# 3.1 Using Data Types (II)



INTRODUCTION TO
**Programming**
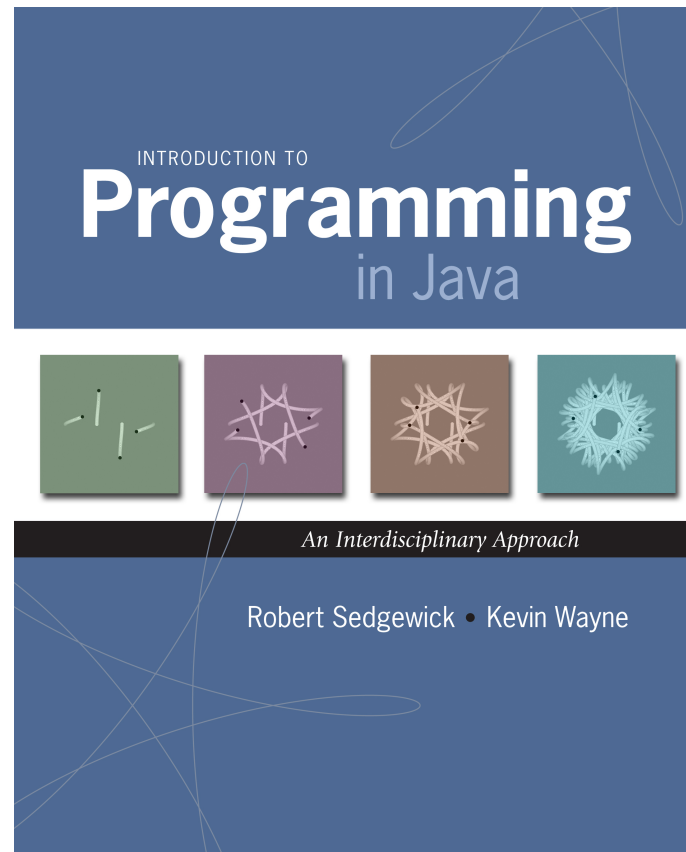in Java

*An Interdisciplinary Approach*

Robert Sedgewick • Kevin Wayne

# Text Processing

# String Data Type

**String data type.** Basis for text processing.
**Set of values.** Sequence of Unicode characters.

**API.**

### public class String (Java string data type)

| | | |
|---|---|---|
| | `String(String s)` | *create a string with the same value as* `s` |
| `int` | `length()` | *string length* |
| `char` | `charAt(int i)` | `i`*th character* |
| `String` | `substring(int i, int j)` | `i`*th through* `(j-1)`*st characters* |
| `boolean` | `contains(String sub)` | *does string contain* `sub` *as a substring?* |
| `boolean` | `startsWith(String pre)` | *does string start with* `pre`*?* |
| `boolean` | `endsWith(String post)` | *does string end with* `post`*?* |
| `int` | `indexOf(String p)` | *index of first occurrence of* `p` |
| `int` | `indexOf(String p, int i)` | *index of first occurrence of* `p` *after* `i` |
| `String` | `concat(String t)` | *this string with* `t` *appended* |
| `int` | `compareTo(String t)` | *string comparison* |
| `String` | `replaceAll(String a, String b)` | *result of changing* `a`*s to* `b`*s* |
| `String[]` | `split(String delim)` | *strings between occurrences of* `delim` |
| `boolean` | `equals(String t)` | *is this string's value the same as* `t`*'s?* |

http://download.oracle.com/javase/6/docs/api/java/lang/String.html

# Typical String Processing Code

| | |
|---|---|
| *is the string a palindrome?* | ```java
public static boolean isPalindrome(String s)
{
   int N = s.length();
   for (int i = 0; i < N/2; i++)
      if (s.charAt(i) != s.charAt(N-1-i))
         return false;
   return true;
}
``` |
| *extract file name and extension from a command-line argument* | ```java
String s = args[0];
int dot = s.indexOf(".");
String base      = s.substring(0, dot);
String extension = s.substring(dot + 1, s.length());
``` |
| *print all lines in standard input that contain a string specified on the command line* | ```java
String query = args[0];
while (!StdIn.isEmpty())
{
   String s = StdIn.readLine();
   if (s.contains(query)) StdOut.println(s);
}
``` |
| *print all the hyperlinks (to educational institutions) in the text file on standard input* | ```java
while (!StdIn.isEmpty())
{
   String s = StdIn.readString();
   if (s.startsWith("http://") && s.endsWith(".edu"))
      StdOut.println(s);
}
``` |

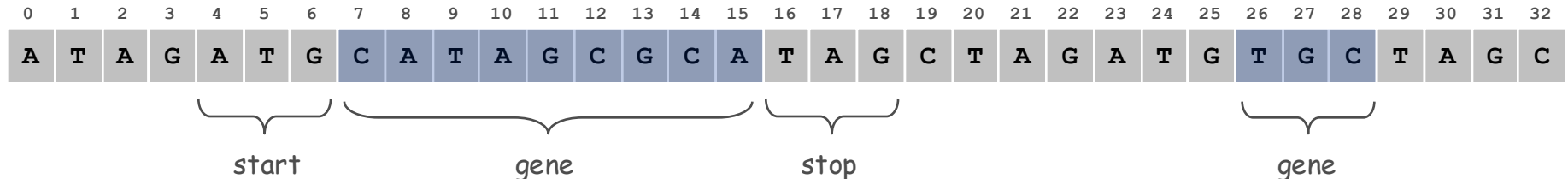# Gene Finding

Pre-genomics era.  Sequence a human genome.
Post-genomics era.  Analyze the data and understand structure.

Genomics.  Represent genome as a string over { A, C, T, G } alphabet.

Gene.  A substring of genome that represents a functional unit.
    Preceded by ATG.                                      [start codon]
    Multiple of 3 nucleotides.                       [codons other than start/stop]
    Succeeded by TAG, TAA, or TGA.          [stop codons]

Goal.  Find all genes.

# Gene Finding:  Algorithm

**Algorithm.**  Scan left-to-right through genome.

If start codon, then set `beg` to index `i`.

If stop codon and substring is a multiple of 3

  - output gene

  - reset `beg` to -1

| i | codon | | beg | gene | remaining portion of input string |
|---|---|---|---|---|---|
| | *start* | *stop* | | | |
| 0 | | | –1 | | ATAGATGCATAGCGCATAGCTAGATGTGCTAGC |
| 1 | | TAG | –1 | | ATAGATGCATAGCGCATAGCTAGATGTGCTAGC |
| 4 | ATG | | 4 | | ATAGATGCATAGCGCATAGCTAGATGTGCTAGC |
| 9 | | TAG | 4 | multiple of 3 | ATAGATGCATAGCGCATAGCTAGATGTGCTAGC |
| 16 | | TAG | 4 | CATAGCGCA | ATAGATGCATAGCGCATAGCTAGATGTGCTAGC |
| 20 | | TAG | –1 | | ATAGATGCATAGCGCATAGCTAGATGTGCTAGC |
| 23 | ATG | | 23 | | ATAGATGCATAGCGCATAGCTAGATGTGCTAGC |
| 29 | | TAG | 23 | TGC | ATAGATGCATAGCGCATAGCTAGATGTGCTAGC |

# Gene Finding: Implementation

```java
public class GeneFind {
    public static void main(String[] args) {
        String start  = args[0];
        String stop   = args[1];
        String genome = StdIn.readAll();

        int beg = -1;
        for (int i = 0; i < genome.length() - 2; i++) {
            String codon = genome.substring(i, i+3);
            if (codon.equals(start)) beg = i;
            if (codon.equals(stop) && beg != -1 && beg+3 < i) {
                String gene = genome.substring(beg+3, i);
                if (gene.length() % 3 == 0) {
                    StdOut.println(gene);
                    beg = -1;
                }
            }
        }
    }
}
```
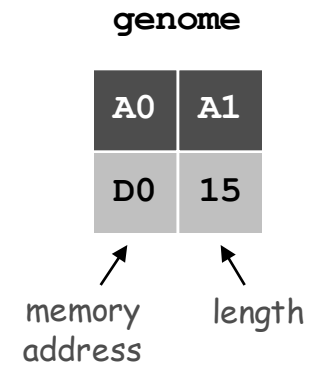
```
% more genomeTiny.txt
ATAGATGCATAGCGCATAGCTAGATGTGCTAGC

% java GeneFind ATG TAG < genomeTiny.txt
CATAGCGCA
TGC
```

# OOP Context for Strings

Possible memory representation of a string.

```
genome = "aacaagtttacaagc";
```

**genome**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | DA | DB | DC | DD | DE |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| a  | a  | c  | a  | a  | g  | t  | t  | t  | a  | c  | a  | a  | g  | c  |

| A0 | A1 |
|----|----|
| D0 | 15 |

memory address     length

# OOP Context for Strings

## Possible memory representation of a string.

```
genome = "aacaagtttacaagc";
```

**genome**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | DA | DB | DC | DD | DE |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| a  | a  | c  | a  | a  | g  | t  | t  | t  | a  | c  | a  | a  | g  | c  |

| A0 | A1 |
|----|----|
| D0 | 15 |

memory address     length

**s**          **t**

```
s = genome.substring(1, 5);
t = genome.substring(9, 13);
```

| B0 | B1 |   | B2 | B3 |
|----|----|---|----|----|
| D1 | 4  |   | D9 | 4  |

s and t refer to different strings that have the same value `"acaa"`

`(s == t)` is `false`, but `(s.equals(t))` is `true`.

compares pointers          compares character sequences

9

# In and Out

# Non-Standard Input

or use OS to redirect from one file

**Standard input.** Read from terminal window.
**Goal.** Read from several different input streams.

**In data type.** Read text from stdin, a file, a web site, or network.

**Ex:** Are two text files identical?

```
public class Diff {
    public static void main(String[] args) {
        In in0 = new In(args[0]);         ← read from one file
        In in1 = new In(args[1]);         ← read from another file
        String s = in0.readAll();
        String t = in1.readAll();
        StdOut.println(s.equals(t));
    }
}
```

# Screen Scraping

**Goal.** Find current stock price of Google.

```
…
<tr>
<td class="yfnc_tablehead1"  width="48%">
Last Trade:
</td>
<td class="yfnc_tabledata1">
<big>
<b>576.50</b>
</big>
</td>
</tr>
<tr>
<td class="yfnc_tablehead1"  width="48%">
Trade Time:
</td>
<td class="yfnc_tabledata1">
11:45AM ET
</td>
</tr>
…
```

`http://finance.yahoo.com/q?s=goog`

NYSE symbol

# Screen Scraping

Goal. Find current stock price of Google.

`s.indexOf(t, i)`: index of first occurrence of pattern `t` in string `s`, starting at offset `i`.

Read raw html from `http://finance.yahoo.com/q?s=goog`.

Find first string delimited by `<b>` and `</b>` after `Last Trade`.

```java
public class StockQuote {
   public static void main(String[] args) {
      String name = "http://finance.yahoo.com/q?s=";
      In in = new In(name + args[0]);
      String input = in.readAll();
      int start    = input.indexOf("Last Trade:", 0);
      int from     = input.indexOf("<b>",  start);
      int to       = input.indexOf("</b>", from);
      String price = input.substring(from + 3, to);
      StdOut.println(price);
   }
}
```

```
% java StockQuote  goog
576.50
```

# Day Trader

Add bells and whistles.

Plot price in real-time.

Notify user if price dips below a certain price.

Embed logic to determine when to buy and sell.

Automatically send buy and sell orders to trading firm.

Warning.  Please, please use at your own financial risk.



*The New Yorker, September 6, 1999*

# OOP Summary

**Object.**  Holds a data type value; variable name refers to object.

**In Java, programs manipulate references to objects.**
     Exception:  primitive types, e.g., `boolean`, `int`, `double`.
     Reference types:  `String`, `Picture`, `Color`, arrays, everything else.
     OOP purist:  language should not have separate primitive types.

**Bottom line.**  We wrote programs that manipulate colors, pictures, and strings.

**Next time.**  We'll write programs that manipulate our own abstractions.

# Extra Slides

## Color Separation

```java
import java.awt.Color;
public class ColorSeparation {
    public static void main(String args[]) {
        Picture pic = new Picture(args[0]);
        int width  = pic.width();
        int height = pic.height();

        Picture R = new Picture(width,  height);
        Picture G = new Picture(width,  height);
        Picture B = new Picture(width,  height);

        for (int x = 0; x < width; x++) {
            for (int y = 0; y < height; y++) {
                Color c = pic.get(x,  y);
                int r = c.getRed();
                int g = c.getGreen();
                int b = c.getBlue();
                R.set(x,  y,  new Color(r,  0,  0));
                G.set(x,  y,  new Color(0,  g,  0));
                B.set(x,  y,  new Color(0,  0,  b));
            }
        }
        R.show();
        G.show();
        B.show();
    }
}
```
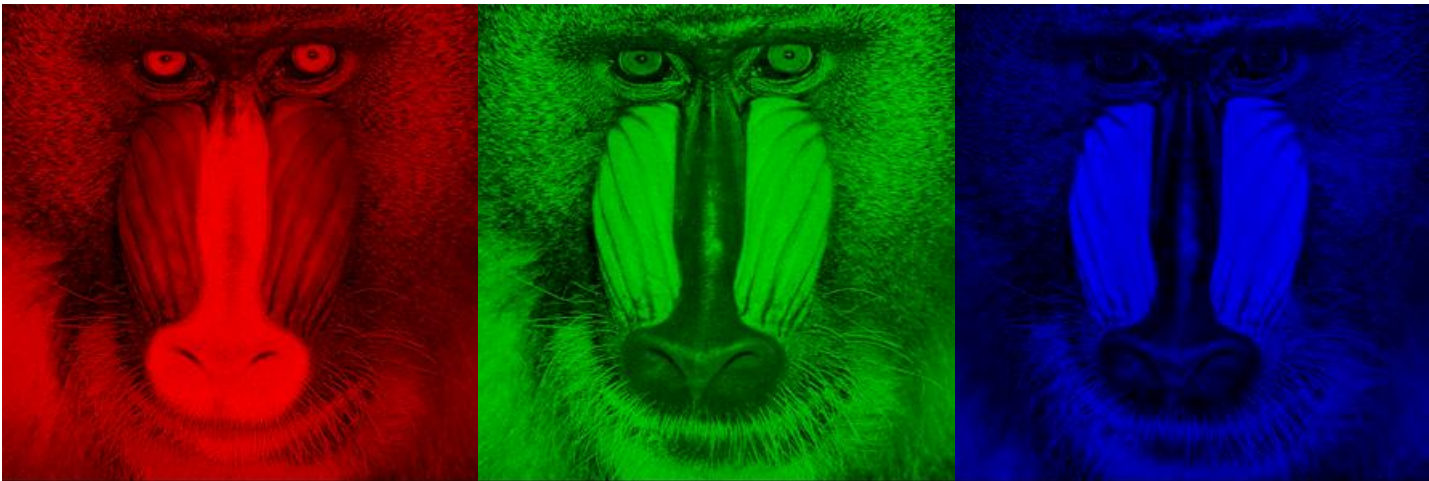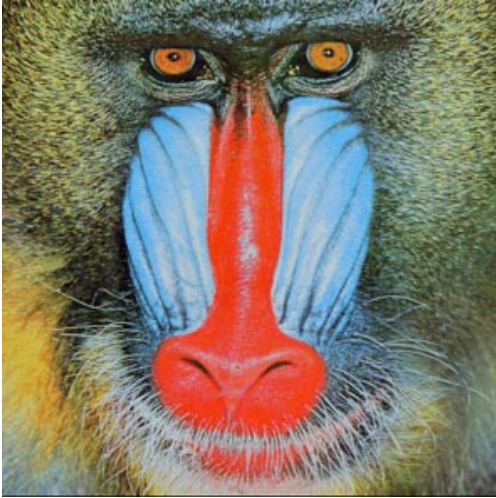
# Color Separation

ColorSeparation.java. Creates three `Picture` objects and windows.

# Memory Management

Value types.
    Allocate memory when variable is declared.
    Can reclaim memory when variable goes out of scope.

Reference types.
    Allocate memory when object is created with `new`.
    Can reclaim memory when last reference goes out of scope.
    Significantly more challenging if several references to same object.

Garbage collector. System automatically reclaims memory;
programmer relieved of tedious and error-prone activity.