# Université Jean Monnet

## Machine Learning & Data Mining

---

# Advanced Algorithms

## Longest Common Subsequence

---

| Edward | Beeching |
| Jérémie | Blanchard |
| Anthony | Deveaux |
| Josselin | Marnat |
| Joseph | Renner |

**UNIVERSITÉ JEAN MONNET**
SAINT-ÉTIENNE

December, 12th 2016

# Table of Contents

# 1 Introduction

The goal of this project is to build a plagiarism detector by using algorithms studied in class to find the Longest Common Subsequence (LCS) between two input texts. The Longest Common Subsequence is the longest subsequence that occurs in both texts. A subsequence is defined as a sequence that can be obtained from another only by deletion of elements, and not reordering. The elements need not be consecutive in the original text.

These are the algorithms implemented to find the LCS between two texts (described in detail in section 2.1):

- Dynamic Programming

- Divide and Conquer (linear space)

- Recursive

- Recursive with Branch and Bound

The project also includes the implementation of algorithms studied in class to solve the printing neatly problem. The printing neatly problem refers to the problem of printing a text while minimizing the number of empty spaces at the end of each line. Formally, the text is composed of n words where the number of characters in each word are represented as $l_1, l_2, \ldots, l_n$. The text is to be printed so each line holds a maximum of M characters. The number of extra spaces printed at the end of a line which contains words $i$ through $j$ is defined as follows:

$$M - j + i - \sum_{k=i}^{j} l_k$$

We want to minimize the sum of the cube of this non-negative value over all lines but the last. This is implemented using the following algorithms (described in detail in section 2.2):

- Dynamic Programming

- Greedy

- Recursive with Branch and Bound

Lastly, a study of the behavior of LCS algorithms and printing neatly algorithms is presented, including runtime analysis, space analysis, experimental runtime analysis and comparisons, and scalability.

# 2 Algorithms presentation

## 2.1 LCS

### 2.1.1 Dynamic programming

The dynamic programming algorithm for computing the LCS between two input texts (X, Y, with lengths m, n, respectively) is found as follows (adapted from Amaury Habrard's slides):

$\rightarrow$ `Algorithm 1 [Appendix]`

$\rightarrow$ `Algorithm 2 [Appendix]`

$\rightarrow$ `Algorithm 3 [Appendix]`

Algorithm for building LCS from $b$ matrix:

$\rightarrow$ `Algorithm 4 [Appendix]`

**LCS Linear Space Forwards:** The following algorithm finds the length of the LCS of two input strings using dynamic programming with linear space complexity.

$\rightarrow$ `Algorithm 5 [Appendix]`

**LCS Linear Space Backwards:** The following algorithm is a linear space implementation of the LCS dynamic programming algorithm, except it starts at the end of texts and works towards the beginning.

$\rightarrow$ `Algorithm 6 [Appendix]`

### 2.1.2 Divide and Conquer

**LCS Divide and Conquer:** The following algorithm uses the LCS linear space forward and backward algorithms to compute the LCS of two input spaces in linear space, with the functionality to return the actual LCS (not just the length).
The algorithm is based on the following 2 properties:

1. The size of the LCS that passes through entry $(i, j)$ is sum of the LCS from $(0, 0)$ to $(i, j)$ (denoted $c[i, j]$) and the LCS from $(i, j)$ to $(m, n)$ (denoted g[i,j])

2. Let $k$ be any number from 0 to $n$, let $q$ be the number that maximizes $c[q, k] + g[q, k]$. There is an optimal solution to the LCS problem (from $0, 0$ to $m, n$) that passes through $q, k$

Combining these properties and the linear space algorithms presented above, we can formulate a divide and conquer algorithm that requires linear space.

$\rightarrow$ `Algorithm 7 [Appendix]`

### 2.1.3  Recursive

**LCS Recursive:**  The recursive algorithm for finding the LCS is a depth first search through the search space, enumerating all possible solutions and returning the largest LCS found.

$$\rightarrow \texttt{Algorithm 8 [Appendix]}$$

### 2.1.4  LCS Branch and bound

**Bounds:**  The bounds used for this implementation were taken from the paper *An Effective Branch-and-Bound Algorithm to Solve the k-Longest Common Subsequence Problem*[1].

- The upper bound is computed as follows:

$$\text{UB}_c = \sum_{\sigma \in \Sigma} \min_{i=1}^{k}(\text{number of characters } \sigma \text{ in sequence } x^{(i)})$$

- The lower bound is computed as follows:

$$\text{LB}_c = \max_{\sigma \in \Sigma} \min_{i=1}^{k}(\text{number of characters } \sigma \text{ in sequence } x^{(i)})$$

To ensure to bound is computed in O(n) time, three hash tables (Python dicts) are used for fast lookup:

1. To count the occurrence of letters in list of words 1

2. To count the occurrence of letters in list of words 2

3. To count the common letters between hash table 1 and 2

For example for lists of letters [A,B,C,A] and [B,C,A,A], the common letters are A:2, B:1, C:1, the upper bound is found to be 2+2+1 = 4 and the lower bound is 2.

**Implementation:**  The algorithm is implemented as a search tree, with a last in first out queue used for a depth first approach. The tree is explored node by node, with the upper bound evaluated for the two children at each step. Children with the higher upper bounds are placed on the queue last, to ensure the solutions with the most potential are explored first. Items are not placed on the queue if their bound is smaller than the best solution found so far. As elements are popped off the queue, if the bound exceeds the best solution found so far the node is explored. The items in the queue are a quadruplet containing (solution, index, tracker_index, upper_bound)
The solution is enumerate in binary as an integer, for example for strings:
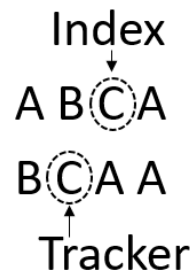
- A,B,C,F

- B,C,D,E

A solution would be 0110 which would be represented as the integer 6, using a binary representation ensures the elements in the queue are as memory efficient as possible.
The "index" and "tracker" index in the quadruplet are the respective indexes of the in the two sets of strings.

---

Index

A B C A

B C A A

Tracker

Example of the index and tracker

## 2.2   Printing Neatly

**Greedy Solution:**   The greedy solution simply prints as many words as possible on a line before going to the next. The solution it gives is not guaranteed to be optimal.

**Dynamic Programming:**   The dynamic programming approach computes an optimal solution in polynomial time. First, we define extras$[i, j]$ to be the number of extra spaces on a line containing words i through j. Using this, we can compute the cost of a line containing the same words:

1. lc$[i, j] = \infty$ if extras$[i, j]$ ¡ 0 (if words do not fit on the line)

2. lc$[i, j] = 0$ if $j = n$ and extras$[i, j] \geq 0$ (if it is the last line and the words fit)

3. lc$[i, j] = ($extras$[i, j])^3$ otherwise

Using this, we can compute the optimal arrangement (minimal cost) of words 1..j = c[j] :

1. $c[j] = 0$ if $j = 0$

2. $c[j] = \min((1 \leq i \leq j), (c[i-1] + lc[i, j]))$ if $j > 0$

**Recursive:**   The recursive solution iterates through all the possible sequences and returns the optimal score and solution.

**Branch and Bound:**   Two branch and bound methods were compared, in both cases an initial estimate of the upper bound was evaluated using the the cost of a greedy approach. The first method then evaluates solutions in the same manner as the recursive, but when the cost of the current branch exceeds the bound given by the greedy solution, the recusion is cut and the algorithm back-tracks. The second algorithm makes an estimate of the cost of the remaining solution by using a greedy approach to find the cost of any remaining words, this enables to algorithm to cut potentially worse solutions early. The second approach is faster than the dynamic approach but unfortunately does not alway produce an optimal solution.

# 3 Preprocessing techniques

**Introduction:** We came up with four different preprocessing techniques to aid in detecting plagiarism. Each technique can be implemented on the two input texts and then the LCS is computed on the results of the preprocessing.

## 3.1 Light Preprocessing

The most basic preprocessing replaces certain symbols to make the texts more universal. In this technique, directional quotes (quotation marks that indicate if a quote is starting or finishing) are replaced by straight quotes, and extra newline characters and spaces are removed.

## 3.2 Advanced Preprocessing

Advanced preprocessing deletes all non-ascii characters and all symbols that are not periods or commas, and makes everything lower case. This is done to make the LCS more robust to small changes in punctuation. For example, the following is an text before preprocessing:

> This kind of relationship can be visualized as a tree structure, where 'student' would be the more general root node and both 'postgraduate' and 'undergraduate' would be more specialized extensions of the 'student' node (or the child nodes).

And here would be the same text after advance preprocessing:

> this kind of relationship can be visualized as a tree structure, where student would be the more general root node and both postgraduate and undergraduate would be more specialized extensions of the student node or the child nodes.

**Stop Word Removal:** For stop word removal preprocessing, an input text is run through the advanced preprocessing technique seen above and then passed to a function that deletes all occurrences of stop words. Stop words are any words that do not carry any significance to the similarity of two texts. We chose to remove them as a possible preprocessing technique to make the LCS robust to students simply adding in meaningless words to cover plagiarism. The stop words are listed in a file. Examples of stop words include "furthermore", "regardless", and "somehow".
Below is the same example sentence from above after stop word removal:

> kind relationship visualized tree structure, student would general root node postgraduate and undergraduate specialized extensions student node child nodes

However, we need to take care about the number of stop words we're removing, because by suppressing too much of them, the sentence can be so simple that finding the LCS don't even make sense anymore.

Stop Word preprocessing delete some 'useless' words (according to a certain dictionary), so the running time is improved. But if you remove too many words, the algorithm is useless. So we have to balance it well, so a lot of test were done.

**Word Ordering**   Word ordering preprocessing starts by feeding an input text to the advanced preprocessing technique from above. Then, each sentence in the text is reordered so that the words in the sentence are arranged alphabetically. This is done so that texts that have copied and rearranged content from an original text are still detected as plagiarism.

Below is the same original sentence from above, after the word ordering preprocessing:

> a and and as be be be both can child extensions general kind of more more node nodes of or postgraduate relationship root specialized structure, student student the the the this tree undergraduate visualized where would.

# 4   Graphical User Interface

## 4.1   Plagiarism detector

The Graphical User Interface (GUI) uses a PYTHON library called `PyQt4`.

In this first tab called 'LCS Unit', we have many separated areas splitted as follows:

1. on the left: LCS inputs:

   (a) option bar:
   - `Task` : task from the corpus
   - `Text` : text from the corpus
   - `LCS preprocess` can be selected among:
     – `Raw`
     – `Advanced preprocessing`
     – `preprocessing`
     – `swr (stop words removal) preprocessing`
     – `WordOrgering preprocessing`
   - `LCS algorithm` can be selected among:
     – `LCS` (classic)
     – `LCS Sentence` (by sentence)

   (b) the text to compute from the corpus

   (c) the 'wiki' text to compare with

2. on the right: LCS results:

   (a) printing neatly section :
   - max-width slider (from 20 to 60)
   - algorithm chooser (Dynamic or Greedy)
   - the LCS printed neatly according to the width and the algorithm

   (b) some numbers resulted from the computation:
   - the original length of the text from the corpus
   - the LCS length

---

- the plagiarism score, defined as follow :
  let $T$ be the original text, and $s_i$ the subsentences of $T$ where words from the LCS are side-by-side in $T$ (words in bold in the GUI). Then we can compute:
  score $= \sum_i (|s_i|)^2$
- the actual running time to find the LCS
- then, a pie chart showing the percentage of plagiarized text *vs.* the original one.

## 4.2   Data table

We used an EXCEL file (`corpus-final09.xls`) to store informations concerning the execution of all the algorithms on all the texts from the corpus.

$\rightarrow$ Figure 3 [Appendix]

- [A] the name of the file

- [B:D] group, person and task

- [E:F] the category of plagiarism, native or not

- [G:H] knowledge and difficulty

- LCS Length and Ratios for:

  - [I:J] no preprocessing
  - [K:L] light preprocessing
  - [M:N] advanced preprocessing
  - [O:P] by sentence, advanced preprocessing
  - [Q:X] without stop words (again with no,advanced,light preprocessing, and alphabetical ordering)

## 4.3   Graphs

**General graph:**   In this third tab, we are simply showing the LCS ratios (LCS length over text length) for all the files, and for the different preprocessing techniques. We can see that the more advanced the preprocessing, the longer the LCS: that is exactly what we wanted.

$\rightarrow$ Figure 4 [Appendix]

**Graph by task:**   The fourth tab shows the LCS Ratio for each task (and for each text of the corpus referring to this task), labeling with different colors the plagiarism category given by the EXCEL file. The labels are confirmed by this graph: we see that the 'cut' category as a much higher LCS ratio than the 'light' and 'heavy' ones, and the 'non' are very low.

$\rightarrow$ Figure 5 [Appendix]

**Repartition by category:**   This last tab is here to display the repartition of each LCS ratio over the plagiarism categories (non plagiarized, heavy modified, light modified and cut from the text)

$\rightarrow$ Figure 6 [Appendix]

---

# 5 Experimental study

## 5.1 Time complexity

**LCS Dynamic Programming:** The dynamic programming algorithm for finding the length of the LCS requires $\mathcal{O}(mn)$ time, where m and $n$ are the lengths (number of words) of the two input texts. This can be observed from the algorithm itself, as there is a for loop from 1 to $n$ nested inside a for loop from 1 to $m$. The recursive algorithm used to build the solution from the b matrix requires $\mathcal{O}(m+n)$ time. Each run of the algorithm takes constant time, and each recursive call to the algorithm decreases at least one of the lengths by one. Combining these two results gives a time complexity of $\mathcal{O}(mn)$ to compute the length and build the LCS.

**LCS Linear Space (Forwards and Backwards):** The linear space versions of the dynamic programming algorithm do the same computations, with an added step at copying the row of the matrix after each row iteration. As a result, the asymptotic time complexity remains $\mathcal{O}(mn)$.

**LCS Divide and Conquer:** Let $T(m, n)$ denote the run time of the algorithm. At each call of the algorithm, it performs two $\mathcal{O}(mn)$ calls to the space-efficient algorithm. Then it makes two recursive calls on strings of size $q$ and $\frac{n}{2}$, and $m - q$ and $\frac{n}{2}$, respectively. Then, for constant $c$, we have:

$$T(m, n) \leq cmn + T(q, \frac{n}{2}) + T(m - q, \frac{n}{2})$$
$$T(m, 2) \leq cm$$
$$T(2, n) \leq cn$$

If we assume $m = n$ and $q = \frac{n}{2}$, then we have:

$$T(m, n) \leq 2T(\frac{n}{2}) + cn^2$$

From the master theorem, we have $a = 2, b = 2, f(n) = n^2$, and $\log b(a) = 1$. Then, since $f(n)$ is $\Omega(n^c)$ where $c = 2 > \log b(a) = 1$, and the regularity condition holds, since $2(\frac{n^2}{4}) \leq kn^2$ with $k = \frac{1}{2}$, then it follows that $T(m, n) = \Theta(f(n)) = \Theta(n^2)$. So, when m=n, the run time is $\Theta(n^2)$. We can assume that when $m \neq n, T(m, n) \leq kmn$, and prove it by induction: assume $T(m', n') \leq km'n'$ for $m' < m$ and $n' < n$:

$$
\begin{aligned}
T(m, n) &\leq cmn + T(q, \frac{n}{2}) + T(m - q, \frac{n}{2}) \\
&\leq cmn + kq(\frac{n}{2}) + k(m - q)(\frac{n}{2}) && \text{(by induction hypothesis)} \\
&= cmn + kq(\frac{n}{2}) + km(\frac{n}{2}) - kn(\frac{n}{2}) \\
&= (c + k/2)mn.
\end{aligned}
$$

So, if $c = k/2$, then the proof works. Thus, the LCS divide and conquer has a time complexity of $\mathcal{O}(mn)$

---

**LCS Recursive:**   The recursion for the recursive solution is as follows:

$$\begin{aligned} T(m,n) &= T(n, m-1) + T(n-1, m) + c \\ &= (T(n, m-2) + T(n-1, m-1) + T(n-1, m-1) + T(n-2, m)) + c' \\ &\geq 2T(n-1, m-1) \end{aligned}$$

Therefore the algorithm is exponential, since it is $\mathcal{O}(2^{\min(n,m)})$

**LCS Branch and Bound:**   In the worst case, the branch and bound solution take the same time as the recursive. Thus, it is exponential. However, in practice, it is faster (see runtimes).

**Printing Neatly Greedy Solution:**   In the greedy algorithm, there is just one iteration over the whole text, as a decision is made to print a word on the same line if possible or move to the next line. As a result, the time complexity is $\mathcal{O}(n)$ where $n$ is the number of words in the text.

**Printing Neatly Dynamic Programming:**   The dynamic programming algorithm can be broken into 3 sections: first, computing the extras matrix; next, computing the line cost matrix; and finally, computing the optimal solution.

1. The extras matrix is an $n \times n$ matrix where each element is computed in constant time. Thus, it is $\mathcal{O}(n^2)$.

2. The line cost matrix is also an $n \times n$ matrix with constant cost for each element, so it is $\mathcal{O}(n^2)$.

3. The optimal cost matrix is an array of length n, but the cost to compute the element at index $i$ is a function of $i$. Thus, it is $\mathcal{O}(n^2)$.

Since all three main parts of the algorithm are $\mathcal{O}(n^2)$, the algorithm is $\mathcal{O}(n^2)$.

**Printing Neatly Recursive:**   The recursive solution iterates through all possible solutions of lines that have at least one word on each line and have 0 or more spaces on the end (that fit). Let $i$ be the number of words that can fit on the first line, then we have the following recursion:

$$T(n) = T(n-1) + T(n-2) + \cdots + T(n-i)$$

Expanding each recursion on the right hand side of the equation gives an exponential growth.

**Printing Neatly Branch and Bound:**   In the worst case, the branch and bound extension behaves just as the plain recursive algorithm, thus it is also exponential. However, the bounding can limit the number of unneeded computations (see runtimes).

## 5.2   Runtime analysis

### 5.2.1   LCS

We have run the different algorithms (Branch & Bound, Recursive, Divide & Conquer and Dynamic) with a different number of words, to study their time complexity. Then, we plotted this data into three different graphs, in the goal to better see what happens.

**Running times on linear scale / |words|:** in this graph, the two scales are linear, and we can already observe the speed of each algorithms, and say that the Branch & Bound and the Recursive are much slower than the Divide & Conquer and the Dynamic algorithms.

$\rightarrow$ Figure 8 [Appendix]

**Running times on logarithmic scale / |words|:** The have a better view of the complexity of the slower algorithms, we draw the same graph, but changing the time scale from linear to logarithmic. Now, we see that the Branch & Bound and the Recursive are linear in this scale, then we can deduce that, graphically, this two algorithms should be exponential.

$\rightarrow$ Figure 9 [Appendix]

**Running times on logarithmic scale / $(|\text{words}|)^2$:** This time, we keep the time scale linear, but the $x$-axis becomes the number of words **squared**. We demonstrate that the running time of the dynamic algorithms is far more efficient than the recursive and branch and bounds methods, and the time complexity is quadratic in $n$ (where $m = n$)

$\rightarrow$ Figure 10 [Appendix]

### 5.2.2 Printing Neatly

**Running times on linear scale / |words|:** Here we see that the more efficient algorithm to print neatly is the Greedy one, and the less efficient is the Recursive one.

$\rightarrow$ Figure 11 [Appendix]

**Running times on logarithmic scale / |words|:** Now, the Branch & Bound and the Dynamic algorithms are close from each other, and the Dynamic one is slightly worse.

$\rightarrow$ Figure 12 [Appendix]

**Running times on linear scale / $(|\text{words}|)^2$:**

$\rightarrow$ Figure 13 [Appendix]

## 5.3 Space complexity

**LCS Dynamic Programming :** In the dynamic programming algorithm, there are two $m \times n$ matrices ($m$ is length of first input text, $n$ is length of second), one for the $c$ matrix denoting optimal solutions, and one for the $b$ matrix that is used to build the solution. These matrices are in addition to the two input word arrays and LCS word array. Thus, the space complexity is $\mathcal{O}(mn)$.

**LCS Linear Space (forward and backward):** In the linear space versions of the dynamic programming algorithm, we don't use the $b$ matrix, and the $c$ matrix is reduced to just 2 rows by $n$ columns. Thus, the space complexity is $\mathcal{O}(n)$.

**LCS Divide and Conquer:**   The divide and conquer algorithm uses the linear space algorithm on each half of the first input text and the full second text. The space complexity for these algorithms alone is $\mathcal{O}(n)$; however, we must save the full column indexed at $j = \frac{n}{2}$ in order to find the optimal row to divide our problem. Since there are m rows, the space needed to save the two columns (one for forwards LCS, one for backwards LCS) is $2m$. Adding this to the $\mathcal{O}(n)$ from the linear space algorithm, we have a space complexity of $\mathcal{O}(m + n)$.

**LCS Recursive:**   Since the algorithm does a depth first search of the search space, a search will go all the way to a leaf node (base case) before backtracking back up the call tree. Each call to the function needs the two input texts that get smaller as the search gets closer to the base case. In terms of space used in addition to the input texts, the space complexity is constant. However, if a new copy of each input text is made for each recursive call, the space complexity becomes $\mathcal{O}(n^2 + m^2)$ in the worst case, which is when a leaf node is found in the search space and both input texts have been reduced to a length of 2 or less:

$$
\begin{aligned}
S(m, n) &= (m + (m - 1) + ... + 1) + (n + (n - 1) + ... + 1) \\
&= \frac{(m^2 + m)}{2} + \frac{(n^2 + n)}{2} \\
&= \mathcal{O}(m^2) + \mathcal{O}(n^2) \\
&= \mathcal{O}(m^2 + n^2)
\end{aligned}
$$

However, if pointers are used and new copies of the input texts are not made for each recursive call, then the space complexity become linear with respect to the input texts lengths.

**LCS Branch and Bound:**   The space complexity for the the branch and bound extension is the same as recursive algorithm, except that a Queue of possible solutions to explore is also kept, which could potentially get quite large.

**Printing Neatly Greedy:**   There are no additional matrices or arrays used to compute the greedy solution, as the input text is the only array stored. Thus, the space needed is linear, $\mathcal{O}(n)$ where $n$ is the number of words in the input text.

**Printing Neatly Dynamic Programming:**   The dynamic programming algorithm uses two matrices (one for extras, one for line cost) of size $n \times n$, where $n$ is the number of words in the input text. Additionally, an array of size $n$ is used to compute the optimal solution. Thus, the asymptotic space complexity is $\mathcal{O}(n^2)$.

**Printing Neatly Recursive:**   There are no additional matrices or arrays used for the plain recursive solution. Thus, as long as pointers are used instead of copying a new input text array for each recursive call, the space complexity is $\mathcal{O}(n)$. If copies are made, then the space complexity would be $\mathcal{O}(n^2)$ (see LCS recursive space complexity analysis).

**Printing Neatly Branch and Bound:**   The space complexity is that of the recursive solution plus an array of possible solutions to explore. This array is created at each level of recursion, and contains at most every word in the input. Since there are at most $n - 1$ levels of recursion (in the case where there is just one word per line), the worst case space complexity would be $\mathcal{O}(n^2)$.

## 5.4    Outliers

If we look again at the Figure 5: Graph by task, we can see that there is three 'cut' texts that have a LCS ratio less than 70% (in task b and c). After looking each files (for example `g4pD_taskb.txt` with `orig_taskb.txt`), we concluded that this is an error by the original classification which was given to us.

## 5.5    Scalability

**LCS:**   For input texts that are very large, it is clear that the LCS Divide and Conquer strategy would scale the best out of all of them. It runs in pseudo-polynomial time which is much faster than the recursive or branch and bound solutions. Also, the space improvement over the classic dynamic programming algorithm would be needed when the input texts are very large.

**Printing Neatly:**   If an input text is very large, the recursive and branch and bound solutions would take too long. Furthermore, since the dynamic programming approach takes polynomial space and time, it would be unsuitable for extremely large input texts. The only solution left is the greedy solution, which takes linear space and time. It offers the best chance of scaling reasonably with large input texts. However, the solution is not guaranteed to be optimal. If an optimal solution was needed, the dynamic programming algorithm would offer the best scalability.

# 6    Plagiarism Detector

Several techniques and methods were explored in order to find the most suitable method of evaluating whether a text is plagiarised or not. Initial attempts were based of a ratio between size of the corpus text and the size of the LCS text, but these were found to be error prone and susceptible to misclassification due to the number of stop words in the texts. Removing the stop words did improve this, but did not take into account how group together words were in the provided text. Our final method of detecting plagiarism uses the LCS that was found between two texts, and combines two types of comparisons between the LCS and the text in question to determine if the text was plagiarized or not.

First, we compute how many words in the LCS appear next to each other in the original text. This is done because texts that are copied are more likely to have copied words next to each other, as opposed to texts that aren't copied and just have similar words. Our function computes the square of the number of words that are adjacent to each other (one word alone is 1, two words next to each other would be 4), divides this number of by the square of the length of the text to obtain a ratio (a ratio of 1 would mean cut and paste) and then returns the square root of this number.

Second, we compute the ratio of sentences in the text that have a percentage of plagiarism higher than a certain threshold. This threshold is an argument to the function and can change. We have found that an optimal threshold is 70%, meaning this function returns the ratio of sentences in the text that are more than 70% plagiarized. However, the user can change the threshold depending on how much time they have to check a text, for example. A lower threshold will result in higher plagiarism scores, and more texts to examine for potential plagiarism.

Finally, we combine these two results to form the suggestion on if the text is plagiarized or not. First, we compute a combined threshold for the sum of both metrics. We have found

through experimentation that all the texts in the corpus that are not plagiarized have combined scores of less than 0.2 out of 1. Thus, we set our threshold sum as 0.2.

$\rightarrow$ `Figure 14 [Appendix]`

Next, we compute the combined plagiarism score for a text by adding the first metric to the second. We then return whether or not this combined score is higher than the combined score threshold.

# 7   Project planning & difficulties

**Planning:**  In general the project went mostly to plan, but there were a few notable changes to the deadlines. The initial plan did not contain either the LCS recursive or the LCS space efficient algorithms. The LCS dynamic backwards took longer than was expected to implement, in order to stay on schedule, another team-member worked on the LCS branch and bound algorithm. The user interface was continually improved beyond the its original dead, as team-members had new ideas of features that could be incorporated into its design. The more advanced pre-processing was done earlier on than planned. Finalised the technique for plagiarism took longer than expected which meant there were delays to both the start of the report and the presentation.

**Difficulties:**  The main goal of coding the LCS algorithm was to implement what we saw during the course.
A hard aspect of the project was translating the pseudo-code into an implementation in Python. It required us to really understand what every part of algorithm was doing, as we had to implement the lower level tasks not included in the slides.
One example of a problem we ran into using Python to implement our algorithms was that Python is zero-indexed, so the first element in a list is indexed at 0. This was a problem because all the algorithms were implemented in the slides using one-indexed pseudo-code. We overcame this by having an empty element in the beginning of every list, to make them line up. Another difficulty we had was combining all of our work into one product. This was not a huge problem, but it did take some time to get all the functions working together in different parts of the software.
Lastly, it was difficult to create preprocessing techniques for the input texts. Our goal was to make the LCS of the preprocessed texts more correlated to the amount of plagiarism. We noticed that some of the corpus texts that were marked as highly plagiarized had short LCS. Thus, we investigated possible reasons for this happening and tried to augment the input texts to make the LCS more revealing. This required a lot of brainstorming and experimenting, as we had to try a lot of different things and collect a lot of data. It turned out to be one of the more time-consuming aspects of the project.

# 8    Conclusion

To achieve the goals set out for this project we decided to split it in different parts for each member of the group. Some had to implement LCS algorithms, others the printing neatly algorithm, some of us worked on the space and time complexity analysis, and finally some on the UI. The precise schedule ensured that we achieved what initially set out to do. Even if small changes were made to it, we tried as much as possible to adhere to the original one to in order to stay on track for the final deadline. To conclude we managed to implement all the LCS algorithms (Classic Forwards and backwards, Linear Space Forwards and Backwards, Divide and Conquer, Recursive and Branch and Bound), and all the Printing Neatly Algorithms, we were able to create a complete and user friendly UI, create different kinds of pre-processing and experiment with and develop techniques for the plagiarism score. Thanks to all of this work, we were able to compare and study them, to improve them and to know which one is the most efficient to use in practice and for our UI. To sum up every member of our group were very involved in the project. For some of us the language Python was unknown, and this gave many of us the opportunity to learn a new language or expand our knowledge of Python.
All of the project team are satisfied with the level of cooperation within the group and the final result that has been achieved.

# 9    References

[1] Gaofeng Huang and Andrew Lim. *An Effective Branch-and-Bound Algorithm to Solve the k-Longest Common Subsequence Problem*

# Appendices

## List of Algorithms

## List of Figures

---

**Algorithm 1:** Optimal substructure

---

**Data:** texts $X, Y$ with lengths $m, n$

1 **for** *an LCS $Z = < z_1, \cdots, z_n >$* **do**
2      **if** $X_m = Y_n$ **then**
3          $X_m = Y_n = z_k$;
4          $z_{k-1}$ is the LCS of $X_{m-1}$ and $Y_{n-1}$;

5 **if** $X_m \neq Y_n$ **then**
6      $z_k \neq X_m \Rightarrow Z$ is the LCS of $X_{m-1}$ and $Y_n$;
7 **if** $X_m \neq Y_n$ **then**
8      $z_k \neq Y_n \Rightarrow Z$ is the LCS of $X_m$ and $Y_{n-1}$;

---

**Algorithm 2:** Recursive solution

---

**Data:** texts $X, Y$ with lengths $m, n$

1 **if** $i = 0$ *or* $j = 0$ **then**
2      $c[i, j] = 0$;
3 **if** $i, j > 0$ *and* $X_i = Y_j$ **then**
4      $c[i, j] = c[i - 1, j - 1] + 1$;
5 **if** $i, j > 0$ *and* $xi \neq yj$ **then**
6      $c[i, j] = \max(c[i - 1, j], c[i, j - 1])$

---

---

**Algorithm 3:** LCS-Length(input: $X, Y$)

---

**1**   $m = |X|$;

**2**   $n = |Y|$;

**3**   $b[1..m, 1..n], c[0..m, 0..n]$;

**4**   **for** $i \in 1..m$ **do** $c[i, 0] = 0$;

**5**   **for** $j \in 0..n$ **do** $c[0, j] = 0$;

**6**   **for** $i \in 1..m$ **do**

**7**     **for** $j \in 1..n$ **do**

**8**       **if** $xi = yj$ **then**

**9**         $c[i, j] = c[i - 1, j - 1] + 1$;

**10**        $b[i, j] =' d'$ #for diagonal;

**11**       **else if** $c[i - 1, j] \geq c[i, j - 1]$ **then**

**12**         $c[i, j] = c[i - 1, j]$;

**13**         $b[i, j] =' u'$ #for up;

**14**       **else**

**15**         $c[i, j] = c[i, j - 1]$;

**16**         $b[i, j] =' l'$ #for left;

**17**   **return** $c$ *and* $b$

---

**Algorithm 4:** Print-LCS(input: $b, X, |X|, |Y|$)

---

**1**   **if** $i = 0$ *or* $j = 0$ **then**

**2**     return;

**3**   **if** $b[i, j] = 'd'$ **then**

**4**     Print-LCS($b, X, i - 1, j - 1$)

**5**   **else if** $b[i, j] = 'u'$ **then**

**6**     Print-LCS($b, X, i - 1, j$)

**7**   **else**

**8**     Print-LCS($b, X, i, j - 1$)

---

---

**Algorithm 5:** LCS_LSF (input: $X, Y$)

---

**1** $m = |X|$;
**2** $n = |Y|$;
**3** $c[1..2, 1..n]$;
**4** col = [] #for use in divide and conquer alg (see next section);
**5** **for** $i \in 1..m$ **do**
**6** $\quad$ $c[1, 0] = 0$;
**7** $\quad$ **for** $j \in 1..n$ **do**
**8** $\quad\quad$ **if** $xi = yj$ **then**
**9** $\quad\quad\quad$ $c[1, j] = c[0, j-1] + 1$;
**10** $\quad\quad$ **else if** $c[1, j-1] > c[0, j]$ **then**
**11** $\quad\quad\quad$ $c[1, j] = c[1, j-1]$;
**12** $\quad\quad$ **else**
**13** $\quad\quad\quad$ $c[1, j] = c[0, j]$;

**14** $\quad$ $c[0, :] = c[1, :]$;
**15** $\quad$ col.append($c[1, n]$);
**16** **return** $c[1, length_Y], col$ #return length of LCS, column

---

---

**Algorithm 6:** LCS_LSB(input: $X, Y$)

---

**1** $m = |X|$;
**2** $n = |Y|$;
**3** $c[1..2, 1..n]$;
**4** col = [] #for use in divide and conquer alg (see next section);
**5** **for** $i \in m - 1..0$ **do**
**6** $\quad$ **for** $j \in n - 1..0$ **do**
**7** $\quad\quad$ **if** $X[i+1] = Y[j+1]$ **then**
**8** $\quad\quad\quad$ $c[1, j] = c[0, j+1] + 1$;
**9** $\quad\quad$ **else if** $c[1, j+1] > c[0, j]$ **then**
**10** $\quad\quad\quad$ $c[1, j] = c[1, j+1]$;
**11** $\quad\quad$ **else**
**12** $\quad\quad\quad$ $c[1, j] = c[0, j]$;

**13** $\quad$ $c[0, :] = c[1, :]$;
**14** $\quad$ col.insert($0, c[1, 0]$);
**15** **return** $c[1, 0], col$ #length of LCS, column

---

---

**Algorithm 7:** LCS_DC(input: $X, Y$)

**1**   $m = |X|$;

**2**   $n = |Y|$;

**3**   **if** $m * n = 0$ **then**

**4**     **return** [] #if one of the texts is empty, return an empty LCS

**5**   **else if** $m \leq 2$ *or* $n \leq 2$ **then**

**6**     **return** *LCS_DyProg(X,Y)* #if one of the texts is small, can do normal
     Dynamic Prog

**7**   **else**

**8**     breakpt = floor($n/2$);

**9**     length, $c$ = LCS_LSF($X, Y[0 : \text{breakpt}]$);

**10**    length2, $c2$ = LCSLSB($X, Y[\text{breakpt} : n]$);

**11**    $q = \text{maxindex}(c, c2)$ #q is the index that maximizes $(c[q] + c2[q])$;

**12**    LCSL = LCS_DC($X[0..q], Y[0..\text{breakpt}]$);

**13**    LCSR = LCS_DC($X[q + 1, m], Y[\text{breakpt} + 1, n]$);

**14**    **return** $LCSL + LCSR$ #return concatenation of sub LCS

---

---

**Algorithm 8:** LCS_Recursive(input: $X, Y$)

**1**   $m = |X|$;

**2**   $n = |Y|$;

**3**   **if** $X[m] = Y[n]$ **then**

**4**     **return** *LCS_Recursive(X[0 : m − 1], Y[0 : n − 1])*;

**5**   **else**

**6**     LCS1 = LCS_Recursive($X[0 : m − 1], Y[0 : n]$);

**7**     LCS2 = LCS_Recursive($X[0 : m], Y[0 : n − 1]$);

**8**     **if** $|LCS1| > |LCS2|$ **then** **return** *LCS1* ;

**9**     **else** **return** *LCS2* ;

---

---

**Algorithm 9:** Brand_and_bound(inputs: $X, Y$)

---

**1** max_lcs_length = min($X$.length, $Y$.length);

**2** best = find_lower_bound($X, Y$) – 1;

**3** starting_solution = (1 << max_lcs_length) – 1;

**4** best_solution = 0;

**5** index = -1;

**6** tracker = -1;

**7** LIFO_queue.put([starting_solution, index, tracker, best]);

**8** **while** *LIFO_queue is not empty* **do**

**9**     node = LIFO_queue.get();

**10**     **if** *node.index = max_lcs_length and node.best > best* **then**

**11**        best = node.best;

**12**        best_solution = node.solution;

**13**     **else**

**14**        **if** *node.bound > best* **then**

**15**           child1 = node.child1 child2 = node.child2

**16**        **if** *find_upper_bound(child1) > find_upper_bound(child2)* **then**

**17**           **if** *child2.bound > best* **then** LIFO_queue.add(child2);

**18**           **if** *child1.bound > best* **then** LIFO_queue.add(child1);

**19**        **else**

**20**           **if** *child1.bound > best* **then** LIFO_queue.add(child1);

**21**           **if** *child2.bound > best* **then** LIFO_queue.add(child2);

**22** **return** *best_solution*;

---

**Algorithm 10:** PN_Greedy(input: X (list of words to be printed), L (number of spaces per line))

---

**1** lines = [];

**2** currentLine = [];

**3** currentLineSpace = 0;

**4** **for** *i in X* **do**

**5**     **if** *currentLineSpace + ($|X[i]| + 1$) ≤ L* **then**

**6**        currentLineSpace += $|X[i]| + 1$ #update space taken on current line (for space between words);

**7**        currentLine += " " + $X[i]$;

**8**     **else**

**9**        lines.append(currentLine);

**10**        currentLineSpace = $|X[i]|$;

**11**        currentLine = $X[i]$;

**12** lines.append(currentLine);

**13** **return** *lines*;

---

---

**Algorithm 11:** PN_DP(input: $X, L$)

---

**1** Compute extras$[i, j], 1 = i, j = n$;

**2** Compute lc$[i, j], 1 = i, j = n$;

**3** $c[0] = 0$;

**4 for** $j \in 1..n$ **do**

**5**      $c[j] = \min((1 \leq i \leq j), (c[i-1] + lc[i, j]))$;

**6**      $p[j] = k$ s.t. $\min(1 \leq i \leq j)(c[i-1] + lc[i, j]) = c[k-1] + lc[k, j]$ #store where the cuts are;

**7 return** $p$

---

---

**Algorithm 12:** PN_Recursive(input: $X, L$)

---

**1** result = $[]$ #where the lines will be stored;

**2** max_words = the maximum number of words that can fit on current line #what the greedy solution would return;

**3 if** $max\_words \geq |X|$ **then**

**4**      **return** $0, X[0 : |X|]$ #base case: where all words inputed will fit on current line;

**5** best_score = $\infty$;

**6 for** $i \in 1..max\_words$ **do**

**7**      subscore = 0 #current line space taken;

**8**      **for** $j \in 0..i-1$ **do**

**9**          subscore += $|X[j]| + 1$ #add current word and space;

**10**      subscore -= 1 #remove last space;

**11**      subscore = $(L - \text{subscore})^3$ #compute cost of space left over;

**12**      score, res = PN_Recursive($X[i : n], L$);

**13**      **if** $score + subscore ¡ best\_score$ **then**

**14**          result.append($X[0..i]$);

**15**          result += res;

**16 return** *best_score, result*

---

Figure 1: Path to the solution of the LCS of [A,B,C,A] and [B,C,A,A]

The path explored first shown as a dashed line, the super-script numbers at the top right show the score and upper bounds computed at each stage.
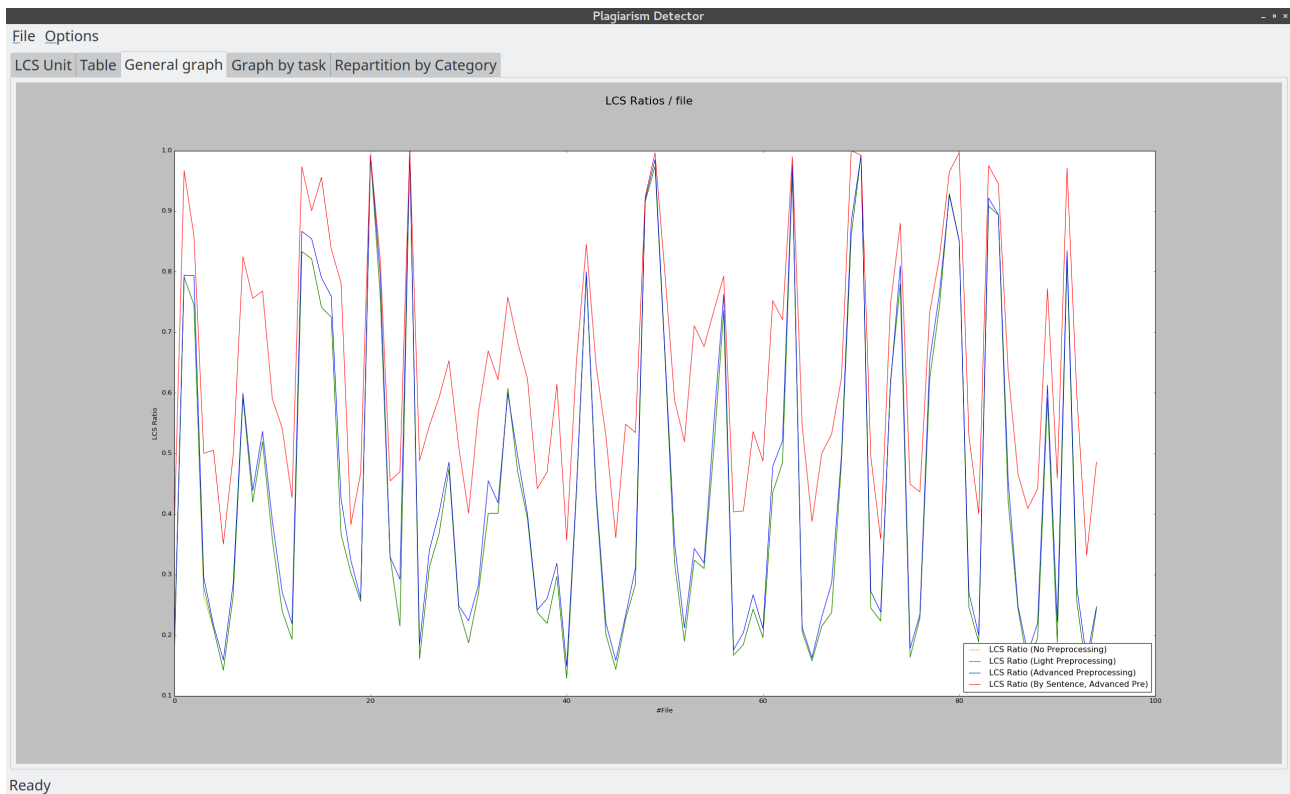
Figure 2: LCS Unit
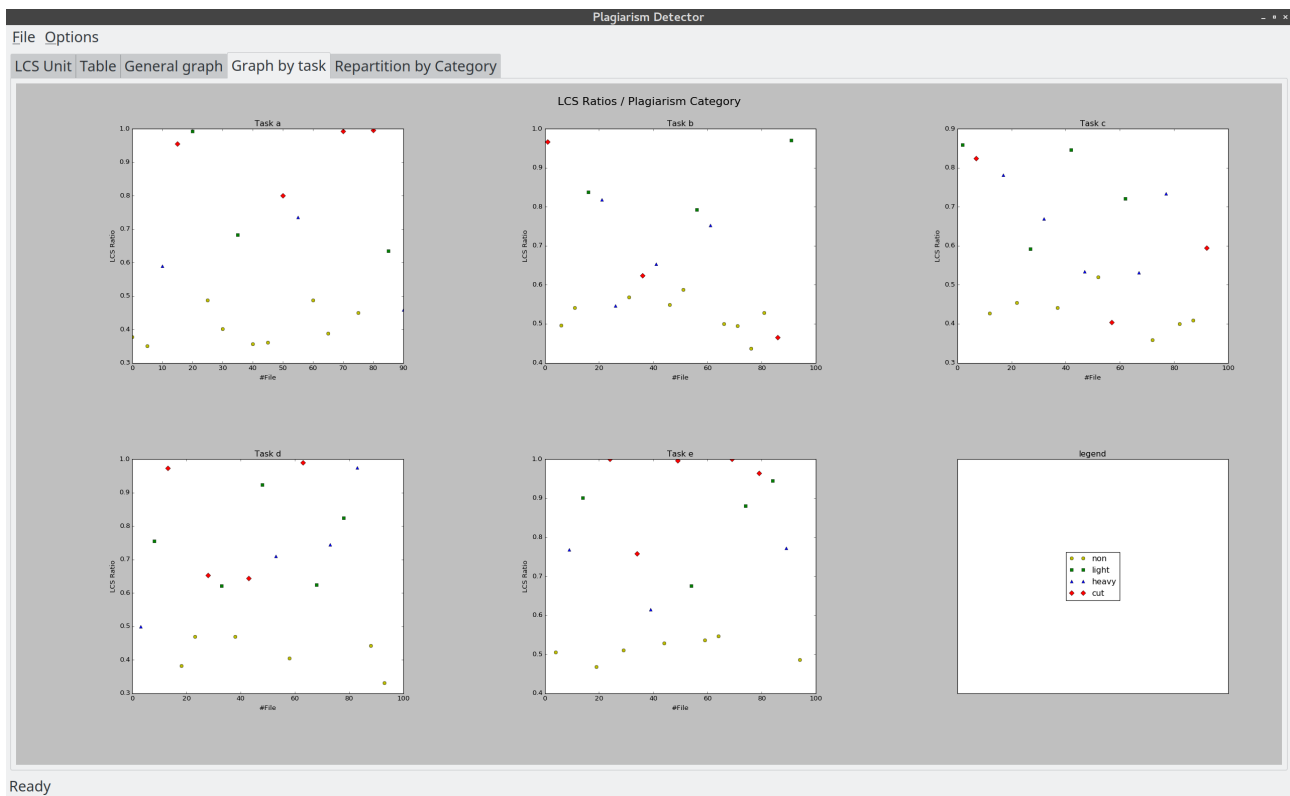
Figure 3: Table

Figure 4: General graph
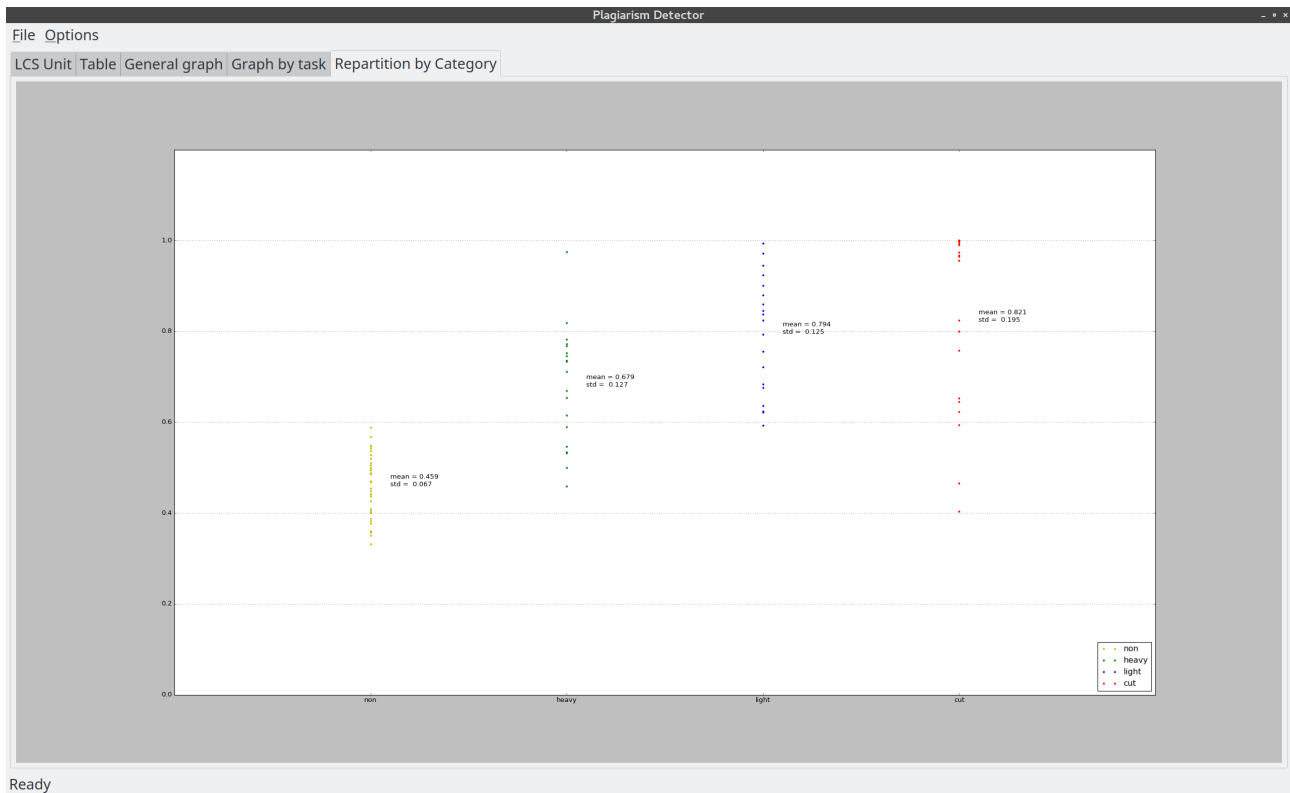


Figure 5: Graph by task

Figure 6: Repartition by Category
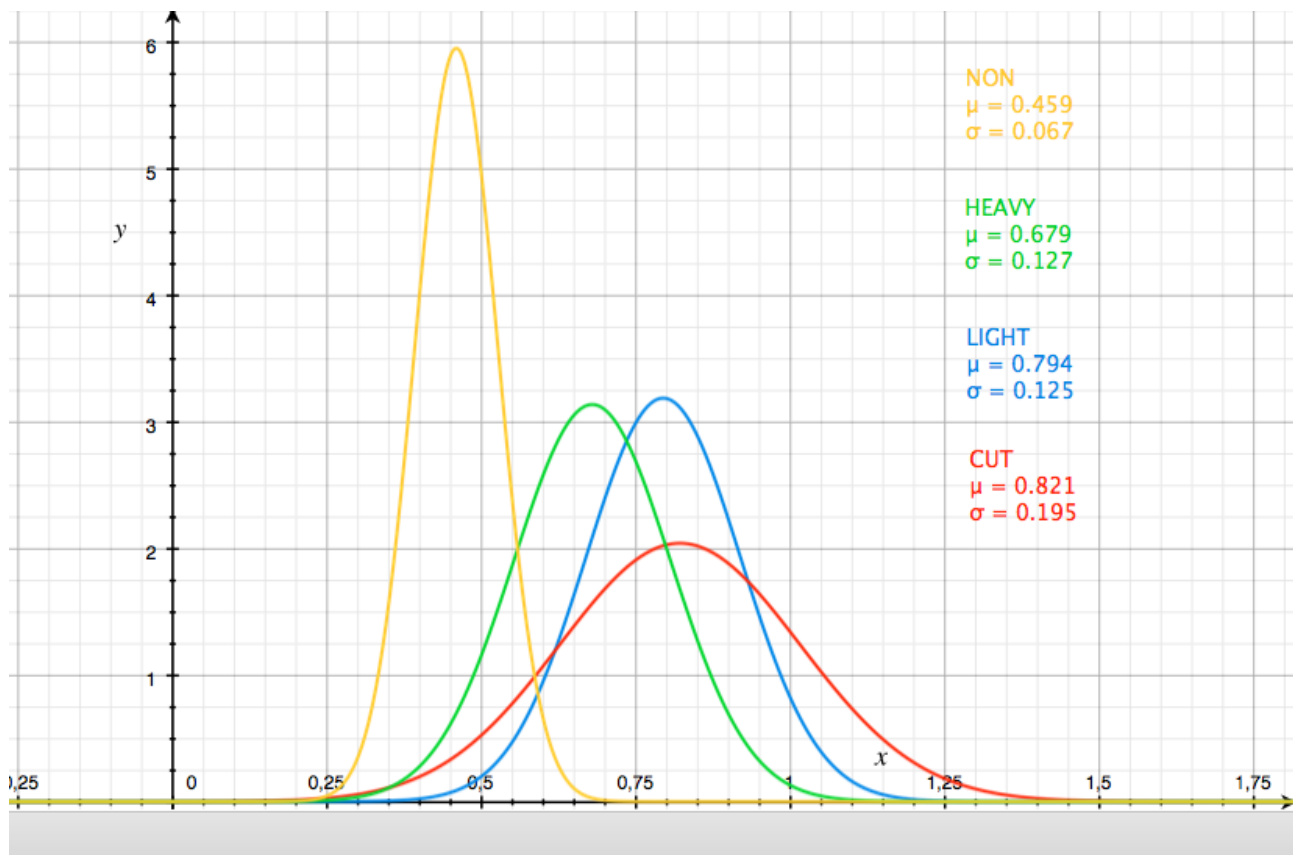


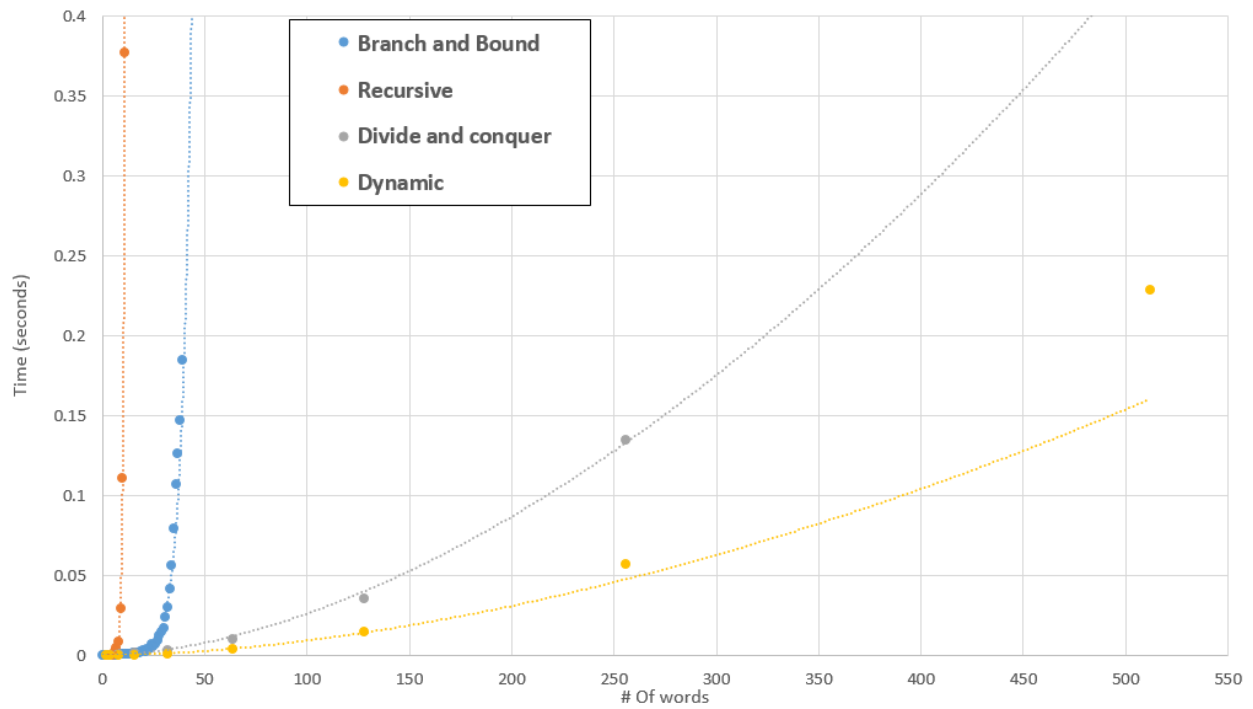Figure 7: Normal distributions of the repartition by category / |words|

Figure 8: Running times for different algorithms on linear scale / |words|



Figure 9: Running times for different algorithms on logarithmic scale / |words|

Figure 10: Running times for different algorithms on logarithmic scale / $(|words|)^2$



Figure 11: Comparison between Printing neatly algorithms, linear scale.

Figure 12: Comparison between Printing neatly algorithms with a log score, showing the exponential runtimes of the recursive and branch and bound (optimal) algorithms.



Figure 13: Comparison between the Printing neatly Dynamic, Branch & Bound (non-optimal) and greedy algorithms

Figure 14: The final scores of the plagiarism detector. Excluding outliers, a good separation is observed between the non-plagiarised texts (yellow).



Figure 15: Planning: the plan at the start of the project.

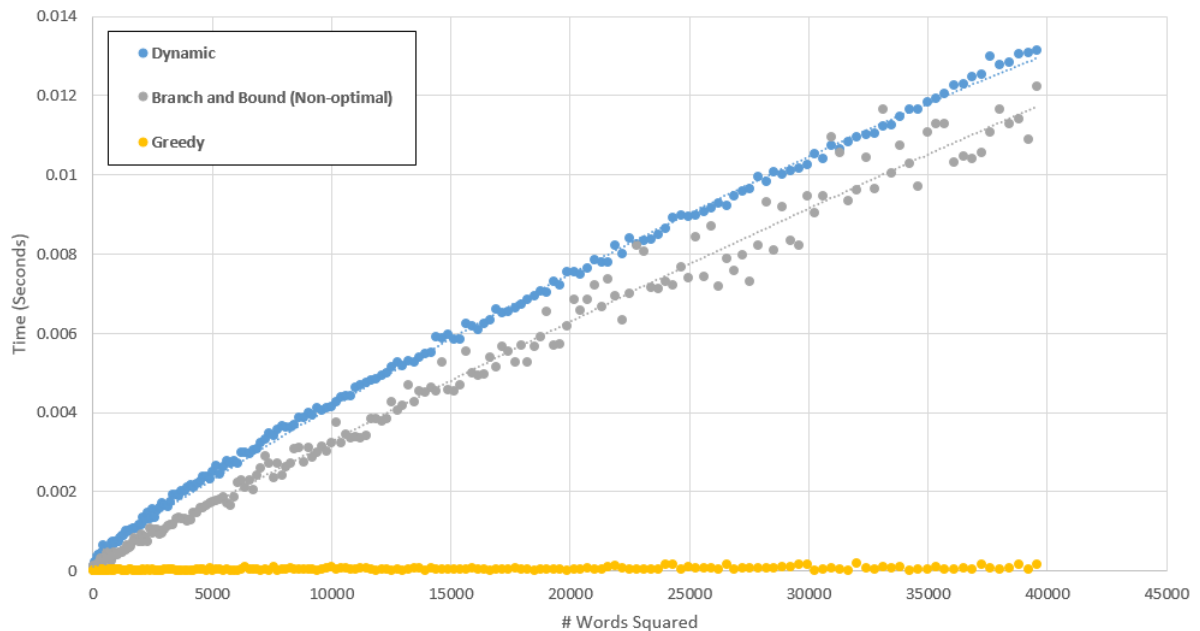| Activity | Person | ACTUAL START | ACTUAL DURATION | | | |
|---|---|---|---|---|---|---|
| Text Pre-Processing | Anthony | 25-Oct | 7 | | | |
| LCS - Classic Forwards | Joe | 18-Oct | 5 | | | |
| LCS - Classic Backwards | Jeremy | 09-Nov | 7 | | | |
| LCS - Space Efficient | Jeremy | 11-Nov | 4 | | | |
| LCS - Divide and Conquer | Joe | 09-Nov | 7 | | | |
| LCS - Recursive | Ed | 10-Nov | 1 | | | |
| LCS - B & B | Ed | 10-Nov | 7 | | Milestone: All LCS Algo. complete | |
| Printing neatly | Ed | 18-Oct | 3 | | | |
| Time and space complexity tests | Ed / Joe | 25-Nov | 12 | | | |
| Comparisons on corpus | Josselin / Joe | 05-Nov | 10 | | | |
| Technique for plagiarism score | Jeremy / Ed | 20-Nov | 20 | | Milestone : Technique for pl | |
| User interface :- Planning | Ed / Anthony | 20-Oct | 10 | | | |
| User interface :- Development | Anthony / Josselin | 07-Nov | 18 | | Milestone : UI complete | |
| Optional: Advanced Pre-Processing | Anthony | 01-Nov | 5 | | | |
| Optional: Advanced Post-Processing | Josselin | 15-Nov | 5 | | | |
| Report | Joe / Josselin / Anthony | 06-Dec | 6 | | | |
| Presentation for defence | Ed / Jeremy | 06-Dec | 6 | | | |

Figure 16: Planning: the actual dates.