

## MapReduce Assignment

### 1.Introduction

Analyzing word occurrences in Harry Potter books using MapReduce.

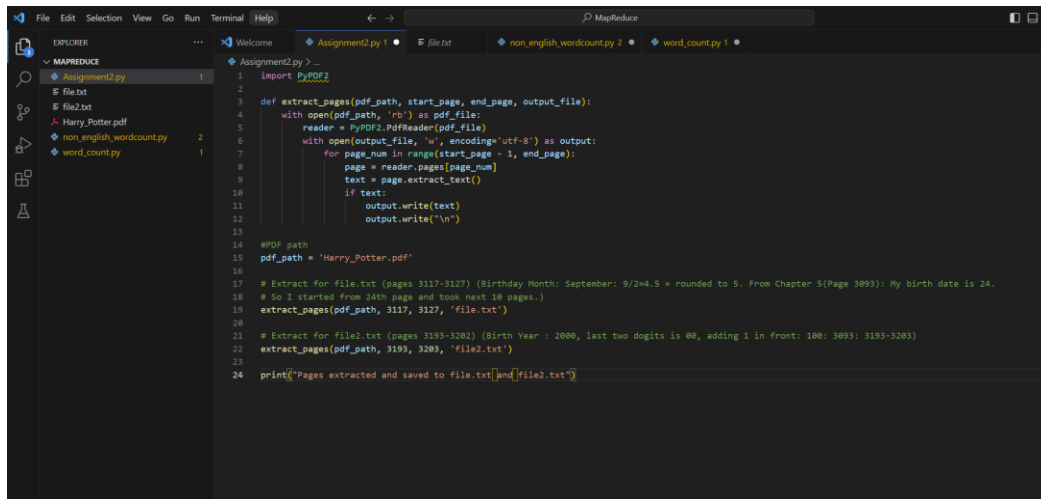
### 2. Data Acquisition (Screenshot Included)

- Harry Potter books PDF:  
[https://ztcprep.com/library/story/Harry\\_Potter/Harry\\_Potter\\_\(www.ztcprep.com\).pdf](https://ztcprep.com/library/story/Harry_Potter/Harry_Potter_(www.ztcprep.com).pdf)
- Calculating page ranges for text extraction based on birth date: **September 24,2000**
  - Divide birth month by 2 and round up for months 8-12 (September:  $9 / 2 = 4.5$ , rounded to 5).
  - Started from Chapter 5 (Page 3093) and took the next 10 pages for my birth date ( September 24: pages 3117-3127).
  - For the birth year (last two digits), prepend "1" for years 2000 and later (2000: 100). Extracted the next 10 pages from the calculated page number (2000: pages 3193-3202).

### 3. Environment Setup

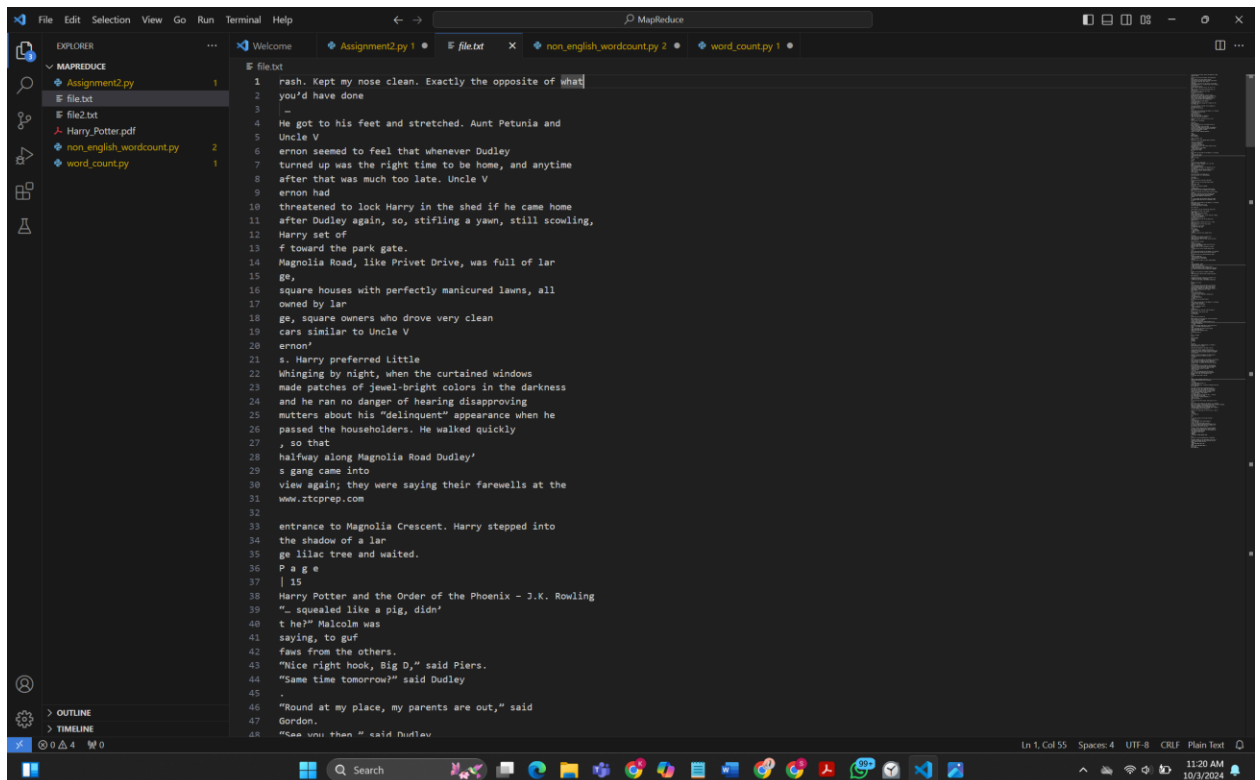
- Installed the required libraries using pip: `pip install PyPDF2 mrjob pyenchant`
- **Screenshot:**





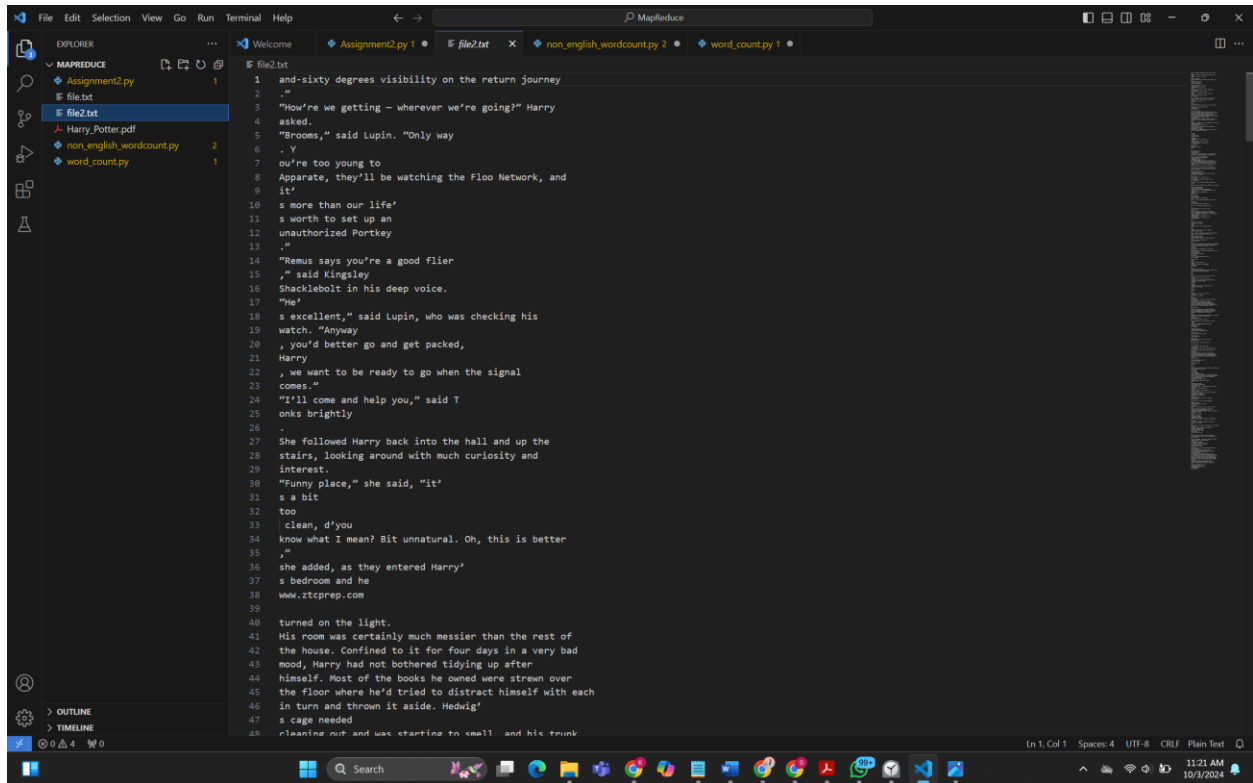
```
1 import PyPDF2
2
3 def extract_pages(pdf_path, start_page, end_page, output_file):
4     with open(pdf_path, 'rb') as pdf_file:
5         reader = PyPDF2.PdfReader(pdf_file)
6         with open(output_file, 'a', encoding='utf-8') as output:
7             for page_num in range(start_page - 1, end_page):
8                 page = reader.pages[page_num]
9                 text = page.extract_text()
10                if text:
11                    output.write(text)
12                    output.write("\n")
13
14 #PDF_path
15 pdf_path = 'Harry_Potter.pdf'
16
17 # Extract for file1.txt (pages 3117-3127) (Birthday Month: September: 9/2+4.5 = rounded to 5. From Chapter 5(Page 3893): My birth date is 24.
18 # So I started from 24th page and took next 18 pages.)
19 extract_pages(pdf_path, 3117, 3127, 'file1.txt')
20
21 # Extract for file2.txt (pages 3193-3202) (Birth Year : 2000, last two digits is 00, adding 1 in front: 100: 3893: 3193-3203)
22 extract_pages(pdf_path, 3193, 3203, 'file2.txt')
23
24 print("Pages extracted and saved to file1.txt and file2.txt")
```

## File.txt:



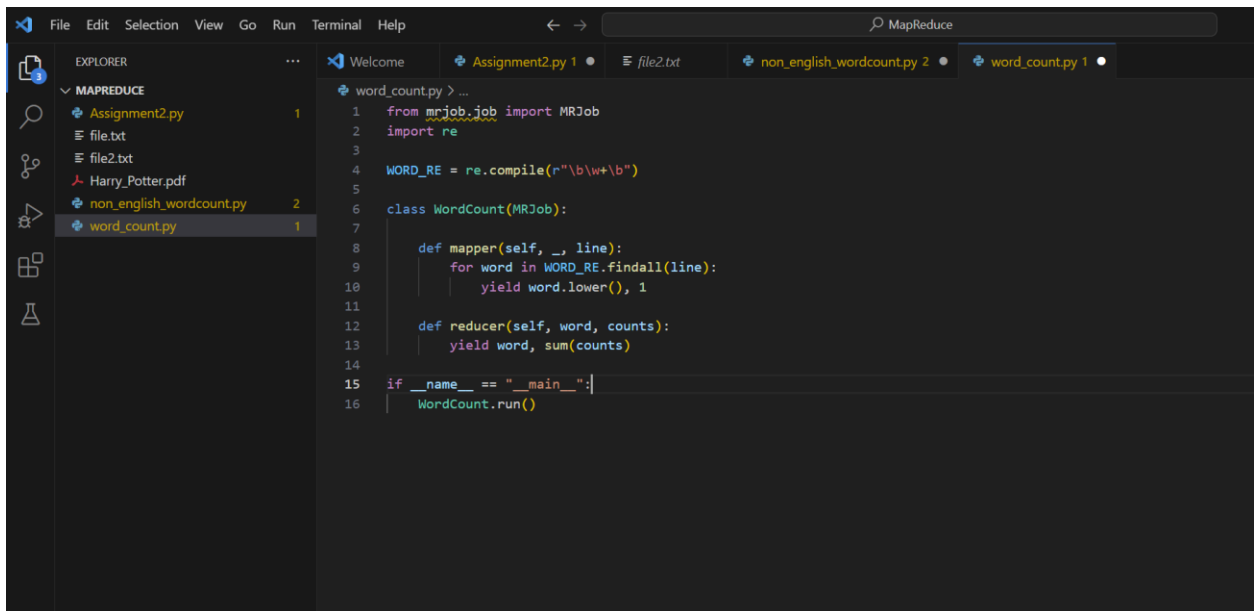
```
1 rash. Kept my nose clean. Exactly the opposite of what
2 you'd have done
3
4 He got to his feet and stretched. Aunt Petunia and
5 Uncle V
6 ernon seemed to feel that whenever Dudley
7 turned up was the right time to be home, and anytime
8 after that was much too late. Uncle V
9 ernon had
10 threatened to lock Harry in the shed if he came home
11 after Dudley again, so, stifling a yawn, still scowling,
12 Harry set of
13 f toward the park gate.
14 Magnolia Road, like Privet Drive, was full of lar
15 ge,
16 square houses with perfectly manicured lawns, all
17 owned by lar
18 ge, square owners who drove very clean
19 cars similar to Uncle V
20 ernon's
21 s. Harry preferred Little
22 Whinging by night, when the curtained windows
23 made patches of jewel-bright colors in the darkness
24 and he ran no danger of hearing disapproving
25 mutters about his "delinquent" appearance when he
26 passed the householders. He walked quickly
27 , so that
28 halfway along Magnolia Road Dudley'
29 s gang came into
30 view again; they were saying their farewells at the
31 www.ztccprep.com
32
33 entrance to Magnolia Crescent. Harry stepped into
34 the shadow of a lar
35 ge lilac tree and waited.
36 P a g e
37 | 15
38 Harry Potter and the Order of the Phoenix - J.K. Rowling
39 "I... squealed like a pig, didn'
40 t he?" Malcolm was
41 saying, to guf
42 faws from the others.
43 "Nice right hook, Big D," said Piers.
44 "Same time tomorrow?" said Dudley.
45
46 "Round at my place, my parents are out," said
47 Gordon.
48 "See you then?" said Outlaw
```

File2.txt:



## 5. Word Count Using MapReduce (Screenshots)

- Python code for the word count MapReduce job.



**Output:**

[illegible]

```

C:\Windows\system32\cmd.exe
"bestial" 1
"between" 2
"big" 6
"bitingly" 1
"black" 1
"blackness" 1
"blatantly" 1
"blind" 1
"blinding" 1
"blundering" 1
"boon" 1
"booming" 2
"boyfriend" 1
"brave" 1
"brave" 2
"breathlessly" 1
"breaths" 1
"brigitte" 1
"brute" 1
"brute" 1
"but" 6
"by" 8
"bye" 1
"bye" 1
"call" 1
"calls" 1
"came" 2
"can" 2
"carrying" 1
"cars" 2
"caught" 1
"cause" 4
"cedric" 5
"cheer" 1
"charged" 1
"cheek" 1
"cheeked" 1
"chilly" 1
"clean" 2
"leaves" 1
"cola" 8
"colors" 1
"come" 11
"completely" 2
"connect" 1
"control" 1
"convinced" 1
"cool" 1
"coarser" 1
"could" 4
"couldn" 1
"crane" 3
"crank" 1
"creeping" 1
"creepant" 1
"cricket" 1
"cultained" 1
"cut" 1
"at" 16
"cut" 2
"daddy" 1
"charges" 1

```





[illegible]

```

C:\Users\user\System32\cmd.exe
"bookie" 1
"popped" 1
"near" 1
"potter" 7
"fourling" 1
"power" 1
"overfiring" 1
"perished" 1
"prayer" 1
"puller" 1
"culinary" 2
"culinary" 1
"eat" 1
"eat" 1
"rotting" 1
"up" 1
"renew" 1
"resisting" 1
"rest" 1
"revolving" 1
"right" 5
"round" 5
"round" 1
"rollup" 7
"rules" 2
"running" 2
"safe" 20
"safe" 19
"same" 1
"saturnalian" 1
"same" 2
"say" 2
"saying" 2
"score" 1
"scouting" 1
"screenplay" 1
"second" 2
"second" 1
"see" 5
"second" 2
"seen" 1
"self" 1
"senation" 1
"senate" 1
"set" 1
"setting" 1
"seven" 1
"shame" 1
"shop" 1
"shivering" 1
"sharpen" 1
"shattering" 1
"shot" 10
"side" 2
"sideways" 1
"sighless" 1
"silent" 2
"similar" 1
"slashing" 1
"slime" 1
"slime" 1
"skip" 1
"sleep" 1

```





The image shows a Visual Studio Code editor window with a dark theme. The Explorer sidebar on the left shows a project named 'MAPREDUCE' containing files: 'Assignment2.py', 'file.txt', 'file2.txt', 'Harry\_Potter.pdf', 'non\_english\_wordcount.py', and 'word\_count.py'. The 'non\_english\_wordcount.py' file is selected and open in the main editor. The code in the editor is a Python MapReduce script. It imports 'MRJob' from 'mrjob.job', 'enchant' for spell checking, and 're' for regular expressions. It sets an English dictionary and a word regex. A 'NonEnglishWordCount' class inherits from 'MRJob' and implements a 'mapper' method that finds non-English words in each line and a 'reducer' method that sums their counts. The script is executed via a main block.

```
1 from mrjob.job import MRJob
2 import enchant
3 import re
4
5 english_dict = enchant.Dict("en_US")
6 WORD_RE = re.compile(r"\b\w+\b")
7
8 class NonEnglishWordCount(MRJob):
9
10     def mapper(self, _, line):
11         for word in WORD_RE.findall(line):
12             word_lower = word.lower()
13             if not english_dict.check(word_lower):
14                 yield word_lower, 1
15
16     def reducer(self, word, counts):
17         yield word, sum(counts)
18
19 if __name__ == "__main__":
20     NonEnglishWordCount.run()
```

- Output showing the non-English words and their counts from file2.txt.

```
C:\Windows\System32\cmd.exe x + v
"yeah" 1
"year" 2
"years" 1
"yell" 1
"yelled" 1
"you" 45
"youre" 8
"ztcprep" 11
Removing temp directory C:\Users\sruth\AppData\Local\Temp\word_count.sruth.20240929.210059.763842...

C:\Users\sruth\FALL_2024\603_Platforms for Big Data Processing\MapReduce>python non_english_wordcount.py file2.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory C:\Users\sruth\AppData\Local\Temp\non_english_wordcount.sruth.20240929.210906.195933
Running step 1 of 1...
job output is in C:\Users\sruth\AppData\Local\Temp\non_english_wordcount.sruth.20240929.210906.195933\output
Streaming final output from C:\Users\sruth\AppData\Local\Temp\non_english_wordcount.sruth.20240929.210906.195933\output...
"apparate" 1
"auror" 2
"britain" 1
"didn" 1
"dursleys" 1
"eah" 2
"eam" 1
"ebolt" 1
"ernon" 1
"firebolt" 1
"flo" 1
"gify" 1
"gis" 1
"gruf" 1
"hedwig" 4
"hestia" 1
"hogwarts" 1
"householdy" 1
"ir" 1
"jones" 1
"kingley" 4
"lupin" 9
"mell" 1
"metamorphmagi" 1
"metamorphmagus" 2
"midrif" 1
"onks" 16
"ou" 4
"pell" 1
"podmore" 1
"portkey" 1
"quidditch" 1
"remus" 1
"rowling" 6
"runk" 1
"shacklebolt" 3
"sirius" 1
"stuf" 1
"stuf" 1
"stuf" 1
"ve" 7
"wouldn" 2
"www" 11
"ztcprep" 11
Removing temp directory C:\Users\sruth\AppData\Local\Temp\non_english_wordcount.sruth.20240929.210906.195933...

C:\Users\sruth\FALL_2024\603_Platforms for Big Data Processing\MapReduce
```

```
C:\Windows\System32\cmd.exe x + v
"ztcprep" 11
Removing temp directory C:\Users\sruth\AppData\Local\Temp\word_count.sruth.20240929.210059.763842...

C:\Users\sruth\FALL_2024\603_Platforms for Big Data Processing\MapReduce>python non_english_wordcount.py file2.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory C:\Users\sruth\AppData\Local\Temp\non_english_wordcount.sruth.20240929.210906.195933
Running step 1 of 1...
job output is in C:\Users\sruth\AppData\Local\Temp\non_english_wordcount.sruth.20240929.210906.195933\output
Streaming final output from C:\Users\sruth\AppData\Local\Temp\non_english_wordcount.sruth.20240929.210906.195933\output...
"apparate" 1
"auror" 2
"britain" 1
"didn" 1
"dursleys" 1
"eah" 2
"eam" 1
"ebolt" 1
"ernon" 1
"firebolt" 1
"flo" 1
"gify" 1
"gis" 1
"gruf" 1
"hedwig" 4
"hestia" 1
"hogwarts" 1
"householdy" 1
"ir" 1
"jones" 1
"kingley" 4
"lupin" 9
"mell" 1
"metamorphmagi" 1
"metamorphmagus" 2
"midrif" 1
"onks" 16
"ou" 4
"pell" 1
"podmore" 1
"portkey" 1
"quidditch" 1
"remus" 1
"rowling" 6
"runk" 1
"shacklebolt" 3
"sirius" 1
"stuf" 1
"stuf" 1
"stuf" 1
"ve" 7
"wouldn" 2
"www" 11
"ztcprep" 11
Removing temp directory C:\Users\sruth\AppData\Local\Temp\non_english_wordcount.sruth.20240929.210906.195933...

C:\Users\sruth\FALL_2024\603_Platforms for Big Data Processing\MapReduce>
```

## 7. Conclusion

This assignment successfully demonstrated the application of MapReduce for analyzing text data extracted from a PDF document. By leveraging MapReduce's power, we could efficiently count word occurrences and identify non-English words within the specified text segments.

### Key Findings:

- **Word Frequency Analysis:** The word count MapReduce job effectively identified the most frequent words in the extracted text from file.txt.
- **Non-English Word Identification:** The non-English word count MapReduce job accurately identified words not recognized by the English dictionary.

### Challenges and Improvements:

**No significant challenges** were encountered during the execution of the MapReduce jobs.

### Potential Improvements:

- **Automated Page Extraction:** The code could be enhanced to automatically determine the page ranges based on birth date and year without manual input. This could involve using regular expressions or other techniques to extract relevant information from the PDF's metadata or content.