

Task 1

Objective: Create a map reduce program to find mutual friends for the following input

A -> B C D

B -> A C D E

C -> A B D E

D -> A B C E

E -> B C D

Algorithm:

I am doing the following steps:

Mapper:

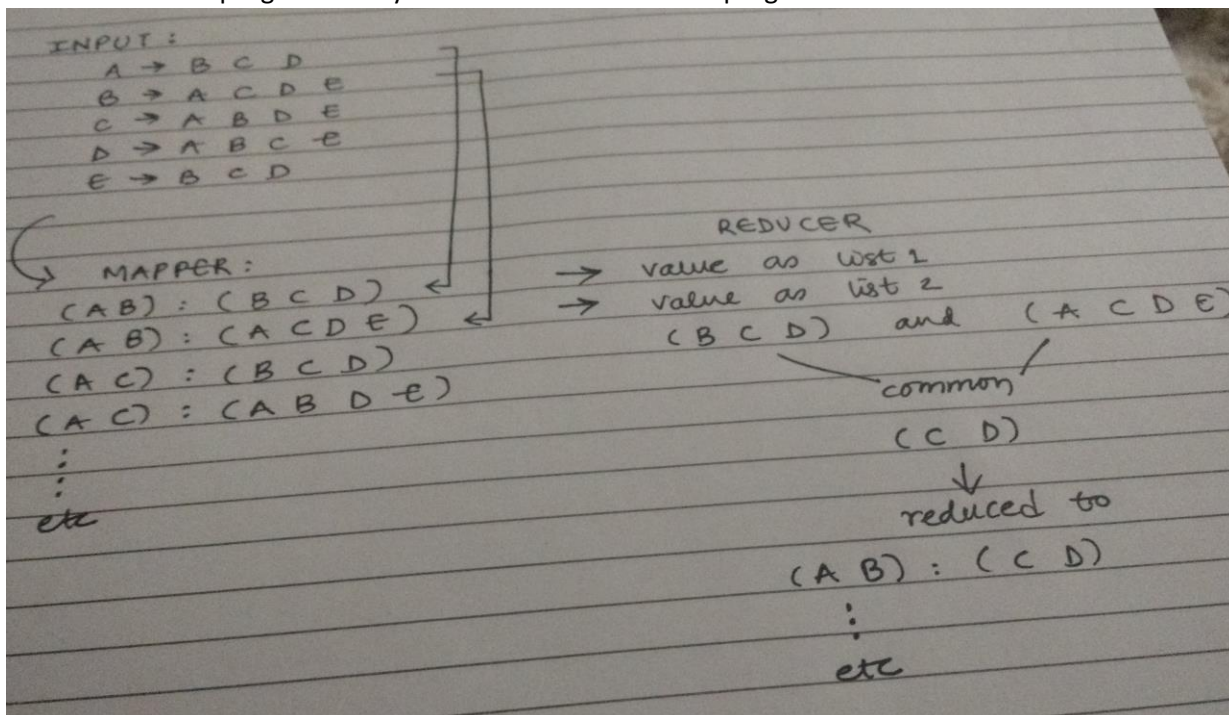
1. Split the input file line by line
2. Split the line by '-'>' format. The left side of the split gives the user name and the right side gives all the friends.
3. The right split can be further split by space to get a list of friends.
4. For each friend f, store the key as sorted (user, f)
5. The mapper program will output (key, value) where value is the original list of friends from the right split.

Reducer:

1. Store the value into lists.
2. Compare the lists to find common elements
3. Output (key, value) where value is the common friends

Main Program:

The main program is very much like the word count program.



Code:

Comments are provided for each important step in the code.

Mapper:

```
MutualFriends.java x
22 public void map(Object key, Text value, Context context
23 ) throws IOException, InterruptedException {
24
25     // Each line separated by newline character
26     StringTokenizer itr = new StringTokenizer(value.toString(), delim: "\n");
27     String line = null;
28     String[] lineSplit = null, friends = null;
29     String f1, f2;
30     int compare;
31     Text keyFriends = new Text();
32
33     while (itr.hasMoreTokens()) {
34         // Get the first line
35         line = itr.nextToken();
36
37         // Split the line into user and friends by "->"
38         // Input: A -> B C D
39         lineSplit = line.split( regex: "-> ");
40
41         // Get friends in array
42         friends = lineSplit[1].split( regex: " ");
43
44         // Get the user as f1 (Need to sort)
45         f1 = lineSplit[0];
46
47         // Get friends in loop
48         for (int i = 0; i < friends.length; i++) {
49
50             // Get friend in variable f2
51             f2 = friends[i];
52
53             // Sort
54             compare = f1.compareTo(f2);
55             if (compare < 0) {
56                 keyFriends.set(f1 + " " + f2);
57             } else {
58                 keyFriends.set(f2 + " " + f1);
59             }
60             context.write(keyFriends, new Text( string: ":" + lineSplit[1]));
61         }
}
```

Reducer:

```
public static class FriendReducer extends Reducer<Text, Text, Text, Text>{
    public void reduce(Text key, Iterator<Text> values,
        OutputCollector<Text, Text> output, Reporter reporter) throws IOException{

        Text[] frndLists = new Text[2];

        int cnt = 0;
        while(values.hasNext()){
            frndLists[cnt++] = new Text(values.next());
        }

        //Get both list values into lists
        String[] list1 = frndLists[0].toString().split( regex: " ");
        String[] list2 = frndLists[1].toString().split( regex: " ");

        //Compare and get common values
        List<String> list = new LinkedList<String>();
        for(String friend1 : list1){
            for(String friend2 : list2){
                if(friend1.equals(friend2)){
                    list.add(friend1);
                }
            }
        }

        // Resulting common friends
        StringBuffer res = new StringBuffer();
        for(int i = 0; i < list.size(); i++){
            res.append(list.get(i));
            if(i != list.size() - 1)
                res.append(" ");
        }
        output.collect(key, new Text(res.toString()));
    }
}
```

Main:

```
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, jobName: "mutual friends");
    job.setJarByClass(MutualFriends.class);
    job.setMapperClass(FriendMapper.class);
    job.setCombinerClass(FriendReducer.class);
    job.setReducerClass(FriendReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(Text.class);
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    System.exit(job.waitForCompletion( verbose: true) ? 0 : 1);
}
}
```

```
Project ▾
  MutualFriends ~\Desktop\5590\Code
    .idea
    .settings
    output
      _SUCCESS.crc
      .part-00000.crc
      _SUCCESS
      part-00000
    src
      com
        friends
        target
      .classpath
      .gitignore
      .project
      input
      pom.xml
    External Libraries
    Scratches and Consoles

MutualFriends.java x
83 //Compare and get common values
84 List<String> list = new LinkedList<String>();
85 for(String friend1 : list1){
86     for(String friend2 : list2){
87         if(friend1.equals(friend2)){
88             list.add(friend1);
89         }
90     }
91 }
92
93 // Resulting common friends
94 StringBuffer res = new StringBuffer();
95 for(int i = 0; i < list.size(); i++){
96     res.append(list.get(i));
97     if(i != list.size() - 1)
98         res.append(" ");
99 }
100 output.collect(key, new Text(res.toString()));
101 }
102
103
104
105 public static void main(String[] args) throws Exception {
106     Configuration conf = new Configuration();
107     Job job = Job.getInstance(conf, "mutual friends");
108     job.setJarByClass(MutualFriends.class);
109     job.setMapperClass(FriendMapper.class);
110     job.setCombinerClass(FriendReducer.class);
111     job.setReducerClass(FriendReducer.class);
112     job.setOutputKeyClass(Text.class);
113     job.setOutputValueClass(Text.class);
114     FileInputFormat.addInputPath(job, new Path(args[0]));
115     FileOutputFormat.setOutputPath(job, new Path(args[1]));
116     System.exit(job.waitForCompletion(verbose: true) ? 0 : 1);
117 }
118
119
120
```

Input:

```
MutualFriends.java x input x
1 A -> B C D
2 B -> A C D E
3 C -> A B D E
4 D -> A B C E
5 E -> B C D
```

Output:

```
Project ▾
  MutualFriends ~\Desktop\5590\Code
    .idea
    .settings
    output
      _SUCCESS.crc
      .part-00000.crc
      _SUCCESS
      part-00000
    src
      com
        friends
        target
      .classpath
      .gitignore
      .project
      input
      pom.xml
    External Libraries
    Scratches and Consoles

MutualFriends.java x part-00000 x
1 A B : C D
2 A C : B D
3 A D : B C
4 B C : A D E
5 B D : A C E
6 B E : C D
7 C D : A B E
8 C E : B D
9 D E : B C
10
```

Task 2

Objective:

a) Consider one of the use case and use a simple dataset. Describe the use case considered based on your assumptions, report the dataset, its fields, datatype etc.

Use Case Considered = Coursera

Problem: Need a high performing and highly scalable database where new features can be easily added.

Solution: They chose Cassandra. But in this task, we will compare both HBase and Cassandra

Dataset:

The table should contain columns like user_id, user_name, course_id, course_name, course_completion_status, course_start_date etc.

The user_id and course_id columns would be numeric and the user_name, course_name and course_completion_status columns would be text. Column course_start_date is date type.

b) Use HBase to implement a Solution for the use case. Report at least 3 queries, their input and output. The query's relevance towards solving the use case is important.

HBase table creation:

```
create 'user_courses', 'user_info', 'course_info'
```

HBase insert rows:

```
put 'user_courses', '1', 'user_info:user_id', '2341'
```

```
put 'user_courses', '1', 'user_info:user_name', 'Avni Mehta'
```

```
put 'user_courses', '1', 'course_info:course_id', '21'
```

```
put 'user_courses', '1', 'course_info:course_name', 'Python'
```

```
put 'user_courses', '1', 'course_info:course_completion_status', 'In Progress'
```

```
put 'user_courses', '1', 'course_info:course_start_date', '04/23/2018'
```

```
put 'user_courses', '2', 'user_info:user_id', '2341'
```

```
put 'user_courses', '2', 'user_info:user_name', 'Avni Mehta'
```

```
put 'user_courses', '2', 'course_info:course_id', '341'
```

```
put 'user_courses', '2', 'course_info:course_name', 'AWS'
```

```
put 'user_courses', '2', 'course_info:course_completion_status', 'Complete'
```

```
put 'user_courses', '2', 'course_info:course_start_date', '06/18/2018'
```

```
put 'user_courses', '3', 'user_info:user_id', '10'
```

```
put 'user_courses', '3', 'user_info:user_name', 'Hardik Mehta'
```

```
put 'user_courses', '3', 'course_info:course_id', '21'
```

```
put 'user_courses', '3', 'course_info:course_name', 'Python'
```

```
put 'user_courses', '3', 'course_info:course_completion_status', 'Complete'
put 'user_courses', '3', 'course_info:course_start_date', '01/23/2018'
```

```
hbase(main):001:0>
hbase(main):002:0* create 'user_courses', 'user_info', 'course_info'
0 row(s) in 1.9190 seconds

=> Hbase::Table - user_courses
hbase(main):003:0>
hbase(main):004:0* put 'user_courses', '1', 'user_info:user_id', '2341'
0 row(s) in 0.5100 seconds

hbase(main):005:0> put 'user_courses', '1', 'user_info:user_name', 'Avni Mehta'
0 row(s) in 0.0210 seconds

hbase(main):006:0>
hbase(main):007:0* put 'user_courses', '1', 'course_info:course_id', '21'
0 row(s) in 0.0220 seconds

hbase(main):008:0> put 'user_courses', '1', 'course_info:course_name', 'Python'
0 row(s) in 0.0130 seconds

hbase(main):009:0> put 'user_courses', '1', 'course_info:course_completion_status', 'In Progress'
0 row(s) in 0.0130 seconds

hbase(main):010:0> put 'user_courses', '1', 'course_info:course_start_date', '04/23/2018'
0 row(s) in 0.0110 seconds

hbase(main):011:0>
hbase(main):012:0* put 'user_courses', '2', 'user_info:user_id', '2341'
0 row(s) in 0.0140 seconds

hbase(main):013:0> put 'user_courses', '2', 'user_info:user_name', 'Avni Mehta'
0 row(s) in 0.0220 seconds

hbase(main):014:0>
hbase(main):015:0* put 'user_courses', '2', 'course_info:course_id', '341'
0 row(s) in 0.0040 seconds

hbase(main):016:0> put 'user_courses', '2', 'course_info:course_name', 'AWS'
0 row(s) in 0.0120 seconds

hbase(main):017:0> put 'user_courses', '2', 'course_info:course_completion_status', 'Complete'
0 row(s) in 0.0090 seconds

hbase(main):018:0> put 'user_courses', '2', 'course_info:course_start_date', '06/18/2018'
0 row(s) in 0.0060 seconds

hbase(main):019:0>
hbase(main):020:0* put 'user_courses', '3', 'user_info:user_id', '10'
0 row(s) in 0.0080 seconds
```

HBase queries:

```
# Scan all rows of table
scan 'user_courses'
```

```
hbase(main):029:0> scan 'user_courses'

COLUMN+CELL
ROW      column=course_info:course_completion_status, timestamp=1529368977311, value=In Progress
1        column=course_info:course_id, timestamp=1529368977144, value=21
1        column=course_info:course_name, timestamp=1529368977235, value=Python
1        column=course_info:course_start_date, timestamp=1529368977393, value=04/23/2018
1        column=user_info:user_id, timestamp=1529368976735, value=2341
1        column=user_info:user_name, timestamp=1529368976904, value=Avni Mehta
2        column=course_info:course_completion_status, timestamp=1529368977786, value=Complete
2        column=course_info:course_id, timestamp=1529368977709, value=341
2        column=course_info:course_name, timestamp=1529368977745, value=AWS
2        column=course_info:course_start_date, timestamp=1529368977842, value=06/18/2018
2        column=user_info:user_id, timestamp=1529368977495, value=2341
2        column=user_info:user_name, timestamp=1529368977583, value=Avni Mehta
3        column=course_info:course_completion_status, timestamp=1529368978141, value=Complete
3        column=course_info:course_id, timestamp=1529368978066, value=21
3        column=course_info:course_name, timestamp=1529368978113, value=Python
3        column=course_info:course_start_date, timestamp=1529368979523, value=01/23/2018
3        column=user_info:user_id, timestamp=1529368977889, value=10
3        column=user_info:user_name, timestamp=1529368977982, value=Hardik Mehta
3 row(s) in 0.2230 seconds
```

```
# Get all rows in a specific a timerange (eg. courses enrolled today)
scan 'user_courses', {TIMERANGE => [1303668804, 1303668904]}
```

```
hbase(main):001:0> scan 'user_courses', {TIMERANGE => [1529368977311, 1539368977311]}
ROW COLUMN+CELL
1 column=course_info:course_completion_status, timestamp=1529368977311, value=In Progress
1 column=course_info:course_start_date, timestamp=1529368977393, value=04/23/2018
2 column=course_info:course_completion_status, timestamp=1529368977786, value=Complete
2 column=course_info:course_id, timestamp=1529368977709, value=341
2 column=course_info:course_name, timestamp=1529368977745, value=AWS
2 column=course_info:course_start_date, timestamp=1529368977842, value=06/18/2018
2 column=user_info:user_id, timestamp=1529368977495, value=2341
2 column=user_info:user_name, timestamp=1529368977583, value=Avni Mehta
```

Get user name and course name

```
scan 'user_courses', {COLUMNS => ['user_name', 'course_name']}
```

```
hbase(main):007:0> scan 'user_courses', {COLUMNS => ['user_info']}
ROW COLUMN+CELL
1 column=user_info:user_id, timestamp=1529368976735, value=2341
1 column=user_info:user_name, timestamp=1529368976904, value=Avni Mehta
2 column=user_info:user_id, timestamp=1529368977495, value=2341
2 column=user_info:user_name, timestamp=1529368977583, value=Avni Mehta
3 column=user_info:user_id, timestamp=1529368977889, value=10
3 column=user_info:user_name, timestamp=1529368977982, value=Hardik Mehta
3 row(s) in 0.0820 seconds
```

c) Use Cassandra to implement a Solution for the use case. Report at least 3 queries, their input and output. The query's relevance towards solving the use case is important.

Table Creation

```
create keyspace coursera with replication={'class':'SimpleStrategy', 'replication_factor':1};
```

USE coursera;

```
CREATE TABLE user_courses (
    user_id int,
    user_name text,
    course_id int,
    course_name text,
    course_completion_status text,
    course_start_date timestamp,
    PRIMARY KEY(user_id, course_id)
);
```

```
cqlsh>
cqlsh> create keyspace coursera with replication={'class':'SimpleStrategy', 'replication_factor':1};
cqlsh>
cqlsh> USE coursera;
cqlsh:coursera>
cqlsh:coursera> CREATE TABLE user_courses ( user_id int, user_name text, course_id int, course_name text, course_completion_status text, course_start_date timestamp, PRIMARY KEY (user_id, course_id));
cqlsh:coursera>
cqlsh:coursera> select * from user_courses;

 user_id | course_id | course_completion_status | course_name | course_start_date | user_name
-----+-----+-----+-----+-----+-----
--+-
(0 rows)
cqlsh:coursera>
```

Insert Records

```
INSERT INTO user_courses (user_id, user_name, course_id, course_name, course_completion_status,
course_start_date)
```

```
VALUES (2341, 'Avni Mehta', 21, 'Python', 'In Progress', '2018-02-05');
```

```
INSERT INTO user_courses (user_id, user_name, course_id, course_name, course_completion_status,
course_start_date)
```

```
VALUES (2341, 'Avni Mehta', 15, 'AWS', 'Complete', '2018-05-05');
```

```
INSERT INTO user_courses (user_id, user_name, course_id, course_name, course_completion_status,
course_start_date)
```

```
VALUES (10, 'Hardik Mehta', 21, 'Python', 'Complete', '2017-05-01');
```

```
cqlsh:coursea> INSERT INTO user_courses (user_id, user_name, course_id, course_
name, course_completion_status, course_start_date) VALUES (2341, 'Avni Mehta', 2
1, 'Python', 'In Progress', '2018-02-05');
cqlsh:coursea>
cqlsh:coursea> INSERT INTO user_courses (user_id, user_name, course_id, course_
name, course_completion_status, course_start_date) VALUES (2341, 'Avni Mehta', 1
5, 'AWS', 'Complete', '2018-05-05');
cqlsh:coursea>
cqlsh:coursea> INSERT INTO user_courses (user_id, user_name, course_id, course_
name, course_completion_status, course_start_date) VALUES (10, 'Hardik Mehta', 2
1, 'Python', 'Complete', '2017-02-05');
cqlsh:coursea>
```

Cassandra Queries

Scan all rows of table

```
select * from user_courses
```

```
cqlsh:coursea> select * from user_courses;

 user_id | course_id | course_completion_status | course_name | course_start_date
         | user_name
-----+-----+-----+-----+-----+-----
-----+-----+-----+-----+-----+-----
      2341 |      15 | Complete | AWS | 2018-05-05 04:00
:00+0000 | Avni Mehta
      2341 |      21 | In Progress | Python | 2018-02-05 05:00
:00+0000 | Avni Mehta
       10 |      21 | Complete | Python | 2017-02-05 05:00
:00+0000 | Hardik Mehta

(3 rows)
cqlsh:coursea>
```

Get all rows in a specific a timerange

```
select * from user_courses where course_start_date > '2017-05-01' ALLOW FILTERING;
```

```
cqlsh:coursera> select * from user_courses where course_start_date > '2017-03-01'
ALLOW FILTERING;
```

user_id	course_id	course_completion_status	course_name	course_start_date	user_name
2341	15	Complete	AWS	2018-05-05 04:00:00+0000	Avni Mehta
2341	21	In Progress	Python	2018-02-05 05:00:00+0000	Avni Mehta

(2 rows)

Get user name and course name

```
select user_name, course_name from user_courses;
```

```
cqlsh:coursera>
cqlsh:coursera> select user_name, course_name from user_courses;
```

user_name	course_name
Avni Mehta	AWS
Avni Mehta	Python
Hardik Mehta	Python

(3 rows)

Demo Video:

Video for the demo can be found at <https://youtu.be/iPJVIpPsUOM>