

## Assignment - B4

- \* Date of Completion: 28/10/2020
- \* Date of Submission: 04/11/2020
- \* Title: Map Reduce
- \* Problem Statement: Write and implement example of Map Reduce operation with suitable collection using MongoDB
- \* Learning Objectives:
  - To understand concept of Map-Reduce as data processing paradigm
  - To compare mapreduce and aggregation pipeline in MongoDB.
- \* Learning Outcomes:
  - Learn and implement Map-Reduce for condensing large volumes of data into useful aggregated results.
  - To know the difference between the use of mapreduce & aggregation pipeline in MongoDB.
- \* Hardware & Software Requirements:
  - MongoDB 4.4.1 Server, Windows 10 64 bit, Intel i5.
- \* Theory:
 

Map Reduce is a data processing paradigm for condensing large volumes of data into useful aggregated results. It is a powerful & flexible



tool for aggregating tools data. It can solve some problems that are too complex to express using the aggregation frameworks query language, & uses Java Script as its 'query language'.

Syntax:

```

db.collection.mapReduce (
  mapFunction () { emit (<key>, <value>) },
  reduceFunction (<key, value>) { // operation & return },
  {
    query: { < query document > },
    out: < collection name >
  }
)

```

Map & Reduce Phases:

- MongoDB applies the "map phase" to each input document (i.e.; the document in the collection that match the query condition). The map emits key-value pairs.
- For those keys that have multiple values, MongoDB applies the "reduce phase", which collects & condenses the aggregated data. optionally the output of reduced data function may pass through a "finalize" function to further condense or process the result of the aggregation.

MapReduce Results:

The map reduce operation can write results to a collection on return the results in life time, if you

if you write mapreduce output to a collection then you can perform subsequent mapreduce operations on the same input collection that merge replace, merge on reduce new results with previous results.

eg: `ab.orders.mapReduce(`  
`function() { emit (this.cust_id, this.amount); }`  
`function (key, value) { return Array.sum (values) } ,`  
`{ query: { status: "A" },`  
`out: "order totals"`  
`}`  
`)`

```
{
  cust_id: "A123"
  amount: 600
  status: "A"
}
```

```
{ cust_id: 'A1234'
  amount: 350
  status: "D"
}
```

```
{
  cust_id: "B212"
  amount: 100
  status: "A"
}
```

```
{ cust_id: 'A123'
  amount: 300
  status: "A"
}
```

Orders



{	cust_id: "A123"
	amount: 600
	status: "A"
}	
{	cust_id: "B212"
	amount: 100
	status: "A"
}	
{	cust_id: "A123"
	amount: 300
	status: "A"
}	

Map

{	"A123" : [ 600 , 300 ]
}	
{	"B212" : [ 100 ]
}	

Reduce

{	"_id": "A123"	order totals
	value: 900	
}		
{	"_id": "B212"	
	value: 100	
}		

- MongoDB: Aggregation framework & MapReduce comparisons
  - MongoDB MapReduce can perform complex and flexible aggregations that can't be done with aggregation framework query language.
  - Average CPU utilization in case of MapReduce is higher in comparison of aggregation.
  - Use of aggregation pipeline is faster as compared to MapReduce.
  - Aggregation framework returns in line result. Thus cannot be used for further processing unlike in MapReduce result can be stored to a collection.

Output from the single pipeline cannot exceed the BSON document size (ie; 16MB) whereas in mapreduce single emit can hold half of BSON size limit (ie; 8MB).

Use cases of MapReduce & Aggregation:

- Queries that are very complex are difficult to handle in Aggregation Framework can be helped out with map reduce.
- The small datasets will take considerable time to load in mapreduce, be it a small data set or big both will take some amount of time to execute. So it's better to handle large data sets in mapreduce.
- Parallelly Aggregation pipeline can be used for handling small data set.

Conclusion:

We studied the use of mapreduce in MongoDB and implemented examples of the same and studied the difference between mapreduce and aggregation Pipeline and their usecases with comparisons.