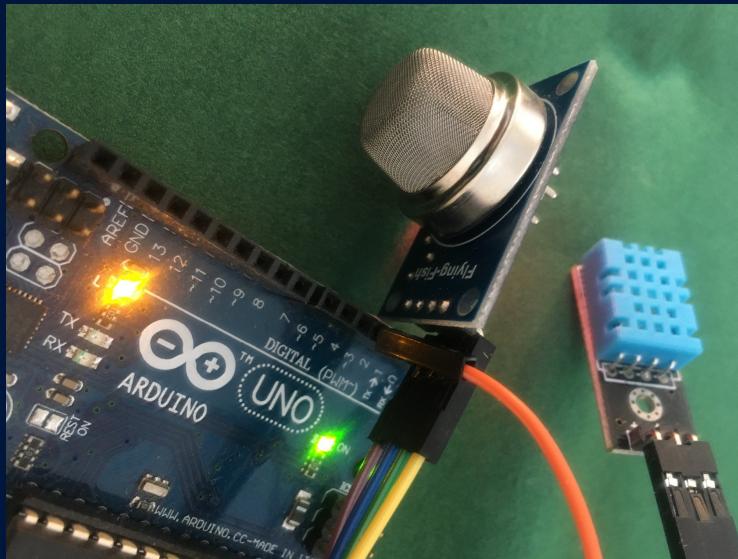


COPENHAGEN BUSINESS ACADEMY



DATA ENGINEERING with MongoDB



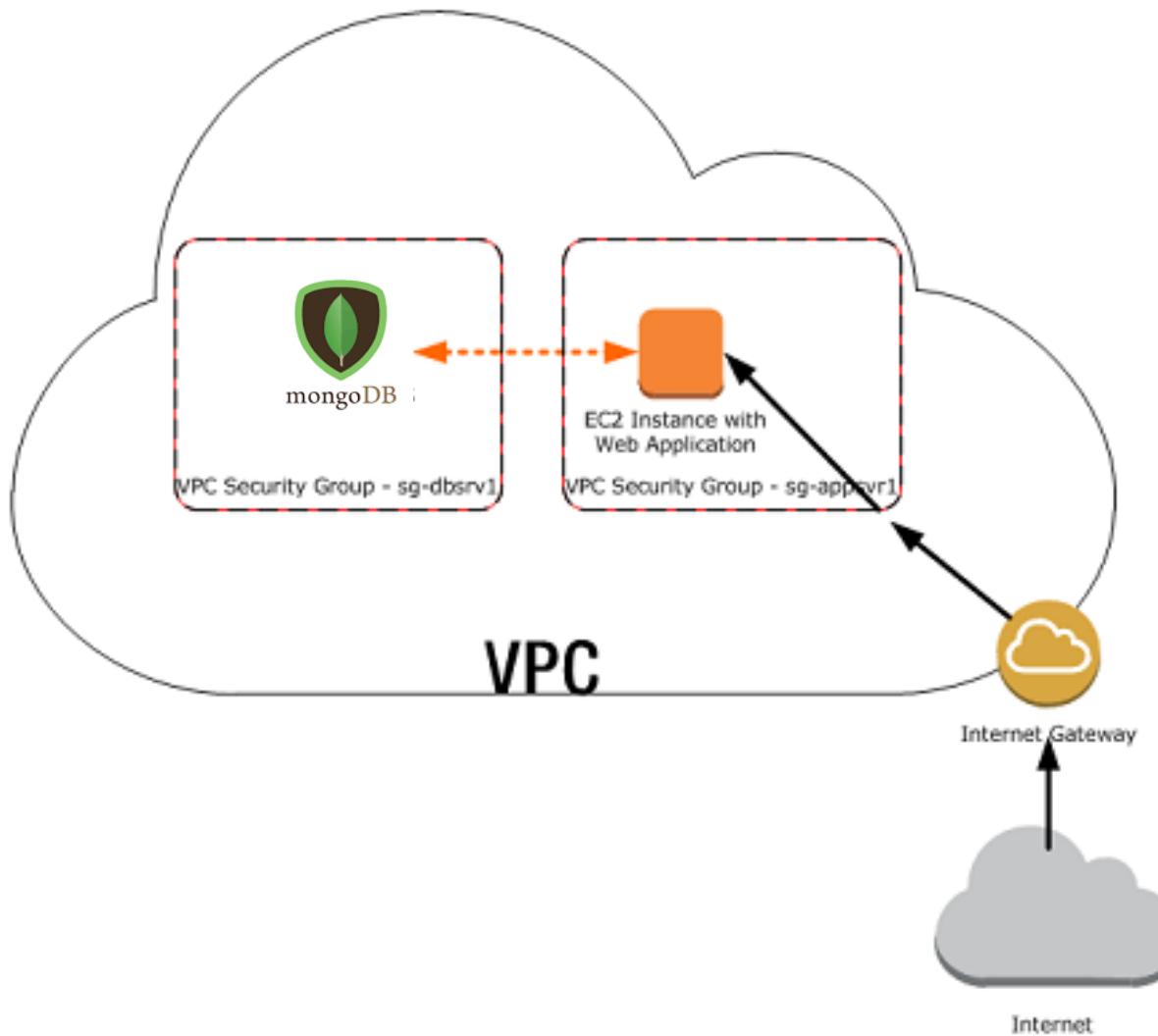
Agenda – MongoDB

- AWS – hvor er vi
 - OLA
 - Databaser
- MongoDB - overview
 - Anvendelse
 - NoSQL vs SQL
- MongoDB
 - CRUD
 - Temrinal eller GUI
 - Øvelser i mongosh
 - Forbind til min database
 - Indlæs restaurent
- Arbejde med aws-opgaven i OLA 5

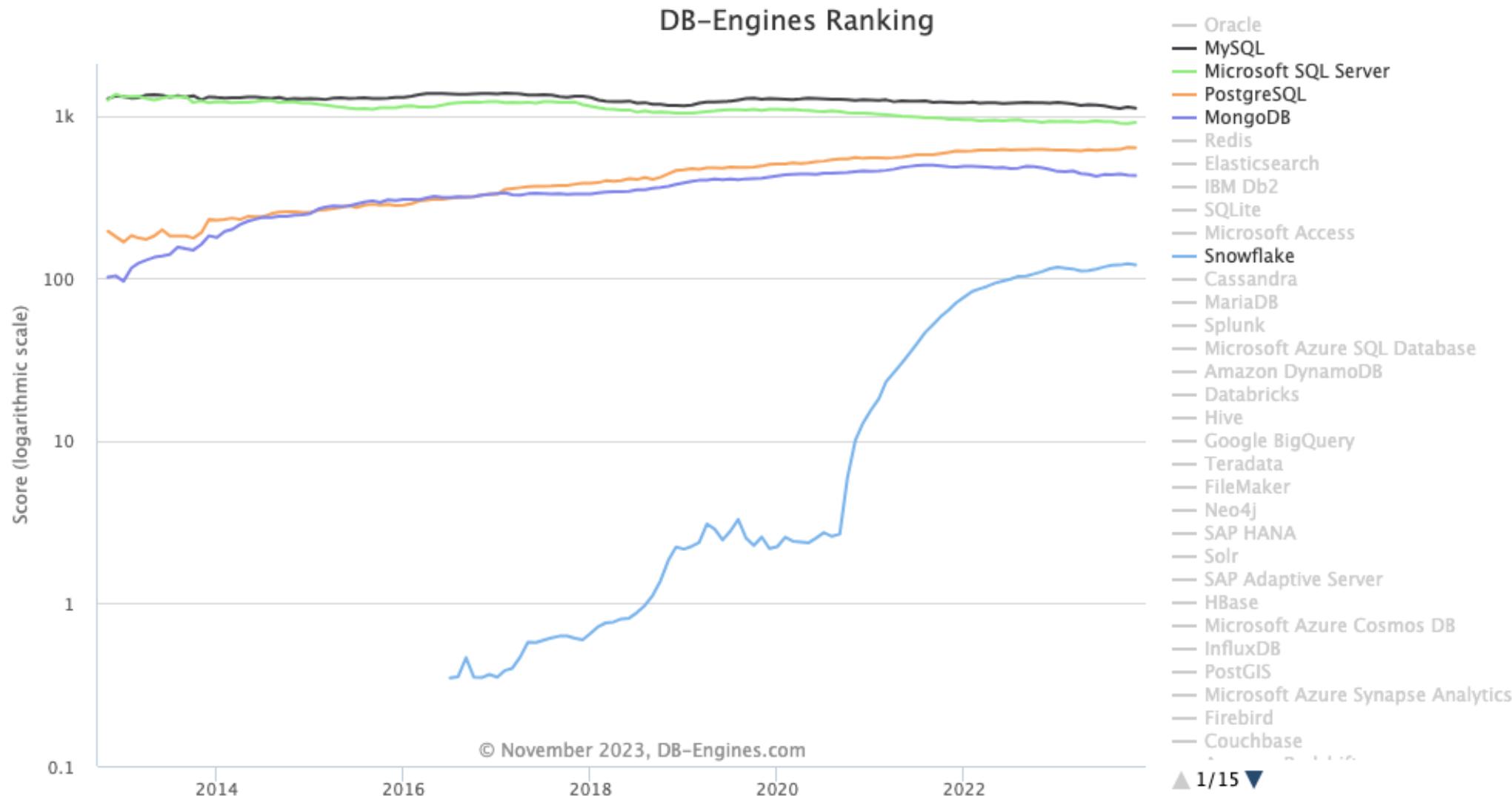
Øvelser – MongoDB

- One House skal modificeres så vardata bliver et array af pris-observationer
- Løs de første opgaver (1-9)
- Løs følgende
 - forbind til min boligsiden-database
 - hvor mange boliger?
 - hvor mange boliger over 4000000?
 - hvor mange over 200 m²?
 - hvor mange over 4000000 og i Charlottenlund?
 - hvor mange er over 200 kvm og bygget før 1950?
 - hvor mange er over 400000 og under 1000000?
 - hvor mange ligger i Odense?
 - hvor mange unikke postnummere?
 - hvor mange ligger i postnumre 8000,8200,8230,8340?
- Løse resten af opgaverne fra Word dokumentet

Bilbasen – slutmålet



- **MongoDB**
 - Create AWS EC2 instance
 - Get rdp-access
 - Client access
 - mongosh
 - Dev access
 - mongolite
- **SQL**
 - MySQL on EC2
 - Client access
 - MySQL-Workbench
 - Dev access
 - RMariaDB
- **Files**
 - Create AWS S3 bucket
 - Dev access
 - aws.s3



MongoDB



mongoDB

HIGH AVAILABILITY

When you need high availability of data with automatic, fast and instant data recovery.

IN-BUILT SHARDING

In future, if you're going to grow big as MongoDB has in-built sharding solution.

UNSTABLE SCHEMA

If you have an unstable schema and you want to reduce your schema migration cost.

NO DBA

If you don't have a Database Administrator (but you'll have to hire one if you're going to go BIG).

CLOUD COMPUTING

If most of your services are cloud-based, MongoDB is best suitable for you.

MySQL



LOW-MAINTENANCE

If you're just starting and your database is not going to scale much, MySQL will help you in easy and low-maintenance setup.

LIMITED BUDGET

If you want high performance on a limited budget.

FIXED SCHEMA

If you've fixed schema and data structure isn't going to change over the time like Wikipedia

HIGH TRANSACTION

If high transaction rate is going to be your requirement (like BBC around 30,000 inserts/minute, 4000 selects/hour)

DATA SECURITY

If data security is your top priority, MySQL is the most secure DBMS.

UNSTABLE SCHEMA

If you have an unstable schema and you want to reduce your schema migration cost.

NO DBA

If you don't have a Database Administrator (but you'll have to hire one if you're going to go BIG).

CLOUD COMPUTING

If most of your services are cloud-based, MongoDB is best suitable for you

FIXED SCHEMA

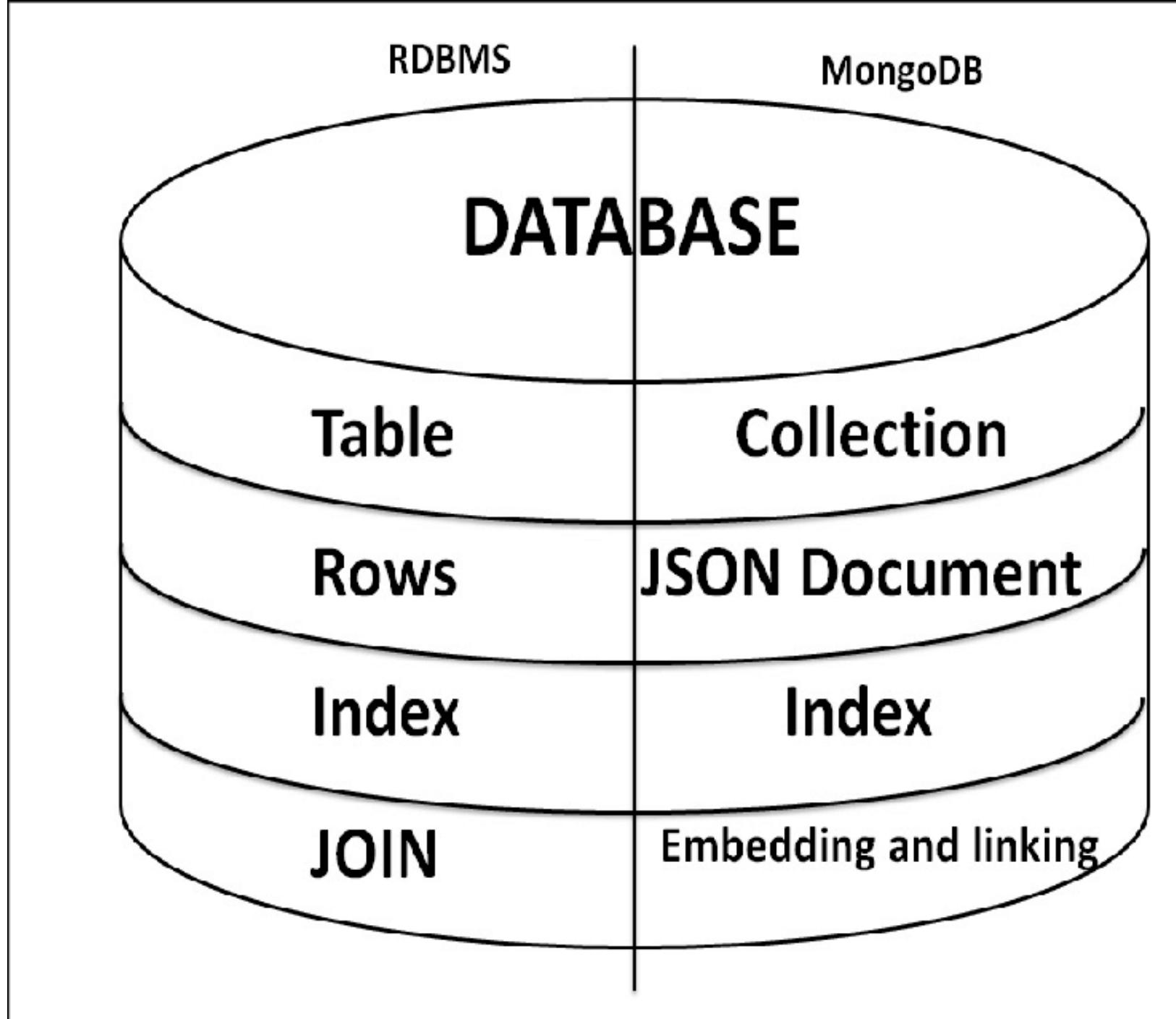
If you've fixed schema and data structure isn't going to change over the time like Wikipedia

HIGH TRANSACTION

If high transaction rate is going to be your requirement (like BBC around 30,000 inserts/minute, 4000 selects/hour)

DATA SECURITY

If data security is your top priority, MySQL is the most secure DBMS.



Data Models: Relational to Document

Relational

Person:

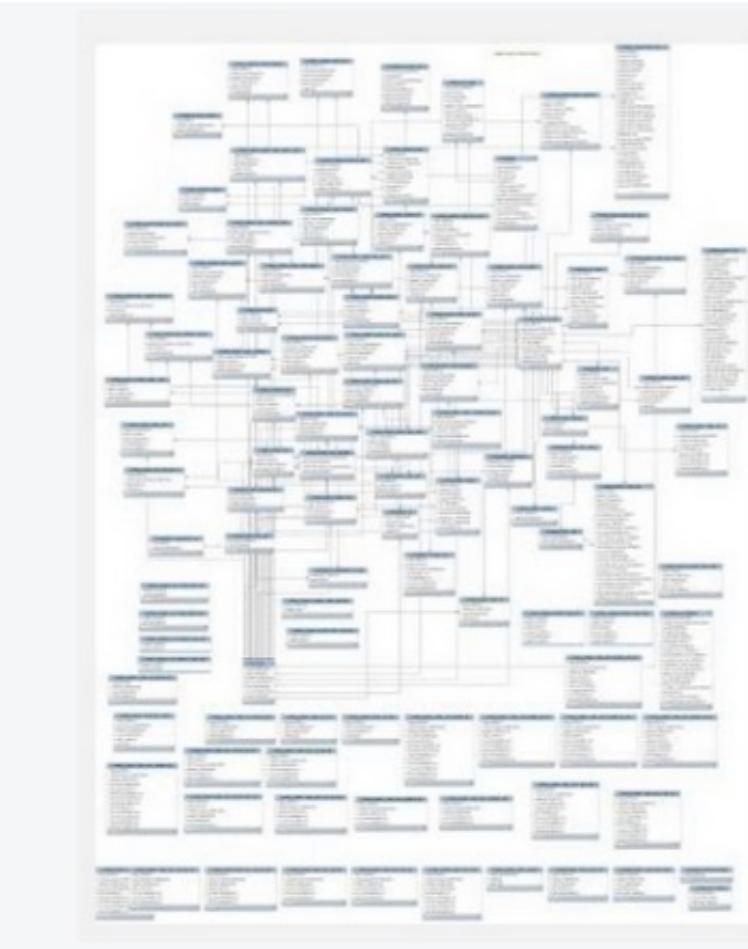
Pers_ID	Surname	First_Name	City
0	Miller	Paul	London
1	Ortega	Avaro	Valencia
2	Huber	Urs	Zurich
3	Blanc	Gaston	Paris
4	Bertolini	Fabrizio	Rom

Car:

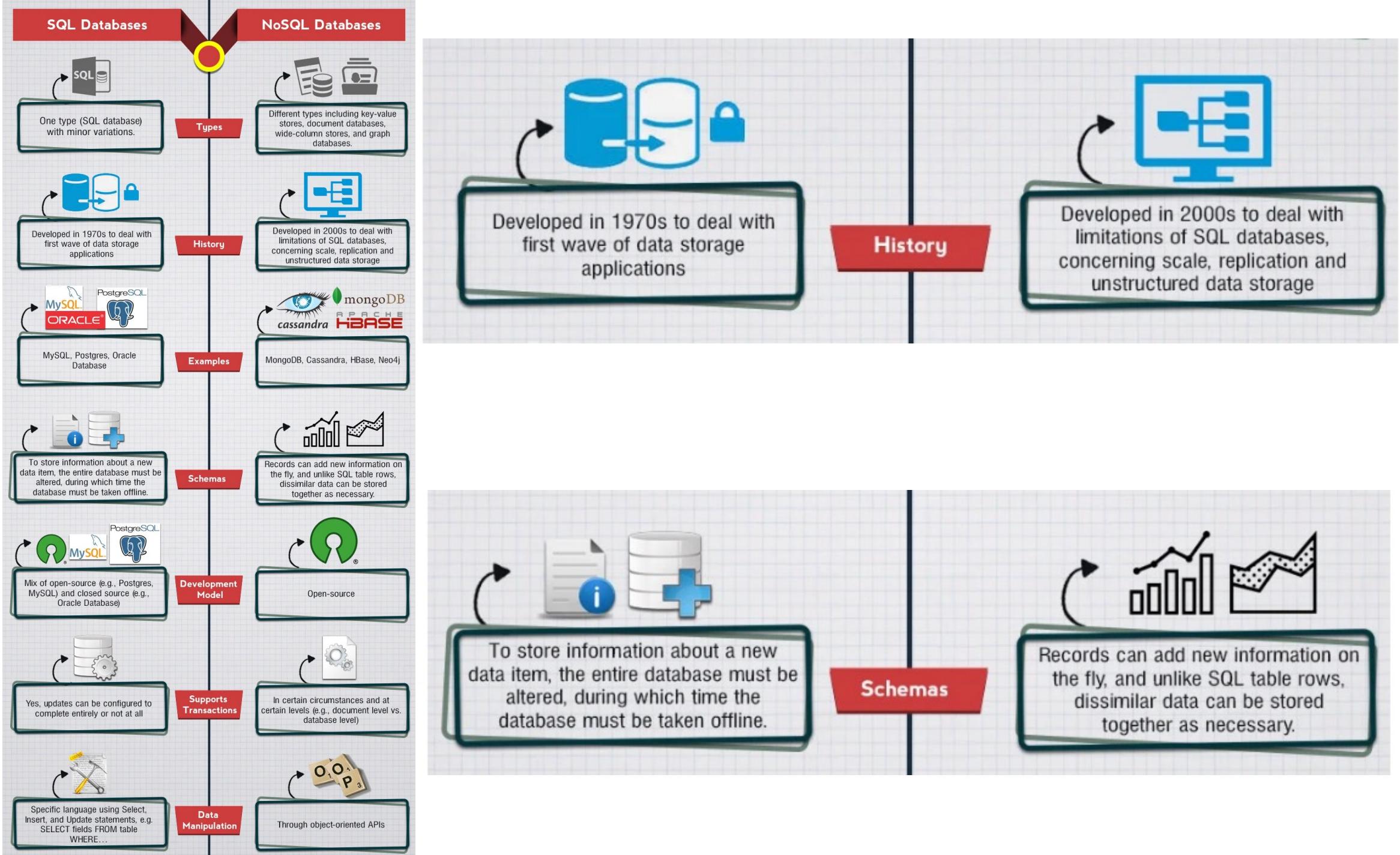
Car_ID	Model	Year	Value	Pers_ID
101	Bentley	1973	100000	0
102	Rolls Royce	1965	330000	0
103	Peugeot	1993	500	3
104	Ferrari	2005	150000	4
105	Renault	1998	2000	3
106	Renault	2001	7000	3
107	Smart	1999	2000	2

MongoDB Document

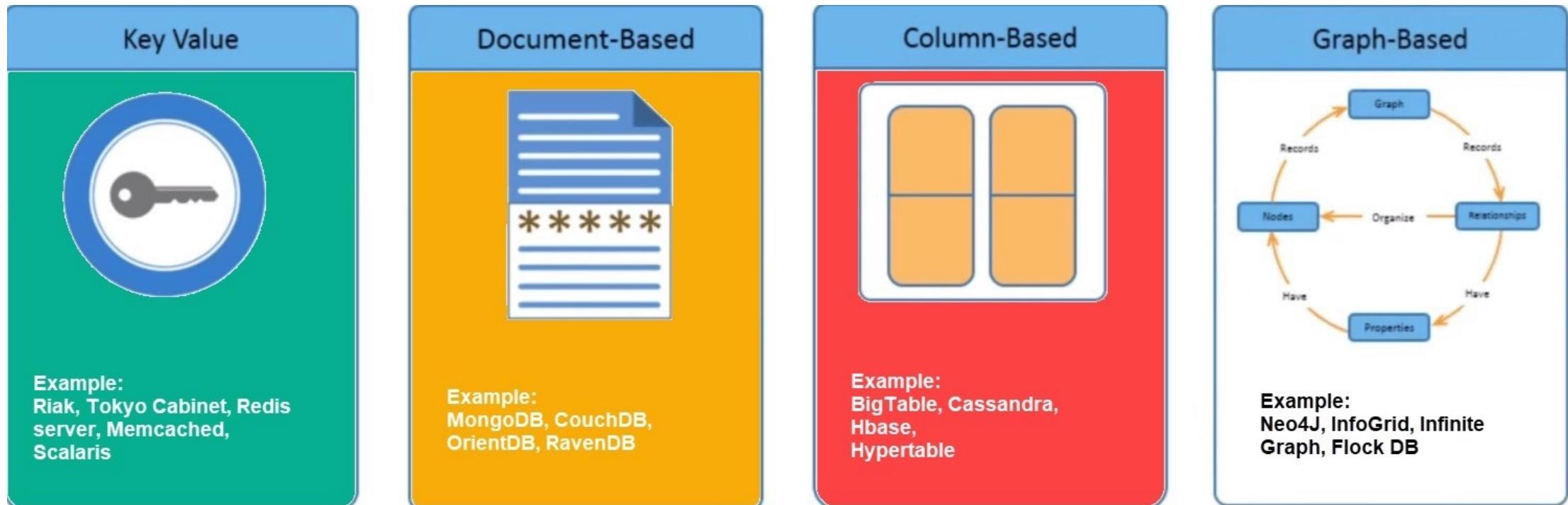
```
{  
    first_name: 'Paul',  
    surname: 'Miller'  
    city: 'London',  
    location: [45.123,47.232],  
    cars: [  
        { model: 'Bentley',  
         year: 1973,  
         value: 100000, ... },  
        { model: 'Rolls Royce',  
         year: 1965,  
         value: 330000, ... }  
    ]  
}
```

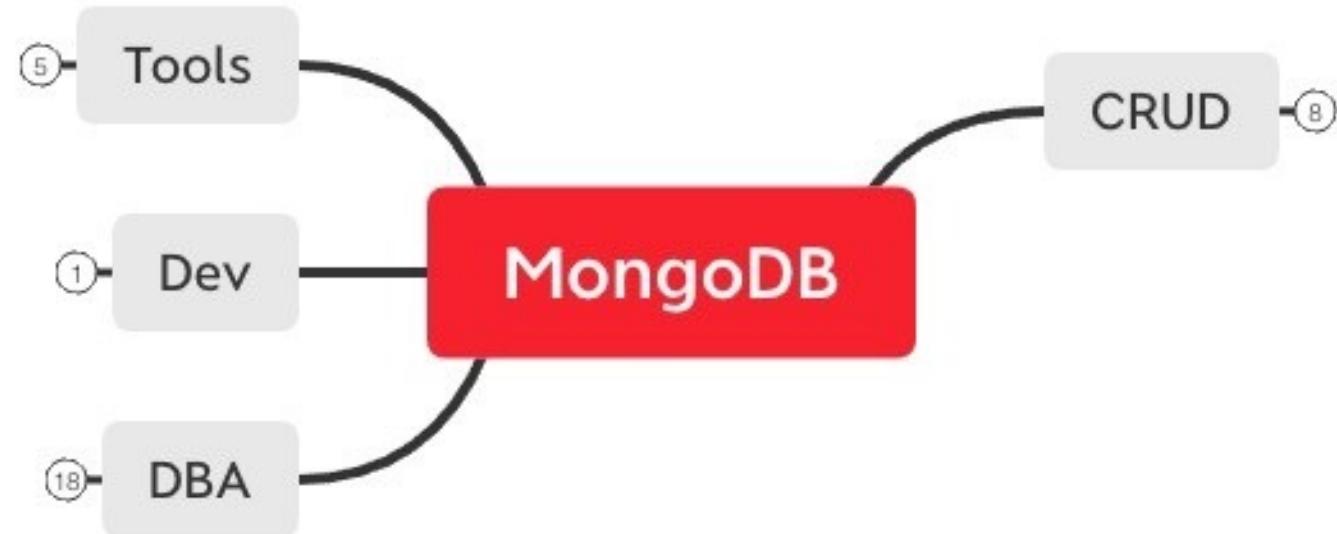


```
{ customer_id : 1,  
  first_name : "Mark",  
  last_name : "Smith",  
  city : "San Francisco",  
  phones: [ {  
      type : "work",  
      number: "1-800-555-1212"  
    },  
    { type : "home",  
      number: "1-800-555-1313",  
      DNC: true  
    },  
    { type : "home",  
      number: "1-800-555-1414",  
      DNC: true  
    }  
  ]  
}
```



NoSQL Landskab







mongoDB

MongoDB – CRUD-operations

C • insert()

R • find()

U • update()

D • remove()

Query Selectors : Comparison

\$gt:Val Greater than Val

\$gte:Val Greater than equals Val

\$lt:Val Lower than Val

\$lte:Val Lower than equals Val

\$all:Array All Array elements are included in field array value

\$in:Array Elements with values contained in Array

\$nin:Array Elements with values Not contained in Array

\$ne:Val Not equal

Val can be any Scalar Integer, String, Date, etc

\$exists

db.users.find({ name: { \$exists: true } })

\$not

db.users.find({ name: { \$not: { \$eq: "Kyle" } } })

db.users.insertOne(← collection

{

 name: "sue", ← field: value

 age: 26, ← field: value

 status: "pending" ← field: value

}

)

db.users.find(← collection

 { age: { \$gt: 18 } }, ← query criteria

 { name: 1, address: 1 } ← projection

).limit(5) ← cursor modifier

db.users.updateMany(← collection

 { age: { \$lt: 18 } }, ← update filter

 { \$set: { status: "reject" } } ← update action

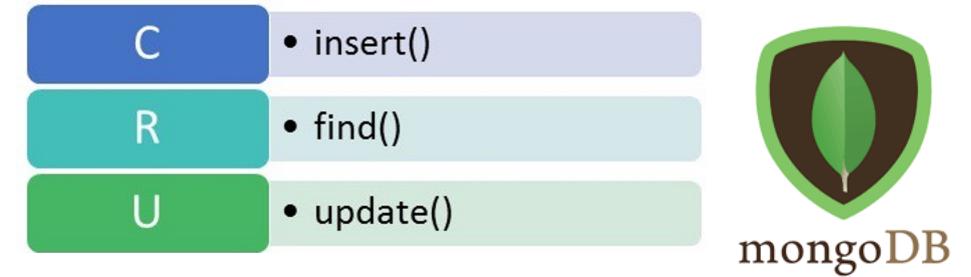
)

db.users.deleteMany(← collection

 { status: "reject" } ← delete filter

)

R Mongolite – CRUD-operations



```
find(  
  query = '{}',  
  fields = '{"_id" : 0}',  
  sort = '{}',  
  skip = 0,  
  limit = 0,  
  handler = NULL,  
  pagesize = 1000  
)
```

Retrieve fields from records matching query. Default handler will return all data as a single dataframe.

```
insert(  
  data,  
  pagesize = 1000,  
  stop_on_error =TRUE,  
  ...)
```

Insert rows into the collection.
Argument 'data' must be a **data-frame**, named list (for single record) or character vector with json strings (one string for each row).

```
update(  
  query = '{}',  
  update = '{"$set":{}}',  
  upsert = FALSE,  
  multiple = FALSE  
)
```

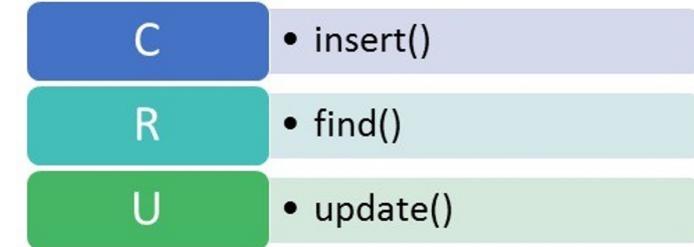
Modify fields of matching record(s) with value of the update argument

MongoDB – CRUD-operations

<https://github.com/jeroen/mongolite>

<https://jeroen.github.io/mongolite/index.html>

<https://cran.r-project.org/web/packages/mongolite/mongolite.pdf>



```
find(  
  query = '{"domain" : { "$regex" : "horton" } }'  
)
```

Retrieve fields from records matching query. Default handler will return all data as a single dataframe.

import data - monogsh tool

```
ls evt_5* | while read x; do mongoimport -d=wyscoutdutch -c=games $x; done  
mongoimport --uri "mongodb+srv://soccer.cukt54z.mongodb.net/myFirstDatabase&authSource=admin" -d=wyscoutdutch
```

Or GUI

The screenshot shows the Studio 3T Pro for MongoDB interface. At the top, there's a toolbar with various icons: Map-Reduce, Compare, Schema, Reschema, Tasks, Export, Import, Data Masking, and SQL Migration. Below the toolbar, a green banner says: "Schedule your first import, export or data comparison by clicking on 'Tasks' in the toolbar." The main window has a tab bar at the top with tabs like Quickstart, matches, Modify Index: textidx, JSON Import*, homes, matches, and restaur. The JSON Import* tab is currently active. Below the tabs, there are buttons for Save task, Save task as..., Schedule, and Run. A section titled "Target connection" shows a dropdown menu set to "mydb (localhost:27017)" with a "Change target" option. There's also a checkbox for "Validate JSON before import". Under "Select JSON sources to import:", there are buttons for Add source, Remove source, and Paste from clipboard. A table below lists the import details:

JSON Source	Target Database	Target Collection
restaurants2.json	Downloads	restaurants2

verify data - mongo shell

```
[dining2> use wyscoutdutch
switched to db wyscoutdutch
[wyscoutdutch> show collections;
games
matches
players
[wyscoutdutch> db.players.findOne()
{
  _id: 869465,
  shortName: 'T. van Bergen',
  firstName: 'Thom',
  middleName: '',
  lastName: 'van Bergen',
  height: 0,
  weight: 0,
  birthDate: '2004-01-06',
  birthArea: { id: 528, alpha2code: 'NL', alpha3code: 'NLD', na
  passportArea: { id: 528, alpha2code: 'NL', alpha3code: 'NLD',
  role: { name: 'Forward', code2: 'FW', code3: 'FWD' },
  foot: null,
  currentTeamId: 23564,
  currentNationalTeamId: null,
  gender: 'male',
  status: 'active',
  imageDataURL: 'https://cdn5.wyscout.com/photos/players/public/ndplayer_1
}
wyscoutdutch>
```

OrgUI

The screenshot shows the Studio 3T Pro interface. At the top, there are tabs for Connect, Collection, IntelliShell, SQL, Aggregate, Map-Reduce, Compare, Schema, and Refresh. Below the tabs, a sidebar lists databases and collections under 'mydb localhost:27017 [direct]'. The 'wyscoutdutch' collection is expanded, showing its sub-collections: games, matches, and players. On the right side, a query editor window is open. It displays a query for the 'players' collection, which returns a single document. The document contains player details such as '_id', 'label', 'date', 'dateutc', 'status', 'competitionId', 'seasonId', 'roundId', and 'gameweek'. The 'Result' tab is selected, showing the JSON output of the query.

```
1  {
2    "_id" : NumberInt(5349369),
3    "label" : "Zagłębie Lubin - Rak
4    "date" : "2022-11-12 20:00:00",
5    "dateutc" : "2022-11-12 19:00:0
6    "status" : "Played",
7    "competitionId" : NumberInt(692
8    "seasonId" : NumberInt(188088),
9    "roundId" : NumberInt(4426775),
10   "gameweek" : NumberInt(17)
11  }
12  {
13    "_id" : NumberInt(5349375),
14    "label" : "Wisła Płock - Cracov
15    "date" : "2022-11-11 20:30:00",
16    "dateutc" : "2022-11-11 19:30:0
17    "status" : "Played",
18    "competitionId" : NumberInt(692
19    "seasonId" : NumberInt(188088),
20    "roundId" : NumberInt(4426775),
21    "gameweek" : NumberInt(17)
22  }
```

work on data - CRUD shell/Gui

```
[sales2> db.sales2.find({"storeLocation":/Lond/}, {"storeLocation":1, "items.quantity":1, "saleDate":1, "_id":0}).limit(2)
[
  {
    saleDate: ISODate("2014-11-11T02:13:51.893Z"),
    items: [
      { quantity: 4 },
      { quantity: 7 },
      { quantity: 5 },
      { quantity: 5 },
      { quantity: 3 },
      { quantity: 5 },
      { quantity: 1 },
      { quantity: 2 },
      { quantity: 3 },
      { quantity: 1 }
    ],
    storeLocation: 'London'
  },
  sales2> db.sales2.find({"storeLocation":/Lond/}, {"storeLocation":1, max:{$max:"$items.quantity"}, saleDate:1, _id:0}).limit(2)
[
  {
    saleDate: ISODate("2014-11-11T02:13:51.893Z"),
    storeLocation: 'London',
    max: 7
  },
  sales2> db.sales2.updateOne({_id:ObjectId("5bd761dcae323e45a93ccff3")}, {$set:{storeLocation:"Nyborg"}})
```

export data - mongodb tool



```
| => mongoexport -v --collection=sales --db=sales -q='{"storeLocation":"London"}' --pretty --jsonArray --out dumplondon.json
2022-03-05T12:44:49.136+0100      will listen for SIGTERM, SIGINT, and SIGKILL
2022-03-05T12:44:49.142+0100      connected to: mongodb://localhost/
2022-03-05T12:44:49.242+0100      exported 794 records
```

```
| => head dumplondon.json
[{
    "_id": {
        "$oid": "5bd761dcae323e45a93ccfed"
    },
    "saleDate": {
        "$date": "2015-09-02T16:11:59.565Z"
    },
    "items": [
        {
            "name": "binder",
            "price": 10.99
        }
    ]
}]
```

export data - mongodb tool

```
=> mongoexport -v --collection=sales --db=sales -q='{"storeLocation":"London"}' --pretty --jsonArray --out dumplondon.json  
2022-03-05T12:44:49.136+0100      will listen for SIGTERM, SIGINT, and SIGKILL  
2022-03-05T12:44:49.142+0100      connected to: mongodb://localhost/  
2022-03-05T12:44:49.242+0100      exported 794 records
```

```
mongoexport --db=statsbomb --collection=matches --query='{"player.name":{ "$in": ["Claire Emslie","Lucy Graham"] }}' --out dump.json
```

Or GUI →

The screenshot shows the MongoDB Compass interface with the title bar "Quickstart × rest × sales × matches × Export* ×". The "Export" tab is selected. Below it, the "Export overview" tab is active, showing "Export unit #1 - JSON". The "Export source" section shows a hierarchy: statsbomb > matches > Find query. The "Change source" button is visible. The "Export target" section has a radio button for "Clipboard" and another for "File" which is selected, pointing to the path "/Users/thor/Documents/mydb/statsbomb/matches". The "Query" field contains the expression "{\$and:[{"player.name":"Claire Emslie"}, {"type.name":"Shot"}]}". The "Projection" field contains {"shot":1,_id:0}. The "Sort" and "Limit" fields are both set to 0. A green checkmark at the bottom indicates "Valid query!".

Update MongoDB from R

<https://jeroen.github.io/mongolite/>

```

g.R x Q ulst x Q ulst x Untitled1* x df3 x london x Mongo_Rest.R x
1 library(mongolite)
2
3 con <- mongo(
4   collection = "rest",
5   db = "dining2",
6   url = "mongodb://localhost",
7   verbose = TRUE
8 )
9 restorig <- con$find(limit=10)
10
11 numt=2
12 for (i in 1:ncol(rests)) {
13   numt<-numt*0.95
14   rests[i,"Viggo"]<-numt[1]
15 }
16
17 for (i in 1:ncol(rests)) {
18   key <- rests[i,6]
19   val <- rests[i,9]
20   con$update(
21     query=paste0('{"restaurant_id": "' ,key, '"}'),
22     update=paste0('{"$set":{"lbl": "' ,val, '"}}')
23   )
24 }
25
26

```

From mongo

restorig	10 obs. of 8 variables
rests	10 obs. of 9 variables

Modif dataframe

name	restaurant_id	age	lbl
Wendy'S	30112340	30112340	kurt
Wilken'S Fine Food	40356483	30112340	kurt
Brunos On The Boulevard	40356151	30112340	kurt
Riviera Caterer	40356018	30112340	kurt
Kosher Island	40356442	30112340	kurt
Taste The Tropics Ice Cream	40356731	30112340	kurt
Regina Caterers	40356649	30112340	kurt
Dj Reynolds Pub And Restaurant	30191841	30112340	NA
Morris Park Bake Shop	30075445	30112340	NA
May May Kitchen	40358429	30112340	NA

name	restaurant_id	age	lbl	Viggo
Wendy'S	30112340	30112340	kurt	1.900000
Wilken'S Fine Food	40356483	30112340	kurt	1.805000
Brunos On The Boulevard	40356151	30112340	kurt	1.714750
Riviera Caterer	40356018	30112340	kurt	1.629012
Kosher Island	40356442	30112340	kurt	1.547562
Taste The Tropics Ice Cream	40356731	30112340	kurt	1.470184
Regina Caterers	40356649	30112340	kurt	1.396675
Dj Reynolds Pub And Restaurant	30191841	30112340	NA	1.326841
Morris Park Bake Shop	30075445	30112340	NA	NA
May May Kitchen	40358429	30112340	NA	NA

Kyle

db.insertOne({})

db.find({}, {})



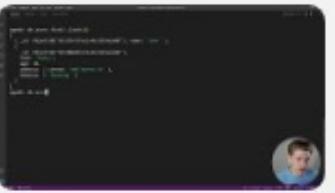
1.23

Basic Commands



4.18

Insert Commands



8.50

Basic Query
Commands

show

db.insertOne({})

db.find({}, {}).limit(2)

use

db.insertMany({})

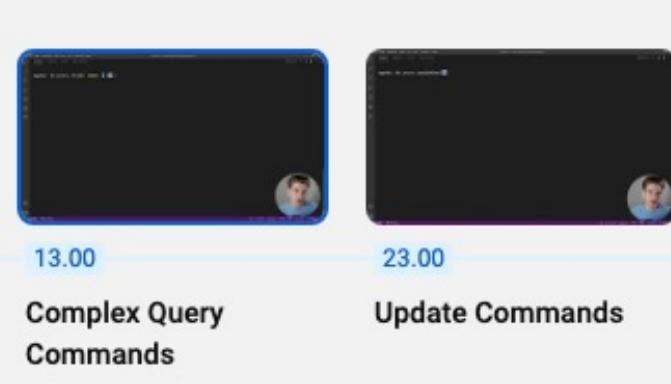
.sort({})
.skip(3)

WHERE

df.find({name: "lb"}, {})

SELECT

df.find({}, {name:1, _id:0})



WHERE

df.find({name: {\$eq:"lb"} })

\$neq

\$gt

\$in:[]

\$nin:[]

\$exists: true

\$gt:20,\$le:2

db.find({\$or: [{ age: {\$lt:1}}, { name:"lb"}]})

db.find({\$expr: {\$gt:["\$d","\$c"]}})

db.find({"vardata.pris": { \$gt:34000}})