

# CSE 1710

## Lecture 16

### *Text, Strings*

#### **Goals/To do:**

Given a string and a character, derive the **frequency of the character within the string**

Given a string, a target character and a replacement character, **implement character substitution.**

Given a numeric value in string format, **parse into numeric type**

#### **Goals/To understand:**

- difference between `char`, `String`, and `StringBuffer`
- The non-primitive `String` masquerades as a primitive type
- Pattern-matching abstractions (regular expressions)
- The difference between raw and formatted text; how to separate content from presentation

# Unicode

- Unicode is a computing industry standard for the consistent encoding, representation and handling of text expressed in most of the world's writing systems.
- Java and other languages use Unicode
- a *unicode character*.
  - is a non-negative numeric value
  - has a corresponding character according to the Unicode character tables (as defined by the Unicode Consortium)

3

The letter J is found in column '004' and row 'A', which makes '004A'

This is a hexadecimal number, denoted \u004A

To convert a hexadecimal number to decimal:

$$d_3d_2d_1d_0 = d_3 \times 16^3 + d_2 \times 16^2 + d_1 \times 16^1 + d_0 \times 16^0$$

where  $d_i$  takes on values [0, 15]  
 the value 10 is denoted by A (or a)  
 the value 11 is denoted by B  
 the value 12 is denoted by C  
 the value 13 is denoted by D  
 the value 14 is denoted by E  
 the value 15 is denoted by F

so to convert \u004A to decimal:

$$= 0 \times 16^3 + 0 \times 16^2 + 4 \times 16^1 + 10 \times 16^0$$

$$= 4 \times 16 + 10 \times 1$$

$$= 64 + 10$$

$$= 74$$

<http://unicode.org/charts/PDF/U0000.pdf>

	000	001	002	003	004	005	006	007
0	NUL	DLE	SP	0	@	P	`	p
1	SOH	DC1	!	1	A	Q	a	q
2	STX	DC2	"	2	B	R	b	r
3	ETX	DC3	#	3	C	S	c	s
4	EOT	DC4	\$	4	D	T	d	t
5	ENG	NAK	%	5	E	U	e	u
6	JACK	SYN	&	6	F	V	f	v
7	BEL	ETB	'	7	G	W	g	w
8	BS	CAN	(	8	H	X	h	x
9	HT	EM	)	9	I	Y	i	y
A	LF	SUB	*	:	J	Z	j	z
B	VT	ESC	+	;	K	[	k	{
C	FF	FS	,	<	L	\	l	
D	CR	GS	-	=	M	]	m	}
E	SO	RS	.	>	N	^	n	~
F	SI	US	/	?	O	_	o	DEL

4

## Unicode

- The Unicode Standard consists of a repertoire of more than 109,000 characters covering 93 scripts
  - Cyrillic, Latin, Bengali, Thai, Greek, ...
  - the basic set is “Controls and Basic Latin”
  - U000.pdf, also see Appendix A of JBA
- Unicode value denoted `\uXXXX`, where `XXXX` is a hexadecimal value
  - the decimal value 15 is represented as `\u000F`
- unicode makes is possible to talk about the *distance* between two character

5

## How Java uses Unicode

- `String` and `StringBuffer` objects encapsulate a string as a *sequence of unicode characters*
- the `char` primitive data type make use of unicode as well

L16App1

6

## About the character sequence...

- the `String` and `StringBuffer` classes encapsulate a string as a *sequence* of characters
- the sequence is indexed L16App2
  - the *first* position is index “0”
  - the *final* position is index “the length of the sequence minus 1”
  
- `String` services to tell us about the sequence
  - `int` : `length()`
  - `char` : `charAt(int)`
  - `String` : `substring(int, int)`
    - first value is start index, inclusive; second value is end index, exclusive
- what if index is out of bounds?
- what if end index is smaller than start index?

7

## What if the sequence has no characters at all?

- this is the *empty string*
- the string has length **zero**
- what if you hear the term “null string”?  
what does this mean?
  - not really a correct-formed phrase, there is no such thing
  - often used to mean a string reference that is set to null.

8

## How can we modify the sequence?

- Once a string object is created, it cannot be changed.
    - This is called *immutability*
    - Strings are *immutable*
  - Instead of modifying the sequence, we just create new strings.
  - It is fast and easy, thanks to the + operator
- L16App2b
- Given this, is it correct to say that String has mutators?
    - not technically; they are actually *generators of new modified objects*

9

## Iterating over a String

- many different ways to iterate
  - several different services to use...
- L16App3
- `theString.toCharArray()`
    - provides array that we can iterate over
    - `for (char c : theString.toCharArray()) { }`
  - `theString.charAt(index)`
    - we can iterate over index values
  - `theString.substring(startIndex, finishIndex)`
    - we can iterate over the starting index values

10

## String matching/comparison (basic)

- does *s1 match s2*?
- does *c1 match c2*?
  
- what does the equality boolean operator `==` tell us?
  - `boolean isMatch = c1==c2;`
  - `boolean isMatch = s1==s2;`
  
- what does `.equals(String)` tell us?
  - `boolean isMatch = s1.equals(s2);` L16App4
  
- what does `.compareTo(String)` tell us?
  - `int differingIndexPos = s1.compareTo(s2);`

11

## Difference between creation of `String` objects

- can construct `String` objects two different ways
  - implicit construction
    - use the pretend “literal” format
    - e.g.,  
`String s1 = "Hello";`
  
  - explicit construction
    - use keyword `new`
    - e.g.,  
`String s3 = new String("Hello");`
  
- explicit construction always creates new object
- if object with the same state already exists, implicit construction will re-use previously created object

12

## Elaboration of “compareTo (String)”

(sort of) “tell me whether the passed string comes before this string in the dictionary”

“aardvark”.compareTo(“anvil”)

- *anvil* does not come before *aardvark* in the dictionary, so the result is no (negative value)

“anvil”.compareTo(“aardvark”)

- *aardvark* **does** come before *anvil* in the dictionary, so the result is yes (positive value)

(better) “tell me whether the passed string comes before this string in the dictionary and, for the first character that is the determining factor, what is the distance”

- the second character is the determining factor (‘a’ vs ‘n’, there is a distance of 13 between them)

13

## Counting Character Frequency

- Given a string and a character, derive the **frequency of the character within the string**
- need to put together *iteration* and *comparision*
  1. iterate over the string to examine each character
  2. for each character, compare to target
  3. conditionally update counter

L16App5

L16App6