

Optimal alignment and plagiarism  
detection software

# Code Execution Manual

---

Diego REYNA REYES & Ryan VAN GREUNEN

January 15th, 2023

## Delivered files

The delivered .zip file contains the following files:

**./TDm2\_REYNAREYES\_VANGREUNEN/documents/report.pdf:** The report containing the answers to questions 1, 2 and 4 of the [instruction website](#).

**./TDm2\_REYNAREYES\_VANGREUNEN/documents/manual.pdf:** This manual explaining how to run and test our code for questions 3 and 4.

**./TDm2\_REYNAREYES\_VANGREUNEN/code/q3/main\_q3.c:** Source code for our implementation of question 3.

**./TDm2\_REYNAREYES\_VANGREUNEN/code/q3/TD2\_q3.c:** A modified version of the [TD2.c](#) file shared on the instructions.

**./TDm2\_REYNAREYES\_VANGREUNEN/code/q3/texte1.txt:** The .txt file containing [texte1](#) as shared with the instructions.

**./TDm2\_REYNAREYES\_VANGREUNEN/code/q3/texte2.txt:** The .txt file containing [texte2](#) as shared with the instructions.

**./TDm2\_REYNAREYES\_VANGREUNEN/code/q4/main\_q4.c:** Source code for our implementation of question 4.

**./TDm2\_REYNAREYES\_VANGREUNEN/code/q4/TD2\_q4.c:** A modified version of the [TD2.c](#) file shared on the instructions. It's important to note that this file is slightly different to the one found in ./TDm2\_REYNAREYES\_VANGREUNEN/code/q3/TD2\_q3.c.

**./TDm2\_REYNAREYES\_VANGREUNEN/code/q4/t1.txt:** The .txt file containing [t1](#) as shared with the instructions.

**./TDm2\_REYNAREYES\_VANGREUNEN/code/q4/t2.txt:** The .txt file containing [t2](#) as shared with the instructions.

## Compiler and directory

**Compiler:** This code was tested and compiled using the following compilers:

- gcc (Ubuntu 7.5.0-3ubuntu1-18.04) 7.5.0

```
gcc (Ubuntu 7.5.0-3ubuntu1-18.04) 7.5.0
Copyright (C) 2017 Free Software Foundation, Inc.
This is free software; see the source for copying conditions. There is NO
warranty; not even for MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.
```

- gcc-10 (Ubuntu 10.3.0-1ubuntu1-18.04-1) 10.3.0

```
gcc-10 (Ubuntu 10.3.0-1ubuntu1-18.04-1) 10.3.0
Copyright (C) 2020 Free Software Foundation, Inc.
This is free software; see the source for copying conditions. There is NO
warranty; not even for MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.
```

Given this is the most recent version of the used compiler that works on Linux, this version will be used in all the following instructions.

- gcc (Rev4, Built by MSYS2 project) 12.2.0

```
gcc (Rev4, Built by MSYS2 project) 12.2.0
Copyright (C) 2022 Free Software Foundation, Inc.
This is free software; see the source for copying conditions. There is NO
warranty; not even for MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.
```

**Directory:** Expecting the provided .zip file was extracted on folder ./

## Compiling and testing question 3

**Compiling command:** To compile the code written for this question we use the following command:

```
gcc-10 ./TDM2_REYNAREYES_VANGREUNEN/code/q3/main_q3.c -o ./q3_exe
```

The structure of the command is:

```
gcc-10 path/to/main_q3.c -o path/to/output/q3_exe
```

**Testing command:** To test the code written for this question we use the following command:

```
./q3_exe ./TDM2_REYNAREYES_VANGREUNEN/code/q3/texte1.txt ./TDM2_REYNAREYES_VANGREUNEN/code/q3/texte2.txt
```

The structure of the command is:

```
path/to/q3_exe path/to/txt_file1 path/to/txt_file2
```

When ran on the provided texte1.txt and texte2.txt we expect the following output:

```
zuramaru@zuramaru-VirtualBox:~/Documents/INF4202B$ gcc-10 ./TDM2_REYNAREYES_VANGREUNEN/code/q3/main_q3.c -o ./q3_exe
zuramaru@zuramaru-VirtualBox:~/Documents/INF4202B$ ./q3_exe ./TDM2_REYNAREYES_VANGREUNEN/code/q3/texte1.txt ./TDM2_REYNAREYES_VANGREUNEN/code/q3/texte2.txt
Score of alignment: 577.
Size of text 1: 1520
Size of text 2: 1799
Similarity score: 82.62 %
-----
Source Wikipedia. La distance de Levenshtein est une distance mathématique donnant une mesure de la similarité entre deux chaînes de caractères. Elle est égale au nombre minimal de caractères qu'il faut supprimer, insérer ou remplacer pour passer d'une chaîne à l'autre. Elle a été proposée par Vladimir Levenshtein en 1965. Elle est également connue sous les noms de distance d'édition ou de distance de déformation dynamique temporelle, notamment en reconnaissance de formes et particulièrement en reconnaissance vocale. Cette distance est d'autant plus grande que le nombre de différences entre les deux chaînes est grand. La distance de Levenshtein peut être considérée comme une généralisation de la distance de Hamming.
-----
Source Wikipedia modifiée par un étudiant du cours IT-4301E, traitement algorithmique de l'information. La distance de Levenshtein est une distance au sens mathématique donnant une mesure de la similarité entre deux séquences de caractères. Elle est égale au nombre minimal de caractères qu'il faut supprimer, insérer ou substituer pour passer d'une séquence à l'autre. Elle a été proposée par Vladimir Levenshtein en 1965. Elle est ainsi également connue sous les noms de distance de Levenshtein ou de distance de déformation dynamique temporelle dans la reconnaissance de formes. Cette distance est une fonction croissante du nombre de différences entre les deux séquences. La distance d'édition peut être considérée comme une généralisation de la distance de Hamming (donnée par le nombre de positions en lesquelles les deux séquences possèdent des caractères différents). On peut montrer en particulier que la distance de Hamming est un majorant de la distance d'édition. Définition formelle : on appelle distance d'édition entre deux mots M et P le coût minimal transformant M en P en effectuant les opérations élémentaires, dites d'édition, suivantes : i) substitution d'un caractère de M par un caractère de P ; ii) ajout d'un caractère de P ; iii) suppression d'un caractère de M. On associe ainsi à chacune de ces opérations un coût. Le coût est toujours égal à 1, sauf dans le cas d'une substitution de caractères identiques. Exemples : si M = "examen" et P = "examen", alors LD(M, P) = 0, parce qu'aucune opération n'a été réalisée. Si M = "examen" et P = "examan", alors LD(M, P) = 1, parce qu'il y a eu un remplacement (changement de e en a), et que l'on ne peut pas en faire un moindre coût.
-----
```

## Compiling and testing question 4

**Compiling command:** To compile the code written for this question we use the following command:

```
gcc-10 ./Tdm2_REYNAREYES_VANGREUNEN/code/q4/main_q4.c -o ./q4_exe
```

The structure of the command is:

```
gcc-10 path/to/main_q4.c -o path/to/output/q4_exe
```

**Testing command:** To test the code written for this question we use the following command:

```
./q4_exe ./Tdm2_REYNAREYES_VANGREUNEN/code/q4/t1.txt ./Tdm2_REYNAREYES_VANGREUNEN/code/q4/t2.txt
```

The structure of the command is:

```
path/to/q4_exe path/to/txt_file1 path/to/txt_file2
```

When ran on the provided t1.txt and t2.txt we expect the following output:

```
zuramaru@zuramaru-VirtualBox:~/Documents/INF4202B$ gcc-10 ./Tdm2_REYNAREYES_VANGREUNEN/code/q4/main_q4.c -o ./q4_exe
zuramaru@zuramaru-VirtualBox:~/Documents/INF4202B$ ./q4_exe ./Tdm2_REYNAREYES_VANGREUNEN/code/q4/t1.txt ./Tdm2_REYNAREYES_VANGREUNEN/code/q4/t2.txt
Score of alignment: 789.
Size of text 1: 1525
Size of text 2: 2020
Similarity score: 77.74 %

-----
Source Wikipedia
-----
La distance de Levenshtein une distance mathématique donnant une mesure de la similarité entre deux chaînes de caractères. Elle est égale au nombre minimal de caractères qu'il faut supprimer, insérer ou remplacer pour passer d'une chaîne à l'autre. Elle a été proposée par Vladimir Levenshtein en 1965. Elle est également connue sous les noms de distance d'édition ou de distance de déformation dynamique temporelle, notamment en reconnaissance de formes et particulièrement en reconnaissance vocale [2]. Cette distance est d'autant plus grande que le nombre de différences entre les deux chaînes est grand. La distance de Levenshtein peut être considérée comme une généralisation de la distance de Hamming.

On peut montrer en particulier que la distance de Hamming est un majorant de la distance de Levenshtein.

-----
Source Wikipedia modifiée par un étudiant du cours IT-4301E, traitement algorithmique de l'information.
-----
La distance d'édition est une distance au sens mathématique donnant une mesure de la similarité entre deux séquences. Elle est égale au nombre minimal de caractères qu'il faut supprimer, insérer ou substituer pour passer d'une séquence à l'autre. Elle a été proposée par Vladimir Levenshtein en 1965. Elle est ainsi également connue sous les noms de distance de Levenshtein ou de distance de déformation dynamique temporelle dans le domaine de la reconnaissance de formes. Cette distance est une fonction croissante du nombre de différences entre les deux séquences. La distance d'édition peut être considérée comme une généralisation de la distance de Hamming (donnée par le nombre de positions dans lesquelles les deux séquences possèdent des caractères différents). On peut montrer en particulier que la distance de Hamming est un majorant de la distance d'édition.

-----
Et pourquoi ne pas raconter des ballivernes entre temps pour détecter les lecteurs attentifs.

-----
Et insérer un paragraphe qui n'a rien à voir avec le chnilitic pour tromper le chaland !

-----
Definition : on appelle distance de Levenshtein entre deux mots M et P le coût minimal pour aller de M à P en effectuant les opérations élémentaires suivantes : i) substitution d'un caractère de M en un caractère de P ; ii) ajout dans M d'un caractère de P ; iii) suppression d'un caractère de M. On associe ainsi à chacune de ces opérations un coût. Le coût est toujours égal à 1, sauf dans le cas d'une substitution de caractères identiques.

Definition formelle : on appelle distance d'édition entre deux mots M et P le coût minimal transformant M en P en effectuant les opérations élémentaires, dites d'édition, suivantes : i) substitution d'un caractère de M par un caractère de P ; ii) insertion dans M d'un caractère de P ; iii) suppression (ou deletion) d'un caractère de M. On associe ainsi à chacune de ces opérations un coût. On choisit souvent un coût égal à 1 pour toutes les opérations excepté la substitution de caractères identiques qui a un coût nul.

Exemples : si M = "examen" et P = "examen", alors LD(M, P) = 0, parce qu'aucune opération n'a été réalisée. Si M = "examen" et P = "examan", alors LD(M, P) = 1, parce qu'il y a eu un remplacement (changement du e en a), et que l'on ne peut pas en faire moins.

Exemples : si M = "examen" et P = "examen", alors Lev(M, P) = 0 parce qu'aucune opération n'a été réalisée. Si M = "examen" et P = "examan", alors Lev(M, P) = 1, parce qu'il y a eu une substitution (changement du e en a), et que l'on ne peut pas en faire une transformation de M en P avec un moindre coût.

-----
Pas super complet ce cours...
```