Anurat Wongbunmak
2488941

**AC50002 - Programming languages for Data Engineering ( SEM 1 22/23 )**

**Python Assignment**

This link will take you to my GitHub repository, where you can access the code for this assignment.

https://github.com/EarnAnurat/PythonAsgmt1_Anurat2488941

The code is made up of five functions that work together to complete the task including readfile(fname), words(name), allabbr(wl) ,cal(wl,abbr) and main().



First, the **readfile()** function reads every line from the file into a list using the .readlines() method, where each line is an item in the list. As a result, it then returns the data "linelist", which is a list of every line.



Second, the **words()** function will get the stored data called "linelist" and turn it into a list of words that can be used as abbreviations. Replace "-" with a space to start, then use.strip() to eliminate the space between them. .split() is used to break apart words, and.upper() is used to change to uppercase. Last step, deleting all special characters with ASCII codes 65-90 that are letters "A" through "Z"



Third, to construct a collection of all potential abbreviations, the **allabbr()** function takes a list of words that are appropriate for abbreviations that we obtain from the previous function. The initial letter, is always be the range(1) which is the first letter of words. The second letters, range(x+1, lgth-1), are not the first and last letters of words. The third letter, range(y+1, lgth), will be one that comes from the rest of the characters that come after the second.



Fourth, the **cal()** function gives the abbreviation a score that represents its effectiveness and then returns the score with the lowest number. An abbreviation's overall score is the result of adding the scores for its second and third letters. These individual letter scores are based on the letter's position in the word and how common or uncommon it is in English.

Finally, **main()** function serves as a starting point for code that performs the primary purpose of the script. It will start execute other functions. Use **readfile()** to get the data name linelist which is a list where each line in .txt file is an item. Then create a list of the best abbreviations to be a final result with the help of **words(), allabbr()** and **cal()**. And finally write the output file.

if __name__ == "__main__": is to ensure that when the module is simply imported, the function is not called, Then if the modules are the main, the module can be run from an ordinary Windows command prompt.
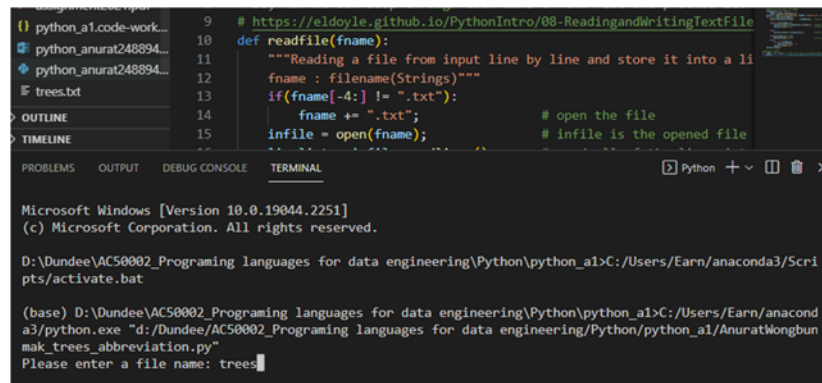
**Using the modules**

Input file



enter a file name: trees



Output file