

## 1. Problem Statement

---

### Case Study: MovieFlix Analytics

---

MovieFlix, a global movie streaming platform, wants to:

- Find the most popular genres and top-rated movies in each country.
- Calculate average watch times, total views per genre, and revenue by region.
- Generate business dashboards and user recommendations, all from millions of viewing records.



**The challenge:**

How can MovieFlix efficiently analyze, summarize, and transform huge volumes of streaming data—without exporting it to another system or writing complex, slow code?

## 2. Learning Objectives

---

By the end of this tutorial, you will:

- Understand what the MongoDB Aggregation Framework is and why it's powerful.
  - Build multi-stage aggregation pipelines using `$match`, `$group`, `$project`, and more.
  - Use pipeline stages to filter, group, reshape, and analyze data.
  - Apply best practices for performance and maintainability.
-

### 3. Concept Introduction with Analogy

#### Analogy: The Data Conveyor Belt

Imagine MovieFlix’s analytics as a high-tech conveyor belt in a mailroom:

- Each letter (document) passes through a series of stations (pipeline stages).
- Some stations filter out junk mail ( `$match` ).
- Others sort mail into bins by city or sender ( `$group` ).
- Some reformat addresses or add labels ( `$project` ).
- At the end, you have exactly the summary or report you need.

The aggregation pipeline is this conveyor belt—each stage transforms or filters the data, step by step, until you get your answer.

### 4. Technical Deep Dive

#### A. What is the Aggregation Framework?

- The Aggregation Framework is MongoDB’s built-in tool for advanced data analysis and transformation.
- It uses a **pipeline**: an array of stages, each performing a specific operation on the documents.
- Each stage’s output becomes the next stage’s input—like a series of assembly line steps.

##### Pipeline Syntax:

```
db.collection.aggregate([
  { $match: { ... } },
  { $group: { ... } },
  { $project: { ... } },
  // ...more stages
])
```

- Pipelines can have any number of stages, in any order

#### B. Key Pipeline Stages

##### 1. \$match: Filtering Documents

- Like a query filter; selects only documents that meet criteria.
- Should be placed early for efficiency.

##### Example:

```
{ $match: { country: "USA", "rating": { $gte: 8 } } }
```

- Filters for US movies with rating 8 or higher.

## 2. \$group: Aggregating Data

- Groups documents by a field (or expression) and computes aggregates (sum, avg, count, etc.).
- The `_id` field specifies the group key.

**Example:**

```
{ $group: {
  _id: "$genre",
  totalViews: { $sum: "$views" },
  avgRating: { $avg: "$rating" }
}
```

- Groups by genre, calculates total views and average rating per genre.
- 

## 3. \$project: Shaping Output

- Reshapes documents, includes/excludes fields, computes new fields.
- Can rename, transform, or format fields.

**Example:**

```
{ $project: {
  _id: 0,
  genre: "$_id",
  totalViews: 1,
  avgRating: { $round: ["$avgRating", 2] }
}
```

- Outputs only genre, totalViews, and a rounded avgRating.
- 

## 4. Other Useful Stages

- `$sort` : Orders documents (e.g., by totalViews descending).
  - `$limit` : Restricts the number of output documents.
  - `$unwind` : Deconstructs arrays into separate documents.
  - `$addFields` : Adds computed fields.
- 

## 5. How the Pipeline Works

- Documents flow through each stage in order.
  - Each stage can filter, group, sort, or reshape data.
  - The output of the last stage is your final result.
- 

## C. Aggregation Pipeline Example

Suppose MovieFlix wants to find the top 3 genres by total views in 2024:

```
db.watchHistory.aggregate([
  // 1. Only 2024 records
  { $match: { year: 2024 } },
  // 2. Group by genre, sum views
  { $group: { _id: "$genre", totalViews: { $sum: "$views" } } },
  // 3. Sort by totalViews descending
  { $sort: { totalViews: -1 } },
  // 4. Limit to top 3
  { $limit: 3 },
  // 5. Project clean output
  { $project: { _id: 0, genre: "$_id", totalViews: 1 } }
])
```

**Explanation:**

- `$match` : Filters for 2024.
- `$group` : Sums up views per genre.
- `$sort` and `$limit` : Gets top 3 genres.
- `$project` : Formats output for reporting.

---

**D. Best Practices for Aggregation Pipelines**

- **Place `$match` early** to reduce data volume for later stages.
- **Use indexes** on fields used in `$match` for performance.
- **Keep documents small**—avoid unnecessary fields with `$project`.
- **Test each stage** separately to debug and optimize.
- **Limit pipeline complexity** for maintainability.

---

## 5. Step-by-Step Data Modeling & Code Walkthrough

Let's build a real MovieFlix aggregation pipeline.

---

**A. Sample Document Structure**

```
{
  "_id": ObjectId("..."),
  "movie": "Edge of Tomorrow",
  "genre": "Sci-Fi",
  "country": "USA",
  "views": 15000,
  "rating": 8.2,
  "year": 2024
}
```

**B. Find the Average Rating and Total Views per Genre in the USA**

```
db.watchHistory.aggregate([
  { $match: { country: "USA" } },
  { $group: {
    _id: "$genre",
    totalViews: { $sum: "$views" },
    avgRating: { $avg: "$rating" }
  } }
```

```
    }
  },
  { $project: {
    _id: 0,
    genre: "$_id",
    totalViews: 1,
    avgRating: { $round: ["$avgRating", 2] }
  }
}
])
```

Explanation:

- Filters for USA records.
- Groups by genre.
- Sums views and averages ratings.
- Projects a clean, rounded output.

C. Find Top 5 Movies by Views in 2024

```
db.watchHistory.aggregate([
  { $match: { year: 2024 } },
  { $sort: { views: -1 } },
  { $limit: 5 },
  { $project: { _id: 0, movie: 1, views: 1 } }
])
```

D. Count Movies per Genre with \$group and \$project

```
db.watchHistory.aggregate([
  { $group: { _id: "$genre", count: { $sum: 1 } } },
  { $project: { _id: 0, genre: "$_id", count: 1 } }
])
```

6. Interactive Challenge / Mini-Project

Your Turn!

- Write an aggregation pipeline to find the average rating for each genre in 2024, but only include genres with more than 10,000 total views.
- Output should show genre, average rating (rounded to 1 decimal), and total views.

7. Common Pitfalls & Best Practices

Pitfall	Best Practice
Placing \$match late	Put \$match early to reduce data volume
Returning too many fields	Use \$project to limit output
Not using indexes	Index fields used in \$match for speed
Complex pipelines in one go	Build and test one stage at a time

## 8. Optional: Programmer’s Workflow Checklist

---

- Define your output and work backwards to design stages.
- Place `$match` as early as possible.
- Use `$group` for aggregation, `$project` for shaping output.
- Test each stage’s output before adding the next.
- Use `$sort` and `$limit` for ranking and pagination.
- Optimize with indexes and by limiting unnecessary fields.