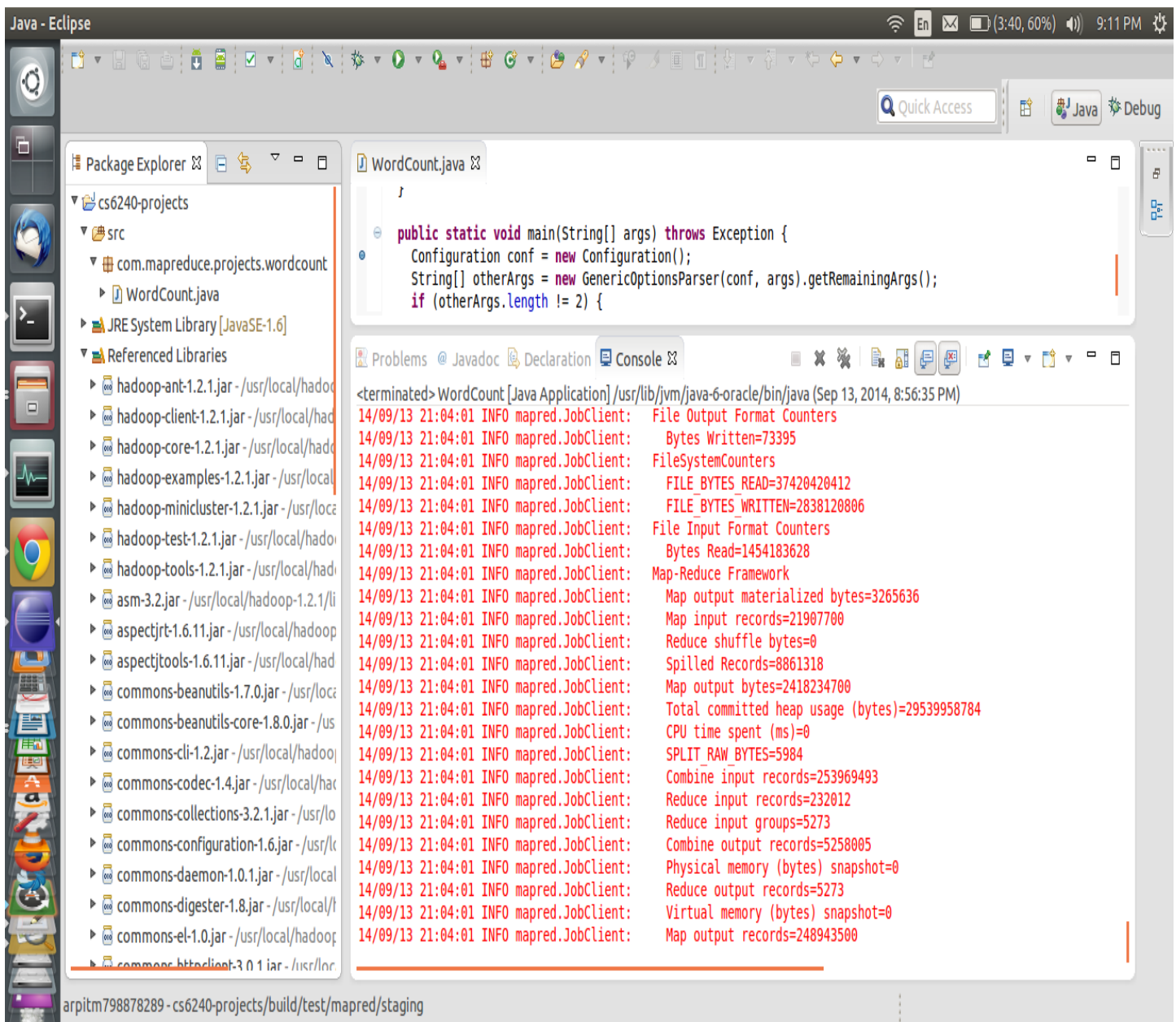# CS6240 – PARALLEL DATA PROCESSING IN MAP-REDUCE
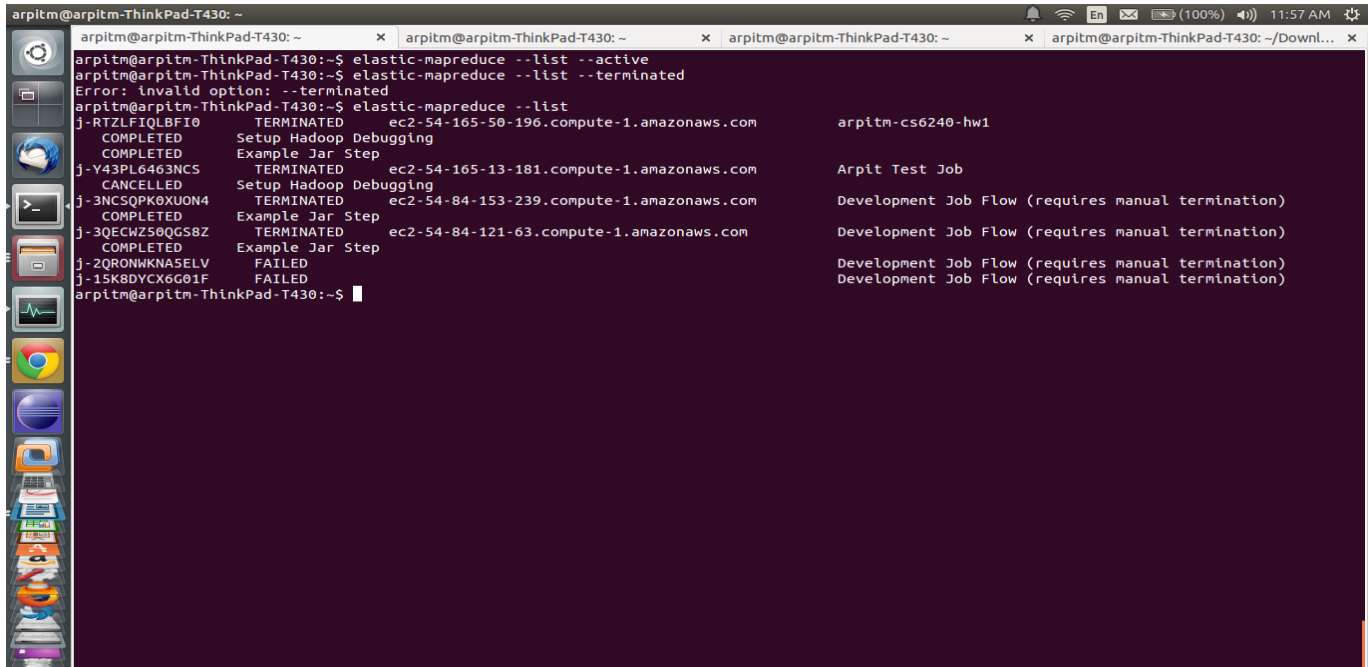
Class: CS6240-02

Homework Number: 1

Name: Arpit Mehta

SCREENSHOT 1 – LOCAL EXECUTION: The below screen shot shows the project directory structure and console output of successful run of WordCount program inside the IDE:

SCREENSHOT 2 – AWS EXECUTION:

The below screen shot shows the list of jobflows that I created on AWS for running the custom wordcount.jar. The job id 'j-RTZLFIQLBFI0' with job name 'arpitm-cs6240-hw1' represents the jobflow of run using 3 small machines.



The below screenshot shows the AWS EMR 'Cluster Details' console detailing the job run.

The below screenshot shows the jobflow steps:



The below screenshot shows my s3 bucket 'arpitm-cs6240' that contains the 'wordcount.jar', the input file 'hw1.txt', the 'logs' directory, the 'hw1-output-3instances' directory contains the output of the jobflow on 3 small machines.

The below screenshot shows 'hw1-output-3instances' containing the output files:



The below screenshot shows the directory content of jobflow 'j-RTZLFIQLBFI0' in the logs directory.

Please find attached with the submission the controller & syslog logs and the final output files part-r-00000, part-r-00001 & part-r-00002.