

MapReduce Programming Model

1) Title of Problem

Implement a MapReduce-style program to count word frequencies from multiple text files while applying the following conditions:

- Ignore common stopwords such as the, and, of, etc.
- Display the Top 5 least frequent words.
- Count only words with more than 5 characters.
- Extend the program to process multiple text files.

2) Code for Problem

The following Python code implements the above requirements:

```
1 from collections import defaultdict
2 import glob
3
4 # Mapper function
5 def mapper(file_path, stopwords):
6     mapped = []
7     with open(file_path, 'r', encoding="utf-8") as f:
8         for line in f:
9             words = line.strip().split()
10            for word in words:
11                w = word.lower().strip(",.?!:;\"'() []{}")
12                if w and w not in stopwords and len(w) > 5:
13                    mapped.append((w, 1))
14    return mapped
15
16 # Reducer function
17 def reducer(mapped_data):
18     reduced = defaultdict(int)
19     for word, count in mapped_data:
20         reduced[word] += count
21     return reduced
22
23 if __name__ == "__main__":
24     # Define stopwords
25     stopwords = {"the", "and", "of", "to", "a", "in",
26                 "is", "it", "that", "for", "on",
27                 "with", "as", "by", "at"}
```

```

29     # Process all text files
30     all_files = glob.glob("*.txt")
31     mapped_data = []
32     for file_path in all_files:
33         mapped_data.extend mapper(file_path, stopwords))
34
35     reduced_data = reducer(mapped_data)
36
37     # Top 5 least frequent
38     least_frequent = sorted(reduced_data.items(), key=lambda x: x[1])
39                          [:5]
40     print("\nTop 5 Least Frequent Words:")
41     for word, freq in least_frequent:
42         print(f"{word}: {freq}")

```

3) Solution Output

Top 5 Least Frequent Words

```

servers: 1
storage: 1
designed: 1
handle: 1
failures: 1

```