

School of Information Technologies
Faculty of Engineering & IT

ASSIGNMENT/PROJECT COVERSHEET - GROUP ASSESSMENT

Unit of Study: COMP 5349 Cloud Computing

Assignment name: Hadoop MapReduce Programming

Tutorial time: Thursday 4p.m-6p.m Tutor name: Waiho Wong

DECLARATION

We the undersigned declare that we have read and understood the [University of Sydney Academic Dishonesty and Plagiarism in Coursework Policy](#), and, except where specifically acknowledged, the work contained in this assignment/project is our own work, and has not been copied from other sources or been previously submitted for award or assessment.

We understand that failure to comply with the *Academic Dishonesty and Plagiarism in Coursework Policy* can lead to severe penalties as outlined under Chapter 8 of the *University of Sydney By-Law 1999* (as amended). These penalties may be imposed in cases where any significant portion of my submitted work has been copied without proper acknowledgement from other sources, including published works, the internet, existing programs, the work of other students, or work previously submitted for other awards or assessments.

We realise that we may be asked to identify those portions of the work contributed by each of us and required to demonstrate our individual knowledge of the relevant material by answering oral questions or by undertaking supplementary work, either written or in the laboratory, in order to arrive at the final assessment mark.

Project team members				
Student name	Student ID	Participated	Agree to share	Signature
1.Xuhong Guo	450489321	<input checked="" type="radio"/> Yes <input type="radio"/> No	Yes <input checked="" type="radio"/> No	Xuhong
2.Yuchen Zhao	460134794	<input checked="" type="radio"/> Yes <input type="radio"/> No	Yes <input checked="" type="radio"/> No	Yuchen Zhao

Assignment 1 Cloud Computing

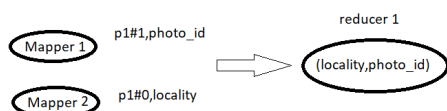
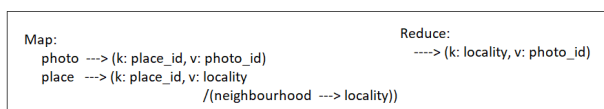
Xuhong Guo 450489321 Yuchen Zhao 460134794

Abstract. In this assignment We will write a series of Hadoop jobs to analyze a data set from flickr.com. The analysis is structured into three separate tasks which build on each other. We will practice basic Hadoop programming features and also observe the ease of use of the Hadoop framework. The data set is the same as those we have used in the labs.

1 Number of photos taken per locality.

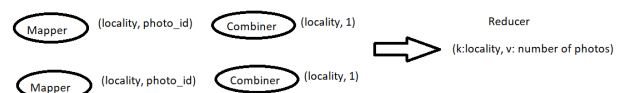
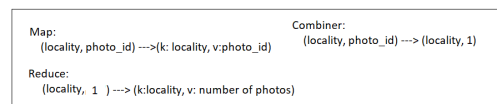
1.1 Mapreduce 1

In order to find every locality photo numbers, we would use two map reduce in this question. In the first part map reduce, we would Implementing and partitioner in Hadoop, because Hadoop supports multiple inputs as well as different mappers, so we use different values from inputs. For the photo record, which comes from n0*.txt. We use place-id as the key followed by #1 to represent the source from and values are photo-id(key:place-id, value: photo-id). For the place record, we set place-id as the key with #0 and value is place-name from place.txt file(key:place-id,value:place-name).Use the reduce-side join by the key place-id, and final output is the values which is locality name and photo-is.We also set place-type-id 22 belongs to the place-type-id7 as the question requested.By the way, the characteristic are shorten of attributes' name in files.



1.2 Mapreduce 2

The second part map reduce in question one, we use the map to read result from previous map reduce. Then we use combiner function to get the photo number for each locality, reduce is used to sum all these photo numbers by each locality which can be shown in diagrams.



1.3 Total running time

The total running time in for question 1 is 5 minutes.

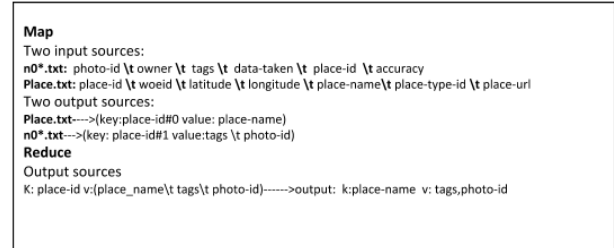
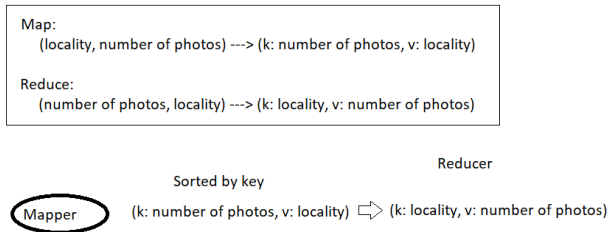
1.4 Result

4	'Ain Taya, Boumerdes, Algeria	7
5	'Ain al 'Awda, Rabat-Sale, Morocco	1
6	'Ali Shah 'avaz, Tehran, Iran	7
7	's-Graveland, North Holland, Netherlands	13
8	's-Gravenbrakel, Hainault, Belgium	33
9	's-Gravendeel, South Holland, Netherlands	4
10	's-Gravenpolder, Zeeland, Netherlands	11
11	's-Gravenweg, South Holland, Netherlands	19
12	's-Gravenwezel, Antwerp, Belgium	8
13	's-Gravenzande, South Holland, Netherlands	56
14	's-Heer Hendrikskinderen, Zeeland, Netherlands	16
15	's-Heer-Abtskerke, Zeeland, Netherlands	3
16	's-Heer-Arendskerke, Zeeland, Netherlands	3

2 The top 50 locality level places based on the number of photos taken in this locality.

2.1 Mapreduce

We could get the number of photos in each locality from the sector 1. Firstly we use map to get these source, which can be shown in table. For the reduce part, we set the a total list sized by 50, using reduce to list the top 50. Moreover, we use a reserve function here by the key. The reduce table could be seen in diagram as well.



2.2 Running Time

This running time is 1 minute.

2.3 Result

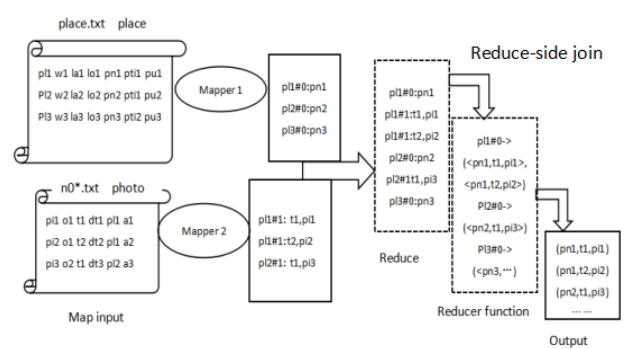
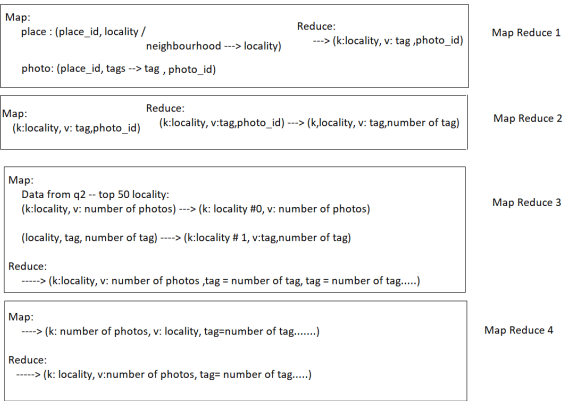
4	New York, NY, US, United States	1188956
5	Seattle, WA, US, United States	698078
6	NY, US, United States	613379
7	Chicago, IL, US, United States	566404
8	Paris, Ile-de-France, FR, France	532534
9	Washington, DC, US, United States	502270
10	Tokyo, Tokyo Prefecture, JP, Japan	498360
11	Taipei City, Taipei City, TW, Taiwan	494227
12	Sydney, NSW, AU, Australia	489161
13	Los Angeles, CA, US, United States	456481
14	Toronto, ON, CA, Canada	445834
15	Vancouver, BC, CA, Canada	429842
16	Melbourne, VIC, AU, Australia	352227
17	Tokyo, Tokyo Prefecture, Japan	342004
18	New York, NY, United States	318895
19	Austin, TX, US, United States	315899
20	Portland, OR, US, United States	281729
21	San Diego, CA, US, United States	272381
22	Barcelona, Catalonia, ES, Spain	265809
23	San Francisco, California, United States	242558
24	Boston, MA, US, United States	235339

In the first map reduce, we use the function partitioner to differ the data comes from two separate input sources based on the original key place-id. The key and value comes from place information is marked as #0, key and value about photo which comes from n0*.txt is marked as #1. Using reduce-side join by the same key which is place-id. we use reduce function to get the same key with different values.We get the key with place-id, value pairs with locality, tags and photo-id. Furthermore,we set the output of reduce is place-name, tags and photo-id. This changing is convenient for further using. The following diagram would show this process design well, however, the diagram shows simple example, due to large data-set, we set 7 map partitoners and 6 reduce partitoners.

The important point is that we filter out place or year tags in the reduce process here in void of useless tags, which including the lower characteristics’ problems, same locality problems and/or data-time problems.

3 Ten most popular tags for each of the top-50 localities.

In this question, we plan to use four map reduce in Hadoop. The graph below shows the simple inputs and outputs of each map reduce process.

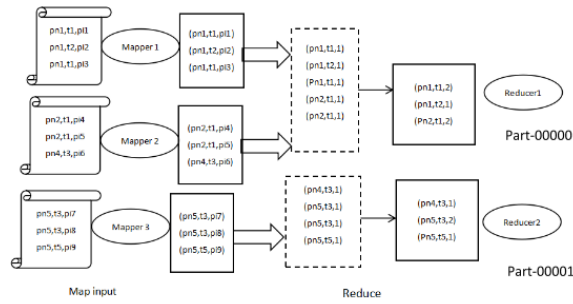


3.2 Mapreduce 2

The second map reduce focuses on tags’ processing. We use this map reduce to get each tags’ number. In this process, map is used to read the information to lines. Reduce is used to count the same tags number under the same key. The output of this map reduce is locality, tags and number of this tag.

3.1 Mapreduce 1

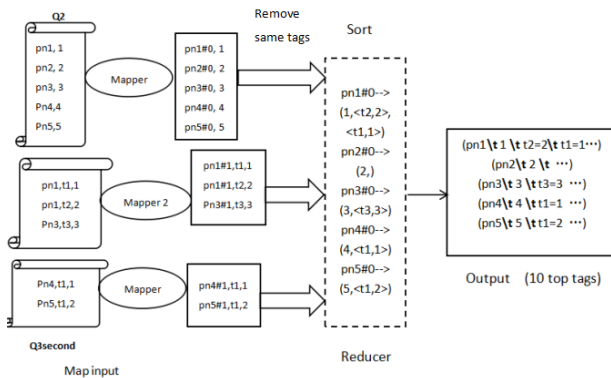
At first, we will explain the first map reduce structure. we could know the input and output of the first map reduce process.



3.3 Mapreduce 3

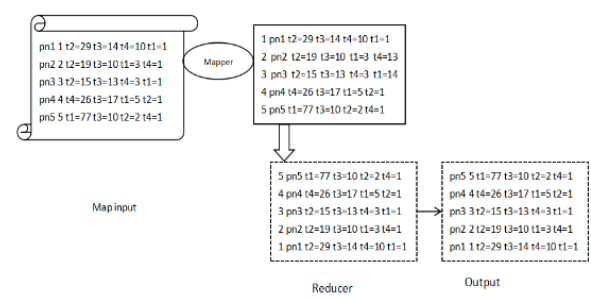
Map is used to get the previous output and input sources. In reduce process, we use keypartitioner to reduce-side join by placename. We set a dictionary to reduce the repeated tags. We set a total list sized by 10 to get the top ten tags and reversed by the tags' number decreasing.

```
Map
Two input sources:
Q2: photo-name \t number of photos
Q3second map reduce: place-name \t tags \t numbers
Two output sources:
Q2---->(key:place-name#0 value: number of photos)
Q3second map reduce---->(key: place-name#1 value:tags \t numbers)
Reduce
Output sources
K: place-name v:(number of photos, tag1\t number1, tag2\t number2...tag10\t number10)
```



3.4 Mapreduce 4

At last, we use a map reduce to process reserve, use -nr in reduce to make the results sorted by each places' photo numbers decreasing. Using map to read lines and using reduce to accomplish sorting. The following diagram shows the structure of fourth map reduce in this question.



3.5 Running Time

This running time is nearly 20 minutes.

3.6 Tags Cleaning

- Tag with year is deleted.
- Tag include london, uk when the place contains London, UK is deleted.
- Tag like sanfrancisco but the place is San Francisco is deleted.
- Tag list new or york, but the place is New York we keep it to prevent not the same meaning.

3.7 Result

```
9 Tokyo 498360 cat=7787, iwalking=7174, asia=7692, sandiegoo=5343, walking=5926, zoo=4725, panda=4596, s
10 Taipei City 494227 台湾=41630, wedding=26390, canon=25259, zodiac=21339, 台北=20034, allkypix=11849, 30d=9707
11 Sydney 489161 wwwintamcom=3670, artchannel=3103, 400d=2457, thepretenda=2331, sculpture=2285, club=2161,
12 Los Angeles 456481 california=59665, hollywood=16846, coop=1265, la=5989, leica=5155, streetart=3131, rays
13 Toronto 445834 ontario=79539, city=13582, canon=11907, tysonwilliams=11399, wwwtysonwilliamscom=10900, ait
14 Vancouver 429842 rolandtanglaophoto=39721, cameraphone=35136, shoru=34146, britishcolumbia=29880, nokia
15 Melbourne 352227 victoria=31009, asianime=18044, anime=16116, animetronics=15139, 7-A=12195, pentash
16 Tokyo 342004 日本=22917, 東京=19873, girl=18956, street=18525, asiana=18509, woman=18285, female=16786, d
17 New York 318895 nyc=76308, newyorkcity=33256, manhattan=26261, brooklyn=16089, usa=15916, art=10003, ci
18 Austin 315899 texas=39889, ate=10455, metal=8955, music=6573, improv=5374, band=5076, austlinimprov=4755,
19 Portland 291729 oregon=46411, pdo=3394, xosbianxoxe=6423, brain=6423, brainlattaphotography=6423, beiat
20 San Diego 272381 california=39581, sandiegoo=20100, zoo=17132, comiccon=9737, panda=8139, giantpanda=
21 Barcelona 26809 catalunya=29029, bcm=14639, cataluna=13043, architecture=9275, art=7209, ashld=6074,
22 San Francisco 242558 usa=10639, sf=8519, city=6388, cameraphone=6115, ca=4410, goldengatepark=4155, tra
23 Boston 235339 massachusetts=13015, dewolf=11577, bw=7560, male=6072, nickdewolf=5663, nick=5663, guys=5657
24 London 224447 uk=26678, canon=8656, travel=1376, art=824, city=4572, eukopa=4385, bw=4373, people=4370,
25 Amsterdam 220168 holland=17325, nederland=12001, amsterdam=1703, city=8666, bicycle=7982, bike=7514,
```

4 Appendix

The HDFS location of our final output files from various executions.

4.1 Question 1

- /user/xguo8788/q1
- /user/yzha3522/q1

4.2 Question 2

- /user/xguo8788/q2
- /user/yzha3522/q2

4.3 Question 3

- /user/xguo8788/q3
- /user/yzha3522/q3