



## Character Encodings (3)

jonkv@ida  
10  
2006-09-12

- **Unicode:** Tens of thousands of characters!
  - Won't even fit into 16 bits...
  - Two **fixed-length encodings:**
    - UCS-2: One 16-bit word per character (does not support all chars!)
    - UCS-4: One 32-bit word per characters
  - Problem: Completely unsupported by most legacy software
    - ASCII: "Hello" coded as 48 65 6c 6c 6f
    - UCS-2: "Hello" coded as 00 48 00 65 00 6c 00 6c 00 6f
    - Contains "00" bytes which C programs interpret as "end of string"
  - Problem: Inefficient for English and English-like languages
    - Twice or four times as many bytes

## Character Encodings (4)

jonkv@ida  
11  
2006-09-12

- Three official **variable-length encodings**
  - UTF-8: Each character stored as 1 to 6 bytes
    - ▶ U-00000000 - U-0000007F: 0xxxxxxx (same as ASCII!)
    - ▶ U-00000080 - U-000007FF: 110xxxxx 10xxxxxx
    - ▶ U-00000800 - U-0000FFFF: 1110xxxx 10xxxxxx 10xxxxxx
    - ▶ U-00010000 - U-001FFFFF: 11110xxx 10xxxxxx 10xxxxxx 10xxxxxx
    - ▶ U-00200000 - U-03FFFFFF: 111110xx 10xxxxxx 10xxxxxx 10xxxxxx 10xxxxxx
    - ▶ U-04000000 - U-7FFFFFFF: 1111110x 10xxxxxx 10xxxxxx 10xxxxxx 10xxxxxx 10xxxxxx
  - UTF-16: Each character stored as one or two 16-bit words
  - UTF-32: Each character stored as one 32-bit word
- Which one to choose?
  - UTF-8 most widely supported
  - Efficient for English and similar languages: Chars fit into 1 byte
- Example for UTF-8:
  - Räksmörgås: 72 c3 a4 6b 73 6d c3 b6 72 67 c3 a5 73
  - In ISO Latin-1, this byte sequence would mean "rÄäksmärgÅs"

## Text in Java

### Characters Strings StringBuilders

2006-09-12

Jonas Kvarnström

12

## Characters 1: Unicode

jonkv@ida  
13  
2006-09-12

- Java uses **Unicode!**
  - Up to Java 1.3: Unicode 2.1, with 38,807 distinct symbols
  - From Java 1.4: Unicode 3.0, with 49,194 distinct symbols
  - From Java 1.5: Unicode 4.0, with 96,248 distinct symbols
- Internally, Java uses 16 bits
  - Java 1.3 and 1.4: Direct mapping (1 char = 1 code point)
  - Java 1.5: UTF-16 (may require two "chars" for one character)

More information: <http://www.unicode.org>

## Characters 2: Character Literals

jonkv@ida  
14  
2006-09-12

- **Character literals** are written using single quotes
  - char c = 'x', c2 = '7';
- But chars are represented as numbers...
  - char c = 120; // Makes c an 'x', since in Unicode, character #120 is an 'x'.
- **Escape codes** represent unprintable characters
  - '\b': Backspace
  - '\t': Tab
  - '\n': Newline
  - '\f': Form feed
  - '\r': Carriage return
  - '\"': Single quote (quote backslash quote quote)
  - '\\': A backslash
  - '\nnn': The character with octal value nnn ('\012' == '\n')
  - '\unnnn': The character with hex value nnnn ('\u03D0' == 'greek beta')

## Characters 3: Class java.lang.Character

jonkv@ida  
15  
2006-09-12

- **Static utility methods** in java.lang.Character
  - Character classes
    - isLowerCase(), isUpperCase(), isTitleCase(), isDigit(), isLetter(), isLetterOrDigit(), isWhitespace(), ...
    - Use these methods to accept greek lower case, kyrillic upper case, non-arabic digits, and so on
    - Character ch; ...; if (ch.isLetter()) { ... }
  - Case conversion
    - toLowerCase(), toUpperCase(), toTitleCase()
  - Conversion from and to decimal/hex/...
    - Character.digit('7', 10) = 7
    - Character.digit('b', 16) = 11
    - Character.forDigit(11, 16) = 'b'
    - Character.getNumericValue('7') = 7
    - Character.getNumericValue('\u216c') = 50 (roman numeral)

The letter 'lj':  
Lower case 'lj'.  
Upper case 'LJ'.  
Title case 'Lj'.

## Strings 1: Strings in Java

jonkv@ida  
16  
2006-09-12

- Two **main classes** representing strings
  - String: Constant strings
  - StringBuilder: Buffers used for building / concatenating strings
    - Temporarily – then you make a String from the StringBuilder
    - Used implicitly for the "+" operator on Strings
- Also other related classes
  - StringWriter
    - I/O class; what you write is appended to a buffer; retrieve as a String
  - CharBuffer
    - I/O class in the "new I/O" system
  - StringBuffer
    - Old version of StringBuilder

## Strings 2: String interfaces in Java

jonkv@ida  
17  
2006-09-12

- Two interfaces characterize **generic functionality**
  - **Appendable:** Anything you can append ("write") text to
    - StringBuilder, StringWriter, CharBuffer
    - PrintWriter (System.out), FileWriter, BufferedWriter,
    - Not String – it is constant
  - **CharSequence:** Anything you can get ("read") text from
    - (Random access required: char charAt(int index))
    - String, StringBuffer, CharBuffer
    - Not StringWriter

## Strings 2: Literals

jonkv@ida  
18  
2006-09-12

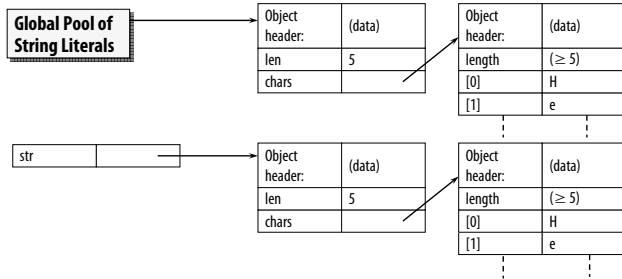
- String literals are written within **double quotes**
  - Escape codes can be used
    - "Hello, World"
    - "This is the first line\nThis is the second line, with a \" quote char"
- String literals are **String objects**
  - They are not char arrays, but can be created from char arrays
    - String str = {'a', 'b', 'c'}; // WRONG
    - char[] array = "abc"; // WRONG
    - char[] array = {'a', 'b', 'c'}; // OK
    - String str = new String(array); // OK
  - No need to assign them to variables; they **are** objects
    - String str = "Hello";
    - System.out.println(str.length()); // 5
    - System.out.println("Hello".length()); // 5

Call a method in  
the string literal!

## Strings 3: Literals and Constructors

jonkv@ida  
25  
2006-09-12

- No need to use a constructor for a literal
  - String str = new String("Hello"); // NOT necessary; inefficient!
  - String str = "Hello"; // Better
- The String literal "Hello" was pre-allocated
  - Calling the constructor creates a copy



## Strings 4: Useful methods

jonkv@ida  
20  
2006-09-12

- Checking the length of a String
  - "foo".length() == 3
- Retrieving part of a String
  - charAt(int pos) returns a specific character
    - "abcdefgabcdefg".charAt(1) == 'b'
  - substring(startindex, endindex) (endindex is exclusive)
  - trim() returns a new String without initial / trailing whitespace
  - split(String regexp) splits around regexp matches
    - "foo:bar::gazonk".split(":") == {"foo", "bar", "", "gazonk"}
    - "foo:bar::gazonk".split("o\*") == {"f", ":bar:gaz", "nk"}

## Strings 5: Useful methods

jonkv@ida  
21  
2006-09-12

- String conversions
  - toUpperCase(Locale), toLowerCase(Locale)
  - toUpperCase(), toLowerCase() (uses default locale)
- Finding matches
  - Finding a matching substring
    - "abcdefgabcdefg".indexOf("cde") == 2
    - "abcdefgabcdefg".indexOf("cde", 5) == 9
    - "abcdefgabcdefg".lastIndexOf("cde") == 9
  - startsWith(), endsWith():
    - "abcdefg".startsWith("abc") == true
  - Does the string match a regular expression? [more later]
    - "aaaabb".matches("a\*b\*") == true

## Strings 6: String Comparisons

jonkv@ida  
22  
2006-09-12

- Comparing Strings
  - equals(String other), equalsIgnoreCase(String other)
    - True if the strings are equal, false otherwise
    - As usual, == tests whether they are stored in the same location!
  - compareTo(String other), compareToIgnoreCase(String other)
    - -1 if this before other, 0 if equal, 1 if this after other
- Different countries have different sorting rules
  - Comparing á to a, sorting ß as ss, sorting ä as a-with-umlaut (German) or after å (Swedish) ...
  - Use [Collator](#), [RuleBasedCollator](#) for searching and sorting
    - Collator collator = Collator.getInstance(); // Locale-dependent
    - Collator collator2 = Collator.getInstance(Locale.CANADA);
    - collator.compare("abc", "ABC") // -1, 0, 1

## Formatting Text

### Converting primitive types to text

### Converting single objects to text

### String concatenation

### Advanced formatting

2006-09-12

Jonas Kvarnström

23

## Formatting 2: Objects to Strings

jonkv@ida  
25  
2006-09-12

- Every object has a **string representation**
  - java.lang.Object has a toString() method
    - This method returns a String representing the object
    - System.out.println(Object o) prints o.toString()
  - Example:
    - public class Circle {  
private double x, y, r;  
private Circle(double x, double y, double r) {  
this.x = x; this.y = y; this.r = r;  
}  
public static void main(String[] args) {  
System.out.println(new Circle(2.0, 4.0, 5.0));  
}  
}  
■ >> java Circle  
Circle@a5e17120

## Formatting 3: Objects to Strings

jonkv@ida  
26  
2006-09-12

- Why "Circle@a5e17120"?
  - Object.toString() returns the class name + "@" + a unique ID
  - The method should be overridden!
- Example:
  - public class Circle {  
...  
public String toString() {  
return "Circle[" + x + ", " + y + ", " + r + "];  
}  
public static void main(String[] args) {  
System.out.println(new Circle(2.0, 4.0, 5.0));  
}  
}  
■ java Circle  
Circle[2.0, 4.0, 5.0]

## Formatting 4: String Concatenation

jonkv@ida  
27  
2006-09-12

- In the previous slide: "Circle[" + x + ", " + y + ", " + r + "]"
    - Java overloads the + operator for string concatenation.
    - Internally, concatenation uses the StringBuilder class
    - The translation is done by Java compiler (javac, jikes, ...)
    - final StringBuilder buf = new StringBuilder("");  
buf.append("Circle[");  
buf.append(x); → String.valueOf(double)  
buf.append(", "); → String.valueOf(double)  
buf.append(y); → String.valueOf(double)  
buf.append(", "); → String.valueOf(double)  
buf.append(r); → String.valueOf(double)  
buf.append("]"); → String.valueOf(char)  
return buf.toString();
- The valueOf methods know how to convert primitive datatypes to Strings!

## Formatting 5: String Concatenation

jonkv@ida  
28  
2006-09-12

- What about constant strings?
  - String str = "hi " + "there"
  - This is evaluated at **compile time!**
    - ▶ Exactly the same bytecode as String str = "hi there"
- What about objects?
  - Circle c = **new** Circle(2.0, 4.0, 5.0);  
String str = "The circle is " + c; // Concatenate String + Object
  - This is translated, and becomes...
  - **final** StringBuilder buf = **new** StringBuilder("");  
buf.append("The circle is ");  
buf.append(c); // → String.valueOf(c) → c.toString()  
String str = buf.toString();

## Formatting 6: StringBuilder Example

jonkv@ida  
29  
2006-09-12

- Using **StringBuilder** in toString():
  - In loops, better to use StringBuilder explicitly
    - **public class** ArrayList {  
**private** Object[] elements;  
**private** int size;  
...  
**public** String toString() {  
**final** StringBuilder buf = **new** StringBuilder();  
buf.append("ArrayList<");  
**for** (int i = 0; i < size; i++) {  
**if** (i != 0) buf.append(", ");  
buf.append(elements[i].toString());  
}  
buf.append(">");  
**return** buf.toString();  
}  
}

Using String and += would be **extremely** inefficient:  
The string is copied twice each concatenation!

## Advanced Text Formatting

jonkv@ida  
30  
2006-09-12

- Use a **Formatter** for advanced text formatting
  - Step 1: Create a formatter (possibly specifying a locale)
    - // Send all output to the Appendable object sb  
StringBuilder sb = **new** StringBuilder();  
Formatter formatter = **new** Formatter(sb, Locale.US);
  - Step 2: Format text using the formatter (similar to printf in C)
    - Each call specifies a format string and any number of arguments
    - Each call appends to the given appendable
    - formatter.format("Hello, %s!\n", person.name);  
formatter.format("You are %d years old.\n", person.age);
    - // Minimum 30, maximum 40 characters (padded with spaces)  
formatter.format("Hello, %30.40s!\n", person.name);
    - // Here we specify an explicit locale...  
formatter2.format(Locale.FRANCE, "e = %+10.4f", Math.E);  
// -> "e = +2,7183"

## Advanced Text Formatting 2

jonkv@ida  
31  
2006-09-12

- Often no need to explicitly create a Formatter
  - Use utility methods instead
  - In PrintStream: printf()
    - System.out.printf("Hello, %s!\n", person.name);
  - In String: format()
    - Calendar c = **new** GregorianCalendar(1995, MAY, 23);  
String s = String.format("Duke's Birthday: %1\$tm %1\$te,%1\$ty", c);  
// -> s == "Duke's Birthday: May 23, 1995"

## Parsing Text

2006-09-12

Jonas Kvarnström

32

## Parsing 1: Strings to Primitive Datatypes

jonkv@ida  
33  
2006-09-12

- Convert **text to data** using a number of methods
  - Depending on data type
  - Every data type (int, float, ...) has a corresponding class
    - In class Integer: **static** int parseInt(String str);
    - In class Float: **static** float parseFloat(String str);
    - In class Double: **static** double parseDouble(String str);
    - ...
    - float val = Float.parseFloat("127"); // 127.0f

## Parsing 2: Tokenizing Formal Languages

jonkv@ida  
34  
2006-09-12

- Tokenizer: **java.util.StringTokenizer**
  - Breaks strings into words given break characters
    - **final** StringTokenizer st;  
st = **new** StringTokenizer("simple test\nLine 2", "\t\r\n");  
**while** (st.hasMoreTokens()) {  
println(st.nextToken());  
}  
}simple  
test  
Line  
2
  - Like String.split(), but more efficient and more restricted
    - String[] tokens = "simple test\nLine 2".split("(\\t|\\r|\\n)");
  - See also:
    - StreamTokenizer – to find tokens in Java-like languages
    - BreakIterator – for natural language text

Space, tab, CR, newline

## Specifying Character Encodings

Different software, different encodings...

How do we specify which encoding is used?

2006-09-12

Jonas Kvarnström

35

## Encodings 1: HTML

jonkv@ida  
36  
2006-09-12

- HTML: **Character Repertoires**
  - HTML 1.0 to 3.2: Only supports ISO Latin-1
  - HTML 4.0 and later: UCS (same chars as Unicode 2.0)
- HTML 4.0: Specifying **Character Encodings**
  - Encoding header in HTTP
    - Content-Type: text/html; charset=ISO-8859-1
    - Content-Type: text/html; charset=utf-8
    - Encoding header comes **before** the HTML file
    - Requires server configuration
  - Encoding meta tag in HTML file
    - <meta http-equiv="Content-Type" content="text/html; charset=iso-8859-1">
    - Up to this header, all characters interpreted as ASCII
  - Regardless of encoding: &#12345; is Unicode char 12345 (dec)

- **File I/O:** Two parallel subsystems in Java
  - Binary I/O: InputStream, OutputStream (pure 8-bit bytes)
    - FileInputStream, DataInputStream, BufferedInputStream, ...
  - Text I/O: Reader, Writer (converts bytes to/from chars)
    - FileReader, StringReader, BufferedReader, ...
    - By default: Uses platform-specific local character set
    - **final** OutputStream os = **new** FileOutputStream("file.txt");
    - **final** Writer wr = **new** OutputStreamWriter(os, "US-ASCII");
    - **final** OutputStream os = **new** FileOutputStream("file.txt");
    - **final** Writer wr = **new** OutputStreamWriter(os, "ISO-8859-1");
    - **final** Reader re = **new** InputStreamReader(is, "UTF-16");
- This is true for standard output too!
  - System.out converts chars to bytes

- **GUI character support** depends on platform/font
  - Regardless of method used
    - myTextField.setString("ä ö ö ö ë ð ð æ Γ ы Ω Α Δ Ж Т");
    - void paintComponent(Graphics g) {
 

```
g.drawString(10, 100, "ä ö ö ö ë ð ð æ Γ ы Ω Α Δ Ж Т");
```
  - Any Java environment can display ISO Latin-1
  - For more characters: May have to modify font.properties
    - Add new fonts that handle parts of the Unicode character set
    - For example:
 

```
sansserif.3=Lucida Sans Unicode,UNICODE_CHARSET
```

- **JDBC:** Depends on database and drivers
  - MySQL: Many character encodings ("sets") and collations
    - SHOW CHARACTER SET; SHOW COLLATION;
  - Default character encoding for DB, table, connection
    - (CREATE|ALTER) DATABASE *db\_name*

```
[[DEFAULT] CHARACTER SET charset_name
[[DEFAULT] COLLATE collation_name]
```
    - CREATE TABLE *tbl\_name* (*column\_list*)
 

```
[DEFAULT CHARACTER SET charset_name [COLLATE collation_name]]
```
    - ALTER TABLE *tbl\_name* [DEFAULT CHARACTER SET *charset\_name*]
 

```
[COLLATE collation_name]
```
    - *col\_name* {CHAR | VARCHAR | TEXT} (*col\_length*)
 

```
[CHARACTER SET charset_name [COLLATE collation_name]]
```
  - SELECT k FROM t1
 

```
ORDER BY k COLLATE latin1_german2_ci;
```

# Patterns and Regular Expressions

See also: [Java Tutorial chapter on regular expressions](#)

<http://java.sun.com/docs/books/tutorial/extra/regexp/>

## Glob Patterns

- Filename pattern matching (**globbing**):
  - "?" matches any single character
  - "\*" matches any sequence of characters
    - \*.java
  - Maybe a few more patterns – but not very powerful
- **Regular Expressions:** Much more powerful patterns
  - Directly supported in Java
  - Syntax as in Perl 5 (mostly)

## Regexp 1: Strings

- **Simplest regular expressions:** Ordinary strings
  - hydro
    - matches "hydro", "hydrogen", "hydrodynamics"
  - top ten
    - matches "top ten", "stop tension"
    - does not match "stop tension" (additional spaces)

## Regexp 2: Period

- **Metacharacters** have special meaning
  - Period (".") matches any single character
  - 2,-dimethylbutane
    - Matches "2,2-dimethylbutane" or "2,3-dimethylbutane"
    - Does not match "2,12-dimethylbutane"
  - 3.14
    - Matches "3.14" or "3x14" or "3714" or "3+14" or "192753514127"
  - 3\14
    - Matches "3.14" or "73.14159" but not "3x14" or "3714"
    - Metacharacter is escaped using a backslash (also metacharacter)
    - Note: Must write test.callMethod("3\\.14") -- two backslashes required in a string literal, so that one remains in the actual string!

## Regexp 3: Character classes

- [...] matches a **character class**
  - Examples
    - [abc] either a or b or c
    - [a-f] a single character between a and f
    - [a-fy-z] a, b, c, d, e, f, y, or z
    - [^abc] anything except a, b, or c
    - [-abc] hyphen or a (hyphen must be first – otherwise marks a range)
    - [fs]inger "finger" or "singer"
  - Some predefined character classes
    - \d digit [0-9]
    - \s whitespace [\t\n\x0B\f\r]
    - \w word character [a-zA-Z\_0-9]
    - \p{Punct} punctuation One of !"#%&'()\*+,-./:;<=>?@[^\_`{|}~
    - \p{XDigit} hex digit [0-9a-fA-F]
    - Pattern.compile("\\s\\w");// white space followed by word char

## Regexp 4: Alternatives

- x|y marks a set of **alternatives**
  - finger|singer
  - html|head|body|title|h[1-6]|p|hr|font|ul|li
    - Some HTML tag names
- May need to place within **parentheses**
  - <html|head|body|title|h[1-6]|p|hr|font|ul|li>
    - matches "<html", "head", ..., "<li>"
  - <(html|head|body|title|h[1-6]|p|hr|font|ul|li)>
    - matches "<html>", "<head>", ..., "<li>"

## Regexps 5: Question mark

jonkv@ida  
46  
2006-09-12

### ■ **Question mark:** Expression should occur 0 or 1 time

- m?ethane
  - matches "methane" or "ethane" or "asdfethane kasd"
- comm?a
  - matches "coma" or "comma"
- [fs]?inger
  - matches "finger", "singer", or "inger"
- (html|body|head)?
- NOTE: Different from "?" in filename matching
  - In filename matching, "?" matches one char – same as "." here

## Regexps 6: Star, Plus, Number of times

jonkv@ida  
47  
2006-09-12

### ■ **Star:** Expression should occur zero or more times

- ab\*c
  - matches "ac", "abc", "abbc"

### ■ **Plus:** Expression should occur one or more times

- ab+c
  - matches "abc", "abbc", "acabbc82", ... (but not "ac")

### ■ **{n,m}:** Expression should occur n to m times

- ab{2,4}c
  - matches "abbc", "abbbc", "abbbbc"
  - does not match "abc"
  - does match "abcabbbcababacab"

## Regexps 7: Start and end of text

jonkv@ida  
48  
2006-09-12

### ■ Regexps can **occur anywhere** in a string

- Use special metacharacters to match start/end
  - ^ matches start (of string or line, depending on settings)
  - \$ matches end (of string or line, depending on settings)
- m?ethane
  - matches "methane" or "ethane" or "asdfethane kasd"
- ^m?ethane
  - matches "methane" or "ethane" or "ethane kasd"
  - does not match "asdfethane kasd"
- ^m?ethane\$
  - matches "methane" or "ethane"
  - does not match anything else

## Regexps 8: Combining metacharacters

jonkv@ida  
49  
2006-09-12

### ■ **Combining metacharacters**

- Metacharacters can be combined
- Common combination: Period and star .\*
  - matches any number of occurrences of any character
  - same as a single "\*" in filename globbing
- cyclo.\*ane
  - cycloane
  - cyclodecane, cyclohexane
  - cyclones drive me insane, did you know that?
- [^>]\*
  - matches any sequence of chars not containing ">"
- <h1[^>]\*>
  - matches "<h1", followed by anything except ">", followed by ">"
  - ^<h1.\*>\$ would have matched "<h1>Hello</h1><p>Some text</p>"

## Regular Expressions in Java (1)

jonkv@ida  
50  
2006-09-12

### ■ In Java: First, **create and compile** a pattern:

- Compiling takes time – but the result can be reused
  - Pattern pat = Pattern.compile("^Hello.\*world\$");
  - For constant patterns: Store as a static constant!
- Compiled to bytecode!

- Creating classes on the fly
- Very efficient pattern matching

### ■ Then, use it to **match** one or more strings:

- Matcher mat = pat.matcher("foo");  
if (mat.matches()) { ... }
- Matcher mat2 = pat.matcher("Hello, wonderful world");  
if (mat2.matches()) { ... }

## Regular Expressions in Java (2)

jonkv@ida  
51  
2006-09-12

### ■ **Example:** Matching lines in a price list

- Something: 100.00  
Something else with a long name: 199.50
- Pattern:
  - One or more non-colon characters
  - Colon
  - Space
  - Any combination of digits/periods
- Some sub-patterns:
  - One non-colon char [^]
  - At least one non-colon char [^:]+
  - Any number of digits/periods [0-9.]+
- Complete pattern: [^:]+: [0-9.]+

## Regular Expressions in Java (3)

jonkv@ida  
52  
2006-09-12

### ■ Parentheses generate **capturing groups**

- What you match inside parentheses can be extracted
- Use this to parse the price list
  - Something: 100.00  
Something else with a long name: 199.50
- Old pattern: [^:]+: [0-9.]+
- Add parentheses for subgroups: ([^:]+): ([0-9.]+)
- Code:

```
final Pattern pat = Pattern.compile("([^:]+): ([0-9.]+)");
final Matcher mat = pat.matcher(myEntireTextTable);
while (mat.matches()) {
    System.out.println(mat.group(1)); // What was within the first ()
    System.out.println(mat.group(2)); // What was within the second ()
}
```

## Regular Expressions in Java (4)

jonkv@ida  
53  
2006-09-12

### ■ **Another example**

- Extract all references from an HTML page:
  - <a href="http://www.ida.liu.se">The IDA Home Page</a>
  - Match <a ...whatever... href="..." ...whatever...>
- Need to match:
  - <a [^>]\* Literally: left-bracket and "a"
  - [^>]\* Any character except end-of-tag ">", as many times as you need
  - href=" The href value...
  - ("^")\* A string of characters, not including quote
  - Closing quote
  - [^>]\*> A sequence of non-">" followed by a ">"
- final Pattern pat = Pattern.compile("<a[>]\*href='\"([^\"]\*)\"' [>]\*>");
- final Matcher mat = pat.matcher(myEntireHTMLDocument);
- while (mat.matches()) System.out.println(mat.group(1));

## Regular Expressions in Java (5)

jonkv@ida  
54  
2006-09-12

### ■ **Replacing subsequences**

- \$1 refers to first capturing group, \$2 to second, etc
- Swap two arguments of a macro:
  - \foo 100,200 → \foo 200,100
  - Pattern to match: \\foo [0-9.]+,[0-9.]+
  - With capturing groups: \\foo ([0-9.]+),([0-9.]+)
  - Replace with: \\foo \$2,\$1
- final Pattern pat = Pattern.compile("\\\\foo ([0-9.]+),([0-9.]+)");
- final Matcher mat = pat.matcher("...");
- final String replaced = mat.replaceAll("\\\\foo \$2,\$1");