

Lecture Outline

Document databases

- Introduction

MongoDB

- Data model
- CRUD operations
 - **Insert**
 - **Update**
 - **Remove**
 - **Find**: projection, selection, modifiers

Document Stores

Data model

- **Documents**
 - Self-describing
 - **Hierarchical tree structures** (JSON, XML, ...)
 - Scalar values, maps, lists, sets, nested documents, ...
 - Identified by a **unique identifier** (key, ...)
- Documents are **organized into collections**

Query patterns

- Create, update or remove a document
- **Retrieve documents according to complex query conditions**

Observation

- Extended key-value stores where the value part is examinable

MongoDB Document Database



MongoDB

JSON document database

- <https://www.mongodb.com/>
- Features
 - Open source, high availability, eventual consistency, automatic sharding, master-slave replication, automatic failover, secondary indices, ...
- Developed by **MongoDB**
- Implemented in C++, C, and JavaScript
- Operating systems: Windows, Linux, MacOS, ...
- Initial release in 2009

RedisJSON vs. MongoDB

1. Architecture & Storage Model

RedisJSON	MongoDB
<ul style="list-style-type: none">• In-memory storage with optional persistence• Part of Redis ecosystem as a module• Native JSON document storage• Single-node by default, clustering through Redis Cluster• Synchronous replication	<ul style="list-style-type: none">• Disk-based storage with memory mapping• Standalone database system• BSON document storage• Built-in support for sharding and replication• Asynchronous replication by default

RedisJSON vs. MongoDB

2. Performance Characteristics

RedisJSON	MongoDB
<ul style="list-style-type: none">• Extremely fast read/write operations (microseconds)• Consistent performance due to in-memory nature• Limited by available RAM• Excellent for caching and real-time operations	<ul style="list-style-type: none">• Fast read/write operations (milliseconds)• Performance can vary based on disk I/O• Can handle datasets larger than available RAM• Good for both real-time and analytical workloads

RedisJSON vs. MongoDB

3. Query Capabilities

RedisJSON	MongoDB
<ul style="list-style-type: none">• JSONPath support for document traversal• Limited query capabilities compared to MongoDB• Basic CRUD operations on JSON documents• No built-in aggregation framework• Atomic operations on JSON values	<ul style="list-style-type: none">• Rich query language• Comprehensive aggregation framework• Full-text search capabilities• Secondary indexes• Complex geospatial queries• Complex array operations and nested queries

RedisJSON vs. MongoDB

4. Use Cases

RedisJSON	MongoDB
<ul style="list-style-type: none">• Real-time applications requiring sub-millisecond response• Session management with JSON data• Caching complex JSON structures• Rate limiting with JSON metadata• Small to medium-sized JSON document storage• Real-time analytics on JSON data	<ul style="list-style-type: none">• Large-scale document storage• Complex data modeling requirements• Applications requiring advanced querying• Data warehousing and analytics• Content management systems• IoT and time-series data

RedisJSON vs. MongoDB

5. Scalability

RedisJSON	MongoDB
<ul style="list-style-type: none">• Vertical scaling (adding more RAM)• Horizontal scaling through Redis Cluster• Limited by total memory across cluster• Best for datasets that can fit in memory• Simple replication model	<ul style="list-style-type: none">• Both vertical and horizontal scaling• Native support for sharding• Automatic load balancing• Can handle petabytes of data• Complex replication with replica sets

RedisJSON vs. MongoDB

When to Choose Which?

RedisJSON	MongoDB
<ul style="list-style-type: none">• You need extremely low latency• Your dataset fits in memory• You're already using Redis• You need atomic operations on JSON documents• Real-time processing is crucial• You want simple JSON document caching	<ul style="list-style-type: none">• You have large datasets• You need complex queries and aggregations• Your data doesn't fit in memory• You need advanced indexing• You want built-in horizontal scaling• You need flexible schema design

Query Example

Collection of movies

```
{  
  _id: ObjectId("1"),  
  title: "Vratné lahve",  
  year: 2006  
}
```

```
{  
  _id: ObjectId("2"),  
  title: "Samotáři",  
  year: 2000  
}
```

```
{  
  _id: ObjectId("3"),  
  title: "Medvídek",  
  year: 2007  
}
```

Query statement

Titles of movies filmed in *2005* and later, sorted by these titles in descending order

```
db.movies.find(  
  { year: { $gt: 2005 } },  
  { _id: false, title: true }  
).sort({ title: -1 })
```

Query result

```
{ title: "Vratné lahve" }
```

```
{ title: "Medvídek" }
```

Data Model

Database system structure

Instance → **databases** → **collections** → **documents**

- Database
- Collection

Collection of documents, usually of a similar structure

- Document

MongoDB **document** = **one JSON object**

- I.e. even a complex JSON object with other recursively nested objects, arrays or values
- Uses BSON, which supports additional data types

Each document has a **unique identifier** (**primary key**)

- Technically realized using a **top-level `_id` field**

Data Model

MongoDB **document**

- Internally stored in **BSON** format (*Binary JSON*)
 - Maximal allowed size 16 MB
 - **GridFS** can be used to split larger files into smaller chunks
 - **MongoDB Atlas Data Lake** is a cloud solution for data storage and analysis

Restrictions on fields

- **Top-level `_id`** is reserved for a **primary key**
- Field names **cannot start with `$`** and **cannot contain `.`**
 - `$` is reserved for query operators
 - `.` is used when accessing nested fields
- The order of fields is preserved
 - Except for `_id` fields that are always moved to the beginning
- **Names of fields must be unique**

Primary Keys

Features of identifiers

- **Unique** within a collection
- **Immutable** (cannot be changed once assigned)
- Can be of **any type** other than a JSON array

Key management

- Natural identifier
- Auto-incrementing number – not recommended
- **UUID** (*Universally Unique Identifier*)
- **ObjectId** – **special 12-byte BSON type** (the default option) Small, likely unique, fast to generate, ordered, based on a timestamp, machine id, process id, and a process-local counter

_id autogeneration

ObjectId Structure (12 bytes):

4 bytes	3 bytes	2 bytes	3 bytes
Timestamp	Machine ID (a hash of the hostname)	Process ID	Counter (Starts from a random value, increments for each new document)

// NOT RECOMMENDED to use incremental _ids:

```
db.collection.insertOne({ _id: 1 })
```

```
db.collection.insertOne({ _id: 2 })
```

//RECOMMENDED to use ObjectId or UUID:

```
db.collection.insertOne({ _id: UUID(), name: "test" })
```

Design Questions

Data modeling (in terms of **collections and documents**)

- Supports **JSON schema validation**
 - However...
 - documents within a collection are similar in practice
 - **implicit schema** is required nevertheless
- Challenge
 - Balancing application requirements, performance aspects, data structure, mutual relationships, query patterns, ...

Two main concepts

- References
- Embedded documents

Denormalized Data Models

Embedded documents

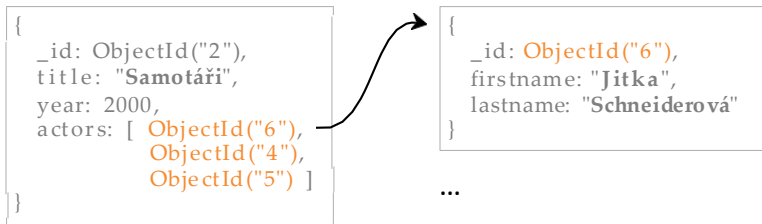
- Related data in a single document
 - with embedded JSON objects, so called **subdocuments**
- Pros: data manipulation (fewer queries need to be issued)
- Cons: possible data redundancies
- Suitable for **one-to-one** or **one-to-many** relationships

```
{
  _id: ObjectId("2"), title: "Samotáři", year: 2000,
  actors: [
    { firstname: "Jitka", lastname: "Schneiderová" },
    { firstname: "Ivan", lastname: "Trojan" },
    { firstname: "Jiří", lastname: "Macháček" }
  ]
}
```

Normalized Data Models

References

- Related data in separate documents
 - These are interconnected via directed links (**references**)
 - Technically expressed using **ordinary values with identifiers of target documents** (i.e. no special construct is provided)
- Features: higher flexibility, follow up queries might be needed
- Suitable for **many-to-many** relationships



Sample Data

Collection of **movies**

```
{
  _id: ObjectId("1"),
  title: "Vratné lahve", year: 2006,
  actors: [ ObjectId("7"), ObjectId("5") ]
}
```

```
{
  _id: ObjectId("2"),
  title: "Samotáři", year: 2000,
  actors: [ ObjectId("6"), ObjectId("4"),
            ObjectId("5") ]
}
```

```
{
  _id: ObjectId("3"),
  title: "Medvídek", year: 2007,
  actors: [ ObjectId("5"), ObjectId("4") ]
}
```

Collection of **actors**

```
{ _id: ObjectId("4"),
  firstname: "Ivan",
  lastname: "Trojan" }
```

```
{ _id: ObjectId("5"),
  firstname: "Jiří",
  lastname: "Macháček" }
```

```
{ _id: ObjectId("6"),
  firstname: "Jitka",
  lastname: "Schneiderová" }
```

```
{ _id: ObjectId("7"),
  firstname: "Zdeněk",
  lastname: "Svěrák" }
```

Application Interfaces

mongo shell

- **Interactive interface to MongoDB**

```
mongosh --port 42222 -u login -p password database
```

GUI Tools

- **MongoDB Compass** - Modern GUI for data visualization and manipulation
- **MongoDB Atlas** - Cloud-based management interface

Drivers

- Java, C, C++, C#, Perl, PHP, Python, Node.js, Go, Rust, Ruby, Scala, ...

Query Language

MongoDB query language is based on **JavaScript**

- **Single command / entire script**
- Read queries return a **cursor**
 - Allows us to iterate over all the selected documents
- While commands are typically applied to a single collection, some operations, like aggregations, can interact with multiple collections.

Query patterns

- Basic **CRUD** operations
 - Accessing documents via identifiers or **conditions on fields**
- Aggregations: **Aggregation Pipeline**
- Atlas Search (full-text search capabilities)
- Time-series operations
- Change streams

CRUD Operations

Overview

- **Create**
 - `insertOne()`, `insertMany()`
 - Insert a new document (or documents) into a collection
- **Update**
 - `replaceOne()`
 - `updateOne()`, `updateMany()`
 - Modify an existing document / documents
- **Delete**
 - `deleteOne()`, `deleteMany()`
 - Delete an existing document / documents
- **Read**
 - `find()`, `findOne()`

Working with Collections in MongoDB

- **db.collection:**

```
db.users.updateOne({ name: "John" }, { $set: { age: 30 } })
```

- **getCollection("collection")**

```
db.getCollection("users").updateOne({ name: "John" },  
{ $set: { age: 30 } });
```

- Reference to collection

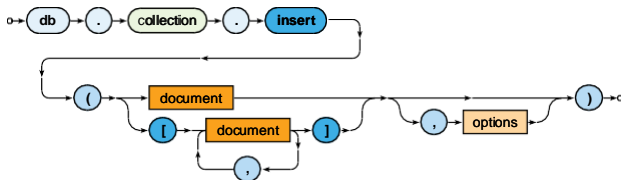
```
const users = db.users;  
users.updateOne({ name: "John" }, { $set: { age: 30 } });
```

- Collections are created automatically when not yet exist

Insert Operation

Insert Operation: insertOne, insertMany

Inserts a new document / documents into a given collection



- Parameters

- Document:** one or more documents to be inserted

- Provided document identifiers (`_id` fields) must be unique
 - When missing, they are generated automatically (**ObjectId**)

- Options**

- ordered:** boolean (default: true)

- **true:** documents are inserted in order, stops on first error
 - **false:** allows parallel processing, continues on error
 - **writeConcern:** level of acknowledgment required

Insert Operation: Examples

Insert a new actor document

```
db.actors.insertOne(  
  {  
    firstname: "Anna",  
    lastname: "Geislerová"  
  }  
)
```

```
{  
  _id: ObjectId("8"),  
  firstname: "Anna",  
  lastname: "Geislerová"  
}
```

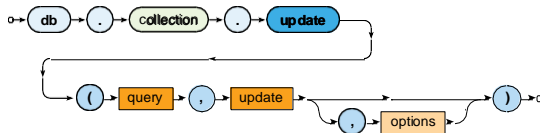
Insert two new movies

```
db.movies.insertMany(  
  [  
    {  
      _id: ObjectId("9"), title: "Želary", year: 2003,  
      actors: [ ObjectId("4"), ObjectId("8") ]  
    },  
    { title: "Anthropoid", year: 2016, actors: [ ObjectId("8") ] },  
  ]  
)
```

Update Operation

Update Operation: replaceOne, updateOne

Modifies / replaces an existing document / documents



- Parameters
 - Query:** description of documents to be updated
 - The same behavior as in find operations
 - Update:** modification actions to be applied
 - Options**
- Use **replaceOne** or **updateOne** to update **one document**
 - Use **updateMany** to update two or more documents

Update Operation: Examples

Replace the whole document of at most one specified actor

```
db.actors.replaceOne(  
  { _id: ObjectId("8") },  
  {  
    firstname: "Aňa",  
    lastname: "Geislerová" }  
)
```

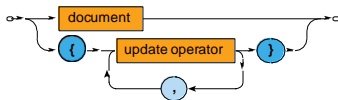
```
{  
  _id: ObjectId("8"),  
  firstname: "Aňa",  
  lastname: "Geislerová"  
}
```

Update all movies filmed in 2015 or later

```
db.movies.updateMany(  
  { year: { $gt: 2015 } },  
  {  
    $set: { new: true },  
    $inc: { rating: 3 }  
  }  
)
```

Update Operation

Update / replace modes

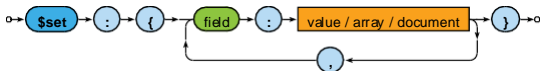


- **Replace (replaceOne method)**
 - **The whole document is replaced** (`_id` is preserved)
- **Update (updateOne, updateMany methods)**
the `update` parameter contains only **update operators**
 - **The current document is updated** using these operators
 - `$set`, `$unset`, `$inc`, `$mul`, ...
 - Each operator can be used at most once

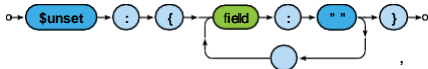
Update Operators

Field operators

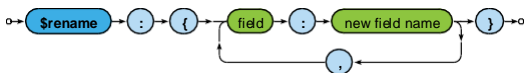
- **\$set** – sets the value of a given field/fields or create a field



- **\$unset** – removes a given field/fields



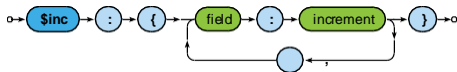
- **\$rename** – renames a given field/fields



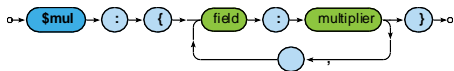
Update Operators

Field operators

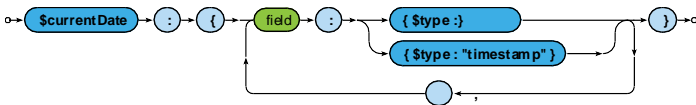
- **\$inc** – increments the value of a given numeric field/fields



- **\$mul** – multiplies the value of a given numeric field/fields



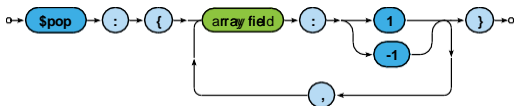
- **\$currentDate** – stores the current date time / timestamp to a given field/fields



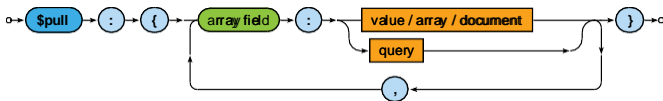
Update Operators

Array operators

- **\$pop** – removes the first / last item of an array



- **\$pull** – removes all array items that match a specified query



Upsert Mode

Upsert behavior of update operation

- When `{ upsert : true }` option is specified, and, at the same time, **no document was updated**
⇒ **new document is inserted**

What this document will contain?

- In case of the **replace** mode...
 - All the fields (i.e. value fields) from the update parameter
- In case of the **update** mode...
 - All the value fields from the query parameter,
 - and the outcome of all the update operators from the update parameter
- `_id` field is preserved, or newly generated if necessary

Upsert Mode: Example

Unsuccessful update of a movie resulting to an insertion

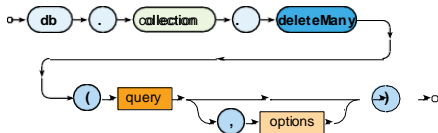
```
db.movies.updateOne(  
  { title: "Tmavomodrý svět", year: { $gt: 2000 } },  
  {  
    $set: {  
      director: { firstname: "Jan", lastname: "Svěrák" },  
      year: 2001  
    },  
    $inc: { rating: 2 }  
  },  
  { upsert: true }  
)
```

```
{ _id: ObjectId("11"),  
  title: "Tmavomodrý svět",  
  director: { firstname: "Jan", lastname: "Svěrák" },  
  year: 2001,  
  rating: 2 }
```

Remove Operation

Delete Operation: deleteOne, deleteMany

Removes a document / documents from a given collection



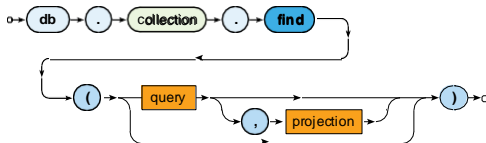
- Parameters
 - Query:** description of documents to be removed
 - The same behavior as in find operations
 - Options:** allows users to specify language-specific rules for string comparison, such as rules for lettercase and accent marks

db.collection.deleteOne() deletes the first document that matches the query

Find Operation

Find Operation: find(), findOne()

Selects documents from a given collection



- Parameters
 - **Query:** description of documents to be selected
 - **Projection:** fields to be included / excluded in the result
- Matching documents are returned via an iterable **cursor**
 - This allows us to chain further `sort`, `skip` or `limit` operations
- **findOne()**
 - Returns a single document directly (not a cursor)
 - Returns null if no document matches
 - More efficient than `find().limit(1)`

Find Operation: Examples

Select all movies from our collection

```
db.movies.find()
```

```
db.movies.find( { } )
```

Select a particular movie based on its document identifier

```
db.movies.findOne( { _id: ObjectId("2") } )
```

Select movies filmed in *2000* with a rating greater than *1*

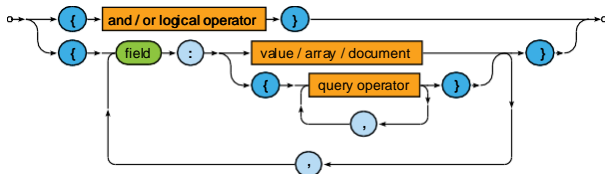
```
db.movies.find( { year: 2000, rating: { $gt: 1 } } )
```

Select movies filmed between *2005* and *2015*

```
db.movies.find( { year: { $gte: 2005, $lte: 2015 } } )
```

Selection

Query parameter describes the documents we are interested in



Boolean expression with one top-level logical operator: $\$and$, $\$or$

Conditions on individual distinct fields

- **Value equality**
 - The actual field value must be identical to the specified value
- **Query operators**
 - The actual field value must satisfy all the provided operators

Selection: Field Conditions

Value equality

- The actual field value must be identical to the specified value
- I.e. identical...
 - including the number, order and names of recursively identical values of all nested **object fields**
 - including the number and order of recursively identical **array items**

Query operators

- The actual field value must satisfy all the provided operators
 - Each operator can be used at most once at one level

Value Equality: Examples

Select movies having a specific director

```
db.movies.find(  
  { director: { firstname: "Jan", lastname: "Svěrák" } }  
)
```

```
db.movies.find(  
  { director: { lastname: "Svěrák", firstname: "Jan" } }  
)
```

Select movies having specific actors

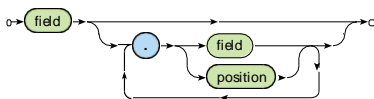
```
db.movies.find( { actors: [ ObjectId("7"), ObjectId("5") ] } )
```

```
db.movies.find( { actors: [ ObjectId("5"), ObjectId("7") ] } )
```

Queries in both the pairs are not equivalent!

Dot Notation

The **dot notation** for field names



- Accessing **fields of embedded documents**
 - **"field.subfield"** (up to 100 b. and up to 100 levels)
 - E.g.: "director.firstname",
"user.address.city.zipcode"
- Accessing **items of arrays**
 - "field.index"**
 - E.g.: "actors.2"
 - Positions start at 0

Value Equality

Example (revisited)

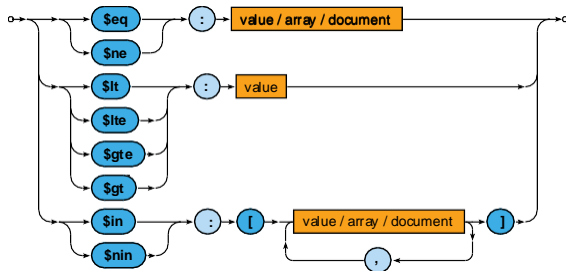
Select movies having a specific director

```
db.movies.find(  
  { director: { firstname: "Jan", lastname: "Svěrák" } }  
)
```

```
db.movies.find(  
  { "director.firstname": "Jan", "director.lastname": "Svěrák" }  
)
```

Query Operators

Comparison operators



- Comparisons take particular **BSON** data types into account
 - Certain numeric conversions are automatically applied

Query Operators

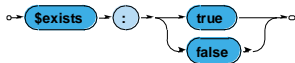
Comparison operators

- **\$eq, \$ne**
 - Tests the actual field value for **equality / inequality**
 - The same behavior as in case of value equality conditions
- **\$lt, \$lte, \$gte, \$gt**
 - Tests whether the actual field value is **less than / less than or equal / greater than or equal / greater than** the provided value
- **\$in**
 - Tests whether the actual field value is equal to **at least one** of the provided values
- **\$nin**
 - Negation of \$in

Query Operators

Element operators

- **\$exists** – tests whether a given field **exists** / **not exists**



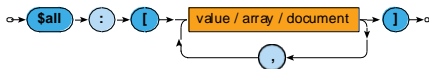
Evaluation operators

- **\$regex** – tests whether a given field value matches a specified **regular expression** (PCRE)
- **\$text** – performs **text search** (text index must exist)
 - Supports for wildcard indexes (`$**`)

Query Operators

Array operators

- **\$all** – tests whether a given array **contains all the specified items** (in any order)



Example (revisited)

Select movies having specific actors

```
db.movies.find(  
  { actors: [ ObjectId("5"), ObjectId("7") ] }  
)
```

```
db.movies.find(  
  { actors: { $all: [ ObjectId("5"), ObjectId("7") ] } }  
)
```

Query Operators

Array operators

- **\$size** – tests the size of a given array against a fixed number (and not, e.g., a range, unfortunately)



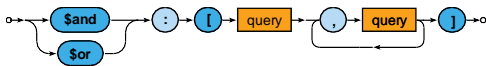
- **\$elemMatch** – tests whether a given array **contains at least one item** that satisfies all the involved query operations



Query Operators

Logical operators

- **\$and, \$or**



- Logical connectives for **conjunction / disjunction**
- At least 2 involved query expressions must be provided
- **Only allowed at the top level** of a query

- **\$not**



- Logical **negation** of exactly one involved query operator
- I.e. **cannot be used at the top level** of a query
- Cannot be used with \$regex

Querying Arrays

Condition based on **value equality** is satisfied when...

- the given field as a whole is identical to the provided value, or
- at least one item of the array is identical to the provided value

```
db.movies.find( { actors: ObjectId("5") } )
```

```
{ actors: ObjectId("5") }
```

```
{ actors: [ ObjectId("5"), ObjectId("7") ] }
```

Querying Arrays

Condition based on **query operators** is satisfied when...

- the given field as a whole satisfies all the involved operators, or
- each of the involved operators is satisfied by at least one item of the given array
 - note, however, that this item may not be the same for all the individual operators

```
db.movies.find( { ratings: { $gte: 2, $lte: 3 } } )
```

```
{ ratings: 3 }
```

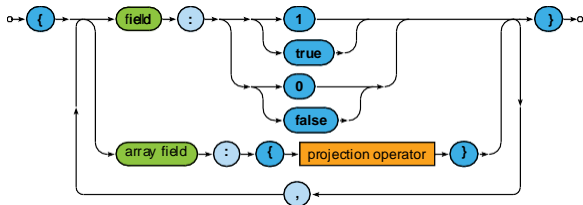
```
{ ratings: [ 3, 7, 5 ] }
```

```
{ ratings: [ 1, 4 ] }
```

Use `$elemMatch` when just a single array item should be found for all the operators

Projection

Projection allows us to determine the fields returned in the result



- **true** or 1 for fields to be **included**
- **false** or 0 for fields to be **excluded**
- Positive and negative enumerations cannot be combined!
 - The only exception is `_id` which is **included by default**
- **Projection operators** – allow to select particular array items

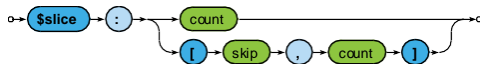
Projection Operators

Array operators

- **\$elemMatch** – selects the first matching item of an array
This item must satisfy all the operators included in `query`
When there is no such item, the field is not returned at all



- **\$slice** – selects the first `count` items of an array (when `count` is positive) / the last `count` items (when negative)
Certain number of items can also be skipped



Projection: Examples

Find a particular movie, select its identifier, title and actors

```
db.movies.find(  
  { _id: ObjectId("2") },  
  { title: true, actors: true }  
)
```

```
{  
  _id: ObjectId("2"),  
  title: "Samotáři",  
  actors: [ ObjectId("6"),  
            ObjectId("4"),  
            ObjectId("5") ]  
}
```

Find movies from 2000, select their titles and the last two actors

```
db.movies.find(  
  { year: 2000 },  
  {  
    title: 1, _id: 0,  
    actors: { $slice: -2 }  
  }  
)
```

```
{  
  title: "Samotáři",  
  actors: [ ObjectId("4"),  
            ObjectId("5") ]  
}
```

Modifiers

Modifiers change the order and number of returned documents

- **sort** – orders the documents in the result
- **limit** – returns at most a certain number of documents



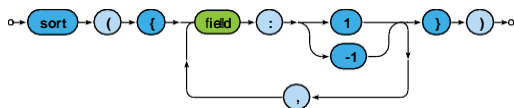
- **skip** – skips a certain number of documents from the beginning



All the modifiers are optional, can be chained in **any order** (without any implications), but **must all be specified before any documents are retrieved** via a given cursor

Modifiers

Sort modifier orders the documents in the result



- **1 (or asc)** for **ascending**, **-1 (or desc)** for **descending** order
- The order of documents is undefined unless explicitly sorted
- Sorting of larger datasets should be supported by indices
- **Sorting happens before the projection phase**
 - I.e. not included fields can be used for sorting purposes as well

Lecture Conclusion

MongoDB

- Document database for **JSON documents**
- **Sharding with master-slave replication architecture**

Query functionality

- CRUD operations
 - **Insert, find, update, remove**
 - Complex filtering conditions
- MapReduce
- Index structures