



SNS COLLEGE OF TECHNOLOGY



Coimbatore-35
An Autonomous Institution

Accredited by NBA – AICTE and Accredited by NAAC – UGC with 'A++' Grade Approved
by AICTE, New Delhi & Affiliated to Anna University, Chennai

DEPARTMENT OF COMPUTER APPLICATIONS

19CAE716 – DATA SCIENCE
II YEAR III SEM

UNIT – V: MAPREDUCE

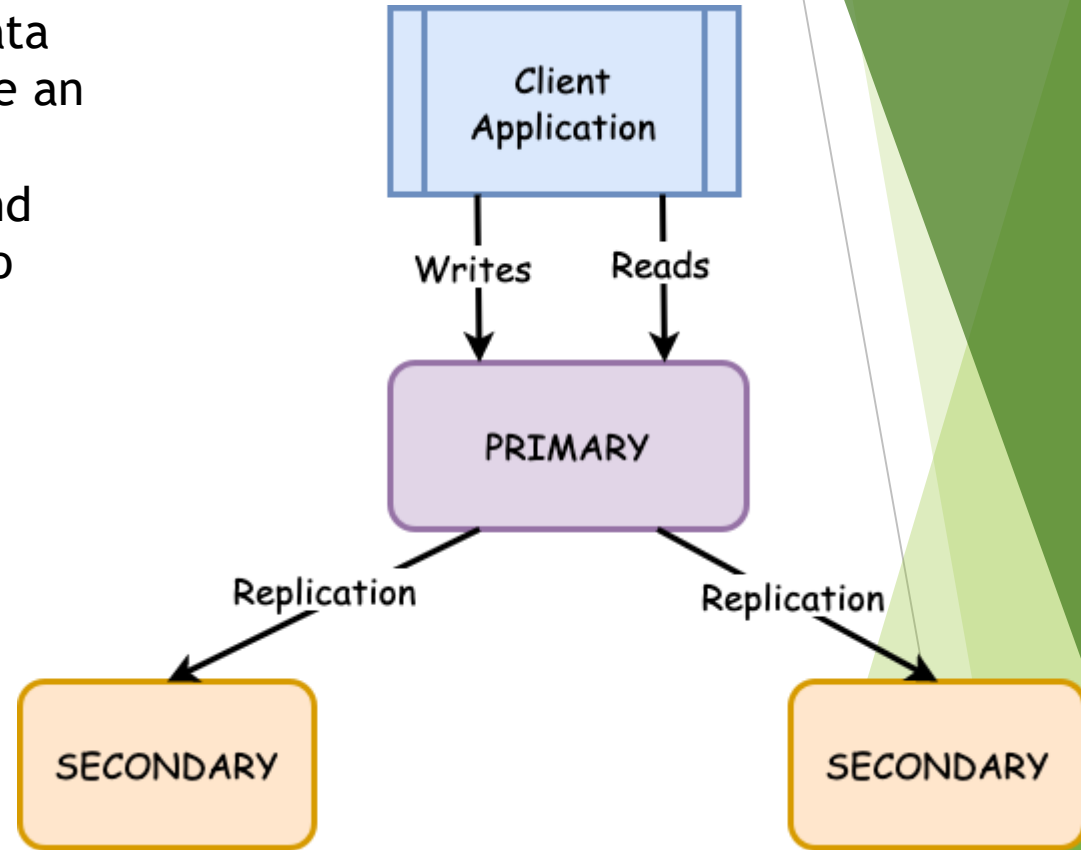


Replication

Replication is the method of duplication of data across multiple servers. For example, we have an application and it reads and writes data to a database and says this server A has a name and balance which will be copied/replicate to two other servers in two different locations.

How replication is formed?

create a ".sh" file create_replicaset.sh and init_mongoreplica.js





Then run the following script :
./create_replicaset.sh

```
EXPLORER
  OPEN EDITORS
    readme.txt
    create_replicaset.sh
  UNTITLED (WORKSPACE)
    CreateReplicaSet
      create_replicaset.sh
      JS init_mongoreplica.js
      readme.txt

CreateReplicaSet > create_replicaset.sh
1  mkdir -p rs1 rs2 rs3 rs4
2  mongod --replSet atique --logpath "1.log" --dbpath rs1 --port 27017 &
3  mongod --replSet atique --logpath "2.log" --dbpath rs2 --port 27018 &
4  mongod --replSet atique --logpath "3.log" --dbpath rs3 --port 27019 &
5  mongod --replSet atique --logpath "4.log" --dbpath rs4 --port 27020 &
```

```
EXPLORER
  OPEN EDITORS
    readme.txt
    JS init_mongoreplica.js
    create_replicaset.sh
  UNTITLED (WORKSPACE)
    CreateReplicaSet
      rs1
      rs2
      rs3
      rs4
      1.log
      2.log

CreateReplicaSet > JS init_mongoreplica.js
1  config = { _id: "atique", members: [
2    { _id: 0, host : "localhost:27017" },
3    { _id: 1, host : "localhost:27018" },
4    { _id: 2, host : "localhost:27019" },
5    { _id: 3, host : "localhost:27020" } ]
6  };
7  rs.initiate(config);
8  rs.status();
```

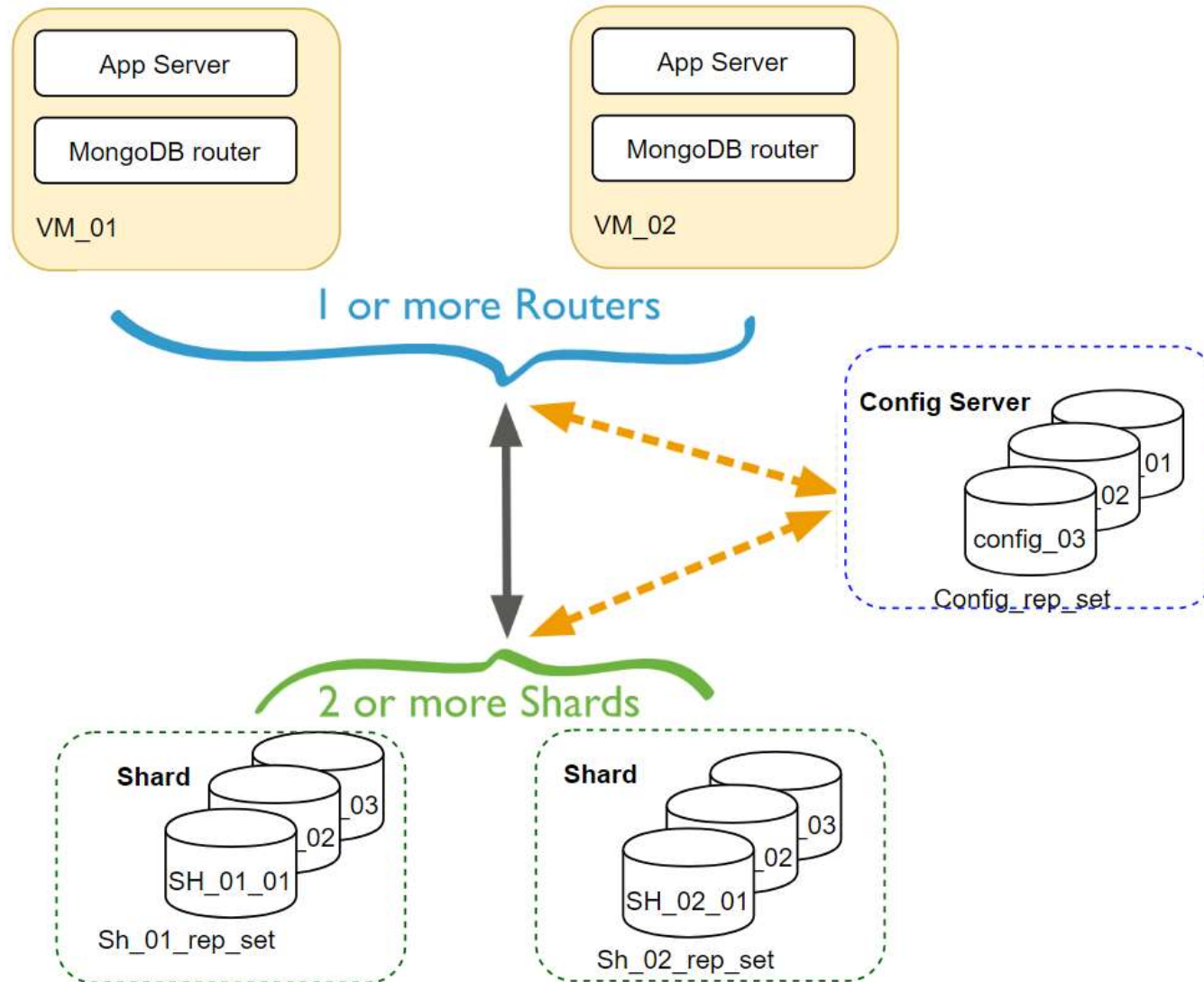


Database Sharding

The need for database sharding arises because as data grows, it becomes more challenging to store it on a single server. It is because the server has limited storage and processing power. As a result, the performance of the database can start to deteriorate.

There are a few things to consider when sharding a database:

1. The data must be divided up in a way that makes sense.
2. The system must handle queries that span multiple servers.
3. The system must be able to recover from failures.



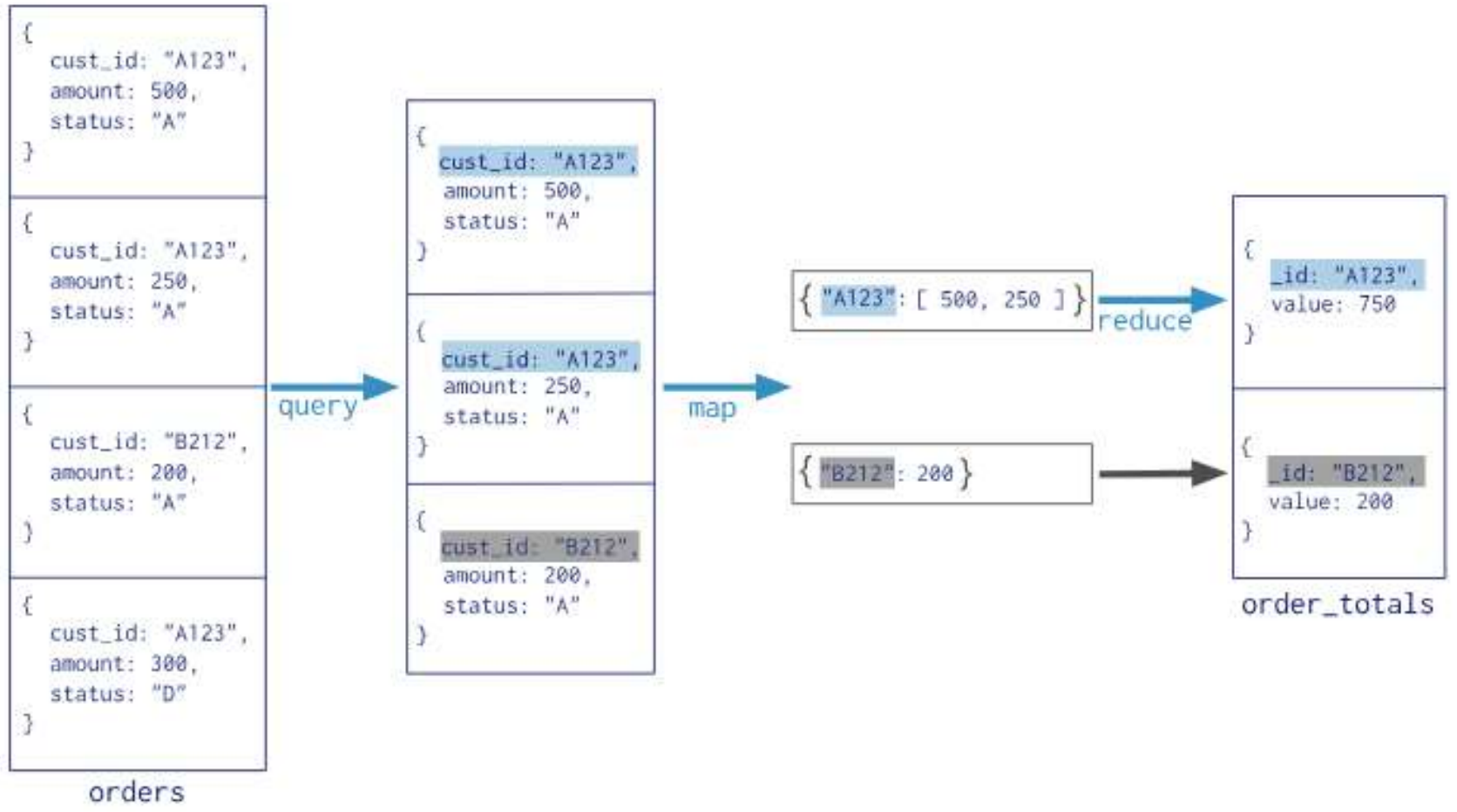


Map-Reduce

Map-reduce is a data processing programming model that helps to perform operations on large data sets and produce aggregated results.

Consider the following map-reduce operation:

```
Collection
↓
db.orders.mapReduce(
  map   → function() { emit( this.cust_id, this.amount ); },
  reduce → function(key, values) { return Array.sum( values ) },
  query → {
  output → { status: "A" },
           out: "order_totals"
  }
)
```





map() function: It uses **emit()** function in which it takes two parameters key and value key.

Here the key is on which we make groups like groups by in MySQL. Example like group by ages or names and the second parameter is on which aggregation is performed like avg(), sum() is calculated on.

reduce() function: It is the step in which we perform our aggregate function like avg(), sum()).

query: Here we will pass the query to filter the resultset.

output: In this, we will specify the collection name where the result will be stored.



```
[> db.employee.find()
{ "_id" : ObjectId("60192cb40cf217478ba935a2"), "name" : "eren", "age" : 25, "rank" : 1 }
{ "_id" : ObjectId("60192cb40cf217478ba935a3"), "name" : "mikasa", "age" : 24, "rank" : 2 }
{ "_id" : ObjectId("60192cb40cf217478ba935a4"), "name" : "jean", "age" : 26, "rank" : 4 }
{ "_id" : ObjectId("60192cb40cf217478ba935a5"), "name" : "conny", "age" : 23, "rank" : 7 }
{ "_id" : ObjectId("60192cb40cf217478ba935a6"), "name" : "sasha", "age" : 25, "rank" : 6 }
{ "_id" : ObjectId("60192cb40cf217478ba935a7"), "name" : "armin", "age" : 24, "rank" : 3 }
> ]
```



**Any
questions?**