

```
1 Problem Statement
2 -----
3 The volume of unstructured data (Text data, log lines, images, binary files) in
4 existence is growing dramatically, and
5 MR, Spark are excellent framework for analyzing this type of data.
6 You will implement a MR application that calculates the most common words from Complete
7 Works of William Shakespeare
8 Please refer a file Complete_Shakespeare.txt
9
10 Here are the brief steps for writing the word counting program:
11
12 1) Create a MR application which is going to calculates the most common words from
13 Complete Works of William Shakespeare -
14 use 'Complete_Shakespeare.txt' file residing in your home on HDFS.
15 2) The most common words will be decided based on the stop words
16 3) Stop words are common words that are often uninteresting. For example "I", "the", "a"
17 etc., are stop words.
18 You can remove many obvious stop words with a list of your own. But for this exercise,
19 you will just remove
20 the stop words from a curated list `stop_words` provided to you in your environment
21 (Refere a file stop_words.txt)
22 4) You will create following classes
23 1. Outer job class acting as a Tool,
24 having name as WordCountCompleteShakespere
25 2. StopWordMapper class representing map task
26 as inner class inside WordCountCompleteShakespere class
27 3. StopWordReducer class representing reduce task
28 as inner class inside WordCountCompleteShakespere class
29 5) Business demands that output should be generated in a folder
30 Shakespere_work on HDFS into 3 files
31 6) Business want you to run the application using following command only
32 yarn jar shakersperework.jar <Input path> <Output path>
33 Input path - Represents a path of a file/dataset on which the
34 Job/application is suppose to run
35 Output path - Represents a path of a output folder on HDFS
36 where the Job/application is suppose to put the final output
37 7) Automate the entire build process by writing down a shell script (This step is
38 optional)
```

Timeline - 27 Oct 2025, EOD