any integer n , bitwise complement of n will be $-(n+1)$



Two's Complement

86

These days, the two's complement of a number represents the negative of that number e.g. -9 is represented as two's complement of 9

Tw
wh If we have n bits, then using one's complement, we can represent numbers between -2^{n-1} and $+(2^{n-1} - 1)$

The first bit acts as a sign bit – if the first bit is 1, it is treated as a negative number, if the first bit is 0, it is treated as a positive number

Easier way of calculating two's complement – take the one's complement and add 1 to the binary string ☺

So for any integer n , bitwise complement of n will be $-(n+1)$



Two's Complement

86

These days, the two's complement of a number represents the negative of that number e.g. -9 is represented as two's complement of 9

Tw
wh If we have n bits, then using one's complement, we can represent numbers between -2^{n-1} and $+(2^{n-1} - 1)$ number

The first bit acts as a sign bit – if the first bit is 1, it is treated as a negative number, if the first bit is 0, it is treated as a positive number

Largest positive number is 01111111 11111111 11111111 11111111
Smallest negative number is 10000000 00000000 00000000 00000000
11111111 11111111 11111111 11111111 now represents

Two's Complement

86

These days, the two's complement of a number represents the negative of that number e.g. -9 is represented as two's complement of 9

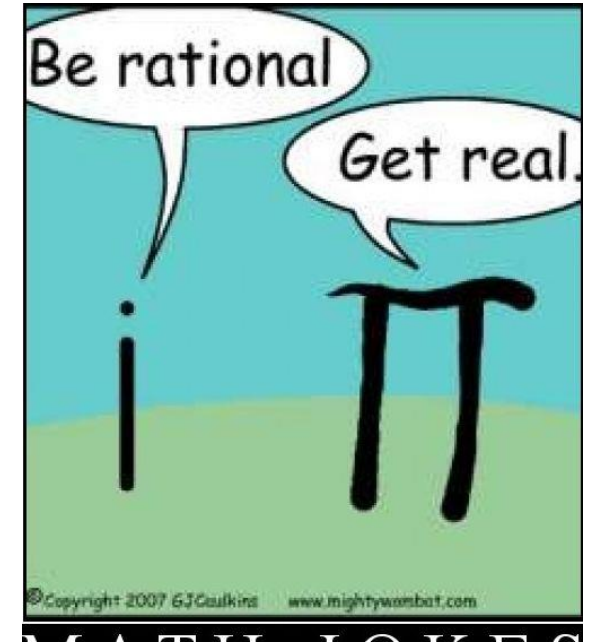
Tw
wh
If we have n bits, then using one's complement, we can represent numbers between -2^{n-1} and $+(2^{n-1} - 1)$

The first bit acts as a sign bit – if the first bit is 1, it is treated as a negative number, if the first bit is 0, it is treated as a positive number

Largest positive number is 01111111 11111111 11111111 11111111
Smallest negative number is 10000000 00000000 00000000 00000000
11111111 11111111 11111111 11111111 now represents -1 😊

Floating Point Representation

14



Floating Point Representation

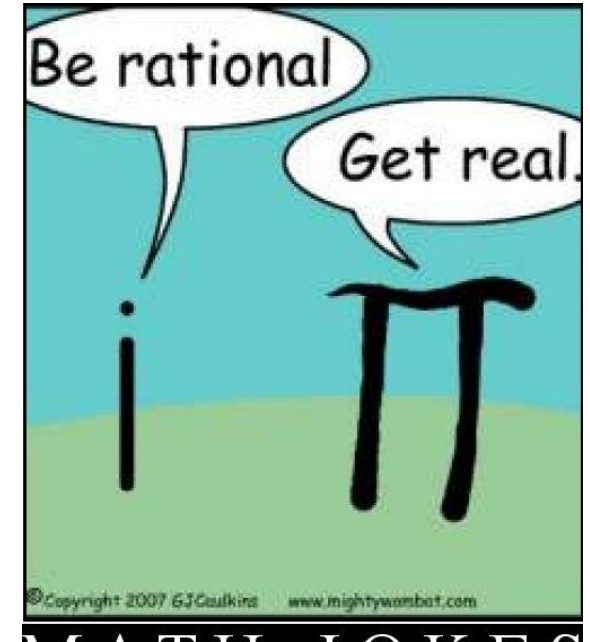
14

Have to represent three things

Sign: uses 1 bit

Mantissa: uses 23 bits

Exponent: uses 8 bits



Floating Point Representation

14

Have to represent three things

Sign: uses 1 bit

Mantissa: uses 23 bits

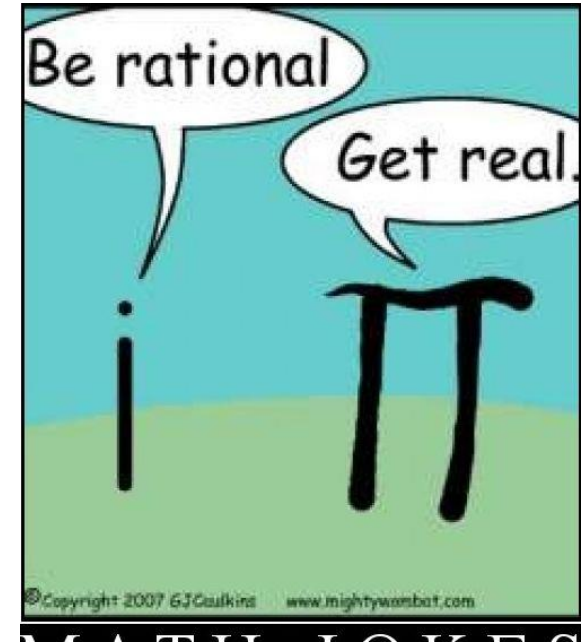
Exponent: uses 8 bits

Base 2: d_i 's are binary digits 0 or 1.

$$\pm d_0.d_1d_2\cdots d_{p-1} \times 2^e$$

Its value is

$$\pm (d_0 + d_1 2^{-1} + d_2 2^{-2} + \cdots + d_{p-1} 2^{-(p-1)}) 2^e$$



Floating Point Representation

14

Have to represent three things

Sign: uses 1 bit

Mantissa: uses 23 bits

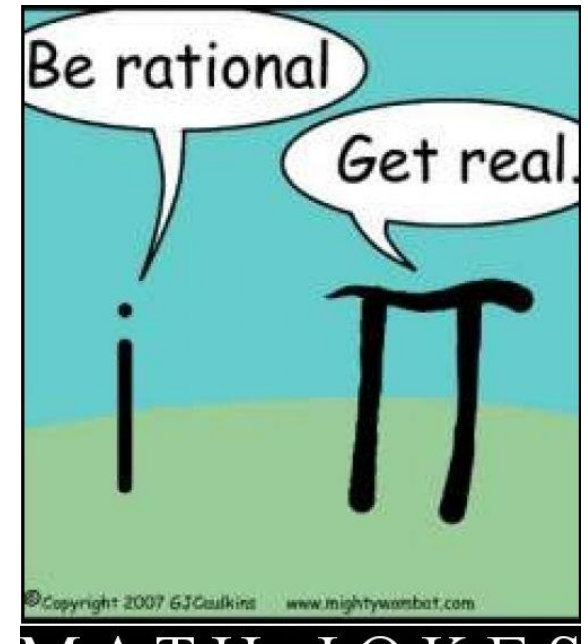
Exponent: uses 8 bits

Base 2: d_i 's are binary digits 0 or 1.

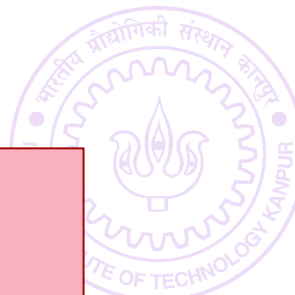
$$\pm d_0.d_1d_2 \cdots d_{p-1} \times 2^e$$

Its value is

$$\pm (d_0 + d_1 2^{-1} + d_2 2^{-2} + \cdots + d_{p-1} 2^{-(p-1)}) 2^e$$



1. p is called the **precision** of the representation.
2. Here base of the representation is 2. In general, the base could be different from 2. e.g., scientific calculators typically use base 10.



IEEE 754 Floating Point Representation

IEEE 754 Floating Point Standard



$$\text{number} = (-1)^s * (1.m) * 2^{e-127}$$



Range of Representation

- Base 2: d_i 's are binary digits 0 or 1.

$$\pm d_0.d_1d_2 \cdots d_{p-1} \times 2^e$$

- E.g., for number $(0.1)_{10}$ (i.e., 0.1 in decimal), $p = 24$, the representation is

$$(1.10011001100110011001100) \times 2^{-4}$$

$$1/16 + 1/32 + 1/256 + 1/512 = 51/512, \dots$$

$(0.1)_{10}$ is not “exactly” represented in binary.

Real numbers cannot be exactly represented with finite precision and limited exponent.



Single-precision (float)

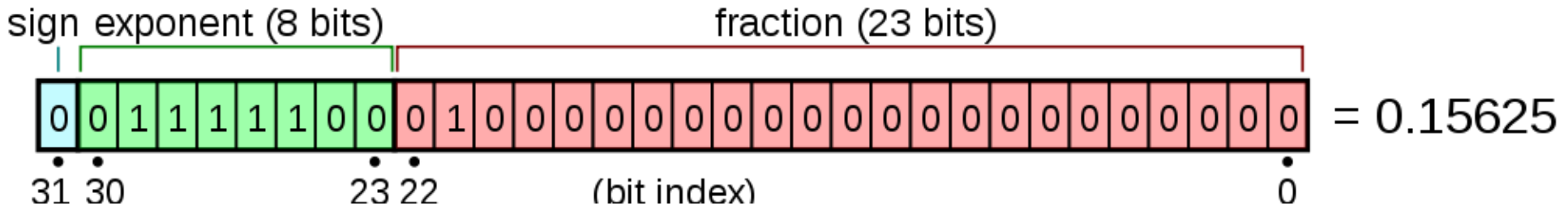
0	0110 1000	101 0101 0100 0011 0100 0010
---	-----------	------------------------------

- Sign: 0 => positive
- Exponent:
 - 0110 1000_{two} = 104_{ten}
 - Bias adjustment: 104 - 127 = -23
- Significand:
 - $1 + 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} + 0 \times 2^{-4} + 1 \times 2^{-5} + \dots$
 - $= 1 + 2^{-1} + 2^{-3} + 2^{-5} + 2^{-7} + 2^{-9} + 2^{-14} + 2^{-15} + 2^{-17} + 2^{-22}$
 - $= 1.0 + 0.666115$
- Represents: $1.666115 \times 2^{-23} \sim 1.986 \times 10^{-7}$

This is what you're using when you are invoking *float*



Imprecise Representation



Next number: $(1 + \frac{1}{4} + 2^{-23})2^{124-127} = 0.15625 + 2^{-26}$

0 0111100 0100000000000000000000001

Previous number: $(1 + 1/8 + 1/16 + \dots 2^{-23})2^{-3}$
 $= (1 + 1/4 - 2^{-23})2^{-3} = 0.15625 - 2^{-26}$

Given number: 0.15625

0 0111100 001111111111111111111111

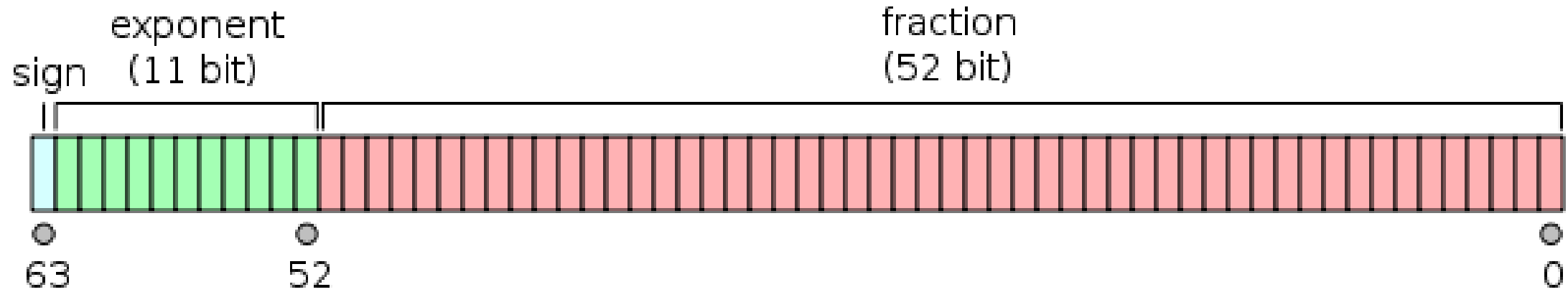
Real numbers in-between 0.15625 and $0.15625 + 2^{-26}$ may be all represented as 0.15625.

Discrete representation: leads to approximation errors.



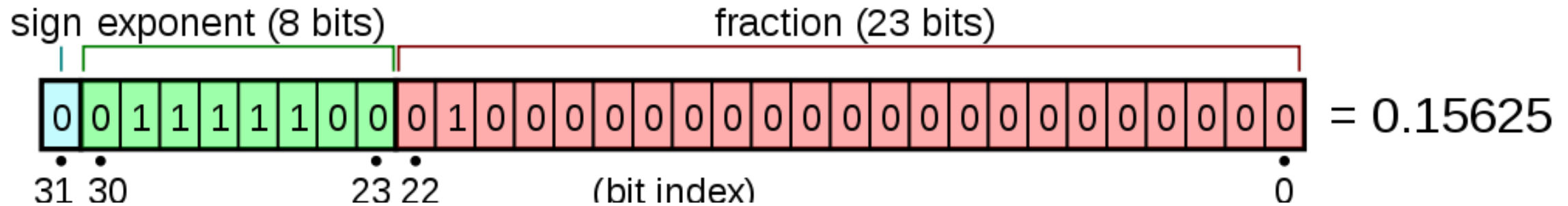
Double Precision (double)

Same logic as single precision, but with 64 bits of memory



Underflow and Overflow

Consider IEEE single precision.



Numbers smaller than 2^{-128} are indistinguishable from 0.

Such numbers occurring in calculations are said to be in **underflow** and are **treated as 0**.

Numbers that are larger than 2^{128} cannot be represented.

Such numbers occurring in calculations are said to have **overflowed**.

Warning: IEEE representation allows for +0.00 and -0.00. Be careful.

