## CS 6240: Assignment 1

### Xiao Wang

GitHub: <a href="https://github.ccs.neu.edu/xiaowang/CS6240-MapReduce">https://github.ccs.neu.edu/xiaowang/CS6240-MapReduce</a>

# 1. Weather Data Results Number of worker threads: 3

#### Running results with no Fibonacci:

(ms)	Sequential	NoLock	CoarseLock	FineLock	NoSharing
1	2849	ERROR	1612	1514	1470
2	3144	ERROR	1545	1593	1486
3	3152	1774	1912	1599	1451
4	2938	1659	1715	1604	1515
5	2916	1527	1536	1473	1449
6	3113	1507	1628	1529	1428
7	2925	1572	1555	1479	1372
8	2923	1611	1626	1529	1530
9	2763	1564	1599	1508	1484
10	2999	1695	1728	1668	1449
Min	2763	1507	1536	1473	1372
Max	3152	1774	1912	1668	1530
Avg	2972.2	1613.625	1645.6	1549.6	1463.4
SpeedUp(avg)		1.841939732	1.806149733	1.918043366	2.031023644

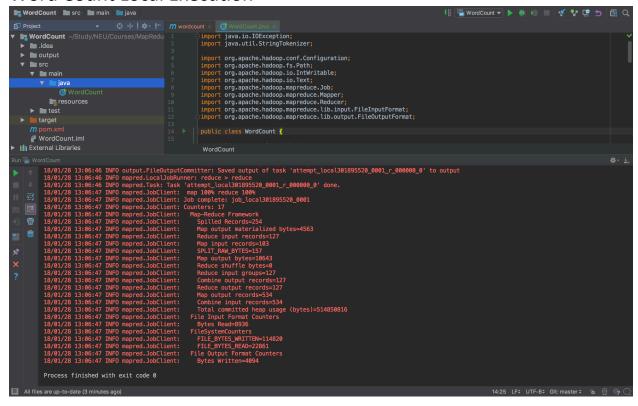
#### Running results with Fibonacci:

(ms)	Sequential	NoLock	CoarseLock	FineLock	NoSharing
1	12312	5503	10566	5695	5481
2	12082	5446	10378	5556	5376
3	11925	5461	10635	5512	5482
4	12169	5604	10587	5608	5735
5	11974	5824	10575	5574	5676
6	11968	5499	10589	5640	5540
7	11969	ERROR	10944	5902	6833
8	11930	5709	10329	5582	5495
9	12010	5742	11095	5711	5594

10	12191	5708	10848	7272	5655
Min	11925	5446	10329	5512	5376
Max	12312	5824	11095	7272	6833
Avg	12053	5610.66667	10654.6	5805.2	5686.7
SpeedUp(avg)		2.148229563	1.131248475	2.07624199	2.11950692

- **Q1.** I expect NoLock to finish fastest because it doesn't consider consistency, therefore it takes the most advantage of parallelization. The experiments mostly confirm my expectation, although NoSharing performs a little better with no Fibonacci. I think it's because NoLock is using a bigger shared data structure and it requires more overhead.
- **Q2.** I expect Sequential to finish slowest because other versions all take advantage of parallelization in some degree. The experiments confirm my expectation.
- **Q3.** The NoLock version crashes with NullPointerError sometimes because of concurrent accesses.
- **Q4.** Sequential is slower than CoarseLock. I think it's because besides locking time, there's some other computation to do and parallelization reduce the time of this part.
- **Q5.** Higher computation cost reduces the difference between Sequential and CoarseLock. I think the reason is higher computation makes it easier for threads to starve.

#### 2. Word Count Local Execution



#### 3. Word Count AWS Execution

```
[INFO] Deleting /Users/Breezen/Study/NEU/Courses/MapReduce/CS6240-MapReduce/HW1/AWS/target
 IINPO]
IINPO]
IINPO]
IINPO]
IINPO]
IINPO]
IINPO]
IINPO]
IINPO]
INPO]
INP
 [INFO] Total time: 4.159 s
[INFO] Finished at: 2018-01-28T19:03:06-05:00
[INFO] Final Memory: 24M/224M
[INFO] Final Memory: 24M/224M
"31" V
--applications Name-Hadoop \
--steps '[474ggs:|[WordCountry,"33://seanxwang-wordcount/inputr,"s3://seanxwang-wordcount/outputr],"Type":"CUSTOM_JAR","Jar":"S3://seanxwang-wordcount/wc-1.0.jar","Actio
nOnFailure":"TEMNINTE_CLUSTEM;"Name":"Custom JAR")! \
--log-uri 33://seanxwang-wordcount/log \
--service-role EMR_DefaultRole \
--ec2-attributes InstanceProfile=EMR_EC2_DefaultRole,SubnetId=subnet-197e4152 \
                             --region us-east-1 \
--enable-debugging \
--auto-terminate
 --auto-terminate
j-1BB9NX8IH8QMA
Seans-MacBook-Pro:AWS Breezen$
               aws
                                        Services v
                                                                    Resource Groups v
                                                                                                                                                                                                                                       xiaowang@ccis.neu *
                                                                                                                                                                                                                                                                                N. Virginia 🕶
                                                                                                                                                                                                                                                                                                            Support *
                                                              Cluster: WordCount Cluster
                                                                                                                                           Terminated Steps completed
     Amazon FMR
                                                                  Summary Monitoring Hardware Events Steps Configurations
  Clusters
     Security configurations
                                                                                                                                                                                                                                                                                                                            G
                                                                Connections:
     VPC subnets
                                                                Master public DNS:
                                                                                                                    ec2-34-236-152-21.compute-1.amazonaws.com SSH
                                                                Tags:
     Events
                                                                                                                                                                      Configuration details
                                                                                                ID: j-1BB9NX8IH8QMA
                                                                                                                                                                                 Release label: emr-5.2.1
                                                                            Creation date: 2018-01-28 19:03 (UTC-5)
                                                                                                                                                                      Hadoop distribution: Amazon 2.7.3
                                                                                    End date: 2018-01-28 19:14 (UTC-5)
                                                                                                                                                                                    Applications: --
                                                                             Elapsed time: 10 minutes
                                                                                                                                                                                           Log URI: s3://seanxwang-wordcount/log/
                                                                         Auto-terminate: Yes
                                                                                                                                                                         EMRFS consistent Disabled
                                                                                Termination Off
                                                                                                                                                                                                view:
                                                                                  protection:
                                                                Network and hardware
                                                                                                                                                                      Security and access
                                                                        Availability zone: us-east-1b
                                                                                                                                                                                       Key name: --
                                                                                   Subnet ID: subnet-197e4152
                                                                                                                                                                      EC2 instance profile: EMR_EC2_DefaultRole
                                                                                        Master: Terminated 1 m4.large
                                                                                                                                                                                          EMR role: EMR_DefaultRole
                                                                                           Core: Terminated 4 m4.large
                                                                                                                                                                          Visible to all users: All Change
                                                                                           Task: --
                                                                                                                                                                         Security groups for sg-d7c001a0 (ElasticMapReduce-
                                                                                                                                                                                             Master: master)
                                                                                                                                                                         Security groups for sg-bab879cd (ElasticMapReduce-
                                                                                                                                                                                     Core & Task: slave)
```

© 2008 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved.