

Cloud Data Infrastructure

Chapter 3 : Data Infrastructure Optimization

Juliette Danel
Polytechnic Institute of Paris
juliette.danel@ip-paris.fr

Godefroy Lambert
Polytechnic Institute of Paris
godefroy.lambert@ip-paris.fr

Marius Ortega
Polytechnic Institute of Paris
marius.ortega@ip-paris.fr

15 March 2024

1 Introduction

Considering the model and needs expressed in chapter 1 ([link](#)) and the denormalization performed in chapter 2 ([link](#)), we will now present how costly each denormalized model is. To do so, we will start by expressing our use cases in MQL (Mongo Query Language) and conclude on the efficiency of each model.

2 MQL queries

2.1 Model 1

In this section, we translate the previously selected queries into MongoDB Queries Language. We base the queries on the first model we produced in the last report (see Fig. 1 1).

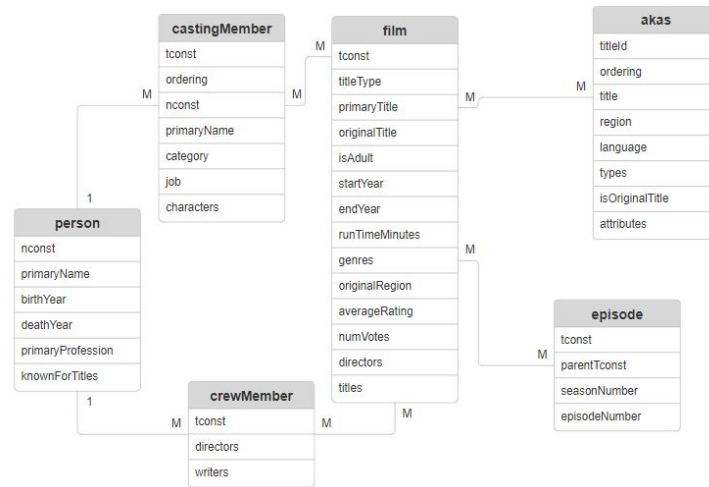


Figure 1: Entity-relationship diagram for model 1

2.1.1 Select the top 10 movies of the year X

To illustrate the query, we set the year to 1999.

```

1    db.film.aggregate([
2      {
3        $match: {
4          "titleType": "movie",
5          "startYear": 1999,
6          "runTimeMinutes": { $lt: 120 }
7        }
8      },
9      {
10       $addFields: {
11         "ratingScore": {
12           $multiply: [
13             "$averageRating",
14             { $log: { $add: ["$numVotes", 1] } }
15           ]
16         }
17       }
18     },
19     { $sort: { "ratingScore": -1 } },
20     { $limit: 10 },
21     {
22       $project: {

```

```

23         "_id": 0,
24         "tconst": 1,
25         "primaryTitle": 1,
26         "startYear": 1,
27         "ratingScore": 1,
28         "runTimeMinutes": 1
29     }
30 }
31 ])

```

2.1.2 Select a title containing the string “XXX” in at least one of its titles

To illustrate the query, we set the string to ”love”.

```

1  db.film.find(
2    {
3      $or: [
4        { "primaryTitle": { $regex: "love", $options: "i" }
5          ↪ },
6        { "titles": { $regex: "love", $options: "i" } }
7      ],
8      {
9        "_id": 0,
10       "tconst": 1,
11       "primaryTitle": 1,
12       "originalTitle": 1,
13       "title": 1
14     }
15   )

```

2.1.3 Actors that played in at least 3 adult films

```

1  db.CastingMember.aggregate([
2    {
3      $lookup: {
4        from: "Person",
5        localField: "nconst",
6        foreignField: "nconst",
7        as: "person"
8      }
9    },

```

```

10     { $unwind: "$person" },
11     {
12         $lookup: {
13             from: "Film",
14             localField: "tconst",
15             foreignField: "tconst",
16             as: "film"
17         }
18     },
19     { $unwind: "$person.job" },
20     { $match: { "film.isAdult": true } },
21     {
22         $group: {
23             "_id": "$person.nconst",
24             "primaryName": { $first: "$person.primaryName" },
25             "totalAdultFilms": { $sum: 1 }
26         }
27     },
28     { $match: { totalAdultFilms: { $gte: 3 } } },
29     {
30         $project: {
31             "_id": 0,
32             "nconst": "$_id",
33             "primaryName": 1,
34             "totalAdultFilms": 1
35         }
36     }
37 ]))

```

2.1.4 Select the director that participated in the most title for each genre

```

1 db.film.aggregate([
2     { $unwind: "$genres" },
3     {
4         $group: {
5             "_id": { "genre": "$genres", "director":
6                 ↳ "$directors" },
7             "titleCount": { $sum: 1 }
8         }
9     },
10    { $sort: { "_id.genre": 1, "titleCount": -1 } },
11    {

```

```

12         "_id": "$_id.genre",
13         "topDirector": { $first: "$_id.director" },
14         "maxTitleCount": { $first: "$titleCount" }
15     }
16 },
17 {
18     $project: {
19         "_id": 0,
20         "genre": "$_id",
21         "topDirector": 1,
22         "maxTitleCount": 1
23     }
24 }
25 ]))

```

2.1.5 Count the number of actors dead before the end of a series

```

1 db.CastingMember.aggregate([
2     {
3         $lookup: {
4             from: "Film",
5             localField: "tconst",
6             foreignField: "tconst",
7             as: "film"
8         }
9     },
10    { $unwind: "$film" },
11    { $match: { "film.typeTitle": "tvSeries" } },
12    { $unwind: "$job" },
13    {
14        $group: {
15            "_id": "$film.tconst",
16            "deadActorsCount": {
17                $sum: {
18                    $cond: [
19                        { $and: [{ $eq: ["$job", "actor"] }, { $gte:
20                            ↳ ["$film.endYear", "$deathYear"] }, {
21                            ↳ $lte: ["$film.startYear", "$deathYear"]
22                            ↳ } ] },
23                        1,
24                        0
25                    ]
26                }
27            }
28        }
29    }
30 ])

```

```

25         }
26     },
27     { $match: { "deadActorsCount": { $gt: 0 } } },
28     { $count: "output" }
29 ]))

```

2.1.6 Select the top Y favorite genre of the year X

To illustrate the query, we set the top to 5 and the year to 2012.

```

1  db.Film.aggregate([
2      {
3          $lookup: {
4              from: "Rating",
5              localField: "tconst",
6              foreignField: "tconst",
7              as: "ratings"
8          }
9      },
10     { $unwind: "$ratings" },
11     { $match: { "startYear": 2012 } },
12     {
13         $addFields: {
14             "score": { $multiply: ["$averageRating", { $ln: {
15                 ↪ $add: ["$ratings.numVotes", 1] } } ] }
16         }
17     },
18     { $sort: { "score": -1 } },
19     { $limit: 5 },
20     {
21         $project: {
22             "_id": 0,
23             "tconst": 1,
24             "startYear": 1,
25             "score": 1
26         }
27     })

```

2.1.7 Select the average number of movies released per decade per region

```

1  db.film.aggregate([

```

```

2      {
3        $match: {
4          "titleType": "movie",
5          "isOriginal": true
6        }
7      },
8      {
9        $addFields: {
10         "decade": { $subtract: [ { $toInt: { $substr:
11           ↪ ["$startYear", 0, 3] } }, { $mod: [ { $toInt: {
12             ↪ $substr: ["$startYear", 0, 3] } }, 10 ] } ] }
13       }
14     },
15     { $unwind: "$region" },
16     {
17       $group: {
18         "_id": { "region": "$region", "decade": "$decade"
19           ↪ },
20         "movieCount": { $sum: 1 }
21       }
22     },
23     {
24       $group: {
25         "_id": "$_id.region",
26         "averageMoviesPerDecade": { $avg: "$movieCount" }
27       }
28     },
29     {
30       $project: {
31         "_id": 0,
32         "region": "$_id",
33         "averageMoviesPerDecade": 1
34       }
35     }
36   ] )

```

2.1.8 Average number of appearance for an actor per series

```

1 db.castingMembers.aggregate([
2   {
3     $match: {
4       "category": "actor"
5     }
6   },

```

```

7   {
8     $group: {
9       "_id": {
10        "actor": "$nconst",
11        "series": "$tconst"
12      },
13      "appearances": { $sum: 1 }
14    }
15  },
16  {
17    $group: {
18      "_id": "$_id.actor",
19      "averageAppearances": { $avg: "$appearances" }
20    }
21  },
22  {
23    $sort: { "averageAppearances": -1 }
24  },
25  {
26    $limit: 3
27  },
28  {
29    $lookup: {
30      from: "personPrimary",
31      localField: "_id",
32      foreignField: "nconst",
33      as: "actorDetails"
34    }
35  },
36  {
37    $project: {
38      "_id": 0,
39      "actorName": { $arrayElemAt:
40        ↪ ["$actorDetails.primaryName", 0] },
41      "averageAppearances": 1
42    }
43  })

```


2.2 Model 2

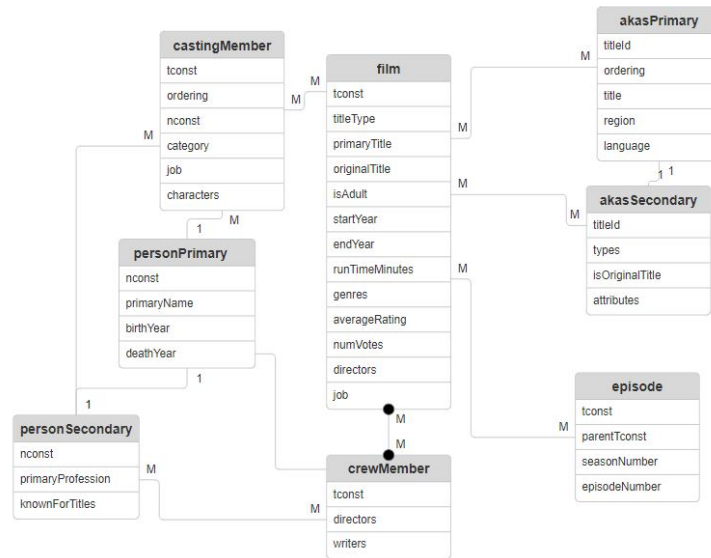


Figure 2: Entity-relationship diagram for model 2

2.2.1 Select the top 10 movies of the year X shorter than 2 hours

```

1  db.film.find([
2      {
3          $match: {
4              "titleType": "movie",
5              "startYear": 1999,
6              "runTimeMinutes": { $lt: 120 }
7          }
8      },
9      {
10         $addFields: {
11             "ratingScore": {
12                 $multiply: [
13                     "$averageRating",
14                     { $log: { $add: ["$numVotes", 1] } }
15                 ]
16             }
17         }
18     },
19     { $sort: { "ratingScore": -1 } },

```

```

20     { $limit: 10 },
21     {
22         $project: {
23             "_id": 0,
24             "tconst": 1,
25             "primaryTitle": 1,
26             "startYear": 1,
27             "ratingScore": 1,
28             "runTimeMinutes": 1
29         }
30     }
31 ]))

```

2.2.2 Select a movie containing the string “XXX” in at least one of its titles

```

1  db.film.aggregate([
2      {
3          $match: {
4              "titleType": "movie"
5          }
6      },
7      {
8          $lookup: {
9              from: "akasPrimary",
10             localField: "tconst",
11             foreignField: "titleId",
12             as: "akas"
13         }
14     },
15     {
16         $unwind: "$akas.title"
17     },
18     {
19         $or: [
20             { "primaryTitle": { $regex: "love", $options: "i" } },
21             { "akas.title": { $regex: "love", $options: "i" } }
22         ]
23     },
24     {
25         $project: {
26             "_id": 0,
27             "tconst": 1,
28             "primaryTitle": 1

```

```

29     }
30   }
31 ])
```

2.2.3 Actors that played in at least 3 adult films

```

1  db.film.aggregate([
2    {
3      $match: {
4        "isAdult": 1,
5        "job.job": "actor"
6      }
7    },
8    {
9      $unwind: "$job"
10   },
11   {
12     $group: {
13       "_id": "$job.job",
14       "count": { $sum: 1 }
15     }
16   },
17   {
18     $match: {
19       "count": { $gt: 3 }
20     }
21   },
22   {
23     $lookup: {
24       from: "personPrimary",
25       localField: "job.nconst",
26       foreignField: "nconst",
27       as: "person"
28     }
29   },
30   {
31     $project: {
32       "_id": 0,
33       "person.name": 1,
34       "job.job": 1,
35       "count": 1
36     }
37   }
38 ])
```

2.2.4 Select the director that participated in the most films for each genre

```
1 db.film.aggregate([
2   {
3     $unwind: "$genres"
4   },
5   {
6     $group: {
7       "_id": { "genre": "$genres", "director": "$directors"
8         ↪ },
9       "count": { $sum: 1 }
10    },
11    {
12      $sort: { "count": -1 }
13    },
14    {
15      $group: {
16        "_id": "$_id.genre",
17        "director": { $first: "$_id.director" },
18        "count": { $first: "$count" }
19      }
20    },
21    {
22      $project: {
23        "_id": 0,
24        "genre": "$_id",
25        "director": 1,
26        "count": 1
27      }
28    }
29  ])
```

2.2.5 Count the number of actors dead before the end of a series

```
1 db.film.aggregate([
2   {
3     $match: {
4       "titleType": "tvSeries"
5     }
6   },
7   {
```

```

8     $unwind: "$job"
9 },
10 {
11     $match: {
12         "job.job": "actor"
13     }
14 },
15 {
16     $lookup: {
17         from: "personPrimary",
18         localField: "job.nconst",
19         foreignField: "nconst",
20         as: "person"
21     }
22 },
23 {
24     $unwind: "$person.deathYear"
25 },
26 {
27     $match: {
28         $or: [
29             { "person.deathYear": { $lt: "$endYear" } },
30             { "person.deathYear": { $exists: true }, "endYear": {
31                 ⇨ $exists: false } }
32         ]
33     },
34     {
35         $group: {
36             "_id": "$person.nconst",
37             "count": { $sum: 1 }
38         }
39     },
40     {
41         $project: {
42             "_id": 0,
43             "count": 1
44         }
45     }
46 ] )
47

```

2.2.6 Select the top Y favorite genres of the year X

```
1 db.film.aggregate([
2   {
3     $match: {
4       "titleType": "movie",
5       "startYear": 1999
6     }
7   },
8   {
9     $group: {
10      "_id": "$genres",
11      "averageRating": { $avg: "$averageRating" }
12    }
13  },
14  { $sort: { "averageRating": -1 } },
15  { $limit: 3 },
16  {
17    $project: {
18      "_id": 0,
19      "genre": "$_id",
20      "averageRating": 1
21    }
22  }
23 ])
```

2.2.7 Select the average number of films released per decade per region

```
1 db.film.aggregate([
2   {
3     $match: {
4       "titleType": "movie"
5     }
6   },
7   {
8     $lookup: {
9       from: "akasPrimary",
10      localField: "tconst",
11      foreignField: "titleId",
12      as: "akas"
13    }
14  },
15  {
```

```

16     $unwind: "$akas.region"
17   },
18   {
19     $group: {
20       "_id": { "region": "$akas.region", "decade": {
21         ↪ $trunc: { $divide: [ "$startYear", 10 ] } } },
22       "count": { $sum: 1 }
23     }
24   },
25   {
26     $group: {
27       "_id": "$_id.region",
28       "average": { $avg: "$count" }
29     }
30   },
31   {
32     $project: {
33       "_id": 0,
34       "region": "$_id",
35       "average": 1
36     }
37   }
38 ] )

```

2.2.8 Average number of appearances for an actor per series

```

1 db.film.aggregate([
2   {
3     $match: {
4       "titleType": "tvSeries"
5     }
6   },
7   {
8     $unwind: "$job"
9   },
10  {
11    $match: {
12      "job.job": "actor"
13    }
14  },
15  {
16    $lookup: {
17      from: "episode",
18      localField: "tconst",

```

```

19     foreignField: "parentTconst",
20     as: "episodes"
21   }
22 },
23 {
24   $unwind: "$episodes"
25 },
26 {
27   $lookup: {
28     from: "personPrimary",
29     localField: "job.nconst",
30     foreignField: "nconst",
31     as: "person"
32   }
33 },
34 {
35   $group: {
36     "_id": "$person.nconst",
37     "count": { $sum: 1 }
38   }
39 }
40 {
41   $project: {
42     "_id": 0,
43     "person.name": 1,
44     "primaryTitle": 1,
45     "count": 1
46   }
47 }
48 ])
49

```

3 Choice of sharding key

For this project, we propose two different types of sharding :

- Sharding on Film.job
- Sharding on Film.startYear

In the first place, we chose to shard on Film.job as our use case often requires to filter on "actor" which is an instance of the "job" attribute. In addition, we chose it because the amount of distinct jobs is known and doesn't vary too much. Consequently, the sharding storage does not scale linearly with the number of films.

For similar reasons, another extremely interesting attribute is Film.startYear which is the release year for films and the first diffusion year for series. Consequently, we decided to propose a second sharding on this attribute.

4 Communications cost model

In the tables bellow (table 1, table 2), we propose a summary of the cost models that we provide interpretation for in the conclusion (section 5). For any further information, you can refer to the [google sheet](#) we used so as to perform these cost computations.

Requête	Coût sharding 1 (job)				Coût sharding 2 (StartYear)			
	S	Shuffle	Output	#msg	S	Shuffle	Output	#msg
Ru1	100	0	1.05E+07	1.05E+11	1	0	3.16E+05	3.16E+09
Ru2	100	7.93E+13	1.03E+04	7.93E+16	100	0	1.03E+04	1.04E+07
Ru3	100	1.39E+14	3061444	6.96E+16	100	1.39E+14	3061444	6.96E+16
Ru4	100	1.39E+14	956976	1.39E+16	100	0	956976	9.57E+07
Rda1	100	1.39E+14	3061444	6.96E+15	100	1.39E+14	3061444	6.96E+15
Rda2	100	0	2360	6.15E+04	1	0	2360	5.90E+04
Rda3	100	7.93E+13	1.05E+07	7.93E+14	15.3	0	1.05E+07	1.05E+08
Rda4	100	1.12E+21	3061444	2.24E+21	100	1.12E+21	3061444	2.24E+21

Table 1: Cost Model summary for Model 1

Requête	Coût sharding 1 (job)				Coût sharding 2 (StartYear)			
	S	Shuffle	Output	#msg	S	Shuffle	Output	#msg
Ru1	100	0	1.05E+07	1.05E+11	3	0	1.05E+07	1.05E+11
Ru2	100	7.93E+13	1.03E+04	7.93E+16	100	7.93E+13	1.03E+04	7.93E+16
Ru3	24	1.39E+14	3061444	6.96E+16	100	1.39E+14	3061444	6.96E+16
Ru4	100	1.39E+14	956976	1.39E+16	100	0	956976	9.57E+07
Rda1	24	0	3061444	1.53E+08	100	306144400	3061444	1.55E+10
Rda2	100	0	2360	6.15E+04	3	0	2360	5.91E+04
Rda3	100	7.93E+13	1.05E+07	7.93E+14	100	7.93E+13	1.05E+07	7.93E+14
Rda4	24	0	3061444	6.12E+06	100	0	3061444	6.12E+06

Table 2: Cost Model summary for Model 2

5 Conclusion

In conclusion, we can see that for both models, aggregation on job is not the best solution. When aggregating on startYear, the cost is drastically lower for most queries which makes it a preferred choice for the sharding. For the model selection, the choice might not be as straight forward. In fact, depending on the more specific use case, one model or the other could be a sensible alternative.

On another aspect, we noticed that our denormalized models produce strategies that appear to be sub-optimal. In fact, while creating our denormalized models

in chapter 2, we wanted to create hybrid strategies using mostly materializations, splits, and overloads while neglecting fusions. As a consequence, queries requiring to lookup information in multiple tables are very costly. If we were to redo the project from the start, we would most certainly emphasize fusion operations to reduce the joining costs despite the fact that it increases the size of documents. An ideal strategy would be to merge the most things possible under two main collection that would be film and person.