

Chapter 1

Document database

In this chapter we will discuss of the organization of the document database and how we handled the replicas. We decided to use MongoDB as DBMS for the document database for the purpose of storing the main information for movies, users, reviews and personnel. In MongoDB we created the following three collection:

- movie
- user

Inside the movie collection there are embedded the documents for the reviews and personnel

1 Collection composition

1.1 Movie collection

The following is the composition of the Movie collection.

```
1 {
2   "_id" : ObjectId(<<id_field>>),
3   "primaryTitle": "The first ever movie",
4   "year": 1989,
5   "runtimeMinutes": 70,
6   "genres": ["Crime", "Drama", "Romantic"],
7   "productionCompany": "Dingo Picture Production" ,
8   "personnel": [
9     {
10      "name": "John Doe",
11      "category": "producer",
12      "job": "writer"
13    },
14    {
15      "name": "Christopher Lee",
16      "category": "actor",
17      "character": ["The one"]
18    },
19    ...
20  ],
21  "top_critic_fresh_count": "4",
22  "top_critic_rotten_count": 0,
23  "user_fresh_count": 14,
```

```

24     "user_rotten_count": 1,
25     "top_critic_rating": 100,
26     "top_critic_status": "Fresh",
27     "user_rating": 93,
28     "user_status": "Fresh",
29     "review":
30     [
31         {
32             "critic_name": "AntonyE",
33             "review_date": "2018-12-07",
34             "review_type": "Rotten",
35             "top_critic": false,
36             "review_content": "I really didn't liked it!",
37             "review_score": "1/10"
38         },
39         {
40             "critic_name": "AntonyE",
41             "review_date": "2015-12-07",
42             "review_type": "Fresh",
43             "top_critic": true,
44             "review_content": "I really liked it!",
45             "review_score": "A-"
46         },
47         ...
48     ]
49 }

```

Listing 1.1: Test

As previously stated the field personnel and review are arrays of embedded documents.

1.2 User collection

The following is the composition of the User collection.

```

1 {
2     "_id"           : ObjectId(<<id_field>>),
3     "username"      : "AntonyE",
4     "password"      : "hashed_password",
5     "firstName"     : "Anton",
6     "lastName"      : "Ego",
7     "registrationreview_date" : "2019-06-29",
8     "review_date of birth"    : "2002-07-16",
9     "last3Reviews":
10    [
11        {
12            "_id"           : ObjectId(<<id_field>>),
13            "primaryTitle"   : "Star Wars: A new Hope",
14            "review_date"    : "2018-12-07",
15            "review_type"    : "fresh",
16            "top_critic"     : false,
17            "review_content"  : "I really liked it!",
18            "vote"           : "8/10"

```

```

19         },
20         {
21             "_id"                : ObjectId(<<id_field>>),
22             "primaryTitle"       : "Ratatouille",
23             "review_date"        : "2019-02-17",
24             "review_type"        : "rotten",
25             "top_critic"         : false,
26             "review_content"     : "The director was also
    ↪ controlled by a rat",
27             "vote"              : "D+"
28         },
29         {
30             "_id"                : ObjectId(<<id_field>>),
31             "primaryTitle"       : "300",
32             "review_date"        : "2020-02-15",
33             "review_type"        : "fresh",
34             "top_critic"         : false,
35             "review_content"     : "Too much slow motion",
36             "vote"              : "3.5/5"
37         }
38     ],
39     "reviews": [
40         {
41             "movie_id" : ObjectId(<<id_field>>),
42             "primaryTitle" : "Evidence",
43             "review_index": 11
44         },
45         ...
46     ]
47
48 }

```

Listing 1.2: Test

Is important to notice that in this collection we have two fields for the reviews, *last3Reviews* and *reviews* which present different ways of storing data for the same underlying entity. As suggested by the name itself *last3Reviews* contains the last three review in the form of an embedded document meanwhile *reviews* contains all the reviews made by a user in the form of document linking towards the movie for which they were written; in particular, the link is formed by the id of the movie, the title and the position in which the review can be found in the array *review* of the movie document. This structure was chosen so that we could avoid a full replica of the reviews in both collections and, at the same time, avoiding to perform join operation for those reviews that are more frequently checked, which are the most recent ones. The idea of creating a separated collection for reviews was immediately discarded because it would have resorted in a design for the document database that resemble a third normal form of relational database.

2 Indexes

In order to provide the best execution speed in search queries we use indexes. In particular we focus on the application part that is available also without registering to the site, like Hall of Fame and search movies functionalities. Without indexes we need a collection scan in user and

movie collections to find the right document.

2.1 Movie collection

- primaryTitle
- year
- genres
- top_critic_rating
- user_rating

```
1 // Dimension is expressed in Kb
2 {
3   _id_: 228,
4   primaryTitle_1: 292,
5   year_1: 92,
6   genres_1: 180,
7   'personnel.primaryName_1': 1680,
8   top_critic_rating_1: 80,
9   user_rating_1: 80
10 }
```

Listing 1.3: Test

2.2 User collection

- username
- date_of_birth

```
1 // Dimension is expressed in Kb
2
3 { _id_: 180, username_1: 168, date_of_birth_1: 100 }
```

Listing 1.4: Test

We decided to use these indexes because they provide a jump in term of speed in search queries without occupying much space. We accept the fact that writes are slower because our application is read-intensive. We consider the idea of using an index on personnel.primaryName but the cost in term of space was higher than potential benefit. For a full analytics report for the indexes you can see the appendix. (Appendix ??)

3 Partition and replicas

To ensure high-availability we deploy a cluster of replica (3). We install and configure MongoDB in all machines, with these priorities

1. 172.16.5.26 (primary)
2. 172.16.5.27 (secondary)
3. 172.16.5.28 (secondary)

We decided to use **nearest read preference** and **W2 write preference**. In fact we can tolerate that users see temporarily an old version of data with a 33% chance.

4 Aggregations

In this section we shall discuss on the different aggregation that the application will implement, the values between «» represent a values passed by an above level.

4.1 Return the best years by top critic and user ratings

```
1 db.movie.aggregate ([
2   { $group:
3     {
4       _id: "$year",
5       topCritic: { $avg: "$top_critic_rating" },
6       rate: { $avg: "$user_rating" },
7       count: { $sum: 1 }
8     }
9   },
10  { $match: { count: { $gte: <<i>> } } },
11  { $sort: { topCritic: -1, rate: -1 } },
12  { $limit: <<j>> }
13 ] )
```

Listing 1.5: Test

4.2 Return the best genres by top critic and user ratings

```
1 db.movie.aggregate ([
2   {
3     $unwind: "$genres"
4   },
5   { $group:
6     {
7       _id: "$genres",
8       topCritic: { $avg: "$top_critic_rating" },
9       rate: { $avg: "$user_rating" },
10      count: { $sum: 1 }
11    }
12  },
13  { $match: { count: { $gte: <<i>> } } },
14  { $sort: { topCritic: -1, rate: -1 } },
15  { $limit: <<j>> }
16 ] )
```

Listing 1.6: Test

4.3 Return the best production houses by top critic and user ratings

```
1 db.movie.aggregate ([
2   { $group:
3     {
4       _id: "$production_company",
5       topCritic: { $avg: "$top_critic_rating" },
6       rate: { $avg: "$user_rating" },
7       count: { $sum: 1 }
8     }
9   },
```

```

10     {$match:{count:{$gte:<<i>>}}},
11     {$sort:{topCritic:-1, rate:-1}},
12     {$limit:<<j>>}}
13 ])

```

Listing 1.7: Test

4.4 Given a movie count the review it has received by each month

```

1 db.movie.aggregate([
2   {
3     $match:{id:<<"id">>}}
4   },
5   {$unwind:"$review"},
6   {
7     $group:
8     {
9       _id:{year:{$year:"$review.review_date"}, month:{$month:"$review.
review_date"}},
10      count:{$sum:1}
11    }
12  },
13  {$sort:{_id:1}}
14 ]
15 ])

```

Listing 1.8: Test

4.5 Given a user count the review he/she has made divided by genres

```

1 db.movie.aggregate([
2   {$match:
3     {"review.critic_name":{$eq:<<"name">>}}
4   },
5   {$unwind:"$genres"},
6   {$group:{_id:"$genres", count:{$sum:1}}},
7   {$sort:{count:-1}},
8   {$limit:<<n>>}}
9 ])

```

Listing 1.9: Test

4.6 Return the number of user divided by an age bucket

```

1 db.user.aggregate([
2   {
3     $match:{ "date_of_birth":{$exists:true}}
4   },
5   {
6     $bucket:
7     {
8       groupBy: {$year:"$date_of_birth"},
9       boundaries: [<<values>>],
10      output:
11      {
12

```

```
13         "population": {$sum:1}
14     }
15 }
16 }
17 ])
```

Listing 1.10: Test