

# 15-110 Principles of Computing – S19

LECTURE 4:

BINARY REPRESENTATIONS, STRING DATA TYPES

**TEACHER:** 

GIANNI A. DI CARO



### Scalar vs. Non-scalar objects

**Indivisible** 

- Scalar type objects:
  - int
  - float
  - complex
  - bool
  - None

- Non-Scalar type literal objects:
  - str: String of characters (text):
    - "Hi", 'Hello!', "Number 5"
  - tuple
  - list
  - set
  - dict

#### Internal structure

- Made of multiple components
- Individual or subsets of components can be addressed for read/write operations

- Scalar vs. Non-scalar terminology, from math
  - ✓ It is termed a scalar any *real number*, or *any quantity* that can be measured using a **single real number**
  - ✓ A vector is made of multiple scalar components (represents a point in a multi-dimensional space)

### Basic examples of using string objects: single and double quotes

str1 = "This is" "Hi" str2 = "spam!" 'Hello!' print( str1, str2) This is spam! "Number 5" Double quotes "abc" name = " " "Hello!" introduction = "My name is" name = "Gianni" 'Z', print( introduction, name) My name is Gianni • 'abc' Single quotes '\_wow\_' "l'm Joe" 'Say "hello!" to her' Double and single quotes together

# Basic examples of using string objects: triple quotes

- long\_str = "'Hi this is a veeeeeeery long string of text that I would like to write over multiple lines "'
- long\_str = """Hi this is a veeeeeeery long string of text that I would like to write over multiple lines """

Triple (single or double) quotes

long\_str = "'Hi this is a veeeeeeery long string of text
that I would like to write over multiple lines "'
print(long\_str)

Hi this is a veeeeeeery long string of text that I would like to write over multiple lines

# String objects

- A string is a *sequence* of <u>characters</u>
  - ✓ Sequence → Ordering, indexing
  - ✓ Characters → Which type of characters are allowed? → Unicode set
  - Sequence:

"Hello Joe"

Н	е	1	- 1	0		J	O	е
0	1	2	3	4	5	6	7	8

Indexing of the positions of the individual characters in the string

→ Access to the individual components of the string type

- Characters: A character is a <u>symbol</u>
  - E.g., the English alphabet has 26 symbols, other alphabets have different sets of symbols, plus we need characters for punctuation, characters for mathematics, characters for ...
  - Computers do not deal with characters, they deal with numbers (binary). Every character is internally stored and manipulated as a combination of 0's and 1's
  - Encoding: Character → Integer number → Binary representation → Python uses Unicode encoding

### Numeric conversions between different bases

- Let's consider an **integer number** x with n = 5 digits, e.g., x = 64523
- This is a base 10 (b = 10) representation of the number, using digits from 0 to 9

$$x = 6 \cdot 10^4 + 4 \cdot 10^3 + 5 \cdot 10^2 + 2 \cdot 10^1 + 3 \cdot 10^0 = 64{,}523$$

Position	4	3	2	1	0
Exponent	$10^{4}$	$10^{3}$	$10^{2}$	$10^{1}$	$10^{0}$
Value	10,000	1,000	100	10	1
Digits	$x_4$	$x_3$	$x_2$	$x_1$	$x_0$

Position	4	3	2	1	0
Exponent	24	$2^3$	2 <sup>2</sup>	$2^1$	$2^0$
Value	16	8	4	2	1
Digits	$x_4$	$x_3$	$x_2$	$x_1$	$x_0$

- Let's consider now a **binary number** x with n=5 digits, e.g., x=11001
- This is a base 2 (b = 2) representation of the number, using digits 0 and 1
- What is the integer value of the number x?

$$x = 1 \cdot 2^4 + 1 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 = 25$$

- How many *unsigned* integer numbers can be represented with 8 bits?  $\rightarrow$  256
- How many *signed* integer numbers can be represented with 8 bits?  $\rightarrow$  128
- Internal non-scalar representation of numbers

Binary	Octal	Decimal	Hexadecimal
0000	0	0	0
0001	1	1	1
0010	2	2	2
0011	3	3	3
0100	4	4	4
0101	5	5	5
0110	6	6	6
0111	7	7	7
1000	10	8	8
1001	11	9	9
1010	12	10	A
1011	13	11	В
1100	14	12	C
1101	15	13	D
1110	16	14	E
1111	17	15	F
Base-2	Base-8	Base-10	Base-16

#### **Bits**

- > One bit (that can take on two values, 0 or 1)
  - We can represent 2 integer numbers: 0 1
  - The max value of an integer that we can represent with 1 bit: 1
- > Two bits
  - We can represent 4 integer numbers: 00 01 10 11, from 0 to 3
  - The max value of an integer that we can represent with 2 bits: 3 (obtained from  $2^2 1$ )
- > Three bits
  - We can represent 8 integer numbers: 000 010 100 110 011 101 001 111, from 0 to 7
  - The max value of an integer that we can represent with 2 bits: 7 (obtained from  $2^3 1$ )

• • • • • • •

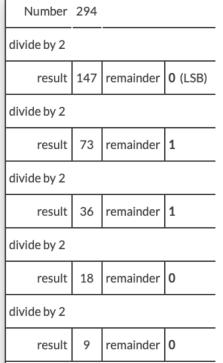
- > 8 bits = 1 byte
  - We can represent 256 (unsigned) integer numbers: from 0 to 255
  - The max value of an integer that we can represent with 8 bits: 255 (obtained from  $2^8 1$ )

#### Numeric conversions between different bases

Position	7	6	5	4	3	2	1	0
Exponent	2 <sup>7</sup>	2 <sup>6</sup>	2 <sup>5</sup>	2 <sup>4</sup>	$2^3$	2 <sup>2</sup>	2 <sup>1</sup>	2 <sup>0</sup>
Value	128	64	32	16	8	4	2	1
Digits	<i>x</i> <sub>7</sub>	$x_6$	$x_5$	$x_4$	$x_3$	$x_2$	$x_1$	$x_0$

MSB LSB

- From base 10 to base 2?
- Keep dividing by 2 and storing the remainder
- → Modulo operation!!!



Dividing each decimal number by "2" as shown will give a result plus a remainder.

If the decimal number being divided is even then the result will be whole and the remainder will be equal to "0". If the decimal number is odd then the result will not divide completely and the remainder will be a "1".

divide by 2						
result	4	remainder	1			
divide by 2						
result	2	remainder	0			
divide by 2						
result	1	remainder	0			
divide by 2						
result	0	remainder	1 (MSB)			

The binary result is obtained by placing all the remainders in order with the least significant bit (LSB) being at the top and the most significant bit (MSB) being at the bottom.

$$(294)_{10} = (100100110)_2$$

### **ASCII** encoding

- Encoding: Character → Integer number → Binary representation
- ASCII (American Standard Code for Information Interchange) standard code, defined in 1968 (and extended later on), assigns a numeric code (that can be hold in 8 bits = 1 byte) to a subset of standard characters
- 1 byte: basic unit of storage in computer memory!

Decimal	Hexadecimal	Binary	0ctal	Char	Decimal	Hexadecimal	Binary	0ctal	Char	Decimal	Hexadecimal	Binary	0ctal	Char
0	0	0	0	[NULL]	48	30	110000	60	0	96	60	1100000	140	`
1	1	1	1	[START OF HEADING]	49	31	110001	61	1	97	61	1100001	141	а
2	2	10	2	[START OF TEXT]	50	32	110010		2	98	62	1100010		b
3	3	11	3	[END OF TEXT]	51	33	110011		3	99	63	1100011		C
4	4	100	4	[END OF TRANSMISSION]	52	34	110100		4	100	64	1100100		d
5	5	101	5	[ENQUIRY]	53	35	110101		5	101	65	1100101		e
6	6	110	6	[ACKNOWLEDGE]	54	36	110110		6	102	66	1100110		f
7	7	111	7	[BELL]	55	37	110111		7	103	67	1100111		g
8	8	1000	10	[BACKSPACE]	56	38	111000		8	104	68	1101000		h
9	9	1001	11	[HORIZONTAL TAB]	57	39	111001		9	105	69	1101001		i
10	Α	1010	12	[LINE FEED]	58	3A	111010		:	106	6A	1101010		i
11	В	1011	13	[VERTICAL TAB]	59	3B	111011		;	107	6B	1101011		k
12	С	1100	14	[FORM FEED]	60	3C	111100		<	108	6C	1101100		ï
13	D	1101	15	[CARRIAGE RETURN]	61	3D	111101		=	109	6D	1101101		m
14	E	1110	16	ISHIFT OUT1	62	3E	111110		>	110	6E	1101110		n
15	F	1111	17	[SHIFT IN]	63	3F	111111		?	111	6F	1101111		0
16	10	10000	20	IDATA LINK ESCAPEI	64	40	1000000		@	112	70	1110000		р
17	11		21	IDEVICE CONTROL 11	65	41	1000001		Ă	113	71	1110001		q
18	12	10010	22	IDEVICE CONTROL 21	66	42	1000010		В	114	72	1110010		r
19	13	10011		IDEVICE CONTROL 31	67	43	1000011		C	115	73	1110011		S
20	14	10100	24	[DEVICE CONTROL 4]	68	44	1000100		D	116	74	1110100		t
21	15	10101		[NEGATIVE ACKNOWLEDGE]	69	45	1000101		E	117	75	1110101		u
22	16	10110	26	[SYNCHRONOUS IDLE]	70	46	1000110		F	118	76	1110110		v
23	17	10111		[ENG OF TRANS. BLOCK]	71	47	1000111		G	119	77	1110111		w
24	18	11000	30	[CANCEL]	72	48	1001000		H	120	78	1111000		x
25	19		31	[END OF MEDIUM]	73	49	1001001		ï	121	79	1111001		у
26	1A	11010	32	[SUBSTITUTE]	74	4A	1001010		i i	122	7A	1111010		z
27	1B		33	[ESCAPE]	75	4B	1001011		ĸ	123	7B	1111011		{
28	1C	11100	34	[FILE SEPARATOR]	76	4C	1001100		Î.	124	7C	1111100		i .
29	1D		35	[GROUP SEPARATOR]	77	4D	1001101		M	125	7D	1111101		}
30	1E	11110	36	IRECORD SEPARATORI	78	4E	1001110		N	126	7E	1111110		~
31	1F	11111		[UNIT SEPARATOR]	79	4F	1001111		0	127	7F	1111111		[DEL]
32	20	100000		[SPACE]	80	50	1010000		Р					
33	21	100001		1	81	51	1010001		Q					
34	22	100010		ii	82	52	1010010		R					
35	23	100011		#	83	53	1010011		S					
36	24	100100		\$	84	54	1010100		Т					
37	25	100101		%	85	55	1010101		Ū					
38	26	100110		&	86	56	1010110		V					
39	27	100111		T.	87	57	1010111		w					
40	28	101000		(	88	58	1011000		X					
41	29	101001		j	89	59	1011001		Y					
42	2A	101010		*	90	5A	1011010		ż					
43	2B	101011		+	91	5B	1011011		ī					
44	2C	101100			92	5C	1011100		Ň					
45	2D	101101		1	93	5D	1011101		i					
46	2E	101110			94	5E	1011110		^					
47	2F	101111		1	95	5F	1011111							
				*					-					

128 8 129 8 130 8 131 8 132 8 133 8 134 8 135 8 136 8 137 8	80h 81h 82h 83h 84h 85h 86h 87h 88h	Simbolo Ç ü é â ä à	160 161 162 163 164 165	A0h A1h A2h A3h A4h	Simbolo á í ó	192 193	HEX C0h C1h	Simbolo L L	224	E0h	Ó
129 8 130 8 131 8 132 8 133 8 134 8 135 8 136 8 137 8	81h 82h 83h 84h 85h 86h 87h	ü é â ä à	161 162 163 164	A1h A2h A3h A4h	í ó	193					_
130 8 131 8 132 8 133 8 134 8 135 8 136 8 137 8	82h 83h 84h 85h 86h 87h	é â ä à	162 163 164	A2h A3h A4h	ó		C1h				
131 8 132 8 133 8 134 8 135 8 136 8 137 8	83h 84h 85h 86h 87h	â ä à	163 164	A3h A4h	_			-	225	E1h	ß Ô Ò
132 8 133 8 134 8 135 8 136 8 137 8	84h 85h 86h 87h	ä à	164	A4h		194	C2h	Ţ	226	E2h	Ó
133 8 134 8 135 8 136 8 137 8	85h 86h 87h	à			ú	195	C3h	F	227 228	E3h E4h	
134 8 135 8 136 8 137 8	86h 87h		105	A5h	ñ Ñ	196 197	C4h C5h		228	E4n E5h	ő Ő
135 8 136 8 137 8	87h		166	A6h	1VI 8	198	C6h	+	230	E6h	
136 8 137 8			167	A7h	0	199	C7h	ä Ä	231	E7h	μ þ
137 8		ç ê	168	A8h	¿	200	C8h	Ë	232	E8h	
	89h	ë	169	A9h	®	201	C9h		233	E9h	Þ Ú Ú Ù
138 8	8Ah	è	170	AAh	7	202	CAh	1	234	EAh	Ď
139 8	BBh	ï	171	ABh	1/2	203	CBh		235	EBh	Ď
140 8	3Ch	î	172	ACh	1/4	204	CCh	Ī	236	ECh	Ý Ý
141 8	BDh	ì	173	ADh	i	205	CDh	=	237	EDh	Ŷ
142 8	BEh	Ä	174	AEh	«	206	CEh	#	238	EEh	_
	8Fh	Ą	175	AFh	<b>&gt;&gt;</b>	207	CFh	п	239	EFh	•
	90h	É	176	B0h	200	208	D0h	ð	240	F0h	
	91h	æ	177	B1h	-000	209	D1h	Ď	241	F1h	±
	92h	Æ	178	B2h		210	D2h	Ė	242	F2h	=
	93h	ô	179	B3h		211	D3h	Ë	243	F3h	3/4
	94h	Ò	180	B4h	-	212	D4h	_	244	F4h	¶
	95h	ò	181	B5h	Á Â	213	D5h	ļ	245	F5h	§
	96h 97h	û	182 183	B6h B7h	À	214 215	D6h D7h	Í Î	246 247	F6h F7h	÷
	9711 98h	ù	184	B8h	©	216	D8h	¦	241	F8h	3
	99h	ÿ Ö	185	B9h		217	D9h	j	249	F9h	
	9Ah	Ü	186	BAh	1	218	DAh		250	FAh	
	9Bh	ø	187	BBh		219	DBh	f	251	FBh	1
	9Ch	£	188	BCh	]	220	DCh	•	252	FCh	3
	9Dh	õ	189	BDh	¢	221	DDh	7	253	FDh	2
	9Eh	×	190	BEh	¥	222	DEh	Ì	254	FEh	
	9Fh	f	191	BFh	,	223	DFh	Ė	255	FFh	-

### Unicode encoding

- Encoding: Character → Integer number → Binary representation
- Developed in recent times to address the widespread use of computers in different countries using different symbols in their alphabet
- Different Unicode codes are around, using encoding larger (and more complex) than the 8 bits of ASCII,
   allowing to index code points (characters) large enough, to represent virtually any language around
- The most commonly used Unicode encoding is the UTF-8, that is fairly compact and includes ASCII codes
- Your Spider makes use of UTF-8!

# String indexing

Indexing:

"Hello Joe"

Н	е	- 1	- 1	0		J	0	е
0	1	2	3	4	5	6	7	8

Indexing of the positions of the individual characters in the string

→ Access to the individual components of the string type

- Index starts from 0 and must be an integer
- Notation to access the n-th component in a string variable my\_string: my\_string[n]

```
greet="Hello Joe"
print(greet[0], greet[4], greet[6])
```

■ We can use *variables* as index:

Н	е	- 1	1	0		J	O	е
-9	-8	-7	-6	-5	-4	-3	-2	-1

• We can also index from the right end of the string (useful to get the last character!) print(greet[x-4])

### String operators

String concatenation, + operator, overloaded: It returns a string consisting of the string operands joined together

```
greet_joe = "Hello Joe"
comma = ","
greet_mary = "hello Mary"
greet = greet_joe + comma + greet_mary
print(greet)
```

Hello Joe, hello Mary

Can I do greet + 1? NO!

String duplication, \* operator, overloaded: It creates multiple copies of a string. If s is a string and n is an integer:

HelloHelloHello

HelloHelloHello

Can I do s\*s? NO!

### String operators

 Part of, in operator, overloaded: Membership operator that returns True if the first operand is contained within the second, and False otherwise

```
s = "Joe"
in_hello = s in "Hello Joe"
in_food = s in "Yummy meal"
print(in_hello, in_food, type(in_hello))
```

True False <class 'bool'>

• **Not Part of, not** in operator, overloaded: Membership operator that returns True if the first operand is not contained within the second, and False otherwise

```
s = "Joe"
in_hello = s not in "Hello Joe"
in_food = s not in "Yummy meal"
print(in_hello, in_food, type(in_hello))
```

False True <class 'bool'>