

# Cloud Computing Practical 4

## Map-Reduce Programming Model

Eoghan Hogan  
17335293  
Eoghan.hogan@ucdconnect.ie

October 23, 2020

## Contents

<b>1</b>	<b>Part 1</b>	<b>3</b>
1.1	Part 1 - Question 1 - Calculate Size MapReduce Signatures . . .	3
1.2	Part 1 - Question 2 - Calculate Size MapReduce . . . . .	4
1.3	Part 1 - Question 3 - Multiplication Signatures . . . . .	5
1.4	Part 1 - Question 4 A - Matrix Multiplication Mapper . . . . .	7
1.5	Part 1 - Question 4 B - Matrix Multiplication Reducer . . . . .	8
1.6	Part 1 - Question 5 A - Matrices Don't Fit . . . . .	9
1.7	Part 1 - Question 5 B - New Mapper . . . . .	10
<b>2</b>	<b>Part 2</b>	<b>11</b>
2.1	Part 2 - Question 1 . . . . .	11
2.2	Part 2 - Question 2 . . . . .	12

## 1 Part 1

### 1.1 Part 1 - Question 1 - Calculate Size MapReduce Signatures

The Mapper Function will take in the Corpus Name and the Corpus (As it expects key, value pairs)

It will output the URL of a page in the Corpus along with its Size

---

**Algorithm 1** The Mapper Function Input and Output

---

```
1: procedure MAP(CorpusName, Corpus)  
2:                                     ▷ Do something to get the url and size...  
3:   emit (HostURL, Size)  
4: end procedure
```

---

The reduce Will take the Urls and Size and emit the sum of the sizes for specific URL's

---

**Algorithm 2** The Reducer Function Input and Output

---

```
1: procedure REDUCE(HostURL, Sizes)          ▷ Sizes will be an iterator  
2:   emit (HostURL, sizeSum)  
3: end procedure
```

---

## 1.2 Part 1 - Question 2 - Calculate Size MapReduce

The implementations of the MapReduce algorithms is fairly simple we simply need to go through the corpus and filter the lines to get the information we are looking for.

---

**Algorithm 3** The Mapper Function

---

```
1: procedure MAPPER(corpusName, Corpus)
2:   for line in Corpus do
3:     HostURL  $\leftarrow$  split(line, ", ")[0]
4:     Size  $\leftarrow$  split(line, ", ")[1]
5:     Size  $\leftarrow$  bytes(Size)
6:     emit (HostURL, Size)
7:   end for
8: end procedure
```

---

In the reducer function we just iterate over the values and do the summation as required.

---

**Algorithm 4** The Reducer Function

---

```
1: procedure REDUCER(HostURL, Sizes)            $\triangleright$  Sizes will be an iterator
2:   sizeSum  $\leftarrow$  0
3:   for bytes in Sizes do
4:     sizeSum  $\leftarrow$  bytes + sizeSum
5:   end for
6:   emit (HostURL, sizeSum)
7: end procedure
```

---

### 1.3 Part 1 - Question 3 - Multiplication Signatures

For matrix Multiplication done via Map reduce things get a bit more Complicated. I found it best to start with the Output Matrix and work my way back. to give a little visual help to the algorithm lets start from the end and work our way backwards.

So we have two Matrices A and B

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}, B = \begin{pmatrix} e & f \\ g & h \end{pmatrix}$$

and the product is this:

$$A * B = C = \begin{pmatrix} a * e + b * g & a * f + b * h \\ c * e + d * g & c * f + d * h \end{pmatrix}$$

But we can rewrite the product by indexing the original Matrices

$$C = \begin{pmatrix} A_{0,0} * B_{0,0} + A_{0,1} * B_{1,0} & A_{0,0} * B_{0,1} + A_{0,1} * B_{1,1} \\ A_{1,0} * B_{0,0} + A_{1,1} * B_{1,0} & A_{1,0} * B_{0,1} + A_{1,1} * B_{1,1} \end{pmatrix}$$

okay now lets say we want to compute

$$A_{0,0} * B_{0,0} + A_{0,1} * B_{1,0}$$

but how do we know which  $A_{0,0}$  or  $A_{0,1}$  we are dealing with? they both appear twice! the row of the final matrix. lets rewrite the final product but include the indexes also

$$C = \begin{pmatrix} (0,0)A_{0,0} * B_{0,0} + A_{0,1} * B_{1,0} & (0,1)A_{0,0} * B_{0,1} + A_{0,1} * B_{1,1} \\ (1,0)A_{1,0} * B_{0,0} + A_{1,1} * B_{1,0} & (1,1)A_{1,0} * B_{0,1} + A_{1,1} * B_{1,1} \end{pmatrix}$$

Now the pattern easily jumps out of us. When we are finding the answer for  $C_{i,j}$  we are going to be using all  $A_{i,k1}$  and all  $B_{k2,j}$ . this means we will need to replicate every value in a the same number of times as columns in B since that's the amount of Columns in C and every value in B the same number of times of Rows in A as that's the number of Rows in C

---

#### Algorithm 5 The Matrix Multiplication Mapper Function Signature

---

```

1:                                     ▷ A is a LxM matrix anf B is an MxN matrix
2: procedure MAP( $A, B$ )
3:                                     ▷ duplicate values in A as many times as columns in B
4:                                     ▷ k = column index B, i = row index A, j = col index A
5:   emit (( $i, k$ ), ("A", j,  $A_{i,j}$ ))
6:                                     ▷ duplicate values in B as many times as rows in A
7:                                     ▷ i = row index A, j = row index B, i = col index B
8:   emit (( $i, k$ ), ("B", j,  $B_{k,j}$ ))
9: end procedure

```

---

For the reducer we now know we are going to take in values that look like  $(M, j, M_{(iork),j})$  and that all our Values that the reduce takes in will represent one value in the final matrix. so we need to sort the Values to so that we end up all the  $A$  values in one array and all the  $B$  values in another

Then we can simple do this with a for loop over  $j$  and indexing the arrays:

$$A_{i,j} * B_{k,j} + A_{i,j} * B_{k,j}$$

---

**Algorithm 6** The Matrix Multiplication Reducer Function Signature

---

```

1:           ▷ positions= (row, col), tuples = (matrix_id, idx, value)
2: procedure REDUCE(positions, tuples)
3:                                     ▷ sort values by matrix
4:                                     ▷ do the summation and multiplication
5:   emit (positions, S)
6: end procedure

```

---

## 1.4 Part 1 - Question 4 A - Matrix Multiplication Mapper

---

**Algorithm 7** The Matrix Multiplication Mapper Function

---

```
1:                                     ▷ A is a LxM matrix and B is an MxN matrix
2: procedure MAP( $A, B$ )
3:    $L, M \leftarrow \text{shape}(A)$ 
4:    $M, N \leftarrow \text{Shape}(B)$ 
5:                                     ▷ duplicate values in A as many times as columns in B
6:   for  $k \leftarrow 0, N$  do
7:     for  $i \leftarrow 0, L$  do
8:       for  $j \leftarrow 0, M$  do
9:         ▷  $k$  = column index B,  $i$  = row index A,  $j$  = col index A
10:        emit  $((i, k), ("A", j, A_{i,j}))$ 
11:      end for
12:    end for
13:  end for
14:                                     ▷ duplicate values in B as many times as rows in A
15:  for  $i \leftarrow 0, L$  do
16:    for  $j \leftarrow 0, M$  do
17:      for  $k \leftarrow 0, N$  do
18:        ▷  $i$  = row index A,  $j$  = row index B,  $k$  = column index B
19:        emit  $((i, k), ("B", j, B_{j,k}))$ 
20:      end for
21:    end for
22:  end for
23: end procedure
```

---

## 1.5 Part 1 - Question 4 B - Matrix Multiplication Reducer

---

**Algorithm 8** The Matrix Multiplication Reducer Function

---

```
1:           ▷ positions= (row, col), tuples = (matrix_id, idx, value)
2: procedure REDUCE(positions, tuples)
3:    $l \leftarrow \text{length}(\text{tuples})/2$ 
4:    $a\_vals \leftarrow \text{array}(0, l)$ 
5:    $b\_vals \leftarrow \text{array}(0, l)$ 
6:                                     ▷ sort values by matrix
7:   for  $t$  in  $\text{tuples}$  do
8:      $idx \leftarrow t[1]$ 
9:     if  $t[0] = "A"$  then
10:       $a\_vals[idx] \leftarrow t[2]$ 
11:    else
12:       $b\_vals[idx] \leftarrow t[2]$ 
13:    end if
14:  end for
15:                                     ▷ do the summation and multiplication
16:   $s \leftarrow 0$ 
17:  for  $i$  in  $0, l$  do
18:     $s \leftarrow s + (a\_vals[i] * b\_vals[i])$ 
19:  end for
20:  emit ( $\text{positions}, s$ )
21: end procedure
```

---



## 1.6 Part 1 - Question 5 A - Matrices Don't Fit

### What if the Matrices do not fit into The Mappers memory?

Well luckily we can keep the Reducer the same we just need to change the input to the mapper and modify what it does. Basically the strategy to deal with this would be to Pre-Processes the Matrices before supplying them to the map functions.

to Pre-Process matrix A we would need to split off the rows such that they are 1xN 2d Arrays. then we would send that to the mapper with the row index and the array "A" id. We also would need The number of columns of B.

To Pre-Process matrix B we would need to split B by its columns so our mapper would take in the Nx1 array representing the column, the index of that column and the number of rows of A.

Its actually quite simple conceptually when we look at the Original Mapper function.

The original Mapper function went through All of A col(B) times and all of B row(A) times, however this time we don't have all of A or B. we only have vectors from them and we are not even sure how many. we Do know however if the vectors are from A or B.

So thinking about the Original Function. while we don't have the whole matrix we do know which column of row the vector belongs to. so we can fix that value and essentially go over it rows(A) or Cols(B) times like would've happened to that specific vector in the original Function.

Thus we can emit pretty much the same information just with a fixed value

---

**Algorithm 9** The Matrix Multiplication Mapper Function Signature

---

```
1: matrixID is "A" or "B"
2: tuple[0] = vector that was split off
3: tuple[1] = col index if B or row index if A
4: tuple[2] = cols(A) if vector is from A
5: tuple[2] = Rows(B) if vector is from B
6: tuple[3] = cols(B) if vector is from A
7: tuple[3] = Rows(A) if vector is from B
8: procedure MAP(matrixID, tuple)
9:   if matrixID = "A" then
10:      $\triangleright$  for each tuple duplicate the values in vector cols(B) times
11:      $\triangleright$  k = column index B, i = row index A, j = col index A
12:     emit ((i, k), ("A", j, Ai,j))
13:   else
14:      $\triangleright$  for each tuple duplicate the values in vector rows(A) times
15:      $\triangleright$  i = row index A, j = row index B, i = col index B
16:     emit ((i, k), ("B", j, Bk,j))
17:   end if
18: end procedure
```

---

## 1.7 Part 1 - Question 5 B - New Mapper

This is the implementation of the Mapper function I discussed on the last page.

---

**Algorithm 10** The Matrix Multiplication Mapper Function

---

```
1: procedure MAP(matrixID, tuple)
2:   if matrixID = "A" then
3:      $\triangleright$  for each tuple duplicate, the values in vector, cols(B) times
4:      $\triangleright$  k = column index B, i = row index A, j = col index A
5:      $i \leftarrow \text{tuple}[1]$ 
6:     for  $t$  in tuple do
7:        $vec = t[0]$ 
8:       for  $k \leftarrow 0, t[3]$  do
9:         for  $j \leftarrow 0, t[2]$  do
10:          emit (( $i, k$ ), ("A",  $j, A_{0,j}$ ))
11:        end for
12:      end for
13:    end for
14:  else
15:     $\triangleright$  for each tuple duplicate, the values in vector, rows(A) times
16:     $\triangleright$  i = row index A, j = row index B, k = column index B
17:     $k \leftarrow \text{tuple}[1]$ 
18:    for  $t$  in tuple do
19:       $vec = t[0]$ 
20:      for  $i \leftarrow 0, t[3]$  do
21:        for  $j \leftarrow 0, t[2]$  do
22:          emit (( $i, k$ ), ("B",  $j, B_{j,0}$ ))
23:        end for
24:      end for
25:    end for
26:  end if
27: end procedure
```

---

## 2 Part 2

### 2.1 Part 2 - Question 1

From the lecture notes We have the following Function Signatures which I have changed a bit to make more sense to me

K means makes alot more sense than matrix multiplication and thus is easier to follow.

The Mapper will take in a data point and the list of K centroids. then all we need to do is compare the distance from each centroid to the data point and return the centroid that is closest to the data point.

---

**Algorithm 11** The K-Means Mapper Signature

---

```
1: procedure MAP(dataPoint, [centroid])
2:   emit (centroid, dataPoint)
3: end procedure
```

---

The Reducer will then take the centroids as the keys so that all the datapoints that were closest to that centroid will be the values. The all we need to do is calculate the new centroid by getting the mean of the datapoints.

---

**Algorithm 12** The K-Means Reducer Signature

---

```
1: procedure REDUCE(centroid, [dataPoint])
2:   emit (centroid or newCentroid, [dataPoints])
3: end procedure
```

---

we Then continue this until there is no change in the centroid position

## 2.2 Part 2 - Question 2

The implementation details of the Mapper and Reducer Functions

---

**Algorithm 13** The K-Means Mapper

---

```
1: procedure MAP(dataPoint, [centroids])
2:   min  $\leftarrow$  centroids[0]
3:   for c  $\leftarrow$  centroids do
4:     if Distance(dataPoint, c)  $\leq$  Distance(dataPoint, min) then
5:       min  $\leftarrow$  c
6:     end if
7:   end for
8:   emit (min, dataPoint)
9: end procedure
```

---

---

**Algorithm 14** The K-Means Reducer

---

```
1: procedure REDUCE(centroid, dataPoints)
2:   newCentroid  $\leftarrow$  mean(dataPoints)
3:   if centroid  $\neq$  newCentroid then
4:     emit (newCentroid, dataPoints)
5:   else
6:     emit (centroid, dataPoints)
7:   end if
8: end procedure
```

---