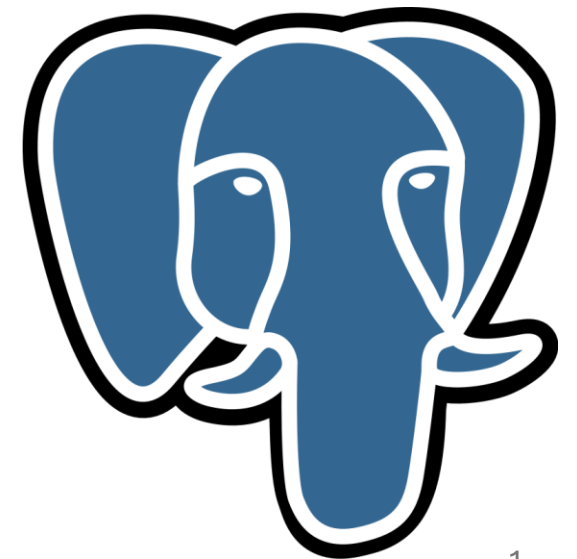# Data Management

Sahar Khanlari - 2107563

Marco Natale – 1929854

A.A. 2024/2025
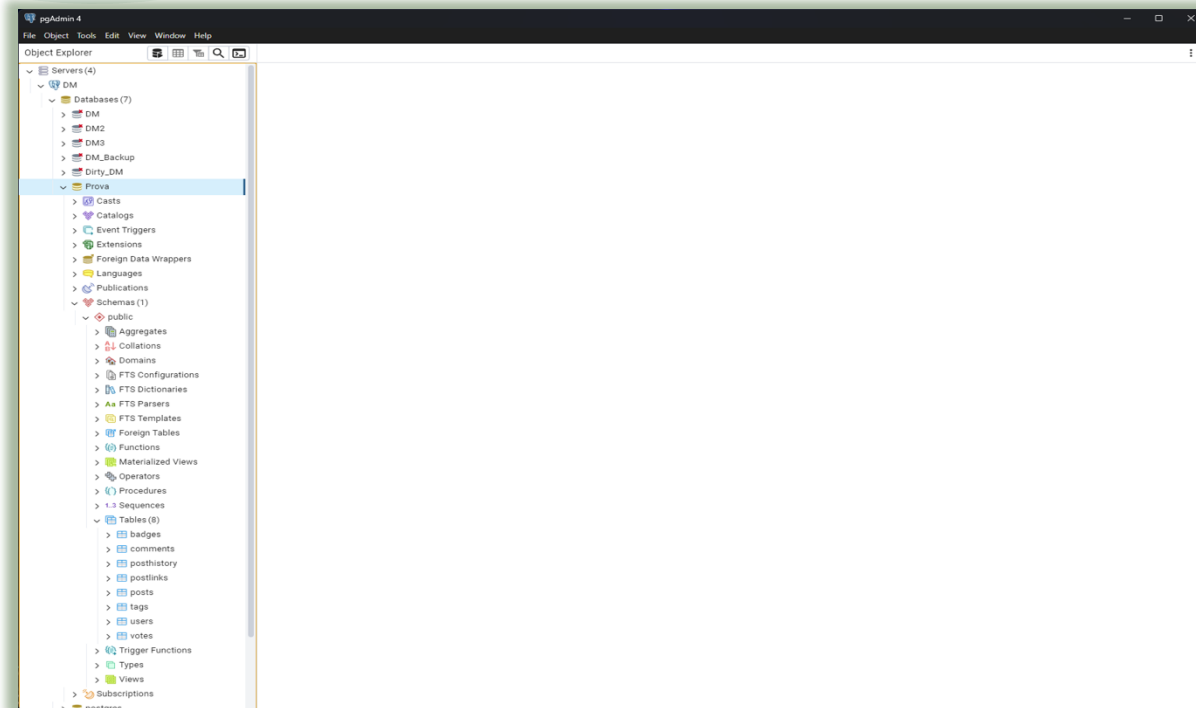
# Introduction

- SQL vs NoSQL (PostgreSQL vs MongoDB)

- Goals:
  - Interface Usability
  - Flexibility
  - Scalability
  - Performance Queries and Comparison
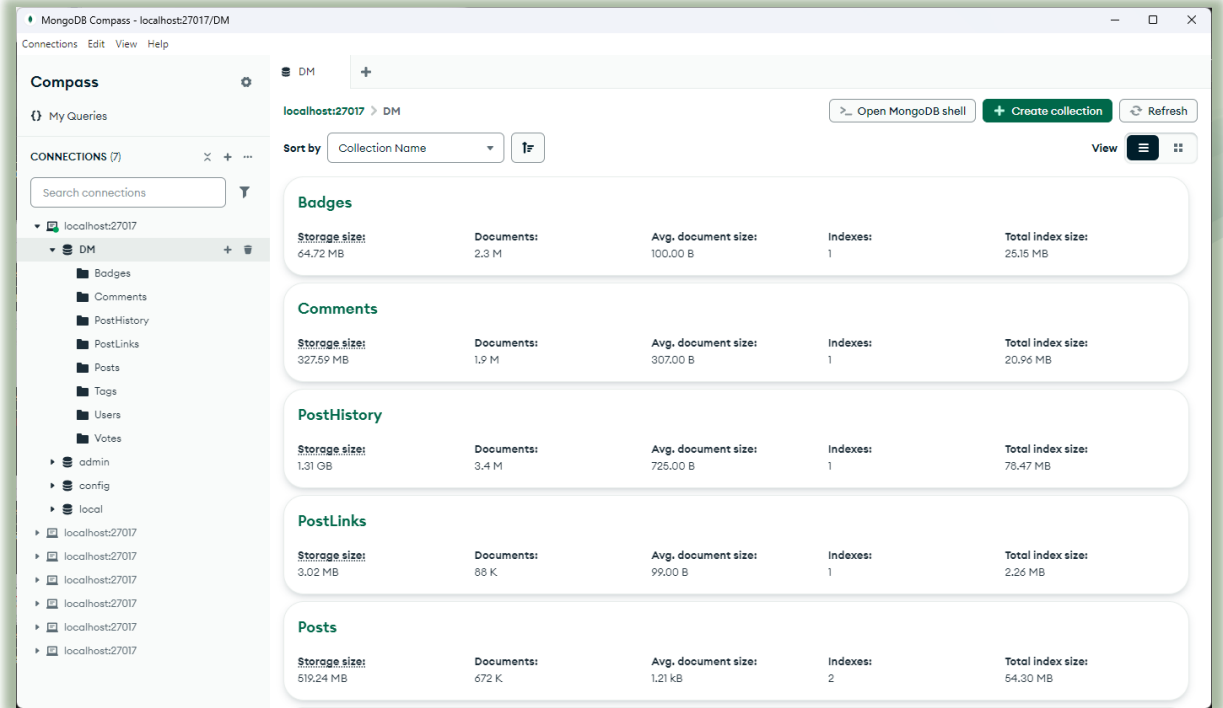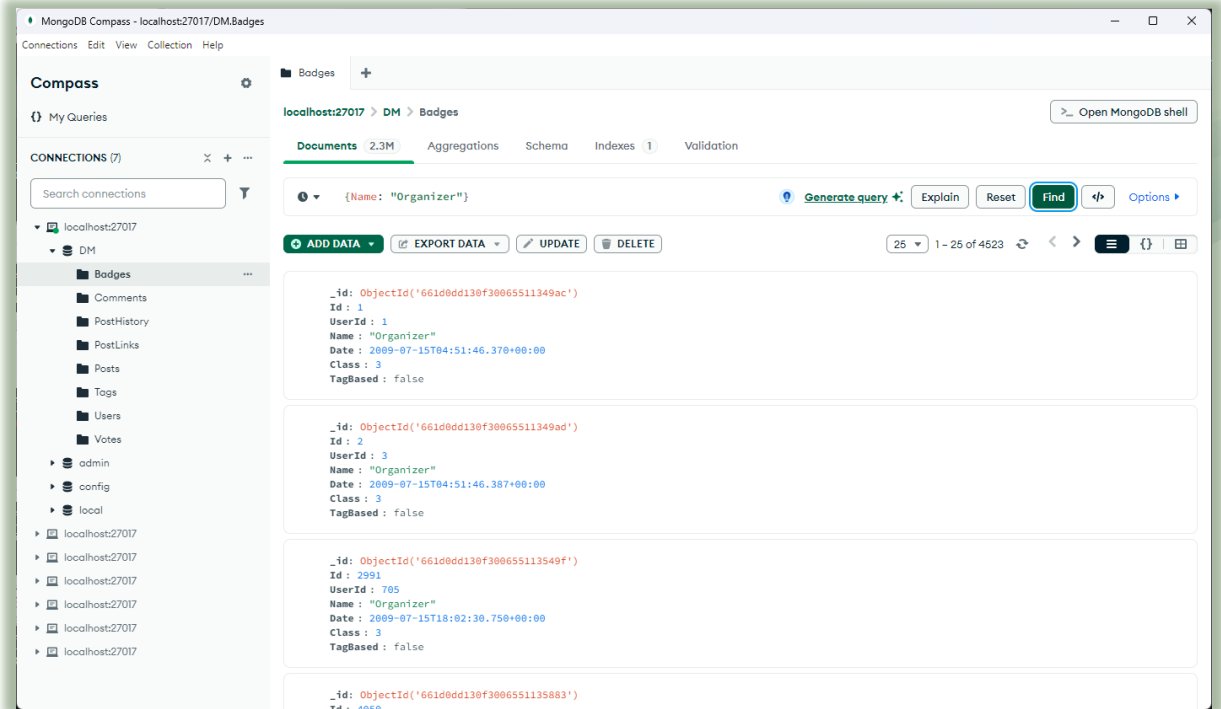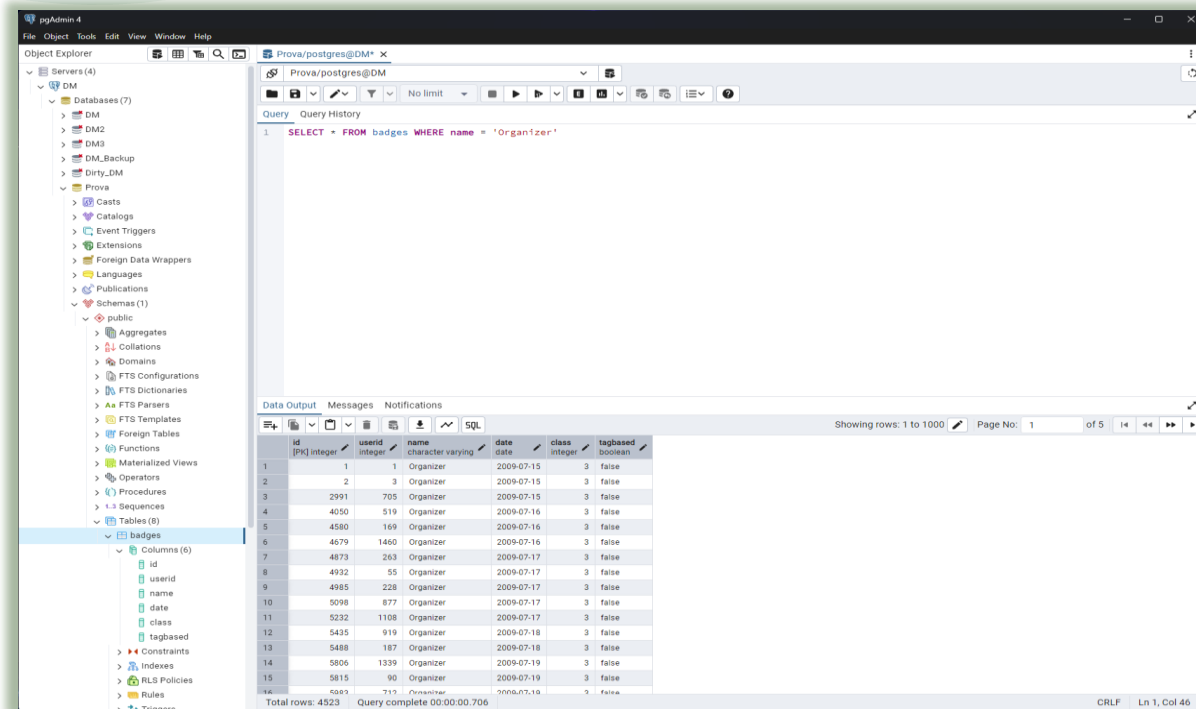
# Interface Usability
## Layout

pgAdmin

MongoDB Compass

# Interface Usability
## Navigation & Usability

pgAdmin                                    MongoDB Compass

# Interface Usability
## Navigation & Usability

pgAdmin

MongoDB Compass

# Interface Usability
## Functionality

➤ **Differences**

**Query Interface:**
- pgAdmin: Uses a simple text editor for writing SQL queries.
- MongoDB Compass: Provides a visual aggregation pipeline builder for NoSQL queries.

**Predefined Queries:**
- pgAdmin: Offers predefined queries (e.g., view data, CRUD scripts) via right-click.
- MongoDB Compass: No predefined SQL-like queries, relies on document-based query building.

➤ **Similarities**

- Both provide a GUI interface for database management.

- Both allow import/export of data.

- Both include a CLI window for command-line interaction.

# Interface Usability
## Functionality

# Flexibility

| Feature | SQL (Relational Databases) | MongoDB (NoSQL) |
|---|---|---|
| Schema | Strict, predefined schema | Schema-less (Flexible documents) |
| Structure | Tables, rows, and columns | Documents (JSON/BSON format) |
| Data Relationships | Normalized, uses joins to relate tables | Denormalized, embeds related data |
| Schema Changes | Requires migrations, downtime | Dynamic, fields can be added/removed |
| Use Case | Best for structured & transactional data | Best for unstructured, evolving data |

# Scalability

| Feature | SQL (Relational Databases) | MongoDB (NoSQL) |
|---|---|---|
| Scaling Strategy | Vertical Scaling (Scale-Up) | Horizontal Scaling (Scale-Out) |
| Data Distribution | Difficult to shard, manual partitioning | Built-in sharding for automatic distribution |
| Read/Write Performance | Efficient for transactions, joins slow large queries | Faster reads, distributed writes |
| Replication | Master-slave or cluster-based replication | Replica sets with automatic failover |
| High Availability | Requires complex setup | Native support for replication & redundancy |
| Global Distribution | Requires additional tools | Geo-distributed databases supported |

# The Dataset

# Importing the tables

➢ **PostgreSQL**:

- Badges: 12,78 sec
- Comments: 16,87 sec
- PostHistory: 62 sec
- PostLinks: 0,72 sec
- Posts: 27,38 sec
- Tags: 0,11 sec
- Users: 14,72 sec
- Votes: 23,34 sec

➢ **MongoDB** doesn't measure time

# Performance

Integrity Checks

Data Cleaning

Foreign Keys
(only for SQL)

Different Query
Types

# Key Observation

Checking if PostLinks references Posts

- **Inefficient query**:
  SELECT pl.id
  FROM postlinks pl
  WHERE pl.postid NOT IN (SELECT p.id FROM posts p);
  **Execution Time: 2 hours 37 minutes**

- **Efficient query**:
  SELECT pl.id
  FROM postlinks pl
  EXCEPT
  SELECT pl.id
  FROM postlinks pl JOIN posts p ON p.id = pl.postid
  ORDER BY id;
  **Execution Time: 258 milliseconds**

# Key Observation #2

- MongoDB requires some indexing for very large data
- The queries couldn't complete/took very long without them

- The following indexes have been created:
  - Posts(Id)
  - Users(Id)

# Integrity Checks

| Among 8 queries | PostgreSQL | MongoDB |
|:---:|:---:|:---:|
| **BEST**<br>PostLinks → Posts | 258 ms | 1.825 sec |
| **AVERAGE** | 4.212 sec | 33.466 sec (~7x) |
| **WORST**<br>Votes → Posts | 7.914 sec | 76.631 sec |

# Data Cleaning

| Among 8 queries | PostgreSQL | MongoDB |
|---|---|---|
| **BEST** PostLinks → Posts | 379 ms | 1.460 sec |
| **AVERAGE** | 9.644 sec | 34.108 sec (~2.5x) |
| **WORST** | 29.958 sec PostHistory → Posts | 87.51 sec Votes → Posts |

- Votes table is **too large** for MongoDB's $lookup.
- PostgreSQL's join is optimized on large tables.

# Foreign Keys

➢ **Posts** references:
- Users(Id) → OwnerUserId
- Posts(Id) → ParentId & AcceptedAnswerId

➢ **Comments** references:
- Posts(Id) → PostId

➢ **Votes** references:
- Users(Id) → UserId
- Posts(Id) → PostId

➢ **Badges** references:
- Users(Id) → UserId

➢ **PostHistory** references:
- Users(Id) → UserId
- Posts(Id) → PostId

➢ **PostLinks** references:
- Posts(Id) → PostId & RelatedPostId

# Performance Queries Overview

| | Read | Write | Joins | Aggregation | Text Search | Nested Updates | Complex Joins | Pagination | Deletion | Count |
|---|---|---|---|---|---|---|---|---|---|---|
| **SQL** | 0.779 | 0.453 | 14.673 | 2.236 | 1.286 | 3.259 | 8.998 | 0.583 | >8h | 0.884 |
| **MongoDB** | 6.867 | 0.003 | 88.112 | 2.518 | 2.3296 | 1.7278 | >6h | 2.470 | 10.869 | 1.144 |

# Simple Read and Write

```sql
SELECT *
FROM Posts
WHERE CreationDate BETWEEN '2009-01-01' AND '2009-
12-31';
```

| | Read | Write |
|---|---|---|
| **SQL** | 0.779 | 0.453 |

```sql
INSERT INTO Posts (Id, PostTypeId, CreationDate, Score, Title, Body, Tags)
VALUES
(2, 1, NOW(), 10, 'Title 1', 'Body 1', '<mac><crash>'),
(3, 2, NOW(), 15, 'Title 2', 'Body 2', '<windows>');
```

- Declarative syntax
- Standard SQL

```
posts.find({"CreationDate":
              {"$gte": datetime(2009, 1, 1),        "$lte":
datetime(2009, 12, 31)}
              })
```

| | Read | Write |
|---|---|---|
| **MongoDB** | 6.867 | 0.003 |

```
posts.insert_many([
{"Id": 2, "PostTypeId": 1, "CreationDate": datetime.now(), "Score": 10,
"Title": "Title 1", "Body": "Body 1", "Tags": "<mac><crash>"},
{"Id": 3, "PostTypeId": 2, "CreationDate": datetime.now(), "Score": 15,
"Title": "Title 2", "Body": "Body 2", "Tags": "<windows>"}
])
```

- JSON-like documents
- Operators and expressions: $gt, $eq, ...

# Simple Joins and Aggregations

```
SELECT p.*, u.DisplayName
FROM Posts p
JOIN Users u ON p.OwnerUserId = u.Id;
```

| | Joins | Aggregation |
|---|---|---|
| SQL | 14.673 | 2.236 |

```
SELECT Tags, COUNT(*)
FROM Posts
GROUP BY Tags;
```

- Optimized for joins
- Relationships are optimized

- Aggregation with group by

```
posts.aggregate([
{"$lookup": {"from": "Users", "localField": "OwnerUserId",
"foreignField": "Id", "as": "owner"}},
{"$unwind": "$owner"},
{"$project": {"_id": 0, "Title": 1, "OwnerDisplayName":
"$owner.DisplayName"}}
])
```

| | Joins | Aggregation |
|---|---|---|
| MongoDB | 88.112 | 2.518 |

```
posts.aggregate([
{"$unwind": "$Tags"},
{"$group": {"_id": "$Tags", "count": {"$sum": 1}}}
])
```

- No native join operation
- $lookup stage is used in aggregations

- Aggregation in stages

# Text Search and Nested Updates

```sql
SELECT *
FROM Posts
WHERE Body LIKE '%virtual machine%';
```

| | Text Search | Nested Updates |
|---|---|---|
| **SQL** | 1.286 | 3.259 |

```sql
UPDATE Users
SET Reputation = Reputation + 10
WHERE Id IN
(SELECT OwnerUserId
FROM Posts
GROUP BY OwnerUserId
HAVING COUNT(*) > 10);
```

```
posts.find({"Body":
{"$regex": "virtual machine", "$options": "i"}
})
```

| | Text Search | Nested Updates |
|---|---|---|
| **MongoDB** | 2.3296 | 1.7278 |

```
owners_with_more_than_10_posts =
posts.aggregate([
{"$group": {"_id": "$OwnerUserId", "postCount": {"$sum": 1}}},
{"$match": {"postCount": {"$gt": 10}}}
])
owner_ids = [owner["_id"] for owner in
owners_with_more_than_10_posts]
users.update_many({"Id": {"$in": owner_ids}}, {"$inc": {"Reputation":
10}})
```

# Complex Joins

```sql
SELECT p.Title, c.Text, u.DisplayName
FROM Posts p
JOIN Comments c ON p.Id = c.PostId
JOIN Users u ON c.UserId = u.Id;
```

| | Complex Joins |
|---|---|
| SQL | 8.998 |
| MongoDB | >6h |

```
posts.aggregate([
{"$lookup": {"from": "Comments", "localField": "Id", "foreignField": "PostId", "as": "comments"}},
{"$unwind": "$comments"},
{"$lookup": {"from": "Users", "localField": "comments.UserId", "foreignField": "Id", "as": "commentUser"}},
{"$unwind": "$commentUser"},
{"$project": {"_id": 0, "Title": 1, "CommentText": "$comments.Text", "CommenterName": "$commentUser.DisplayName"}}
])
```

# Pagination, Deletion and Count

| | Pagination | Deletion | Count |
|---|---|---|---|
| **SQL** | 0.583 | >8h | 0.884 |
| **MongoDB** | 2.470 | 10.869 | 1.144 |

```sql
SELECT *
FROM Posts
ORDER BY CreationDate DESC LIMIT 10
OFFSET 20;
```

```sql
DELETE FROM Posts
WHERE Score < 5 AND Id NOT IN (SELECT
PostId FROM Comments);
```

```sql
SELECT COUNT(*)
FROM Users
WHERE Reputation > 100;
```

```python
posts.find().sort("CreationDate", -
1).skip(20).limit(10)
```

```python
post_ids_with_comments = comments.distinct("PostId")
posts.delete_many({"Score": {"$lt": 5}, "Id": {"$nin":
post_ids_with_comments}})
```

```python
users.count_documents({"Reputation": {"$gt": 100}})
```

# Conclusions

| | Read | Write | Joins | Aggregation | Text Search | Nested Updates | Complex Joins | Pagination | Deletion | Count |
|---|---|---|---|---|---|---|---|---|---|---|
| **SQL** | 0.779 | 0.453 | 14.673 | 2.236 | 1.286 | 3.259 | 8.998 | 0.583 | >8h | 0.884 |
| **MongoDB** | 6.867 | 0.003 | 88.112 | 2.518 | 2.3296 | 1.7278 | >6h | 2.470 | 10.869 | 1.144 |

- MongoDB better suited for changing data (faster write, update and delete operations)
- PostgreSQL faster in structured data

Use case of this dataset: large forum
→ Dynamic data, but frequent join requests (loading posts and comments)
→ Optimizations are put in place (e.g., on delete set null)

# Thank you

Sahar Khanlari - 2107563

Marco Natale – 1929854

A.A. 2024/2025