

Bases de Données Avancées,

Rapport Phase 2 :

Migration MongoDB

Guilhem ALBALADEJO
Info 4A
Polytech Marseille

TABLE DES MATIÈRES

I	Introduction	1
II	Migration des données	1
III	Traduction des requêtes	1
III-A	Indexes	1
III-B	Pipelines des requêtes	1
III-B1	Filmographie d'un acteur	1
III-B2	Top N films	2
III-B3	Acteurs multi-rôles	3
III-B4	Collaborations	3
III-B5	Genres populaires	4
III-B6	Évolution de carrière	4
III-B7	Classement par genre	5
III-B8	Carrière propulsée	5
III-B9	Enfants stars	6
III-C	Comparaison avec SQLite	6
IV	Documents structurées	7
V	Conclusion	7

Bases de Données Avancées,

Rapport Phase 2 :

Migration MongoDB

I. INTRODUCTION

Ce rapport correspond au compte-rendu de la deuxième phase du cours de bases de données avancées. Nous verrons au sein de celui-ci comment créer une base de données mongoDB, gérer la migration de nos données depuis SQLite vers des collections plates dans une premier temps, puis vers des documents structurés. Nous verrons comment reconstruire nos requêtes sous mongoDB à l'aide de pipelines d'agrégation. Enfin nous verrons les avantages et inconvénients des bases de données noSQL comparées aux bases SQL.

II. MIGRATION DES DONNÉES

On cherche dans cette partie à migrer les données contenues dans notre base de données `imdb.db` créée dans la phase 1, vers une base MongoDB. On commence par installer la dernière version de MongoDB Community Edition depuis le site officiel, ainsi que la bibliothèque `pymongo` pour pouvoir interagir avec notre base de données via des scripts python. On lance notre base de données à l'aide de la commande :

```
mongod --dbpath data/mongo/standalone
```

La migration s'effectue grâce à un script `migrate_flat.py`, qui convertit nos tables telles quelles en collections plates. Comme on effectue aucune transformation sur nos tables, ce script est très court et indépendant de notre schéma. La seule exception concerne la conversion des booléens, qui sont stockés comme nombres entiers dans SQLite et doivent donc être convertis pendant la migration. Heureusement, les colonnes booléennes ont été créées avec le type "BOOLEAN" dans le schéma de notre base SQL, on peut donc vérifier le type des colonnes pour correctement convertir les champs booléens dans notre script d'import.

III. TRADUCTION DES REQUÊTES

Dans MongoDB, les requêtes s'effectuent à l'aide de pipelines d'agrégation. Nous ne pourrons donc pas réutiliser nos requêtes SQL écrites précédemment. Il faudra toutes les adapter pour pouvoir les lancer sur nos collections.

A. Indexes

Pour cela, on crée un script `queries_mongo.py` qui contiendra nos requêtes sous la forme de fonctions python. Un autre point à considérer avant de commencer à écrire nos requêtes concerne les indexes, en effet MongoDB étant un système

de gestion de base de donnée non relationnelle, l'exécution de requêtes relationnelles risque d'être très longue, de plus, MongoDB ne crée pas d'indexes automatiques comme SQLite sur nos clés primaires, mais sur les clés "`_id`", que nous n'utiliserons pas dans nos requêtes.

On crée donc les indexes suivants :

- Sur la collection `persons` : `person_id`.
- Sur la collection `titles` : `title_id`.
- Sur la collection `movies` : `movie_id`.
- Sur la collection `ratings` : `movie_id`.
- Sur la collection `principals` : `movie_id`, (`person_id`, `profession_id`, `movie_id`).
- Sur la collection `movie_titles` : `movie_id`, `title_id`.
- Sur la collection `person_profession` : `person_id`, `profession_id`.
- Sur la collection `cast` : `movie_id`, `person_id`.
- Sur la collection `movie_genres` : `movie_id`, `genre_id`.
- Sur la collection `title_ordering` : `movie_id`, `title_id`

B. Pipelines des requêtes

Nous définissons les requêtes suivantes dans un fichier `queries_mongo.py` :

1) Filmographie d'un acteur: On lance le pipeline suivant sur la collection "persons" :

```
[  
  {"$match": {"name": {  
    "$regex": actor_name,  
    "$options": "i"  
  }}},  
  {"$lookup": {  
    "from": "principals",  
    "localField": "person_id",  
    "foreignField": "person_id",  
    "as": "pr"  
  }},  
  {"$unwind": "$pr"},  
  {"$lookup": {  
    "from": "movie_titles",  
    "localField": "pr.movie_id",  
    "foreignField": "movie_id",  
    "as": "mt"  
  }},  
  {"$unwind": "$mt"},  
  {"$match": {  
    "mt.is_primary": True  
  }},  
  {"$lookup": {  
    "from": "movies",  
    "localField": "mt.movie_id",  
    "foreignField": "id",  
    "as": "m"  
  }},  
  {"$unwind": "$m"},  
  {"$group": {  
    "actor": "$name",  
    "titles": {"$push": "$m.title"},  
    "count": {"$sum": 1}  
  }},  
  {"$sort": {"titles": -1}}]
```

```

25     "localField": "mt.movie_id",
26     "foreignField": "movie_id",
27     "as": "m"
28   },
29   { "$unwind": "$m" },
30   { "$lookup": {
31     "from": "titles",
32     "localField": "mt.title_id",
33     "foreignField": "title_id",
34     "as": "t"
35   }},
36   { "$unwind": "$t" },
37   { "$lookup": {
38     "from": "cast",
39       "let": { "person_id": "$t._id" },
40       "from": "person",
41       "let": { "person_id": "$person_id" },
42       "pipeline": [
43         { "$match": {
44           "$expr": {
45             "$and": [
46               { "$eq": [
47                 "$person_id",
48                 "$$person_id"
49               ]},
50               { "$eq": [
51                 "$movie_id",
52                 "$$movie_id"
53               ]}
54             ]
55           }
56         }
57       }
58     ],
59     "as": "ca"
60   }},
61   { "$lookup": {
62     "from": "characters",
63     "localField": "ca.character_id",
64     "foreignField": "character_id",
65     "as": "c"
66   }},
67   { "$lookup": {
68     "from": "ratings",
69     "localField": "mt.movie_id",
70     "foreignField": "movie_id",
71     "as": "r"
72   }},
73   { "$unwind": "$r" },
74   { "$sort": {
75     "m.year": -1
76   }},
77   { "$project": {
78     "_id": 0,
79     "title": "$t.title_name",
80     "year": "$m.year",
81     "character_name": "$c.name",
82     "average_rating": "$r.average_rating"
83   }}
84 ]

```

On commence par filtrer la collection persons avec le nom⁵¹ de l'acteur passé dans la variable `actor_name`, puis on⁵² joint les collections contenant les informations à récupérer. On⁵³ pense bien à filtrer pour garder uniquement les titres principaux⁵⁴ et à joindre cast à la fois sur `person_id` et `movie_id`. Finalement⁵⁵ on trie par année décroissante et on projette uniquement les⁵⁶

champs qui nous intéressent.

2) *Top N films*: On lance le pipeline suivant sur la collection "movies" :

```

1  [
2    { "$match": {
3      "year": { "$gte": start_year, "$lte": end_year}
4    }},
5    { "$lookup": {
6      "from": "movie_genres",
7      "localField": "movie_id",
8      "foreignField": "movie_id",
9      "as": "mg"
10   }},
11   { "$unwind": "$mg" },
12   { "$lookup": {
13     "from": "genres",
14     "let": { "genre_id": "$mg.genre_id" },
15     "pipeline": [
16       { "$match": {
17         "$expr": {
18           "$eq": [
19             "$genre_id",
20             "$$genre_id"
21           ]
22         }
23       }},
24       { "$match": {
25         "genre_name": {
26           "$regex": genre,
27           "$options": "i"
28         }
29       }},
30       "as": "g"
31     }},
32     { "$unwind": "$g" },
33     { "$lookup": {
34       "from": "ratings",
35       "localField": "movie_id",
36       "foreignField": "movie_id",
37       "as": "r"
38     }},
39     { "$unwind": "$r" },
40     { "$sort": { "r.average_rating": -1 } },
41     { "$limit": N },
42     { "$lookup": {
43       "from": "movie_titles",
44       "let": { "movie_id": "$movie_id" },
45       "pipeline": [
46         { "$match": { "$expr": { "$eq": [
47           "$movie_id",
48           "$$movie_id"
49         ] } } },
50         { "$match": { "is_primary": True } }
51       ],
52       "as": "mt"
53     }},
54     { "$lookup": {
55       "from": "titles",
56       "localField": "mt.title_id",
57       "foreignField": "title_id",
58       "as": "t"
59     }},
60     { "$unwind": "$t" },
61     { "$project": {
62       "title": "$t.title_name",
63       "year": "$t.year",
64       "primary": "$mt.is_primary"
65     }}
66   }
67 ]

```

```

57     "_id": 0,
58     "movie_id": 1,
59     "title": "$t.title_name",
60     "year": 1,
61     "average_rating":
62       → "$r.average_rating",
63     "num_votes": "$r.num_votes"
64   } }
]

```

On commence par filtrer movies pour garder uniquement les films sortis dans l'intervalle [start_year, end_year]. On joint "movie_genres" et "genres" pour garder uniquement les films du genre qui nous intéresse à l'aide d'une pipeline imbriquée dans le lookup sur "genres" pour ne joindre que sur les films du bon genre. Finalement on joint les collections contenant les informations restantes, on ne garde que les titres principaux, on trie par note moyenne décroissante et on projette les champs voulus.

3) *Acteurs multi-rôles*: On lance le pipeline suivant sur la collection "cast" :

```

1 [
2   { "$group": {
3     "_id": {
4       "person_id": "$person_id",
5       "movie_id": "$movie_id"
6     },
7     "charactersPlayed": {"$addToSet":
8       → "$character_id"
9   }},
10   { "$addFields": {
11     "charactersCnt": {"$size":
12       → "$charactersPlayed"
13   }},
14   { "$match": {
15     "charactersCnt": {"$gt": 1}
16   }},
17   { "$sort": {"charactersCnt": -1}},
18   { "$lookup": {
19     "from": "persons",
20     "localField": "_id.person_id",
21     "foreignField": "person_id",
22     "as": "pe"
23   }},
24   { "$unwind": "$pe"},
25   { "$lookup": {
26     "from": "movie_titles",
27     "localField": "_id.movie_id",
28     "foreignField": "movie_id",
29     "as": "mt"
30   }},
31   { "$unwind": "$mt"}, // mt.is_primary: True
32   { "$match": {"mt.is_primary": True}},
33   { "$lookup": {
34     "from": "titles",
35     "localField": "mt.title_id",
36     "foreignField": "title_id",
37     "as": "t"
38   }},
39   { "$unwind": "$t"}, // t._id: 0
40   { "$project": {
41     "_id": 0,
42     "movie_id": 1,
43     "title": "$t.title_name",
44     "charactersCnt": 1,
45   } }
46 ]

```

```

40   "person_id": "$_id.person_id",
41   "actor_name": "$pe.name",
42   "movie_id": 1,
43   "title": "$t.title_name",
44   "charactersCnt": 1,
45 }
]

```

On regroupe la collection par films et acteurs, pour garder uniquement les films dans lequel un même acteur est associé à plusieurs personnages, on compte le nombre de personnages pour chaque groupe en les stockant dans un tableau puis en lisant la taille de ce tableau. On filtre pour ne garder que les acteurs multi-rôles. Finalement on récupère le reste des informations à renvoyer.

4) *Collaborations*: On commence par récupérer tous les films joués par l'acteur dans un tableau avec :

```

pmPipeline = [
  { "$match": {
    "name": {
      "$regex": actor_name,
      "$options": "i"
    }
  }},
  { "$lookup": {
    "from": "principals",
    "localField": "person_id",
    "foreignField": "person_id",
    "as": "pr"
  }},
  { "$lookup": {
    "from": "professions",
    "localField": "pr.profession_id",
    "foreignField": "profession_id",
    "as": "p"
  }},
  { "$unwind": "$p" },
  { "$match": {
    "p.job_name": {
      "$regex": "actor",
      "$options": "i"
    }
  }},
  { "$project": {
    "_id": 0,
    "movie_id": "$pr.movie_id"
  }}
]
# array contenant les films joués par
→ l'acteur
playedMovies =
→ next(db["persons"].aggregate(pmPipeline))
[ "movie_id" ]

```

Puis on utilise ce tableau pour filtrer toutes les personnes ayant travaillé dans ces films, avant de filtrer pour ne garder que les directeurs, en lançant la pipeline suivante sur la collection "persons" :

```

1 [
2   { "$match": {
3     "movie_id": {"$in": playedMovies}
4   } },

```

```

5     { "$lookup": {
6         "from": "professions",
7         "localField": "profession_id",
8         "foreignField": "profession_id",
9         "as": "p"
10    },
11    { "$unwind": "$p",
12    { "$match": {
13        "p.job_name": {
14            "$regex": "director",
15            "$options": "i"
16        }
17    },
18    { "$lookup": {
19        "from": "persons",
20        "localField": "person_id",
21        "foreignField": "person_id",
22        "as": "pe"
23    },
24    { "$unwind": "$pe",
25    { "$group": {
26        "_id": {"person_id": "$person_id"},
27        "name": {"$first": "$pe.name"},
28        "movies": {"$addToSet": "$movie_id"}
29    },
30    { "$addFields": {"movieCnt": {"$size":
31        "$movies"}},,
32    { "$sort": {"movieCnt": -1}},
33    { "$project": {
34        "_id": 0,
35        "name": 1,
36        "movieCnt": 1
37    } }
38  ]

```

5) *Genres populaires:* On lance le pipeline suivant sur la collection "movie_genres" :

```

1 [
2     { "$lookup": {
3         "from": "genres",
4         "localField": "genre_id",
5         "foreignField": "genre_id",
6         "as": "g"
7    },
8     { "$unwind": "$g"},,
9     { "$lookup": {
10        "from": "movies",
11        "localField": "movie_id",
12        "foreignField": "movie_id",
13        "as": "m"
14    },
15    { "$unwind": "$m"},,
16    { "$lookup": {
17        "from": "ratings",
18        "localField": "m.movie_id",
19        "foreignField": "movie_id",
20        "as": "r"
21    },
22    { "$unwind": "$r"},,
23    { "$group": {
24        "_id": {"genre_id": "$genre_id"},,
25        "genre_name": {"$first": -->
26        "$g.genre_name"},,
27        "avg_rating": {"$avg": -->
28        "$r.average_rating"},,
29        "movies": {"$addToSet": "$movie_id"}}
30    },
31    { "$addFields": {
32        "movieCnt": {"$size": "$movies"}}
33    },
34    { "$match": {
35        "avg_rating": {"$gt": 7.0},
36        "movieCnt": {"$gt": 50}
37    },
38    { "$project": {
39        "_id": 0,
40        "genre_name": 1,
41        "avg_rating": 1,
42        "movieCnt": 1
43    } }
44  ]

```

```

26     "avg_rating": {"$avg": -->
27     "$r.average_rating"},,
28     "movies": {"$addToSet": "$movie_id"}}
29   },
30   { "$addFields": {
31     "movieCnt": {"$size": "$movies"}}
32   },
33   { "$match": {
34     "avg_rating": {"$gt": 7.0},
35     "movieCnt": {"$gt": 50}
36   },
37   { "$project": {
38     "_id": 0,
39     "genre_name": 1,
40     "avg_rating": 1,
41     "movieCnt": 1
42   } }
43 ]

```

On joint les collections concernées, on groupe par genre et on filtre par genre ayant une note moyenne > 7.0 et un nombre de films > 50. Finalement on projette les champs qui nous intéressent.

6) *Évolution de carrière:* On lance le pipeline suivant sur la collection "movies" :

```

1 [
2     { "$lookup": {
3         "from": "cast",
4         "localField": "movie_id",
5         "foreignField": "movie_id",
6         "as": "c"
7    },
8     { "$unwind": "$c"},,
9     { "$lookup": {
10        "from": "persons",
11        "localField": "c.person_id",
12        "foreignField": "person_id",
13        "as": "pe"
14    },
15    { "$unwind": "$pe"},,
16    { "$match": {
17        "pe.name": {
18            "$regex": actor_name,
19            "$options": "i"
20        }
21    },
22    { "$lookup": {
23        "from": "ratings",
24        "localField": "movie_id",
25        "foreignField": "movie_id",
26        "as": "r"
27    },
28    { "$unwind": "$r"},,
29    { "$addFields": {
30        "decade": {
31            "$subtract": [
32                "$year",
33                {"$mod": [
34                    "$year",
35                    10
36                ] }
37            ] }
38        } }
39  ]

```

```

39     } },
40     { "$group": {
41         "_id": "$decade",
42         "movie_list": { "$addToSet":
43             { "$movie_id": ,
44             "avg_rating": { "$avg":
45                 "r.average_rating" }
46         } },
47         { "$addFields": {
48             "movieCnt": { "$size": "movie_list" }
49         } },
50         { "$project": {
51             "_id": 0,
52             "decade": "$_id",
53             "movieCnt": 1,
54             "avg_rating": 1
55         } },
56         { "$sort": {
57             "decade": -1
58         } }
59     } ]
60 
```

On commence par filtrer par le nom de l'acteur après deux lookups, puis on calcule la décennie de chaque film selon décennie = year - (year % 10). Finalement on regroupe par décennie, sur lesquelles on calcule la note moyenne de tous les films et le nombre de films.

7) Classement par genre: On lance le pipeline suivant sur la collection "movie_genres" :

```

1 [ {
2     { "$lookup": {
3         "from": "ratings",
4         "localField": "movie_id",
5         "foreignField": "movie_id",
6         "as": "r"
7     } },
8     { "$unwind": "$r" },
9     { "$setWindowFields": {
10         "partitionBy": "$genre_id",
11         "sortBy": { "r.average_rating": -1 },
12         "output": {
13             "rank": { "$rank": {} }
14         }
15     } },
16     { "$match": {
17         "rank": { "$lte": 3 }
18     } },
19     { "$lookup": {
20         "from": "genres",
21         "localField": "genre_id",
22         "foreignField": "genre_id",
23         "as": "g"
24     } },
25     { "$unwind": "$g" },
26     { "$lookup": {
27         "from": "movie_titles",
28         "localField": "movie_id",
29         "foreignField": "movie_id",
30         "as": "mt"
31     } },
32     { "$unwind": "$mt" },
33     { "$match": { "mt.is_primary": True } },
34     { "$lookup": {
35         "from": "titles",
36         "localField": "mt.title_id",
37         "foreignField": "title_id",
38         "as": "t"
39     } },
40     { "$unwind": "$t" },
41     { "$sort": { "g.genre_name": 1 } },
42     { "$project": {
43         "_id": 0,
44         "genre": "$g.genre_name",
45         "title": "$t.title_name",
46         "rank": 1
47     } }
48 ] 
```

```

35         "from": "titles",
36         "localField": "mt.title_id",
37         "foreignField": "title_id",
38         "as": "t"
39     } },
40     { "$unwind": "$t" },
41     { "$sort": { "g.genre_name": 1 } },
42     { "$project": {
43         "_id": 0,
44         "genre": "$g.genre_name",
45         "title": "$t.title_name",
46         "rank": 1
47     } }
48 ] 
```

On utilise \$setWindowFields pour partitionner par genre et trier par note moyenne décroissante, on garde les 3 notes les plus hautes à l'aide d'un champ "rank" que l'on ajoute lors de la partition. Finalement on récupère et projette les champs que l'on veut renvoyer.

8) Carrière propulsée: On commence par créer deux pipelines distinctes "fh_pipeline" et "lnh_pipeline" qui récupèrent respectivement les premiers "hit" et derniers "non hit" de chaque personne, puis on lance la pipeline suivante sur la collection "persons" :

```

[ {
1     { "$lookup": {
2         "from": "principals",
3         "let": { "person_id": "$person_id" },
4         "pipeline": fh_pipeline,
5         "as": "fh"
6     } },
7     { "$unwind": "$fh" },
8     { "$lookup": {
9         "from": "principals",
10        "let": { "person_id": "$person_id" },
11        "pipeline": lnh_pipeline,
12        "as": "lnh"
13    } },
14     { "$unwind": {
15         "path": "$lnh",
16         "preserveNullAndEmptyArrays": True
17     } },
18     { "$match": {
19         "$or": [
20             { "lnh": None },
21             { "$expr": { "$gt": [ "$fh.year",
22             "lnh.year" ] } }
23         ]
24     } },
25     { "$lookup": {
26         "from": "movie_titles",
27         "localField": "fh.movie_id",
28         "foreignField": "movie_id",
29         "as": "mt"
30     } },
31     { "$unwind": "$mt" },
32     { "$match": { "mt.is_primary": True } },
33     { "$lookup": {
34         "from": "titles",
35         "localField": "mt.title_id",
36         "foreignField": "title_id",
37         "as": "t"
38     } }
39 ] 
```

```

37     "as": "t"
38   },
39   {"$unwind": "$t"},
40   {"$project": {
41     "_id": 0,
42     "name": 1,
43     "first_hit_title": "$t.title_name",
44     "first_hit_year": "$fh.year"
45   }
46 }

```

On utilise les deux pipelines définies précédemment sur principaux lors de l'agrégation pour connaître l'année du first et last non hit pour chaque personne. On filtre ensuite pour garder les personnes n'ayant eu aucun non hit après leur premier hit. Finalement on récupère les champs que l'on veut renvoyer.

9) *Enfants stars*: On lance le pipeline suivant sur la collection "cast" :

```

1 [
2   {"$lookup": {
3     "from": "ratings",
4     "localField": "movie_id",
5     "foreignField": "movie_id",
6     "as": "r"
7   },
8   {"$unwind": "$r"},
9   {"$match": {
10     "r.num_votes": {"$gt": 200000}
11   }},
12   {"$lookup": {
13     "from": "movies",
14     "localField": "movie_id",
15     "foreignField": "movie_id",
16     "as": "m"
17   },
18   {"$unwind": "$m"},  

19   {"$lookup": {
20     "from": "persons",
21     "localField": "person_id",
22     "foreignField": "person_id",
23     "as": "pe"
24   },
25   {"$unwind": "$pe"},  

26   {"$addFields": {
27     "age": {
28       "$subtract": ["$m.year",
29                     "$pe.birth_year"]
30     }
31   },
32   {"$match": {
33     "age": {"$lt": 18}
34   }},
35   {"$lookup": {
36     "from": "movie_titles",
37     "localField": "movie_id",
38     "foreignField": "movie_id",
39     "as": "mt"
40   },
41   {"$unwind": "$mt"},  

42   {"$match": {"mt.is_primary": True}},
43   {"$lookup": {
44     "from": "titles",
45     "localField": "mt.title_id",
46   }
47 }

```

```

45   "foreignField": "title_id",
46   "as": "t"
47 },
48   {"$unwind": "$t"},  

49   {"$sort": {
50     "pe.name": 1,
51     "age": 1
52   }},
53   {"$project": {
54     "_id": 0,
55     "name": "$pe.name",
56     "movie_title": "$t.title_name",
57     "votes": "$r.num_votes",
58     "age": 1
59   }
60 }
61 ]

```

On commence par garder uniquement les films ayant plus de 200000 votes, puis on ajoute un champ contenant l'âge des acteurs au moment où ils ont joués dans un film donné, finalement on filtre pour ne garder que ceux ayant moins de 18 ans et on projette les champs que l'on veut retourner.

C. Comparaison avec SQLite

On mesure à présent le temps d'exécution de chaque requête dans un fichier test_queries.py pour comparer leurs performances avec celles que l'on avait définies sur SQLite (on prends ici les temps que l'on avait mesuré après la création d'indexées, pour refléter la performance réelle de SQLite) :

Requête	SQLite (ms)	mongoDB (ms)	Augmentation (%)
Q1 - Filmographie	1 850	2 027	9
Q2 - Top N films	72	2 492	3 361
Q3 - Multi-rôles	931	2 012	116
Q4 - Collaborations	204	196	-4
Q5 - Genres populaires	111	16 275	14 562
Q6 - Évolution de carrière	108	17 006	15 646
Q7 - Classement par genre	260	4 686	1 702
Q8 - Carrières propulsées	1765	78 427	4 343
Q9 - Enfants star	78	9 312	11 838

TABLE I. TABLEAU DE COMPARAISON DES TEMPS D'EXÉCUTION

On voit que si l'augmentation du temps d'exécution est non négligeable pour les requêtes les plus simples, il devient extrêmement important pour les requêtes plus complexes. Certaines d'entre elles atteignent des temps d'exécutions beaucoup trop élevés pour pouvoir être utilisé dans le contexte d'une application réactive sans impacter l'expérience utilisateur.

En comparant l'écriture de nos requêtes à celle que l'on avait écrit dans la phase 1 du projet, on remarque également que nos pipelines d'agrégation sont beaucoup plus longs et complexes que nos requêtes SQL. Cela s'explique en partie par la forme que prennent ces pipelines : des documents JSON, ou dict python, mais également par l'écriture de celles-ci, très longues et complexes.

Cela s'explique par la forme que prend notre base de données : elle suit encore le schéma relationnel que l'on avait défini pour notre base de données SQLite, avec des collections très normalisées et beaucoup de tables d'associations. Si MongoDB permet de traiter ce genre de base de données,

on remarque vite que cela impacte fortement la complexité des requêtes : la majeure partie de nos pipelines est occupée par des lookups et des unwinds, qui en plus d'augmenter la complexité et réduire la lisibilité de nos requêtes, impacte fortement la performance de celles-ci, en effet MongoDB n'est pas optimisé pour cela.

IV. DOCUMENTS STRUCTURÉS

Afin de tirer pleinement partie des capacités noSQL de mongoDB, on va créer une collection de documents structurés movies_complete contenant toutes les informations associées à un film pour pouvoir les récupérer à partir d'un seul document sans lookup.

Les documents de notre collection suivront le modèle :

```

1  {
2      "_id": "tt0111161",
3      "title": "The Shawshank Redemption",
4      "year": 1994,
5      "runtime": 142,
6      "genres": [ "Drama" ],
7      "rating": {
8          "average": 9.3,
9          "votes": 2500000
10     },
11     "directors": [
12         { "person_id": "nm0001104", "name": "Frank Darabont" }
13     ],
14     "cast": [
15         {
16             "person_id": "nm0000209",
17             "name": "Tim Robbins",
18             "characters": [ "Andy Dufresne" ],
19             "ordering": 1
20         },
21         {
22             "person_id": "nm0000151",
23             "name": "Morgan Freeman",
24             "characters": [ "Red" ],
25             "ordering": 2
26         }
27     ],
28     "writers": [
29         { "person_id": "nm0175840", "name": "Stephen King", "category": "novel" },
30         { "person_id": "nm0001104", "name": "Frank Darabont", "category": "screenplay" }
31     ],
32     "titles": [
33         { "region": "US", "title": "The Shawshank Redemption" },
34         { "region": "FR", "title": "Les Évadés" }
35     ]
36 }
```

Pour construire ces documents, on se base sur les données déjà présentes dans notre base de données, pour construire les documents structurés à l'aide d'une pipeline d'agrégation dans le fichier migrate_structured.py.

Cette pipeline, servant à récupérer un film complet est une pipeline très complexe est très longue, récupérant des informations dans de nombreuses collections, en effet si l'on voulait se passer de la collection movies_complete, et reconstruire un film à l'aide de cette pipeline, cela prends : 146 ms, alors que simplement aller chercher le même document dans movies complete prends : 0,1 ms.

De plus, en utilisant la méthode storageSize() sur les collections dans mongosh, on remarque que la taille totale de nos collections plates avec index est de : 82,69 Mo, alors que celle de notre collection structurée est de : 36,32 soit plus de deux fois moins grande.

Créer une collection structurée représente donc une manière de stocker l'information sous une forme plus condensée et de la récupérer de manière plus performante, on voit ici clairement l'avantage de construire notre base de données en suivant le modèle de document dénormalisé classique de MongoDB plutôt que d'utiliser un schéma relationnel.

V. CONCLUSION

Si les bases de données noSQL comme mongoDB peuvent être manipulées comme des bases relationnelles, elles pour cet usage là bien moins performantes que des bases de données relationnelles classiques. Leur avantage se trouve dans leurs capacité à gérer des données structurées autrement que sous forme de tableau rigide, et ainsi gagner en flexibilité, en espace de stockage, et en performance, à condition de stocker sous forme pré-agrégée les résultat de nos requêtes.