

ITI 1121. Introduction to Computing II

Winter 2013

Assignment 1

(Last modified on February 4, 2013)

Deadline: Friday February 1, 2013, 18:00

[[PDF](#)]

Solution

- [a1-solution.jar](#)

Learning objectives

- Using arrays and matrices to implement complex types
- Recognize the role of reference variable to create class associations
- Applying basic of object oriented programming principles to solve problems
- Constructing software systems using one-to-many class associations
- Editing, compiling, debugging, and running Java programs
- Raising awareness concerning the university policies for academic fraud

Background information

Analytics¹ is one of “the four key growth areas of Information and Communications Technology (ICT)²”. A Google search for “analytics jobs” returns over 50 million Web documents! Analytics aims to extract meaningful information (knowledge) from large and complex data-sets. One of its challenges is to effectively and efficiently compare diverse documents for the purpose of finding similarities. This has applications for searching the Web, searching the literature, detecting plagiarism, detecting intrusions in computer systems, searching genomic data bases, etc. For simplicity, we will limit ourselves to strings representing genomic data (genes). However, the methodology can also be applied to other forms of documents, including text, music, video, and stock exchange data.

Problem statement

How does one compare documents whose length and content are different? One approach, which we explore here, is to transform documents into feature vectors.

In “text mining”, a ***k*-gram** (or word) is a contiguous sequence of k symbols over some alphabet. Given a four letter alphabet, A,C,G,T, there is a total $4^2 = 16$ 2-grams. A **profile** gives the number of occurrences of k -grams in a document. Given an input String of length n , there are $n - k + 1$ overlapping strings (k -grams) of length k .

As an illustration, the complete mitochondrial genome of Neanderthal is a string of length 16,565 over a four-letter alphabet (A,C,G,T):

```
>gi|196123578|ref|NC_011137.1| Homo sapiens neanderthalensis mitochondrion
GATCACAGGTCTATCACCTATTAACCACTCACGGGAGCTCTCCATGCATTTGGTATTTTCGTCTGGGGG...(16565)...T
TGGGGGTAGCTAAAGTGAAGTGTATCCGACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCTTA
AATAAGACATCACGATG
```

¹A term akin to “Big Data” and “Data Mining”

²[IT worker shortage has serious implications for Canada by Dan Oysey, Financial Post, December 11, 2012.](#)

Its profile is a vector of counts of size 16:

Profile: [words = AA(1589), AC(1490), AG(799), AT(1229), CA(1533), CC(1769), ..., TT(1016)]

The above profile shows that the 2-gram AA was found 1,589 times in the mitochondrial genome of Neanderthal, whilst AC was found 1,490 times. Any genomic sequence can be transformed into a profile, and compared to that of Neanderthal.

Implementation

This assignment considers two implementations of a profile. In the first implementation, all possible k -grams are enumerated and stored into memory. Accordingly, this implementation will be called **WordProfile**. For the second implementation, each k -gram is interpreted as a base 4 number. These base 4 numbers represent the index of the k -grams in the profile. Accordingly, this implementation will be called **HashProfile**.

1 Rules and regulation (10 marks)

Follow all the directives available on the [assignment directives web page](#), and submit your assignment through the on-line submission system [uottawa.blackboard.com](#).

You must preferably do the assignment in teams of two, but you can also do the assignment individually. Pay attention to the directives and answer all the following questions.

2 Word (10 marks)

The class **Word** will be used for the implementation of **WordProfile**. Here are the characteristics of the class **Word**.

1. A **Word** object stores a k -gram (a String of size k) as well as a count. The initial value of the count is zero.
2. The class **Word** has a single constructor. It has a single formal parameter, **word** (of type String).
3. **void count()**: A **Word** object has a method **count** that increments the counter of this object by one.
4. **int getCount()**: The method **getCount** returns the current value of the counter.
5. **String getWord()**: The method **getWord** returns the k -gram represented by this object.
6. Finally, each **Word** has a **toString()** method that returns a String object consisting of the k -gram represented by this object, as well as the count, in between parentheses.

3 WordProfile (30 marks)

A **WordProfile** stores the number of occurrences of possible k -grams in a given input string. The k -grams are stored and returned in lexicographic order.

1. **WordProfile** uses an array of objects of the class **Word** to tally the counts of all k -grams for a given input String.
2. Has a constructor with the following signature, **WordProfile(int k, String input)**, where k is the length of the k -grams and **input** is a reference to a **String** object.
3. **int getSize()**: returns the length of the profile.
4. **String getWord(int index)**: returns the word (k -gram) at the designated position.
5. **int getCount(int i)**: returns the count the i th word.
6. **int getCounts()**: returns the total number of k -grams examined from the input.

7. **double getDistance(WordProfile other)**: returns the distance between *this* and the **other** profile. For this assignment, we measure the Euclidean distance of the two profiles. Some care must be taken so that documents of different lengths can be compared. Indeed, longer documents will display higher count values. We therefore normalize each count of a profile by dividing it by the total number of counts. Each value of the profile is now seen as a frequency. Here is the definition of the Euclidean distance for vectors of size d :

$$((x_0 - y_0)^2 + (x_1 - y_1)^2 + \dots + (x_{d-1} - y_{d-1})^2)^{\frac{1}{2}}$$

where x_0 is the frequency of the first k -gram in this profile, that is the count of the first k -gram divided by the total number of k -grams in the input, y_0 has the same definition for the other profile. Given an alphabet of size 4, there are $d = 4^k$ distinct k -grams.

8. **String getMostFrequentWord()**: returns the word that has the highest count.
9. **String toString()**: returns a String representation of the project. Specifically, it displays all the k -grams together with their counts.

4 HashProfile (30 marks)

A **HashProfile** stores the number of occurrences of k -grams in a given input string. Unlike **WordProfile**, the k -grams are not explicitly stored! Instead, internally, each k -gram is interpreted as a base 4 number, which represents the index of the k -gram in the array of counts. Be careful, the counts for k -grams must be stored in k -grams lexicographic order. Therefore, **getWord(0)** will return AA for both, a **WordProfile** object and a **HashProfile** object, assuming $k = 2$ and the alphabet is A,C,G,T.

1. **HashProfile** uses an array of integers to tally the counts of all k -grams for a given input.
2. Has a constructor with the following signature, **HashProfile(int k, String input)**, where k is the length of the k -grams and **input** is a reference to a **String**.
3. **int getSize()**: returns the length of the profile.
4. **String getWord(int index)**: returns the word (k -gram) at the designated position.
5. **int getCount(int i)**: returns the count of the i th word.
6. **int getCounts()**: returns the total number of k -grams in the input.
7. **double getDistance(HashProfile other)**: returns the distance between *this* and the **other** profile. For this assignment, we measure the Euclidean distance of the two profiles. Some care must be taken so that documents of different lengths can be compared. Indeed, longer documents will display higher count values. We therefore normalize each count of a profile by dividing it by the total number of counts. Each value of the profile is now seen as a frequency. Here is the definition of the Euclidean distance for vectors of size d :

$$((x_0 - y_0)^2 + (x_1 - y_1)^2 + \dots + (x_{d-1} - y_{d-1})^2)^{\frac{1}{2}}$$

where x_0 is the frequency of the first k -gram in this profile, that is the count of the first k -gram divided by the total number of k -grams in the input, y_0 has the same definition for the other profile. Given an alphabet of size 4, there are $d = 4^k$ distinct k -grams.

8. **String getMostFrequentWord()**: returns the word that has the highest count.
9. **String toString()**: returns a String representation of the project. Specifically, it displays all the k -grams together with their counts.

5 WordCompareAll and HashCompareAll (15 marks)

WordCompareAll and **HashCompareAll** are two test programs. Their content will be identical, except that **WordCompareAll** uses **WordProfile** objects, whereas **HashCompareAll** uses **HashProfile** objects. The programs are reading the arguments from the command line:

```
java WordCompareAll 2 toxin.txt unknown_a.txt unknown_b.txt
```

where k is the length of the k -grams. This is followed by a list of file names. Here is the behaviour of the test programs.

1. Display the student information.
2. Read and create a profile for each text file given as an argument to the program.
3. Display all the profiles.
4. Display the distance between all pairs of profiles.
5. Display the most frequent word for each profile

6 Application (5 marks)

Cases of food poisoning are unfortunately quite common — “Two people are paralyzed after drinking botulism-contaminated carrot juice s(. . .)” from Botulism-tainted juice paralyzes two in Canada, www.reuters.com, 2006/10/17 — we have been asked to help the investigators.

Background information. Botulism is a rare disease caused by a toxin (a protein) expressed by the bacterium *Clostridium botulinum*. From the above source of information: “Botulism can cause nausea, fatigue, double-vision, paralysis and respiratory failure. In severe cases, the toxin can be fatal.”

The investigators suspect that at least one of the following proteins, Unknown A or Unknown B, has a similar domain to that of *Clostridium botulinum*’s toxin. Which protein(s) is (are) related to this toxin? Why?

```
>gi|27867582 (fragment of the known Clostridium botulinum toxin gene)
```

```
GTGAATCAGCACCTGGACTTTTCAGATGAAAAATTAAATTTAACTATCCAAAATGATGCTT
ATATACCAAAATATGATTCTAATGGAACAAGTGATATAGAACAACATGATGTTAATGAAC
TTAATGTATTTTTCTATTTAGATGCACAGAAAGTGCCCGAAGGTGAAAATAATGTCAATC
TCACCTCTTCAATTGATACAGCATTATTAGAACAACCTAAAATATATACATTTTTTTCAT
CAGAATTTATTAATAATGTCAATAAACCTGTGCAAGCAGC
```

```
> Unknown A
```

```
TCTATCAAGTAGATTATTAATACTACTGCACAAAATGATTCTTACGTTCCAAAATATGA
TTCTAATGGTACAAGTGAATAAAGGAATATACTGTTGATAAACTAAATGTATTTTTCTA
TTTATATGCACAAAAAGCTCCTGAAGGTGAAAGTGCTATAAGTTTAACTTCTTCAGTTAA
TACAGCATTATTAGATGCATCTAAAGTTTATACGTTTTTTTCTTCAGATTTTATTAATAC
```

```
> Unknown B
```

```
TCCTGGCTCAGGACGAACGCTGGCGGCGTGTGCTTAACACATGCAAGTCGAGCGATGAAG
CTTCCTTCGGAAGTGGATTAGCGGCGGACGGGTGAGTAACACGTGGGTAACCTGCCTCA
AAGTGGGGGATAGCCTTCCGAAAGGAAGATTAATACCGCATAACATAAGAGAATCGCATG
ATTTTCTTATCAAAGATTTATTGCTTTGAGATGGACCCGCGCGCATTAGCTAGTTGGTA
```

7 Academic fraud

This part of the assignment is meant to raise awareness concerning plagiarism and academic fraud. Please read the following documents:

- <http://web5.uottawa.ca/admingov/regulations.html#r71>
- <http://web5.uottawa.ca/mcs-smc/academicintegrity/>
- <http://www.uottawa.ca/plagiarism.pdf>

Cases of plagiarism will be dealt with according to the university regulations.

By submitting this assignment, you acknowledge:

1. having read the above documents, and
2. understanding the consequences of plagiarism

Files

You must hand in the following files (in .jar file).

- README.txt
- Word.java
- WordProfile.java
- WordCompareAll.java
- HashProfile.java
- HashCompareAll.java
- Utilities.java
- StudentInfo.java
- toxin.txt
- unknown_a.txt
- unknown_b.txt
- alq6.txt (the answer to question 6)

The JavaDoc of the classes can be found here: [Documentation](#).

Here are some tests, but remember that “testing shows the presence, not the absence of bugs”, Edsger W. Dijkstra.

- Assignment1AllTest.java
- HashProfileTest.java
- WordProfileTest.java
- WordHashProfileTest.java
- WordTest.java

A Frequently Asked Questions (FAQ)

1. “Will the value of k always be 2?”

No. You can assume that the value of k will always be greater than 0, but there is no upper bound.

2. “Is the alphabet fixed?”

Yes, for this assignment, you can assume that the alphabet will always be “A,C,G,T”.

3. “Is a profile storing words with non-zero counts or all of them?”

It is the latter, a profile stores all the words, even those having a count of zero.

4. “Could you give us a worked out example?”

The figure below shows the resulting profiles for the same input but different values of k , $k = 2, 3, 4$.

k=2

ATGTACAGCTTAGTAGCTACA**AT**ACTGTGGACTAC**GA**ATG

AA(2), AC(5), AG(3), **AT(3)**, CA(2), CC(0), **CG(1)**, CT(4), GA(2), GC(2), GG(1), GT(3), TA(6), TC(0), TG(4), TT(1)

16 2-grams

k=3

ATGTACAGCTTAGTAGCTACA**AAT**ACTGTGGACTAC**GA**ATG

AAA(0), AAC(0), AAG(0), **AAT(2)**, ACA(2), ACC(0), ..., CAA(1), ..., TTT(0)

64 3-grams

k=4

ATGTACAGCTTAGTAGCTACA**AATA**CTGTGGACTAC**GA**ATG

AAAA(0), AAAC(0), AAAG(0), AAAT(0), AACA(0), AACC(0), AACG(0), AACT(0), AAGA(0), AAGC(0), AAGG(0), AAGT(0), **AATA(1)**, ..., TTTT(0)

256 4-grams

5. “Do we have to handle error situations?”

No. Until we have seen exceptions in class, which is the mechanism for handling error situations in Java, you can assume that all the input values are valid.

6. “Can we add new methods?”

You cannot add **public** methods, but you can add (and are encouraged to) as many **private** methods as you want.

7. “Could you please post test cases?”

Here are some tests, but remember that “testing shows the presence, not the absence of bugs”, Edsger W. Dijkstra.

- [Assignment1AllTest.java](#)
- [WordProfileTest.java](#)
- [HashProfileTest.java](#)
- [WordHashProfileTest.java](#)
- [WordTest.java](#)

8. “I think I found an error in your tests. I think the Euclidean distance is ill-defined for profiles of different lengths.”

Indeed, the Euclidean distance between two profiles has a meaning only if the two profiles being compared were produced using the same alphabet and k . This has now been fixed.

9. “What word should the method `getMostFrequentWord()` return when there are more than one possibility?”

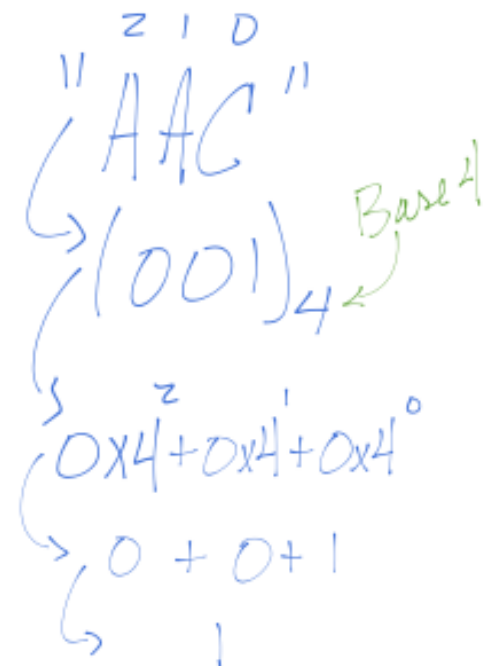
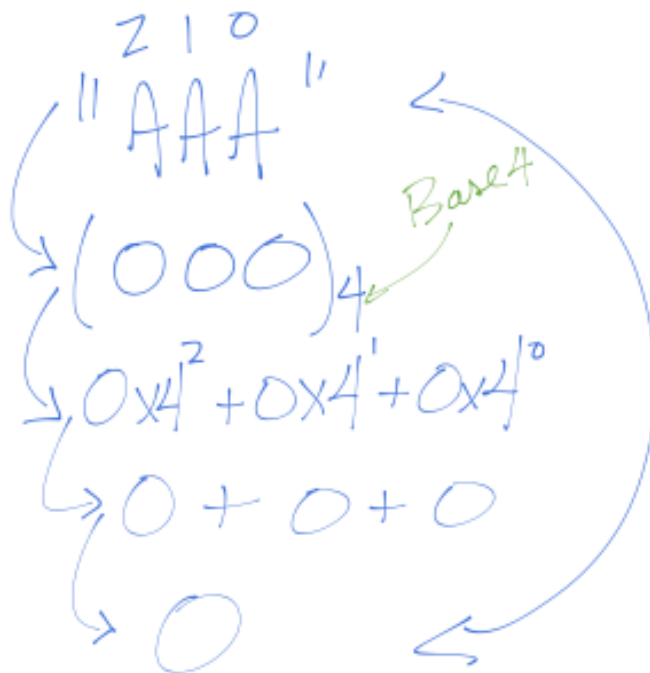
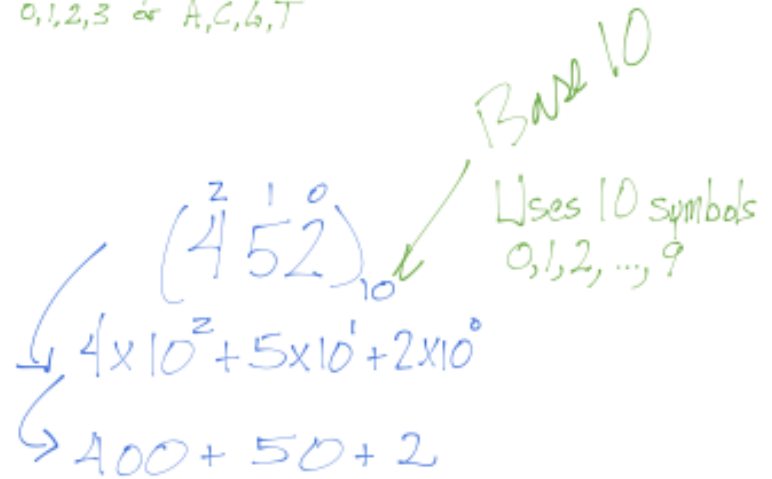
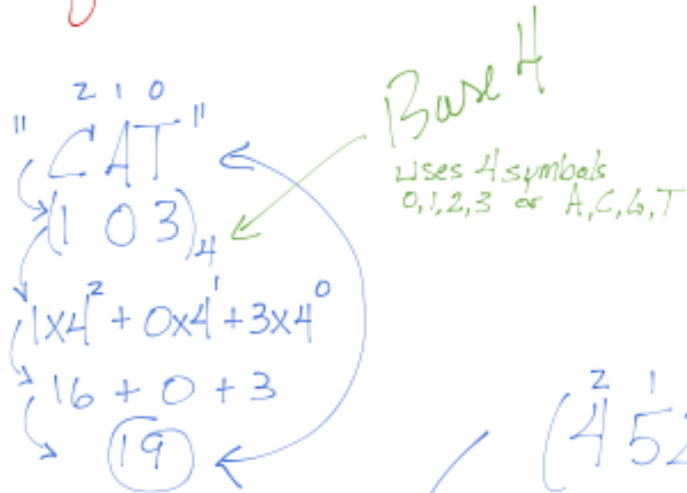
Since, this is not specified in the description, any word whose count is maximal is a valid solution. I would suggest returning the one that comes first in lexicographic order.

10. “In class, you showed how to transform a k -gram into an index. Could you post these notes.”

Here are my notes:

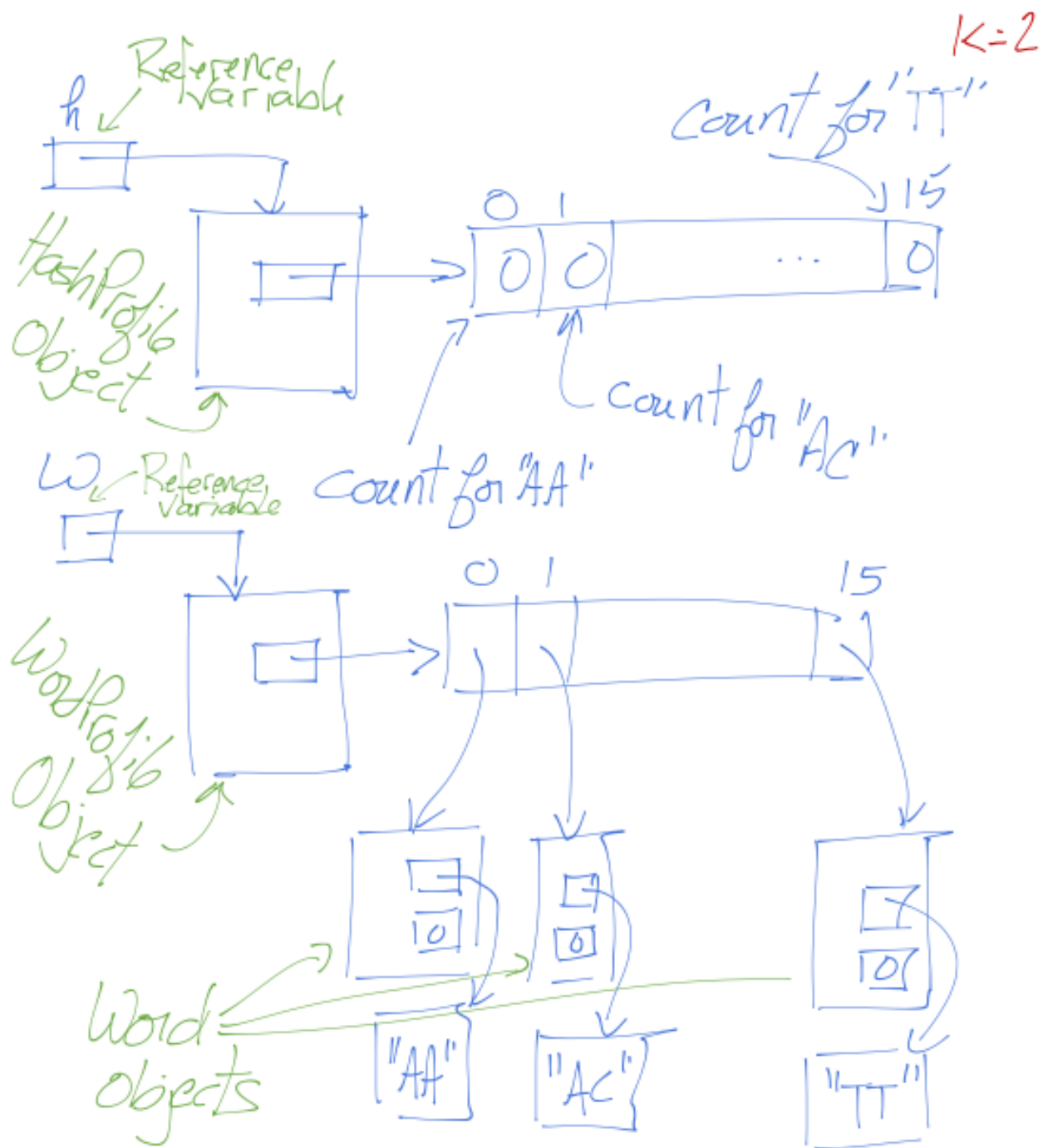
HashProfile
K=3

A	≡ 0
C	≡ 1
G	≡ 2
T	≡ 3



11. "In class, you showed memory diagrams for WordProfile and HashProfile, can I get a copy of those diagrams?"

Here are my notes:



Memory diagrams for HashProfile & WordProfile objects

12. "What approach do you suggest for reading the input files?"

This part of the assignment has already been solved for you. I am providing with this assignment a class called `Utilities.java` that has a method `getLine(String fname)`. The method reads a line from the input file. For this assignment, the input files contain a single line. Here is how it is used.

```
String input;
input = Utilities.getLine("toxin.txt");
```

13. "Can I use the methods of the class String?"

Indeed, you have to use the methods of the class String in order to extract the k-grams out of the input String.

Last Modified: February 4, 2013