

Sravani Beeram

CS6240 Parallel Data Processing using Mapreduce

Homework 1

Running Time without Fibonacci delay

Program	Minimum Time(ms)	Maximum Time(ms)	Average Time(ms)
No Lock	417	512	446
Fine Lock	429	530	452
No Sharing	409	639	477
Coarse Lock	766	804	784
Sequential	755	1901	885

Running Time with Fibonacci delay

Program	Minimum Time(ms)	Maximum Time(ms)	Average Time(ms)
No Lock	439	689	524
Fine Lock	427	770	610
No Sharing	444	766	658
Coarse Lock	777	1041	901
Sequential	785	2377	1134

Number of Threads used for running the program : 2

Speed up calculation for Fibonacci delay :

No Lock :2.16

Fine Lock:1.86

No Sharing:1.72

Coarse Lock:1.26

1. Which program version (SEQ, NO-LOCK, COARSE-LOCK, FINE-LOCK, NO-SHARING) would you normally expect to finish fastest and why? Do the experiments confirm your expectation? If not, try to explain the reasons.

No-Lock is expected to finish fast as data is shared between the threads and threads will act independently to update the data structure. Run times of my programs confirm the same

2. Which program version (SEQ, NO-LOCK, COARSE-LOCK, FINE-LOCK, NO-SHARING) would you normally expect to finish slowest and why? Do the experiments confirm your expectation? If not, try to explain the reasons.

Sequential is expected to finish slowest as there is only one thread in the system to complete the execution. Hence it will be slower than the multithreaded programs. My programs results confirm that SEQ does take maximum time for execution for case B.

3. Compare the temperature averages returned by each program version. Report if any of them is incorrect or if any of the programs crashed because of concurrent accesses.

Except No_Lock program, temperature averages returned by the other program version were same. Sometimes No_Lock program returned Null Pointer Exception. This was because as one thread had already access to object ID but had not updated in the common data structure.

And second thread already tries to access the Id and doesn't get any value from the data structure.

4. Compare the running times of SEQ and COARSE-LOCK. Try to explain why one is slower than the other. (Make sure to consider the results of both B and C—this might support or refute a possible hypothesis.)

In general Sequential should be slowest without Fibonacci and Coarse-lock slowest with Fibonacci. Sequential has only one thread to execute the program and is slower compared to multiple thread programs. Hence, without Fibonacci, Sequential is slowest and program run times confirm the same.

With Fibonacci, Coarse lock has a thread that locks the object and no other threads have access to the object and fibonacci slows it down further. My coarse-lock code with fibonacci is

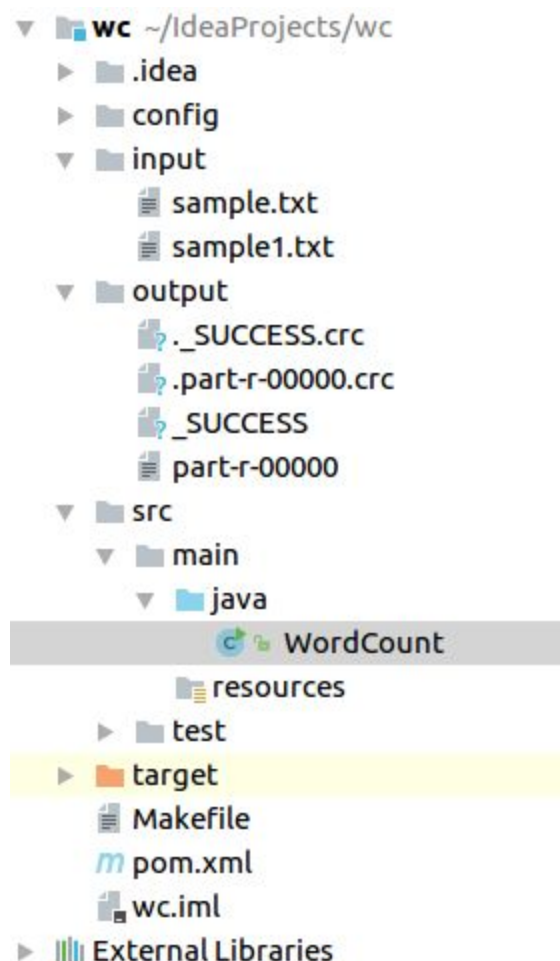
60% slower than sequential but other times sequential is slower. It could be because the lock on the object is released quickly and hence runs faster than sequential. It would have been more evident if the task was bigger.

5. How does the higher computation cost in part C (additional Fibonacci computation) affect the difference between COARSE-LOCK and FINE-LOCK? Try to explain the reason.

Higher computation cost doesn't affect the difference between Coarse-lock and Fine-lock since in both of them before we update the values in the accumulated data structure we make the threads wait. Code results depict the same.

Word Count Local Execution

1. Project directory structure, showing that the WordCount.java file is somewhere in the src directory



2. The console output for a successful run of the WordCount program inside the IDE. The console output refers to the job summary information Hadoop produces, not the output your job emits.

```
Terminal
+
X
File System Counters
  FILE: Number of bytes read=36285390488
  FILE: Number of bytes written=338020493
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
Map-Reduce Framework
  Map input records=21907703
  Map output records=248943503
  Map output bytes=2418234724
  Map output materialized bytes=6456811
  Input split bytes=5221
  Combine input records=248943503
  Combine output records=458752
  Reduce input groups=5274
  Reduce shuffle bytes=6456811
  Reduce input records=458752
  Reduce output records=5274
  Spilled Records=1370982
  Shuffled Maps =45
  Failed Shuffles=0
  Merged Map outputs=45
  GC time elapsed (ms)=2520
  CPU time spent (ms)=0
  Physical memory (bytes) snapshot=0
  Virtual memory (bytes) snapshot=0
  Total committed heap usage (bytes)=21492137984
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=1454183639
File Output Format Counters
  Bytes Written=73401
sravani@sravani:~/IdeaProjects/wc$
```

Word Count AWS Execution

1. Show a similar screenshot that provides convincing evidence of a successful run of the Word Count program on AWS. Make sure you run the program using at least three machines, i.e., one master node and two workers

Amazon EMR

Cluster list

Security configurations

VPC subnets

Help

Add step

Resize

Clone

Terminate

AWS CLI export

Cluster: WordCount Cluster

Terminated

Steps completed

Connections:

--

Master public DNS:

ec2-52-11-154-80.us-west-2.compute.amazonaws.com

SSH

Tags:

--

Summary

ID: j-1SJMJAYSLROF1

Creation date: 2017-01-29 15:28 (UTC-5)

End date: 2017-01-29 15:49 (UTC-5)

Elapsed time: 20 minutes

Auto-terminate: Yes

Termination protection: Off

Configuration Details

Release label: emr-5.2.1

Hadoop distribution: Amazon 2.7.3

Applications: --

Log URI: s3://cs6240sraav/logs/

EMRFS consistent view: Disabled

Network and Hardware

Availability zone: us-west-2a

Subnet ID: subnet-115ba158

Master: Terminated 1 m1.medium

Core: Terminated 2 m1.medium

Task: --

Security and Access

Key name: --

EC2 instance profile: EMR_EC2_DefaultRole

EMR role: EMR_DefaultRole

Visible to all users: All [Change](#)

Security groups for sg-49028031 (ElasticMapReduce-Master: master)

Security groups for sg-4d028035 (ElasticMapReduce-Core & Task: slave)

Monitoring

Hardware

Steps

Configurations

Events

Bootstrap Actions