

1. Write Python code and use MapReduct to count occurrences of each word in the first text file (file.txt). How many times each word is repeated?

MY DOB IS 24/SEP/1998

the month is september(09), it is divided by 2 and rounded to 5. the 5th book is Harry Potter and the order of phonex

the date is 24 so I have extracted the 10 pages starting from the 24th page (file1)

. the year is 98, so i have extracted the 10 pages starting from 98th (page file2)

```
In [1]: import re                                # importing the required packages
import string
from functools import reduce
from collections import Counter
```

```
In [2]: (text):
e.findall(r'\w+', text.lower()) # regular expression '\w+' takes whole words
                                # returning the list of words
```

```
In [3]: a open('file1.txt', 'rb') as file:        # Opening the 'file1.txt'
lines = file.readlines()                        # Reading all lines from the file
lines = [line.decode('utf-8') for line in lines]
text = ''.join(lines)                          # Joining the list of lines into a single string
words = map_words(text)                        # calling the 'map_words' function
```

```
In [4]: rd_counts_list = list(map(Counter, [words])) # Using the map function to create a list of Counters
total_word_counts = reduce(lambda x, y: x + y, word_counts_list) # Using 'reduce' to sum up the Counters

for word, count in total_word_counts.items(): # Iterating through the total word counts
    print(f'{word}: {count}')
```

```
of: 66
phoenix: 11
j: 11
k: 11
rowling: 11
barely: 2
an: 9
inch: 1
from: 13
```

```
In [6]: def map_words(text):
        return re.findall(r'\w+', text.lower()) # regular expression '\w+'

In [7]: def is_not_english_word(word, spellchecker):
        s: 3
        return not word.isdigit() and word not in spellchecker # Checking if the word is not a digit and not in the spellchecker
when: 5
silver: 2
antlers: 1
caught: 1
it: 29
thing: 2
was: 45
thrown: 1
checked: 1
```

```
In [8]: # Initializing a SpellChecker instance
spell_checker = SpellChecker.from_dict(word_counts_list) # Initializing a dictionary to store the word counts
# Iterating through each word in the list and checking if it is a non-English word
for word, count in word_counts_list.items():
    if not is_not_english_word(word, spell_checker):
        word_counts_list[word] = word_counts_list.get(word, 0) + 1 # If the word is not in the dictionary, add it with a count of 1
    else:
        word_counts_list[word] = word_counts_list.get(word, 0) # If the word is in the dictionary, keep the current count
```

```
In [5]: # Iterating through each word in the list and checking if it is a non-English word
for word, count in word_counts_list.items():
    if not is_not_english_word(word, spell_checker):
        word_counts_list[word] = word_counts_list.get(word, 0) + 1 # If the word is not in the dictionary, add it with a count of 1
    else:
        word_counts_list[word] = word_counts_list.get(word, 0) # If the word is in the dictionary, keep the current count
```

From the second text file (file2.txt), write Python code and use MapReduct to count how many times non-English words (names, places, spells etc.) were used. List those words and how many times each was repeated.

```
In [5]: # Iterating through each word in the list and checking if it is a non-English word
for word, count in word_counts_list.items():
    if not is_not_english_word(word, spell_checker):
        word_counts_list[word] = word_counts_list.get(word, 0) + 1 # If the word is not in the dictionary, add it with a count of 1
    else:
        word_counts_list[word] = word_counts_list.get(word, 0) # If the word is in the dictionary, keep the current count
```

```
p: 11
g: 11
...
```

```

In [6]: def map_words(text):
        an: 9
        return re.findall(r'\w+', text.lower()) # regular expression '\w+'
        inch: 1
        from: 13
In [7]: def is_non_english_word(word, spellchecker):
        s: 3
        return not word.isdigit() and word not in spellchecker # Checking
        when: 5
        silver: 2
        antlers: 1
        caught: 1
        it: 29

```

```

In [8]: thing: 2
        was: 45
        hyphen: 1
        checker: 1
        # Initializing a SpellChecker instance
        # Initializing a dictionary to store word counts

```

From the second text file (file2.txt), write Python code and use MapReduce to count how many times non-English words (names, places, spells etc.) were used. List those words and how many times each was repeated.

```

In [5]: # Iterating through each word in the list
        is_non_english_word(word, spellchecker) # Checking if the word is non-English
        word_counts[word] = word_counts.get(word, 0) + 1 # If the word is not in the dictionary,
        from functools import reduce
        from spellchecker import SpellChecker # Importing the required packages
        # Iterating through the 'word_counts' dictionary
        {count}')

```

```

p: 11
g: 11
j: 11
k: 11
rowling: 11
s: 40
irritably: 1
t: 12
heatedly: 1
kneacher: 1
c: 2
m: 5
mrs: 17
weasley: 21
tonks: 8
ll: 2
exasperatedly: 1
earsplitting: 1
yoooooooo: 1
mr: 4
greenland: 1
evanesco: 1
mundungus: 13
n: 1
sleepily: 1
gree: 1
balefully: 1
arry: 1
pology: 1
crookshanks: 2
absentmindedly: 1
incredulously: 1
d: 2
animagus: 1
wormtail: 1
bracingly: 1
butterbeer: 3
breadboard: 1

```

In []:

In []: