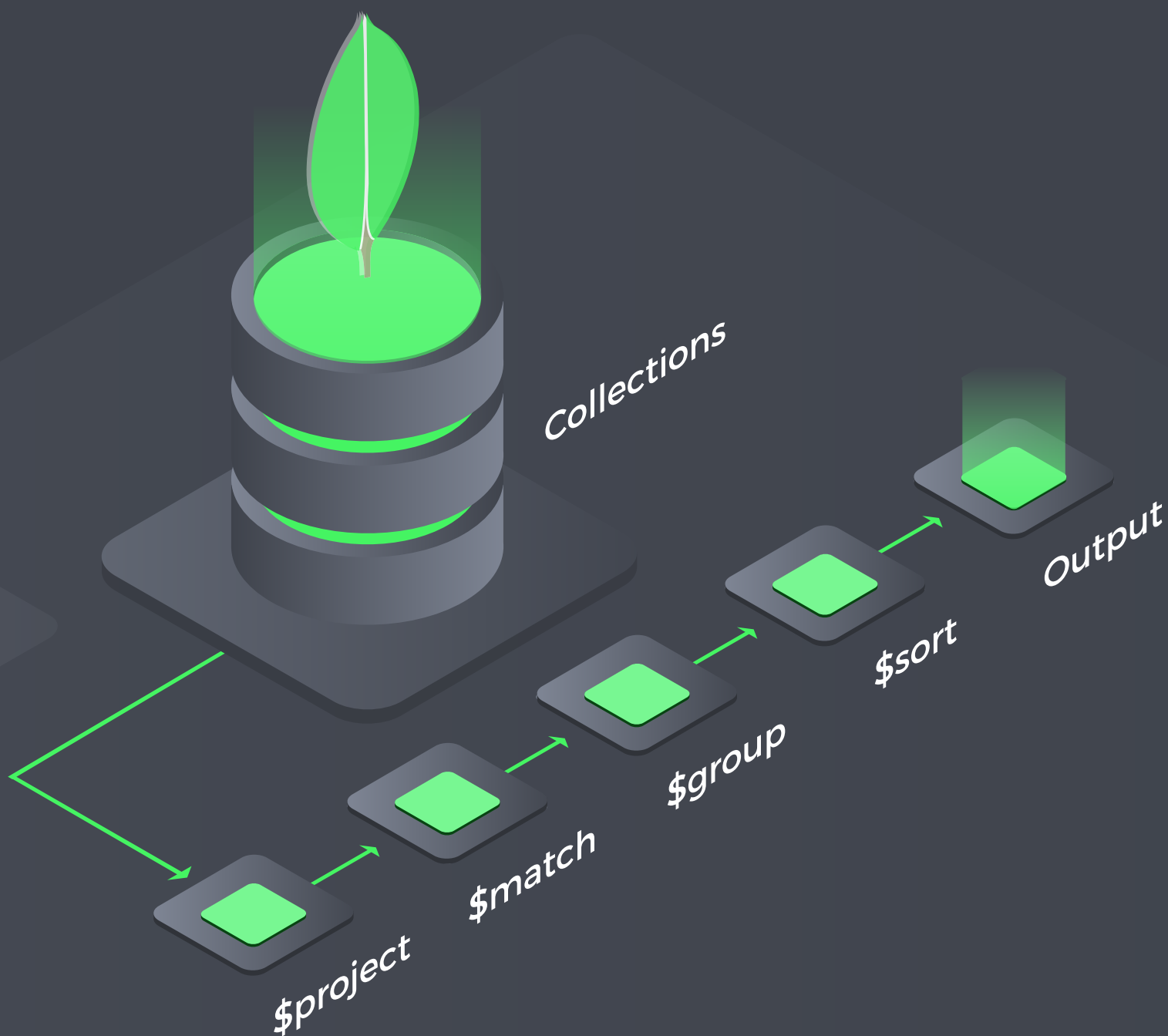


severalnines

MongoDB Aggregation Framework Stages and Pipelining

by Onyancha Brian Henry



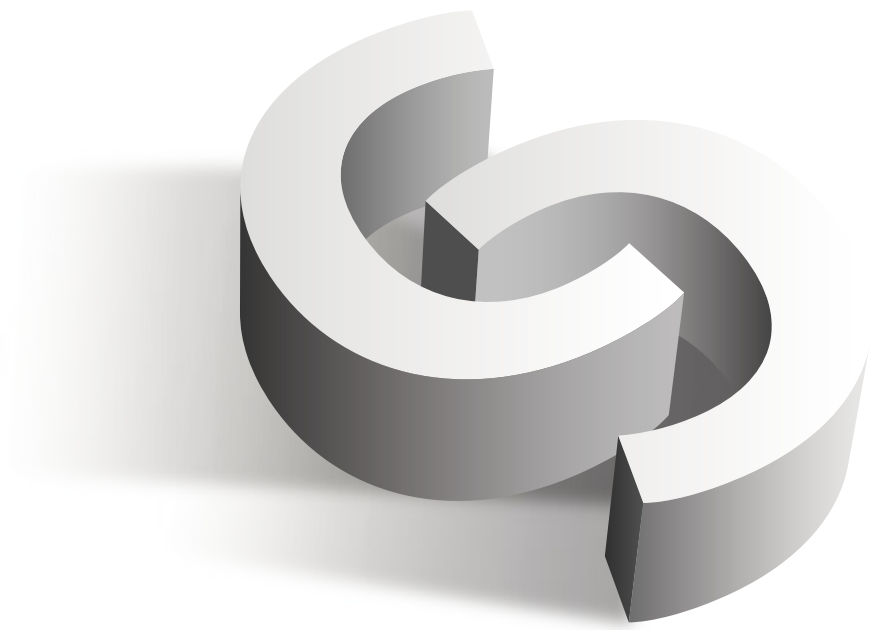


Table of Contents

Introduction	4
What is the Aggregation Framework?	5
Aggregation Pipeline	6
Basic Stages of Aggregation Pipeline	6
\$match	6
\$group	9
\$unwind	10
\$project	12
Points to note	12
\$sort	13
\$sample	15
\$limit	15
\$lookup	16
Aggregation Process	18
Accumulator Operators	21
\$sum	21
\$avg	22
\$max and \$min	23
\$push	24
Similarity of the Aggregation Process in MongoDB with SQL	25
Aggregation Pipeline Optimization	27
Projection Optimization	27
Pipeline Sequence Optimization	27
MapReduce in MongoDB	28
MapReduce JavaScript Functions	30
Incremental MapReduce	31
Comparison Between MapReduce and Aggregation Pipeline in MongoDB	34
Summary	35
About ClusterControl	36
About Severalnines	36
Related Resources	37



Introduction

Using the CRUD find operation while fetching data in MongoDB may sometimes become tedious. For instance, you may want to fetch some embedded documents in a given field but the find operation will always fetch the main document and then it will be upon you to filter this data and select a field with all the embedded documents, scan through it to get ones that match your criteria. Since there is no simple way to do this, you will be forced to use something like a loop to go through all these subdocuments until you get the matching results. However, what if you have a million embedded documents? You will unfortunately get frustrated with how long it will take. Besides, the process will take a lot of your server's random memory and maybe terminate the process before you get all the documents you wanted, as the server document size may be surpassed.

In this paper, we will deep dive into MongoDB's Aggregation Framework and look into the different stages of the Aggregation Pipeline. We'll see how we make use of these stages in an aggregation process. We'll then look at the operators that can assist in the analysis process of input documents. Finally, we'll compare the aggregation process in MongoDB with SQL, as well as the differences between the aggregation process and MapReduce in MongoDB.

What is the Aggregation Framework?

Regarding the limitations associated especially with a large number of documents, there is the need to group them to enhance the scanning process. Aggregation framework is therefore an operation process which manipulates documents in different stages, processes them in accordance with the provided criteria and then return the computed results. Values from multiple documents are grouped together, on which more operations can be performed to return matching results.

Aggregation Pipeline

To scan the documents one by one in order to apply some operation to them will obviously outdo the purpose of aggregation framework because it will consequently take much time to do this. Therefore, the data processing is done at the same time from different stages using the UNIX pipelining technique. Documents from a collection are channeled into a multistage pipeline from which they are converted into aggregated data.

Contrary to the map-reduce functions in MongoDB which is Javascript code interpreted, the aggregation pipeline runs compiled C++ code. MongoDB data is stored in BSON format and for this reason, the map-reduce takes longer to convert this BSON data into JSON format for processing. On the other hand, the aggregation pipeline does not need to perform any conversion.

We can show the aggregation process using a simple flow chart as below:



Basic Stages of Aggregation Pipeline

As mentioned before, documents pass a number of defined stages in order to be filtered to the desired result. The stages that may be involved are:

\$match

Like other MongoDB operations, this uses the standard MongoDB queries to filter documents without any modification and then passes them to the next stage. A document has to match the provided criteria in the query for it to pass to the next stage.

Example:

Let's create a simple collection named `users` and populate it with the data below.

```
1 | {
2 |   "_id" : ObjectId("5b43cbe2106c21d21c776e81"),
3 |   "userId" : NumberLong("1530442083133"),
4 |   "name" : "George", "eyeColor" : "blue",
5 |   "connections" : [
6 |     {
```

```

7         "status" : "Disconnect",
8         "userName" : "Valencia",
9         "name" : "Derrick Clinton",
10        "id" : NumberLong("1530444522597")
11    },
12    {
13        "status" : "Connect",
14        "userName" : "Carliston",
15        "name" : "James Good",
16        "id" : NumberLong("1530444522597")
17    }
18 ]
19 }
20 {
21     "_id" : ObjectId("5b43cbe2106c21d21c776e82"),
22     "userId" : NumberLong("153044201111"),
23     "name" : "Monica",
24     "eyeColor" : "normal",
25     "connections" : [
26         {
27             "status" : "Disconnect",
28             "userName" : "JohnDoh",
29             "name" : "Alex Xian",
30             "id" : NumberLong("153044445903")
31         },
32         {
33             "status" : "Connect",
34             "userName" : "MaryCartie",
35             "name" : "Mary Carey",
36             "id" : NumberLong("1530444522597")
37         }
38     ]
39 }
40 {
41     "_id" : ObjectId("5b43cbe2106c21d21c776e83"),
42     "userId" : NumberLong("15304420836758"),
43     "name" : "Harrison",
44     "eyeColor" : "blue",
45     "connections" : [
46         {
47             "status" : "Connect",
48             "userName" : "Kevin",
49             "name" : "Keni Sems",
50             "id" : NumberLong("34435345345343")
51         },
52         {
53             "status" : "Disconnect",
54             "userName" : "Mayaka",
55             "name" : "Andrew Fake",

```

```

56         "id" : NumberLong("15304445224357")
57     }
58 ]
59 }

```

We now have 3 documents in `users` collection with 2 more embedded documents in the `connections` field in each. Using the `$match` stage, let us return the document with name equal to `Harrison`.

Query

```

1 db.getCollection('users').aggregate([
2   {$match: {'name': 'Harrison'}}
3 ])

```

Result

```

1 {
2   "_id" : ObjectId("5b43cbe2106c21d21c776e83"),
3   "userId" : NumberLong("15304420836758"),
4   "eyeColor" : "blue",
5   "name" : "Harrison",
6   "connections" : [
7     {
8       "status" : "Connect",
9       "userName" : "Kevin",
10      "name" : "Keni Sems",
11      "id" : NumberLong("34435345345343")
12    },
13    {
14      "status" : "Disconnect",
15      "userName" : "Mayaka",
16      "name" : "Andrew Fake",
17      "id" : NumberLong("15304445224357")
18    }
19  ]
20 }

```

In this case, the document with name equal to `Harrison` will be passed to the next stage since it matched the criteria.

In order to achieve the best performance of the `$match` stage, use it early in the aggregation process since it will:

1. Take advantage of the indexes hence become much faster
2. Limit the number of documents that will be passed to the next stage.

However, you must not use the `$where` clause in this `$match` stage since it is catered for within the match condition.

\$group

For a specified expression, data is grouped accordingly in this stage. For every distinct group that is formed, it is passed to the next stage as a document with a unique `_id` field.

The syntax for this group operation is:

```
1 | { $group: { _id <expression>, <field>: {<accumulator>: <expression>}}}
|
```

The accumulator operations that may be involved include: `$sum`, `$avg`, `$max`, `$last`, `$push`. For our `users` collection above, we will group the documents using the `eyeColor` field and see how many groups we will get.

```
1 | db.getCollection('users').aggregate([
2 |     { $group: {
3 |         _id: "$eyeColor",
4 |         }
5 |     }
6 | ])
|
```

Using the `_id` field here we are specifying which criteria we are using to group and in this case we use the `eyeColor` field. The result from this operation is:

```
1 | { "_id" : "blue" }
2 | { "_id" : "normal" }
|
```

We can go further and sum the number of people in each of this group with an accumulator expression of `sum`. I.e.

```
1 | db.getCollection('users').aggregate([
2 |     { $group: {
3 |         _id: "$eyeColor",
4 |         numberOfPeople: { $sum: 1 }
5 |         }
6 |     }
7 | ])
|
```

The result for this query will be

```
1 | { "_id" : "blue", "numberOfPeople" : 2 }
2 | { "_id" : "normal", "numberOfPeople" : 1 }
|
```

The number of people with blue `eyeColor` is 2 and for normal color is 1. Besides, you can fetch the names of people in this groups as an array using the `push` operator and field name in the expression as:

```

1 | db.getCollection('users').aggregate([
2 |   {$group: {
3 |     _id: "$eyeColor",
4 |     names: {$push: "$name"}
5 |   }
6 | }
7 | ])

```

With this operation the result is

```

1 | { "_id" : "blue", "names" : [ "George", "Harrison" ] }
2 | { "_id" : "normal", "names" : [ "Monica" ] }

```

This is the goodness of the aggregation framework. Otherwise you could use a loop to group this data and as mentioned above this will be tedious besides taking plenty of your time.

\$unwind

More often you will employ embedding of documents and would like to fetch those documents as separate entities from the main document. The unwind stage will help us get these documents with its simple syntax of

```

1 | {$unwind: <field path>}

```

Using our `users` collection we can fetch the `connections` for each user using the simple operation below and also the position of each subdocument in the array.

```

1 | db.getCollection('users').aggregate([
2 |   {$unwind:
3 |     {
4 |       path: "$connections",
5 |       includeArrayIndex: "arrayIndex"
6 |     }
7 |   }
8 | ])

```

```

1 | {
2 |   "_id" : ObjectId("5b43cbe2106c21d21c776e81"),
3 |   "userId" : NumberLong("1530442083133"),
4 |   "eyeColor" : "blue",
5 |   "name" : "George",
6 |   "connections" : {
7 |     "status" : "Disconnect",
8 |     "userName" : "Valencia",
9 |     "name" : "Derrick Clinton",
10 |    "id" : NumberLong("1530444522597")

```

```

11     },
12     "arrayIndex" : NumberLong(0)
13 }
14 {
15     "_id" : ObjectId("5b43cbe2106c21d21c776e81"),
16     "userId" : NumberLong("1530442083133"),
17     "eyeColor" : "blue",
18     "name" : "George",
19     "connections" : {
20         "status" : "Connect",
21         "userName" : "Carliston",
22         "name" : "James Good",
23         "id" : NumberLong("1530444522597")
24     },
25     "arrayIndex" : NumberLong(1)
26 }
27 {
28     "_id" : ObjectId("5b43cbe2106c21d21c776e82"),
29     "userId" : NumberLong("153044201111"),
30     "name" : "Monica",
31     "eyeColor" : "normal",
32     "connections" : {
33         "status" : "Disconnect",
34         "userName" : "JohnDoh",
35         "name" : "Alex Xian",
36         "id" : NumberLong("153044445903")
37     },
38     "arrayIndex" : NumberLong(0)
39 }
40 {
41     "_id" : ObjectId("5b43cbe2106c21d21c776e82"),
42     "userId" : NumberLong("153044201111"),
43     "name" : "Monica",
44     "eyeColor" : "normal",
45     "connections" : {
46         "status" : "Connect",
47         "userName" : "MaryCartie",
48         "name" : "Mary Carey",
49         "id" : NumberLong("1530444522597")
50     },
51     "arrayIndex" : NumberLong(1)
52 }
53 {
54     "_id" : ObjectId("5b43cbe2106c21d21c776e83"),
55     "userId" : NumberLong("15304420836758"),
56     "name" : "Harrison",
57     "eyeColor" : "blue",
58     "connections" : {
59         "status" : "Connect",
60         "userName" : "Kevin",

```

```

61     "name" : "Keni Sems",
62     "id" : NumberLong("34435345345343")
63   },
64   "arrayIndex" : NumberLong(0)
65 }
66 {
67   "_id" : ObjectId("5b43cbe2106c21d21c776e83"),
68   "userId" : NumberLong("15304420836758"),
69   "name" : "Harrison",
70   "eyeColor" : "blue",
71   "connections" : {
72     "status" : "Disconnect",
73     "userName" : "Mayaka",
74     "name" : "Andrew Fake",
75     "id" : NumberLong("15304445224357")
76   },
77   "arrayIndex" : NumberLong(1)
78 }

```

\$project

In this stage, the documents are modified either to add or remove some fields that will be returned. In a nutshell, this stage passes the documents with only specified fields.

The syntax for project is

```
1 | {$project: {<specifications>}}
```

Points to note

1. If a field is described with a value of 1 or true, the document that is to be returned will have that field.
2. You can suppress the `_id` field so that it cannot be returned by describing it with 0 or false value.
3. You can add a new field or reset the field by describing it with a value of some expression.
4. The `$project` operation will basically treat a numeric or boolean values as flags. For this reason, you will need to use the `$literal` operator for you to set a field value numeric or boolean.

From the `users` collection we can fetch the username of people in the connections documents without necessarily getting the main document information and also suppressing the `_id` field using this stage.

```

1 db.getCollection('users').aggregate([
2   {$unwind:
3     {
4       path: "$connections",
5       includeArrayIndex: "arrayIndex"
6     }
7   },
8   {$project: {"connections.userName": 1, _id:0}}
9 ])

```

The resulting documents will be.

```

1 { "connections" : { "userName" : "Valencia" } }
2 { "connections" : { "userName" : "Carliston" } }
3 { "connections" : { "userName" : "JohnDoh" } }
4 { "connections" : { "userName" : "MaryCartie" } }
5 { "connections" : { "userName" : "Kevin" } }
6 { "connections" : { "userName" : "Mayaka" } }

```

We can also add a new field of connections status and fetch the correspondent data as:

```

1 db.getCollection('users').aggregate([
2   {$unwind:
3     {
4       path: "$connections",
5       includeArrayIndex: "arrayIndex"
6     }
7   },
8   {$project: {"connections.userName": 1, _id:0,
9     "ConnectionStatus": "$connections.status"}}
10 ])

```

The result for this will be:

```

1 {"connections": {"userName": "Valencia"},"ConnectionStatus":
  "Disconnect"}
2 {"connections": {"userName" : "Carliston"},"ConnectionSta-
  tus" : "Connect"}
3 {"connections": {"userName": "JohnDoh"},"ConnectionStatus" :
  "Disconnect"}
4 {"connections": {"userName": "MaryCartie"},"ConnectionSta-
  tus" : "Connect"}
5 {"connections" : {"userName" : "Kevin"},"ConnectionStatus" :
  "Connect"}
6 {"connections": {"userName" : "Mayaka"},"ConnectionStatus" :
  "Disconnect"}

```

\$sort

The sort stage arranges the returned documents in relation to some specified order in the sort key parameter. The documents are never modified. Only the order changes. Let's consider this simple collection of students

```
1  {"_id": ObjectId("5b47c275106c21d21c776e84"), "name": "Gadaf-  
2  fy", "age": 20}  
3  {"_id": ObjectId("5b47c275106c21d21c776e85"), "name" :  
4  "John", "age" : 18}  
5  {"_id": ObjectId("5b47c275106c21d21c776e86"), "name": "Da-  
6  vid", "age": 30}  
7  {"_id": ObjectId("5b47c275106c21d21c776e87"), "name": "Emi-  
8  ly", "age": 16}  
9  {"_id": ObjectId("5b47c275106c21d21c776e88"), "name": "Cyn-  
10  thia", "age": 14}  
11 {"_id": ObjectId("5b47c275106c21d21c776e89"), "name":  
12 "Mary", "age": 28}
```

This aggregation stage will return documents which are sorted using the age as the sort key. If it is set to 1 then the arrangement is in ascending order otherwise if set to -1, then the arrangement will be in a descending order.

```
1  db.getCollection('students').aggregate([  
2      {$project: {'name':1, 'age':1, '_id':0}},  
3      {$sort: {age: 1}}  
4  ])
```

The resulting documents will be:

```
1  { "name" : "Cynthia", "age" : 14 }  
2  { "name" : "Emily", "age" : 16 }  
3  { "name" : "John", "age" : 18 }  
4  { "name" : "Gadaffy", "age" : 20 }  
5  { "name" : "Mary", "age" : 28 }  
6  { "name" : "David", "age" : 30 }
```

We can also change the sort key to name and check on the alphabetical order besides getting the results in a descending order.

```
1  db.getCollection('students').aggregate([  
2      {$project: {'name':1, 'age':1, '_id':0}},  
3      {$sort: {name: -1}}  
4  ])
```

The resulting documents will be:

```
1 | { "name" : "Mary", "age" : 28 }
2 | { "name" : "John", "age" : 18 }
3 | { "name" : "Gadaffy", "age" : 20 }
4 | { "name" : "Emily", "age" : 16 }
5 | { "name" : "David", "age" : 30 }
6 | { "name" : "Cynthia", "age" : 14 }
```

\$sample

This stage randomly selects and returns a number of documents that have been specified. For example from the `students` collection, fetch 2 random documents as

```
1 | db.getCollection('students').aggregate([
2 |     {$project: {'name':1, 'age':1, '_id':0}},
3 |     {$sample: {"size":2}}
4 | ])
```

The resulting documents

```
1 | { "name" : "Gadaffy", "age" : 20 }
2 | { "name" : "Mary", "age" : 28 }
```

If you run the operation a couple of times you will be getting different documents.

\$limit

As opposed to the `$sample` stage that returns documents randomly, `$limit` returns the first N documents and N as the specified limit. An example from the `students` collection:

```
1 | db.getCollection('students').aggregate([
2 |     {$project: {'name':1, 'age':1, '_id':0}},
3 |     {$limit: 2}
4 | ])
```

The resulting documents will be

```
1 | { "name" : "Gadaffy", "age" : 20 }
2 | { "name" : "John", "age" : 18 }
```

\$lookup

The `$lookup` stage is basically like doing a left outer join from one collection to another but in the same database. It filters documents from the joined collection. A new array field is added with elements that are matching documents from the joined collection.

The syntax for the `$lookup` stage is:

```
1  {
2  $lookup:
3      {
4          from: <collection to join>,
5          localField: <input document field>,
6          foreignField: <field from documents of the from
7          collection>,
8          as: <output array field>
9      }
}
```

- **from**: this takes the value name of the collection you want to perform the join and it has to be in the same database as the collection you are querying.
- **localField**: this is a field in the current collection which you want to perform equality match on to the **foreignField**.
- **foreignField**: is a field in the collection you are joining with that you are to use in performing an equality match on the **localField**.
- **as**: is a new array field to add to the input documents. It contains matching documents from the foreign collection.

As an example, we will make another collection named `address` as shown below:

```
1  { "_id" : ObjectId("5b485eb5106c21d21c776e8b"), "name" : "John", "town" : "Nairobi" }
2  { "_id" : ObjectId("5b485eb5106c21d21c776e8c"), "name" : "David", "town" : "Beijing" }
3  { "_id" : ObjectId("5b485eb5106c21d21c776e8d"), "name" : "Emily", "town" : "Juba" }
4  { "_id" : ObjectId("5b485eb5106c21d21c776e8e"), "name" : "Cynthia", "town" : "London" }
5  { "_id" : ObjectId("5b485eb5106c21d21c776e8f"), "name" : "Mary", "town" : "California" }
```

Let's for example get the town for each student in the `students` collection from the `address` collection. In this case, we will use the `name` field in the `students` collection as the `localField` and the `name` field in the `address` collection as the `foreignField`.

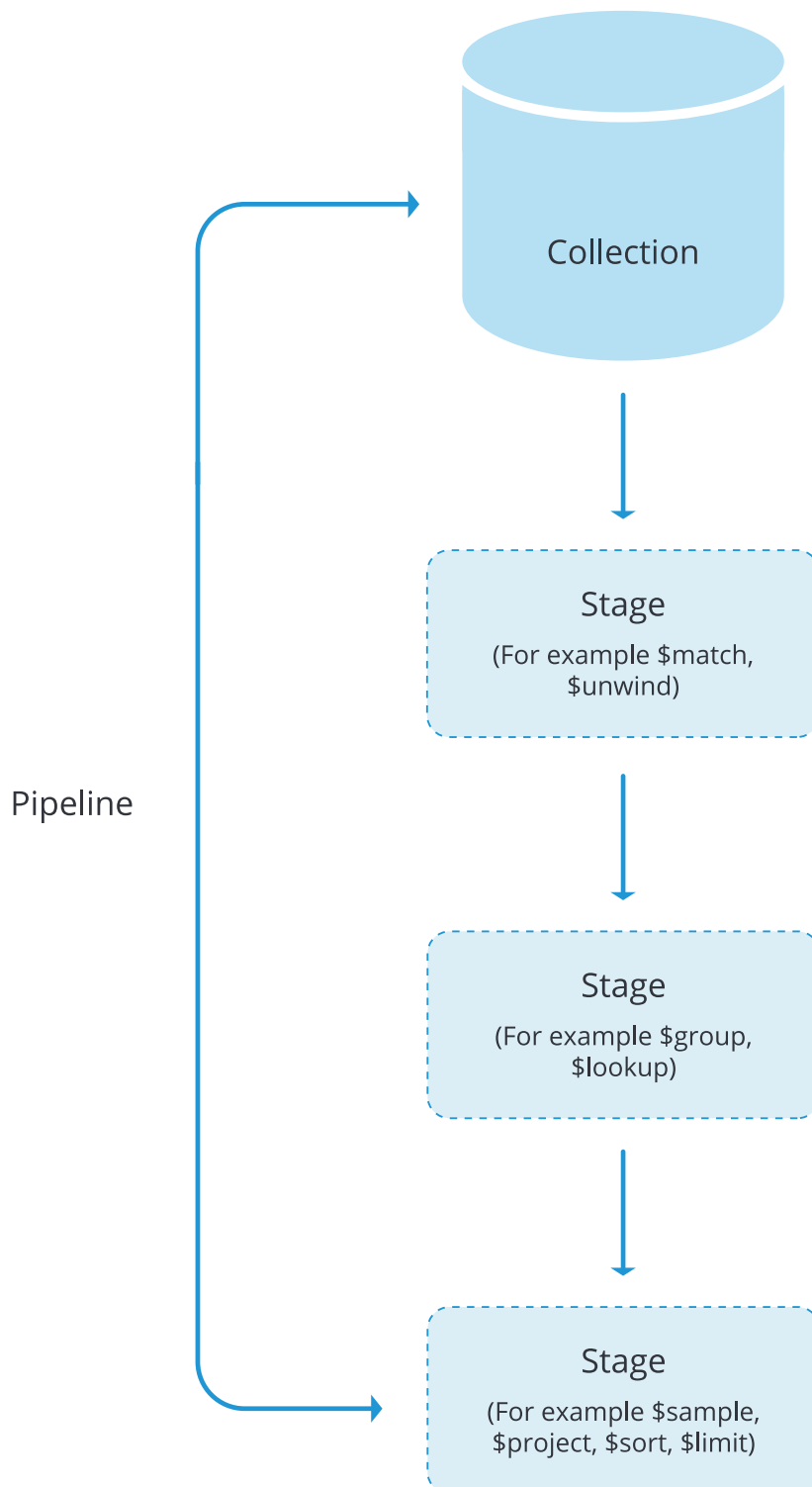
```
1  db.getCollection('students').aggregate([
2      {$lookup: {from: 'address', localField: 'name', foreignField: 'name', as: 'address'}},
3      {$project: {'address.town': 1, age: 1, name: 1, _id: 0}}
4  ])
```


The resulting documents will be

```
1  {"name" : "Gadaffy", "age" : 20, "address" : [{"town" : "Georgia"}]}
2  {"name" : "John", "age" : 18, "address" : [{"town" : "Nairobi"}]}
3  {"name" : "David", "age" : 30, "address" : [{"town" : "Beijing"}]}
4  {"name" : "Emily", "age" : 16, "address" : [{"town" : "Juba"}]}
5  {"name" : "Cynthia", "age" : 14, "address" : [{"town" : "London"}]}
6  {"name" : "Mary", "age" : 28, "address" : [{"town" : "California"}]}
```

Aggregation Process

As we have discussed all these stages, we need to understand how to associate them in our querying process so that we get the desired results. The process starts with inputting documents from the selected collection into the first stages. These documents can pass through 1 or more stages with each stage involving different operations. A simple diagram of the pipeline process is shown below.



The output of each stage becomes the input of the next stage and any stage can be repeated in order to filter the documents further.

As an example, we are going to use the students collection to fetch only the name and age of each student whose age is greater than 20 and then sort the results using the age field as our key value.

This is our data:

```
{ "_id": ObjectId("5b47c275106c21d21c776e84"), "name": "Gadaffy", "age": 20 }
{ "_id": ObjectId("5b47c275106c21d21c776e85"), "name": "John", "age": 18 }
{ "_id": ObjectId("5b47c275106c21d21c776e86"), "name": "David", "age": 30 }
{ "_id": ObjectId("5b47c275106c21d21c776e87"), "name": "Emily", "age": 16 }
{ "_id": ObjectId("5b47c275106c21d21c776e88"), "name": "Cynthia", "age": 14 }
{ "_id": ObjectId("5b47c275106c21d21c776e89"), "name": "Mary", "age": 28 }
```

We are going to use the `$match` stage to filter out documents whose age value is less than 20, then using `$sort` stage, arrange the documents in relation to the age and finally return the name and age of each student only using the `$project` stage.

```
1 db.getCollection('students').aggregate([
2   {$match: {age: {$gt:20}}},
3   {$project: {_id: 0, name:1, age:1} },
4   {$sort: {age: 1}}
5 ])
```

The resulting documents will be:

```
1 { "name" : "Mary", "age" : 28 }
2 { "name" : "David", "age" : 30 }
```

Basically, we can show the result of each stage in the output in the diagram on next page.



↓ All documents in the collection are passed

\$match stage returns

```
{ "_id": ObjectId("5b47c275106c21d21c776e86"), "name": "David", "age": 30 }  
{ "_id": ObjectId("5b47c275106c21d21c776e89"), "name": "Mary", "age": 28 }
```

↓ Only documents with age greater than 20 are passed

\$project stage returns

```
{ "name": "David", "age": 30 }  
{ "name": "Mary", "age": 28 }
```

↓ Documents are arranged in ascending order using age field

\$sort stage returns

```
{ "name": "Mary", "age": 28 }  
{ "name": "David", "age": 30 }
```

Accumulator Operators

These are operators which assist in the analysis process of the input documents. All of them are found in the `$group` stage but as from version 3.2, some are also found in the `$project` stage.

We are going to use this data to elaborate on the usage of these operators (add it to a collection and name it `students`)

```
1  {"_id" : ObjectId("5b485eb5106c21d21c776e8a"), "name" :
   "Gadaffy", "town" : "Georgia", "unit" : "A", "age" : 18, "marks"
   : [20, 50, 38] }
2  {"_id" : ObjectId("5b485eb5106c21d21c776e8b"), "name" :
   "John", "town" : "Nairobi", "unit" : "B", "age" : 24, "marks" :
   [38, 60, 70]}
3  {"_id" : ObjectId("5b485eb5106c21d21c776e8c"), "name" : "Da-
   vid", "town" : "Beijing", "unit" : "A", "age" : 28}
4  {"_id" : ObjectId("5b485eb5106c21d21c776e8d"), "name" : "Emi-
   ly", "town" : "Juba", "unit" : "C", "age" : 30, "marks" : [40,
   87, 34]}
5  {"_id" : ObjectId("5b485eb5106c21d21c776e8e"), "name" : "Cyn-
   thia", "town" : "London", "unit" : "B", "age" : 16, "marks" :
   [60, 90, 98]}
6  {"_id" : ObjectId("5b485eb5106c21d21c776e8f"), "name" :
   "Mary", "town" : "California", "unit" : "B", "age" : 22,
   "marks" : [52, 50, 56]}
```

The accumulator operators include;

\$sum

This operator returns the sum of numeric value while ignoring non numeric values. It is found both in `$group` and `$project` stages as from version 3.2

In the `$group` stage, it will return the collective sum of all numeric values in accordance to some applied expression to each document in a group who share the same key name. The syntax for the `$sum` operator is

```
1  | {$sum: <expression> }
```

However, for the `$project` stage we can add the number of expressions and make an array thereby the syntax becomes:

```
1  | {$sum: [<expression1>, <expression2> ...]}
```

Using the data above, let's group the students in relation to their unit and sum the number students in each group.

```
1 | db.getCollection(students).aggregate([{$group: {"_id": "$unit", sum: {$sum: 1}}}]])
```

And the result for this is:

```
1 | { "_id" : "A", "NumberOfStudents" : 2 }
2 | { "_id" : "B", "NumberOfStudents" : 3 }
3 | { "_id" : "C", "NumberOfStudents" : 1 }
```

If the field specified in the expression does not exist, then the operation will return a value of 0.

We can also sum the marks of each student using the \$project stage. Remember, in this case we have used only 1 field that is the marks field, you can add many field as your data involves.

```
1 | db.getCollection('students').aggregate([{$project: {Total-Marks: {$sum: ["$marks"]}}}]])
```

The operation will give the following result:

```
1 | { "_id" : ObjectId("5b485eb5106c21d21c776e8b"), "Total-Marks" : 168 }
2 | { "_id" : ObjectId("5b485eb5106c21d21c776e8c"), "Total-Marks" : 182 }
3 | { "_id" : ObjectId("5b485eb5106c21d21c776e8d"), "Total-Marks" : 161 }
4 | { "_id" : ObjectId("5b485eb5106c21d21c776e8e"), "Total-Marks" : 848 }
5 | { "_id" : ObjectId("5b485eb5106c21d21c776e8f"), "Total-Marks" : 158 }
```

\$avg

Give the average value of numeric values while ignoring non-numeric values. From version 3.2, this operator is available in both \$group and \$project stages.

In the \$group stage, it will return the collective average of all numeric values in accordance to some applied expression to each document in a group who share the same key name. The syntax for the \$avg operator is

```
1 | {$avg: <expression> }
```

However, for the \$project stage we can add the number of expressions and make an array thereby the syntax becomes:

```
1 | {$avg: [<expression1>, <expression2> ...]}
```

In the group stage, we can get the average age of every group as:

```
1 | db.getCollection('students').aggregate([{$group: {"_id": "$unit", AverageAge: {$avg: "$age"}}}])
```

And the result will be:

```
1 | { "_id" : "A", "AverageAge" : 23 }
2 | { "_id" : "B", "AverageAge" : 20.666666666666668 }
3 | { "_id" : "C", "AverageAge" : 30 }
```

In the project stage, we can get the average marks of each student as:

```
1 | db.getCollection('students').aggregate([{$project: {"_id": 0, "name": 1, AverageMarks: {$avg: ["$marks"]}}}])
```

The result for this operation will be:

```
1 | { "name" : "Gadaffy", "AverageMarks" : 36 }
2 | { "name" : "John", "AverageMarks" : 56 }
3 | { "name" : "David", "AverageMarks" : 60.666666666666664 }
4 | { "name" : "Emily", "AverageMarks" : 53.666666666666664 }
5 | { "name" : "Cynthia", "AverageMarks" : 82.66666666666667 }
6 | { "name" : "Mary", "AverageMarks" : 52.666666666666664 }
```

\$max and \$min

Return the maximum and minimum values respectively from a given numeric array. Returns a 0 for field that does not exist. In the group stage we can check for the maximum and minimum ages for each group formed as:

```
1 | db.getCollection('students').aggregate([{$group: {"_id": "$unit", "Maximum age": {$max: "$age"}, "Minimum age": {$min: "$age"}}}])
```

The result from the operation is:

```
1 | { "_id" : "A", "Maximum age" : 28, "Minimum age" : 18 }
2 | { "_id" : "B", "Maximum age" : 24, "Minimum age" : 16 }
3 | { "_id" : "C", "Maximum age" : 30, "Minimum age" : 30 }
```

In the project stage, we can find the maximum marks and minimum marks of each student as:

```
1 db.getCollection('students').aggregate([{$project: {"_id": 0,"name": 1,"Maximum marks": {$max: "$marks"},"Minimum marks": {$min: "$marks"}}}])
```

The result for this operation is:

```
1 { "name" : "Gadaffy", "Maximum marks" : 50, "Minimum marks" : 20 }
2 { "name" : "John", "Maximum marks" : 70, "Minimum marks" : 38 }
3 { "name" : "David", "Maximum marks" : 76, "Minimum marks" : 40 }
4 { "name" : "Emily", "Maximum marks" : 87, "Minimum marks" : 34 }
5 { "name" : "Cynthia", "Maximum marks" : 98, "Minimum marks" : 60 }
6 { "name" : "Mary", "Maximum marks" : 56, "Minimum marks" : 50 }
```

\$push

This is available only in the **\$group** stage. It is used to return an array of the expression values. For our example above, after grouping the students according to their unit, what if we want to get the names of students in each group? We will have to push their names into an array using the **\$push** operator as:

```
1 db.getCollection('students').aggregate([{$group: {"_id": "$unit","students": {$push: "$name"}}}])
```

The result from this operation will be

```
1 { "_id" : "A", "students" : [ "Gadaffy", "David" ] }
2 { "_id" : "B", "students" : [ "John", "Cynthia", "Mary" ] }
3 { "_id" : "C", "students" : [ "Emily" ] }
```


Similarity of the Aggregation Process in MongoDB with SQL

With the introduction of the aggregation process in MongoDB, there is a greater capability of doing most of the data processing just like in SQL. We can contrast the operations hand in hand from MongoDB to SQL as depicted below.

We will use our **student** collection above as a table in SQL and as a collection in MongoDB.

SQL	MongoDB
SELECT	\$project
WHERE/HAVING	\$match
JOIN	\$lookup
LIMIT	\$limit
GROUP BY	\$group
ORDER BY	\$sort
COUNT()	\$sum/ \$sortByCount
SUM()	\$sum
AVG()	\$avg

Table: Operators in MongoDB that offer equivalent SQL functions

SQL	MongoDB	Explanation
<pre>SELECT town, age,unit FROM `students` WHERE name = 'Mary' LIMIT 1</pre>	<pre>db.students.aggregate([{ \$match: {name: 'Mary'} }, { \$project: { "town": 1, "age": 1, "Unit": 1 } }, { \$limit: 1 }])</pre>	<p>Fetches the ages, towns and units of the students whose names correspond to Mary and return only 1 result.</p>

SQL	MongoDB	Explanation
SELECT * FROM `students` WHERE unit = 'B' ORDER BY name ASC	db.students.aggregate([{ \$match: {unit: 'B'} } { \$sort: {name: 1} }])	Fetches all students who belong to unit B and arranges the results in accordance to their names in an ascending order.
SELECT * FROM `students` GROUP BY unit COUNT(*) as numberOfStudents	db.students.aggregate([{ \$group: { _id: 'unit', numberOfStudents: {"\$sum": 1} } }])	Groups the students that share the same unit value and then counts the number of students in each group formed.
SELECT AVG(age) AS averageAge FROM `students`	db.students.aggregate([{\$group: { _id: null, averageAge : {"\$avg": "\$age"}}})	Calculates the average age of the students.
SELECT * FROM `students` HAVING age > 20	db.students.aggregate([{\$match: {'age':{\$gt: 20}}})	Return all students whose age is greater than 20.

Table: Examples of similarity in MongoDB and SQL

Aggregation Pipeline Optimization

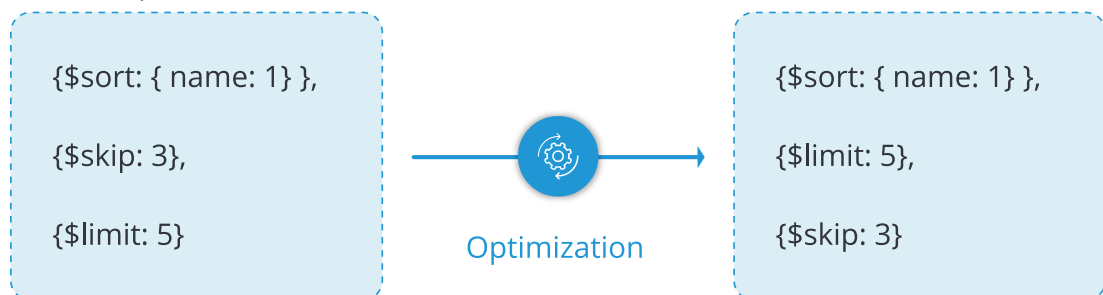
As much as the aggregation concept improves performance more than the CRUD find operation, there are further techniques you can involve to improve the performance. This is achieved through a reshaping of the pipeline. The query optimizer in MongoDB is most effective at picking the best of multiple plans but with the aggregation framework, it is always under your complete control and design how the steps should be executed.

Projection Optimization

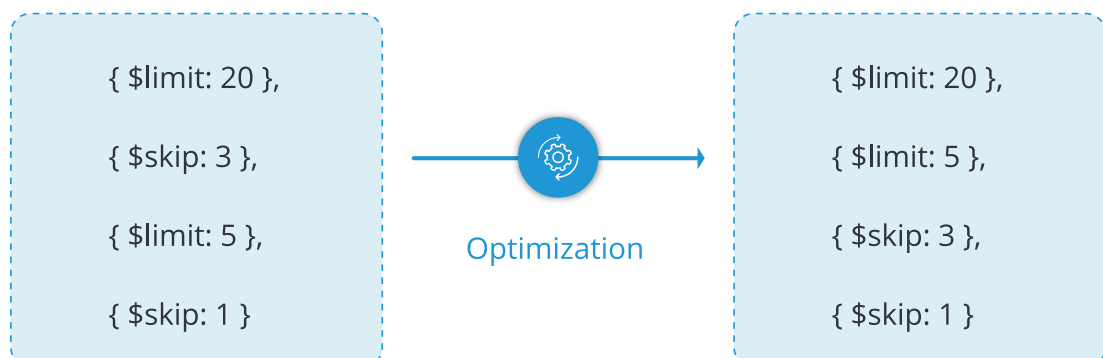
In this case you select the only fields you want to be returned. For example if you want embedded documents in a field, then there will be no need to return other fields in the main document as this will reduce data amount passing through the pipeline hence save of time. The \$project stage is used in this case

Pipeline Sequence Optimization

1. **\$sort + \$skip + \$limit** Sequence Optimization
If you have a sequence of sort followed by a skip and then a limit, an optimization phase will occur to bring the **limit** stage before the **skip** stage. For example



2. **\$limit + \$skip + \$limit + \$skip** Sequence Optimization
For a continuous sequence of skip and limit, the optimization phase will attempt to group the limit stages together and then the skip stages together.



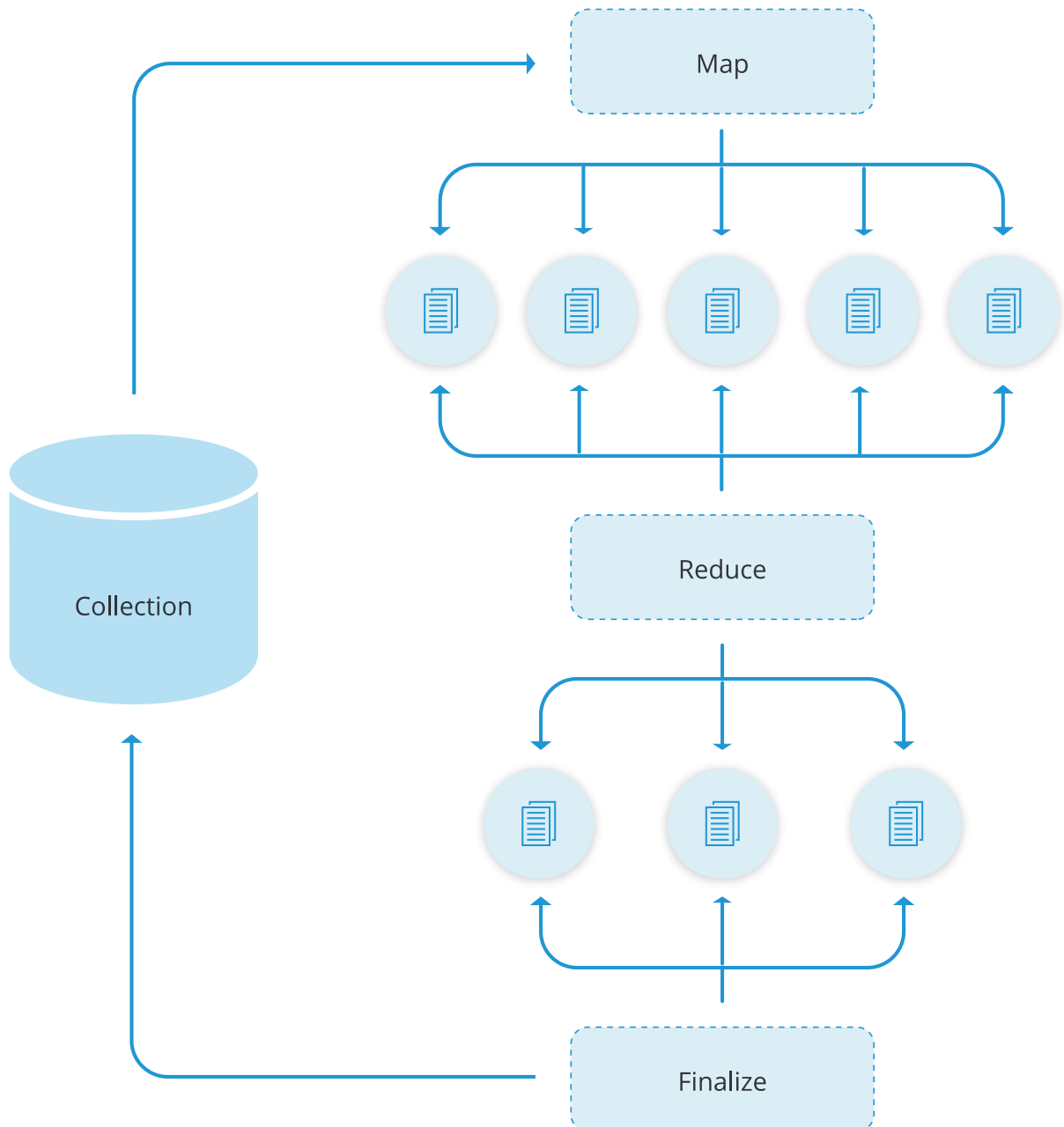
MapReduce in MongoDB

This is also an aggregation process which condenses large volumes of data into an aggregated form.

Before doing a map-reduce on your collection, you need to

1. Understand your data structure and how are you going to analyze it.
2. Design your end result by having a structure how it should look like.
3. Manipulate a sample of the data to know which equations and modifications that need to be integrated.

The MapReduce structure can be summarized with the diagram below



Majorly, there are 2 stages involved, that is, the mapping process and then reducing the mapped results.

For every input document, arbitrary sorting and limiting is done and then the map phase is applied with an end result of producing key-value pairs.

If there are keys with multiple values, they are passed to the reduce phase to condense the aggregated data. A query is also used to limit the number of documents entering into the map phase.

The final result is stored in a new collection whose name is specified in the `out` property of the map-reduce operation.

In MongoDB, we use the `mapReduce` function to achieve this operation. The syntax for MapReduce operation in MongoDB is:

```
1 db.collection.mapReduce({
2     /*map*/ function() { emit (fields to be returned);},
3     /*reduce and sum some field value*/ function(key,
4     values){ return Array.sum( values) },
5     /*query*/ {
6         query: { field: value},
7     },
8     /*collection to store data*/
9     out: "collection name"
10 })
```

For example, for our students collection we can calculate the average age of each group using this simple map reduce function:

```
1 db.getCollection('students').mapReduce(
2     function(){emit(this._id, this.age)},
3     function(key, values){return Array.avg(values)},
4     {
5         query: {},
6         out: 'results'
7     }
8 )
```

This operation will result in a new collection named `results` with the average age value for each group i.e.

```
1 { "_id" : "A", "value" : 23 }
2 { "_id" : "B", "value" : 20.666666666666668 }
3 { "_id" : "C", "value" : 30 }
```

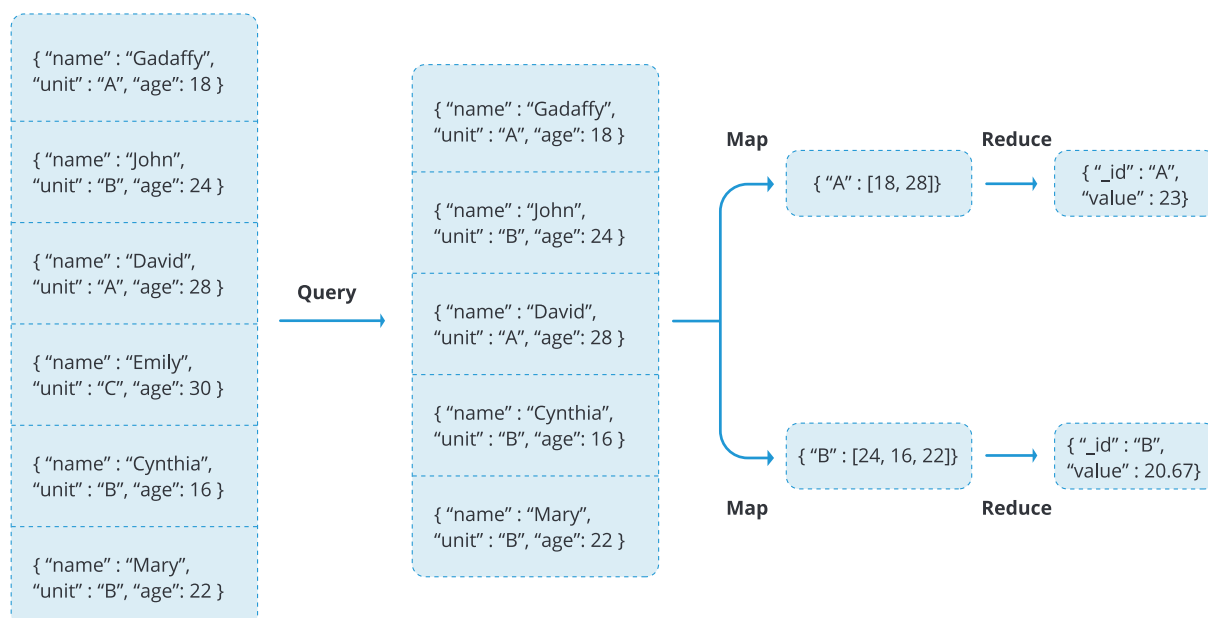
We can also supply some expression to the query attribute like getting the group A and B only

```
1 db.getCollection('students').mapReduce(  
2   function(){emit(this.unit, this.age)},  
3   function(key, values){return Array.avg(values)},  
4   {  
5     query: {unit: {$ne: 'C'}},  
6     out: 'results'  
7   }  
8 )
```

The result for this operation will result in a collection with this data:

```
1 { "_id" : "A", "value" : 23 }  
2 { "_id" : "B", "value" : 20.666666666666668 }
```

In simple representation, this is how the process takes place



MapReduce JavaScript Functions

In order to associate or map values to a key, the map-reduce function uses the custom JavaScript functions. The **reduce** operation will then identify keys with multiple values and reduce them to single objects.

An advantage with this JavaScript operation is flexibility which can allow further modifications and calculations. For instance, in the reduce part of our operation above, we are able to calculate the average age in each group because the data is more flexible.

Incremental MapReduce

As mentioned above, the map-reduce operation returns data which is stored in a new collection, this is not always the case. Sometimes the data is returned inline for you to carry out more aggregation operations. Sometimes the map-reduce data set may constantly be growing resulting into some final return issues. The result documents are always expected to be within the limit size defined by the BSON Document Size of 16 megabytes. Because of this reason, it is advisable to carry out an incremental MapReduce function for a large set of data rather than a single map-reduce operation on the entire data set.

The query parameter will help you to specify the conditions for which new documents will be passed while the out parameter will specify the reduce action with which new results will be merged into the existing output collection.

Let's consider data in a sessions collection as shown below.

```
1  {_id: ObjectId("5b4b36a726d09bd6a0953731"),userid: "a",ts:
   ISODate("2011-11-03T14:17:00Z"),length: 131}
2  {_id: ObjectId("5b4b36a726d09bd6a0953732"),userid: "b",ts:
   ISODate("2011-11-03T14:23:00Z"),length: 128}
3  {_id: ObjectId("5b4b36a726d09bd6a0953733"),userid: "c",ts:
   ISODate("2011-11-03T15:02:00Z"),length: 138}
4  {_id: ObjectId("5b4b36a726d09bd6a0953734"),userid: "d",ts:
   ISODate("2011-11-03T16:45:00Z"),length: 63}
5
6  {_id: ObjectId("5b4b36a726d09bd6a0953735"),userid: "a",ts:
   ISODate("2011-11-04T11:05:00Z"),length: 123}
7  {_id: ObjectId("5b4b36a726d09bd6a0953736"),userid: "b",ts:
   ISODate("2011-11-04T13:14:00Z"),length: 138}
8  {_id: ObjectId("5b4b36a726d09bd6a0953737"),userid: "c",ts:
   ISODate("2011-11-04T17:00:00Z"),length: 148}
9  {_id: ObjectId("5b4b36a726d09bd6a0953738"),userid: "d",ts:
   ISODate("2011-11-04T15:37:00Z"),length: 83}
```

To run the initial MapReduce on the current collection, we define some functions we are going to use in our MapReduce function.

i.e.

1. A map function which will ideally map the `userid` to an object which contains the fields `userid`, `total_time`, `count` and `avg_time`.

```
1  var mapping = function(){
2      emit(this.userid, {
3          userid: this.userid,
4          total_time: this.length,
5          count: 1
6          avg_time: 0
7      });
8  }
```

2. Since we have defined we have the total time and count, we need to define a reduce function that will do this calculation.

```
1 | var reducing = function(key, values){
2 |     var objectReturn = {
3 |         userid: key,
4 |         total_time: 0,
5 |         count: 0,
6 |         avg_time: 0
7 |     };
8 |     values.map(function(value){
9 |         objectReturn.total_time += value.total_time;
10 |        objectReturn.count += value.count;
11 |    });
12 |    return objectReturn;
13 | }
```

3. Regarding the `avg_time`, we need to define a finalize function with 2 arguments that is the key and `reduceValue` to add the average and return a modified document.

```
1 | var finalizing = function(key, reduceValue){
2 |
3 |     if(reduceValue.count > 0){
4 |         reduceValue.avg_time = reduceValue.total_time /
5 |         reduceValue.count
6 |     }
7 |     return reduceValue;
8 | }
```

4. We then use these functions to do an incremental map-reduce on our data set in the `sessions` collection as:

```
1 | db.sessions.mapReduce(mapping, reducing,
2 | {
3 |     out: 'sessionsInfo',
4 |     finalize: finalizing
5 | })
6 | )
```

If you add new documents to the `sessions` collection, you will need to modify the query to determine which documents will be passed. For example if we had the last document for a given day as

```
1 | {_id: ObjectId("5b4b36a726d09bd6a0953738"),userid: "d",ts:
  |   ISODate("2011-11-04T15:37:00Z"),length: 83}
```


For the new documents we can restrict them to have a timestamp greater than for this last document in order to be passed to the next step. .i.e.

```
1 db.sessions.mapReduce(mapping, reducing,  
2 {  
3   query: {ts: {$gt: ISODate("2011-11-04T15:37:00Z")}}  
4   out: {reduce: 'sessionsInfo'},  
5   finalize: finalizing  
6 }  
7 )
```

Comparison Between MapReduce and Aggregation Pipeline in MongoDB

As much as the MapReduce operation try to provide almost equivalent operations as in the aggregation pipeline, there are some distinctive features that can lead a user to prefer one to the other. We will discuss this in the below table.

MapReduce	Aggregation Pipeline
<p>Relatively slower process.</p> <p>The fact that the MapReduce function is based on a JavaScript interpreter, the data in MongoDB is in a BSON format hence has to be converted into a JSON format before application of the Map-Reduce operation. It tends to take more time for a correlated function as the aggregation pipeline.</p>	<p>The process is quite faster.</p> <p>Pipeline concept is based on the process of parallel operations. It implies that several operations are carried out almost at the same time. Also, the data is not converted into any other format. For this reason, the results are generated at a faster rate</p>
<p>More flexibility on aggregated data.</p> <p>If one would wish to pick the results at some point, there are several options like: inline, merge, reduce and a new collection.</p>	<p>Limited data flexibility.</p> <p>The results are only available in the inline-block. To get them into another collection you will therefore be required to write more CRUD queries.</p>
<p>Incremental Aggregation.</p> <p>Due to the restricted document size of 16 megabytes, incremental aggregation provides an opportunity for one to compute for more results by inputting new documents that match a supplied query and update the initial results using the reduce operation.</p>	<p>Incremental aggregation is not supported and since the documents are returned inline, the size of the document is always restricted to 16MB. Besides, there is only 1 supported output option therefore impossible to update values if incremental aggregation was to be applied.</p>
<p>Customizability.</p> <p>Data that is available within the functions can be manipulated to suit own specifications.</p>	<p>Limited to operators and expressions supported by the aggregation framework and therefore it is impossible for one to write custom functions.</p>
<p>Supports non-sharded and sharded input collections.</p>	<p>Support non-sharded and sharded input collections.</p>

Summary

Aggregation is the process of manipulating large data sets with some specified procedures to return calculated results. These results are provided in a simplified format to enhance analysis of the associated data.

The aggregation process can be done by either MapReduce operation or the aggregation pipeline concept in MongoDB. This process is run on the mongod instance to simplify the application code, beside the need to limit resource requirements.

The input to an aggregation process is the documents in collections and the results is also a document or a number of documents.

Aggregation stages involve operators such as addition, averaging values for given fields, finding the maximum and minimum values among many more operators. This makes the analysis of data even more simplified.

About ClusterControl

ClusterControl is the all-inclusive open source database management system for users with mixed environments that removes the need for multiple management tools. ClusterControl provides advanced deployment, management, monitoring, and scaling functionality to get your MySQL, MongoDB, and PostgreSQL databases up-and-running using proven methodologies that you can depend on to work. At the core of ClusterControl is its automation functionality that lets you automate many of the database tasks you have to perform regularly like deploying new databases, adding and scaling new nodes, running backups and upgrades, and more. Severalnines provides automation and management software for database clusters. We help companies deploy their databases in any environment, and manage all operational aspects to achieve high-scale availability.

About Severalnines

Severalnines provides automation and management software for database clusters. We help companies deploy their databases in any environment, and manage all operational aspects to achieve high-scale availability.

Severalnines' products are used by developers and administrators of all skills levels to provide the full 'deploy, manage, monitor, scale' database cycle, thus freeing them from the complexity and learning curves that are typically associated with highly available database clusters. Severalnines is often called the "anti-startup" as it is entirely self-funded by its founders. The company has enabled over 12,000 deployments to date via its popular product ClusterControl. Currently counting BT, Orange, Cisco, CNRS, Technicolor, AVG, Ping Identity and Paytrail as customers. Severalnines is a private company headquartered in Stockholm, Sweden with offices in Singapore, Japan and the United States. To see who is using Severalnines today visit:

<https://www.severalnines.com/company>



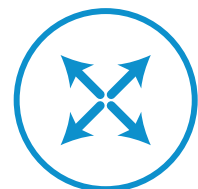
Deploy



Manage

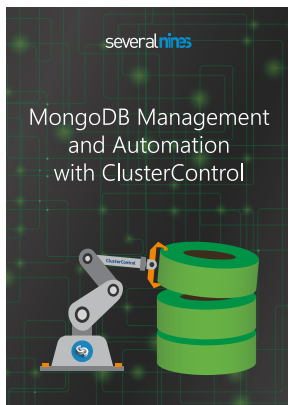


Monitor



Scale

Related Resources



MongoDB Management and Automation with ClusterControl

This white paper reviews the challenges involved in managing MongoDB at scale and introduces mitigating features of ClusterControl from Severalnines. As a best of breed database management solution, ClusterControl brings consistency and reliability to your database environment, and simplifies your database operations at scale.

[Download whitepaper](#)



Become a MongoDB DBA: Bringing MongoDB to Production

Learn from our MongoDB experts what it takes to ensure your MongoDB stacks are production-ready. This whitepaper includes tips and tricks that we have collected from our best resources to help you deploy, monitor, manage and scale MongoDB in your environment.

[Download whitepaper](#)



Become a MongoDB DBA Blog Series

Read our popular blog series on how to become a MongoDB DBA: we cover everything from deployment and monitoring via management through to scaling your MongoDB database setups.

[Read the blog](#)

MongoDB supports rich queries through its powerful aggregation framework, and allows developers to manipulate data in a similar way to SQL. Effectively, it allows developers to perform advanced data analysis on MongoDB data.

This whitepaper provides a foundation of essential aggregation concepts - how multiple documents can be efficiently queried, grouped, sorted and results presented in appropriate ways for reports and dashboards.

