

DATABASE RAPPORT

Hassan Raza Hussain, Hassan Mohdi, Gordan Mikkelsen & Kristoffer Rath
hansen



29. MAJ 2020

CPHBUSINESS

GITHUB LINK: <https://github.com/lakridserne/db-exam>

Github link: <https://github.com/lakridserne/db-exam>

INDHOLD

| | |
|--|----|
| Introduction | 2 |
| Problem statement..... | 2 |
| Design and architecture | 3 |
| Requirements and analysis..... | 3 |
| Business case..... | 3 |
| User stories..... | 4 |
| Database fields | 5 |
| Database structure..... | 5 |
| Design of databases | 5 |
| Choice of database technologies | 6 |
| Implementation of databases..... | 8 |
| Logging..... | 8 |
| Session..... | 9 |
| Transactions | 10 |
| Installation instructions..... | 11 |
| Conclusion..... | 14 |
| Bilag: | 15 |

Github link: <https://github.com/lakridserne/db-exam>

..

INTRODUCTION

We are a group of 4 members who have been given an assignment for the database exam project. The objective of this assignment is to provide database support to an application, as a proof of concept. In the assignment, multiple diverse database models are built and implemented for different purposes.

We have chosen to deal with data from two sources, which has to do with houses and apartments in Denmark. One source, DAWA, contains a list on every address in Denmark and additional information regarding those addresses that is public to the general public. The other source is from a real estate agent containing data about houses and apartments sold in Denmark.

Based on these data sources, we have made the following problem statement.

PROBLEM STATEMENT

Getting a quote on your house or apartment can take time. If you want to quickly find out roughly where you stand can be very challenging. Therefore we want to give an immediate rough estimate on your house or apartment based on available public data

DESIGN AND ARCHITECTURE

Requirements and analysis

Business case

| | |
|--|---|
| Contribution to the business strategy | Our strategy is to provide homeowners with a tool to quickly estimate the value of their house or apartment. This is done through analysis of the public data regarding the address and the area. |
| Options considered | <ol style="list-style-type: none">1. More and better advertising to get more customers to find out how much value their house or apartment has.2. Live with slow quotation time and thus less customers buying and selling3. The new system |
| Benefits | <ol style="list-style-type: none">1. Immediate quotes2. Increased sales based on the immediate quotes |
| Timescales | Our current analysis shows us that the system will take around 6 weeks to develop |
| Cost | The new system is estimated to cost around \$20,000. |

Github link: <https://github.com/lakridserne/db-exam>

| | |
|---------------------------------------|--|
| Expected return on investments | It is expected that the system will recoup the costs and start giving a return on investment through increased sales due to immediate quotes. |
| Risks | There are various risks such as inaccurate quotes that we are trying to mitigate through sampling different house prices in different areas and compare it with real listings. |

User stories

User story 1 - House price

As customer on the website, I want to find out what my house is worth, so I go to our website where I enter my house or apartments address. Then the website will respond with an approximate price with how much the house is worth. That way, I can better decide if I want to sell my house.

User story 2 - number of rooms on a specific house

As customer on the website, I want to find out how many rooms a specific house has that was sold by Boliga in the past. Therefore, I go to our website, where I enter address on the specific house, the website then respond back with the number of rooms in the house.

So I can see if the specific house fits my needs

User story 3 - Graphs

As salesman on the on database, I want to see statistics and graphs on the different houses and apartments in question.

Github link: <https://github.com/lakridserne/db-exam>

This way, I can better target the right customers.

Database fields

Our Database fields is designed on the way to give you access to accurate the information you want, since a proper design on our database is essential for reaching our goals, even though we spent a lots of time on it, it is worth it since we end up with a database that meets our needs and can easily be adapted to changes. In our database We don't really have a lots of different fields, because we are mostly working with basic data, where most of our data type are either varchar or int since we feel like they fits best with data we have in our boliga database.

Database structure

In the following, we'll discuss how we designed our database structure and chose which technologies were appropriate to use taking the user stories and business case into consideration.

Design of databases

We have chosen 4 different databases to work with (per the requirements of the assignment). The databases were selected based on criterias we'll delve deeper into below, but first, I'd like to go through what data we have and how we are planning on using those data.

In our business case, we have some data from sales of houses and apartments. Based on these sales, and based on public data from DAWA (Danmarks Adresse Web API) we will be giving some very rough estimates on how much you could get for selling your house or apartment in your area.

In order to fulfill that we need some databases to store our data in. We have thus decided to have databases with these purposes.

- Database for DAWA data (general information about addresses in Denmark)
- Database for boliga data (data showing which sales boliga have made)

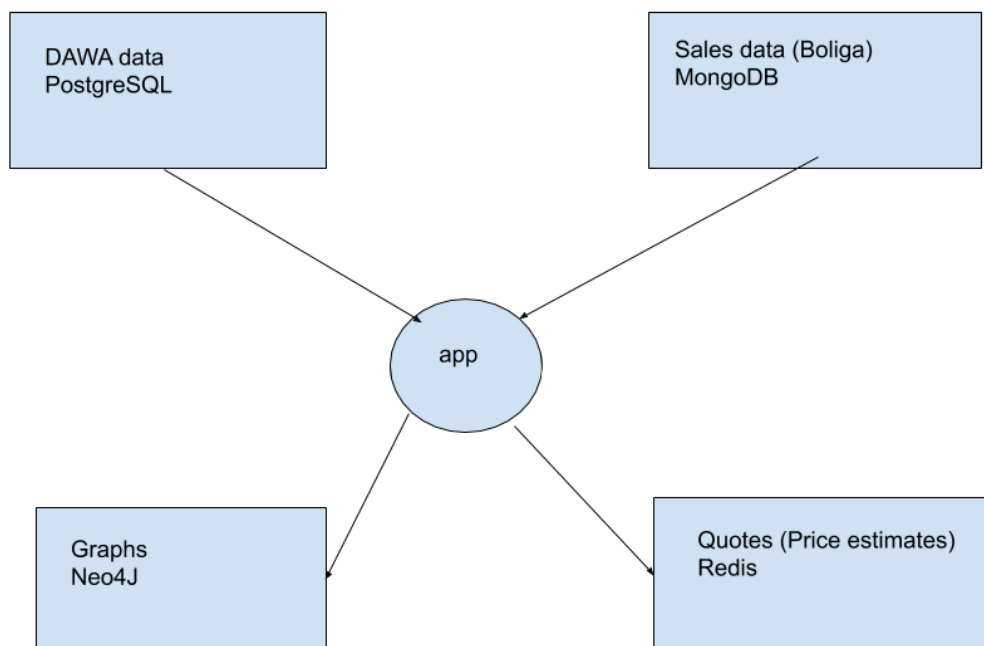
Github link: <https://github.com/lakridserne/db-exam>

- Database used for generating graphs for statistics purposes
- Database for storing the data we generate

Could we have stored the data in the same database? Yes, we could theoretically have stored a large part of the data in the same database if we wanted to. However, if we do that, the underlying databases will not be optimized for the data we store and process.

By using these 4 different databases, we can choose each one from what things we expect to store, how popular we expect it to, how resilient it should be, and many more parameters.

Below you can see a diagram detailing what databases we will be using for what data. We'll delve deeper into the choices of the different databases in "Choice of database technologies".



(IMAGE 1)

Choice of database technologies

We have evaluated the different databases on a number of criteria to see which database types would be good to use for our data.

Github link: <https://github.com/lakridserne/db-exam>

One of the things we thought about when deciding which databases we wanted to use was whether we should use relational databases or NoSQL databases. Relational databases are pretty good if you have data where you need to relate to it with some other specific data.

However, NoSQL databases can be faster than relational databases in certain circumstances (note not in all) and allows you more flexibility to store different kinds of data (where relational databases are more strict). This can be good because it allows you to stretch your models a bit, but it can also be dangerous if you do it too much because you can end up having inconsistent data or data structured in a way so it'll take a very long time to find the data you are looking for.

We chose to have our data from DAWA in a PostgreSQL as that is a relational database. This is good for the kind of data we have in DAWA, where we want to quickly be able to find specific addresses from different informations.

For our sales data, we have chosen MongoDB to store those. The reason we chose to store our sales data in a NoSQL database is that it is highly scalable (you can just launch more mongod processes).

We have a lot less data to search through with the sales data and can with indexes optimize things a lot so performance is optimal.

Neo4J was without a question the choice for making graphs for us, since it's made for making graphs primarily.

Lastly, we have Redis that are driving our quotes database. Redis is like MongoDB (and for a part also PostgreSQL) a NoSQL database. The advantage of using Redis is that it's very fast as it runs in-memory.

We also looked at things like CAP and ACID for the different databases.

| | CAP qualities | ACID |
|---------|-----------------------------|-------------------------|
| MongoDB | MongoDB is highly available | From MongoDB 4.0 multi- |

Github link: <https://github.com/lakridserne/db-exam>

| | | |
|------------|--|---|
| | through automatic failover and strongly consistent by default. Therefore the CAP theorem is fulfilled. | document ACID transactions is possible. |
| Redis | Since Redis is not designed to be run on several servers, it is not as highly available as the other databases. However, it is consistent and tolerates partitions. Therefore it follows CP and not CAP. | Redis is single threaded allowing it to be ACID compliant by default. |
| Neo4J | Neo4J only follows CA and is not partition tolerant. | Neo4J is fully ACID compliant. |
| PostgreSQL | PostgreSQL is not partition tolerant - so it's CA. | PostgreSQL is ACID compliant from the ground up. |

Implementation of databases

Logging

When we build applications and services, it's being more and more important to be able to discover information about the running system. More precisely being able to analyse the information to diagnose problems, discover information so we can avoid problems before they become severe, or just to discover information about users. All this requires logging.

Redis:

Redis has built in logging capabilities. In the config file you can set up where it should send the log lines to - for example to stdout (the command prompt) - or a specific file. In Linux it also have native support for syslog. Syslog accepts log messages from every program that sends it a message and it routes those message to various on-disk log files, handling rotation and deletion of old logs. With configuration, it can also forward messages to other server for

Github link: <https://github.com/lakridserne/db-exam>

further processing. When we talk about it as a service, it's much more convenient than logging to files directly, because things like all of the special log file rotation and deflection is handled for us already.

PostgreSQL

By having PostgreSQL run like a cloud service at Amazon AWS, it gives us an opportunity to see the logging on Amazon's website. It gave us a overview of which events there has been in our PostgreSQL database. In the overview for logging it shows the time for the recent event in the database and a system note for what the event was.

MongoDB:

Before we could enable logging for MongoDB, we had to set a single verbosity for log messages. That had the following form: `db.setLogLevel(<level>, <component>)`.

In this form, `<level>` take int as its type, and `<component>` take string as its type. `db.setLogLevel()` sets a single verbosity level. If we had to set more than one verbosity level in a single operation, we had to use a `setParameter` command to set the `LogComponentVerbosity` parameter. Another way we could do it was to specify verbosity settings in the configuration file.

Neo4J:

Neo4J does contain some possibilities for logging. Thus it does save some log information in some log files. However, most of the log functions is only for the enterprise version of Neo4J, and it was not that relevant for our project.

Session

Redis:

To store the sessions in Redis what we do is that we use Redis hashes which is a pretty standard way for session state. "Redis Hashes are maps between string fields and string

Github link: <https://github.com/lakridserne/db-exam>

values" ⁽¹⁾ . So we use the session ID as the key and to use the hash fields to describe the state. ⁽²⁾

MongoDB:

Regarding sessions storage we implemented an API called Spring Session MongoDB. This allows us to manage a users session information and storing it in MongoDB. The dependencies is added in the xml file so we can import the framework. We do this because we want to get an overview of interaction with the application in case something goes wrong. ⁽³⁾

PostgreSQL:

When we store sessions in PostgreSQL, we gather the database information that is important. After that we create a session store schema. Then follow a tutorial on how to point the policy server to the database. Lastly, we restart the server so the changes goes into effect ⁽⁴⁾

Transactions

Redis:

Redis is different compared to most of the other database transactions, where it can use MULTI, EXEC, DISCARD and WATCH as the foundation of transactions. Redis also allows you to execute a group of command on a single step.

Compared to our other databases, Redis does not support roll backs, which means it can fail during a transaction, where it will keep trying to execute the rest of the transaction.

¹ <https://redis.io/topics/data-types>

² <https://stackoverflow.com/questions/54798523/how-to-configure-redis-cache-for-session-management-with-uuids>

³ <https://spring.io/projects/spring-session-data-mongodb#learn>

⁴ <https://techdocs.broadcom.com/content/broadcom/techdocs/us/en/ca-enterprise-software/layer7-identity-and-access-management/single-sign-on/12-7/installing/install-a-policy-server/configure-odbc-databases-as-policy-session-key-and-audit-stores/configure-odbc-databases-as-session-store/store-session-information-in-postgresql.html>

Github link: <https://github.com/lakridserne/db-exam>

PostgreSQL;

PostgreSQL has been having transaction support for long time, where you have to use the BEGIN, COMMIT or ROLLBACK queries through a client to control the transaction, because node-postgres is a low level and un-opinionated language. PostgreSQL has support for ACID per default.

MongoDB:

After MongoDB updated to version 4.0, they now have support for multi-document ACID transactions, and bit later in version 4.2 they also added support for distributed transaction, which make you able to do multi-document transaction on a sharded clusters and incorporates the existing support for multi-document transaction on replica sets.

Installation instructions

Mongo:

We installed MongoDB community Edition on Windows. We used the official MongoDB documentation to install it. We installed MongoDB version 4.2. Before we could install MongoDB and work with the database we had to be sure that mongoDB 4.2 community edition supports our version of Windows, and since we met the requirements, we could use MongoDB without any problems. Users that use earlier Windows versions than Windows 10 must install an update before installing mongoDB.⁵

Neo4J:

When we were installing Neo4J, the first thing we did was, to download it from their website along with going to the official documentation to get installation instructions. It was pretty easy work. Before we could start the Neo4J download on our computers we had to fill out some information about us on Neo4J's website. Then we could install the Neo4J desktop on our computers. The important thing here was we had to get the "Activation key", which you

⁵ <https://docs.mongodb.com/manual/tutorial/install-mongodb-on-windows/>

can see after we filled out the information as mentioned earlier. Another thing you also can do is if you not get the activation key for a reason, then you can also generate a key from within the app filling out the form on the right side of the app screen. ⁶

(IMAGE 2)

In the start we downloaded and ran PostgreSQL on the computer. But afterwards we decided to run it like a cloud service on Amazon AWS. The reason we did that was our DAWA data was too big to run it on PostgreSQL on our local machine. We decided to run it on Amazon, because one group member knew we could do it that way. We experienced that doing it in that way it gave us a lot features and we don't needed to think about if the server is running. As we can see on the picture (IMAGE 3) below it shows us how our PostgreSQL dashboard looked like in Amazon.

6 <https://neo4j.com/download-thanks-desktop/?edition=desktop&flavour=windows&release=1.2.8&offline=true#installation-guide>

Github link: <https://github.com/lakridserne/db-exam>

The screenshot displays the Amazon RDS console interface. On the left is a navigation sidebar with options like Dashboard, Databases, Performance Insights, Snapshots, Automated backups, Reserved instances, Subnet groups, Parameter groups, Option groups, Events, Event subscriptions, Notifications, and Certificate update. The main content area shows the 'house-quotes' database instance. At the top, there's a breadcrumb trail: RDS > Databases > house-quotes. Below this is a 'Summary' section with a table of instance details:

| Summary | | | |
|-------------------------------|-----------------------------------|----------------------|----------------------------|
| DB identifier house-quotes | CPU 1.00% | Info Available | Class db.t3.micro |
| Role Instance | Current activity 0 Connections | Engine PostgreSQL | Region & AZ eu-north-1b |

Below the summary is a tabbed interface with 'Connectivity & security' selected. This section contains three sub-sections:

- Endpoint & port:** Endpoint: house-quotes.ctyfhkfvzw5.eu-north-1.rds.amazonaws.com, Port: 5432
- Networking:** Availability zone: eu-north-1b, VPC: vpc-e26ace8b
- Security:** VPC security groups: default (sg-c04378ac) (active), Public accessibility: Yes

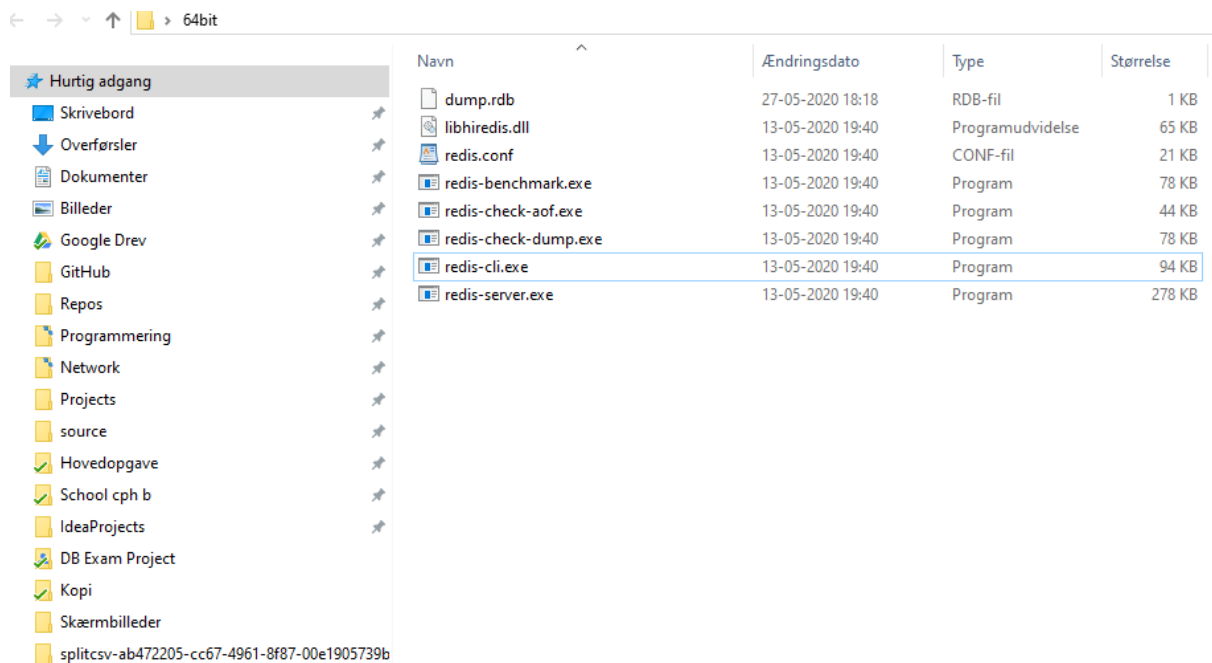
The footer of the console shows 'Feedback', 'English (US)', and copyright information for Amazon Web Services, Inc. (© 2008 - 2020).

(IMAGE 3)

Redis:

When we installed Redis, we downloaded a zip folder which we unpacked. This folder contained a Redis server and Redis CLI. Then if we wanted to start the redis server we just have to click on redis-server.exe. If we wanted to be sure about if the Redis server is running, we could check it by click on redis-cli.exe which also is a command prompt and in that we can write "ping" and then if the server is running it'll answer with "pong". In the image which is showed below (IMAGE 4), we can see how the Redis file looked.

Github link: <https://github.com/lakridserne/db-exam>



(IMAGE 4)

CONCLUSION

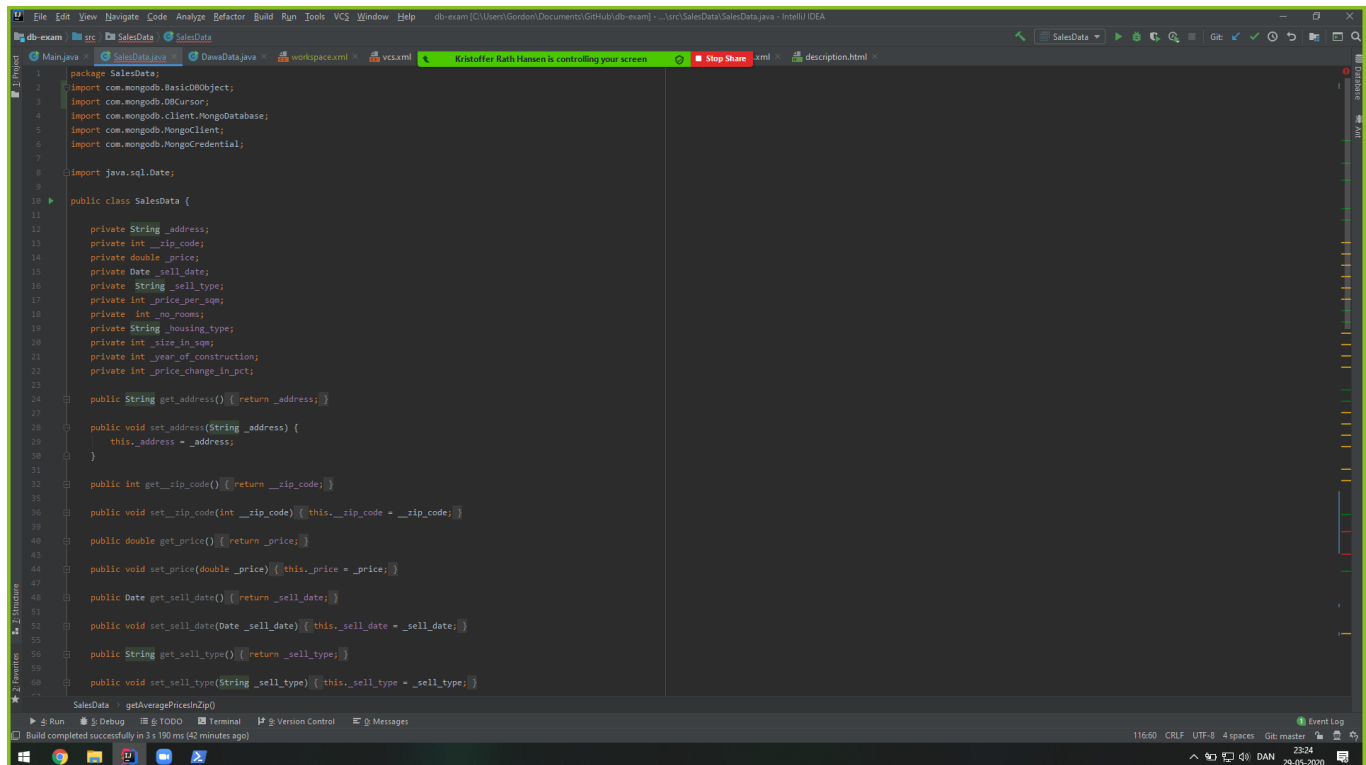
Based on our analysis and design we chose to use 4 databases, namely Redis, MongoDB, Neo4J and PostgreSQL and have them connected to a single application. For creating the application , we did some research and chose Java since it seems to have a large user community for the databases we chose. After having connected the different databases to our application, we imported our two data sources, the first one into PostgreSQL and the other MongoDB.

Then by retrieving the different data from multiple sources the system will calculate an average price. These prices that gets generated will then get stored in our Redis database for future use.

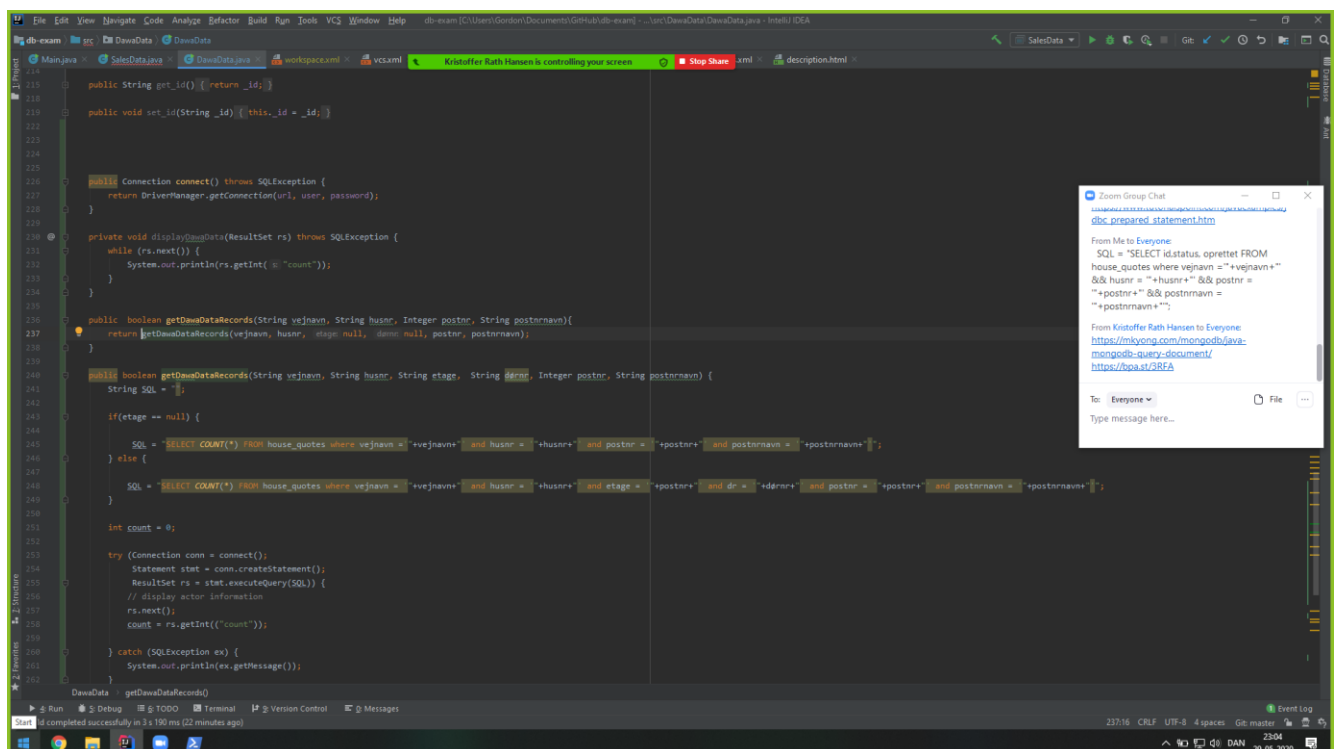
What our last database will do then is that Neo4j will collect the data that is necessary from MongoDB and PostgreSQL. From these data Neo4J will generate a graph which will display statistics.

Github link: <https://github.com/lakridserne/db-exam>

BILAG:

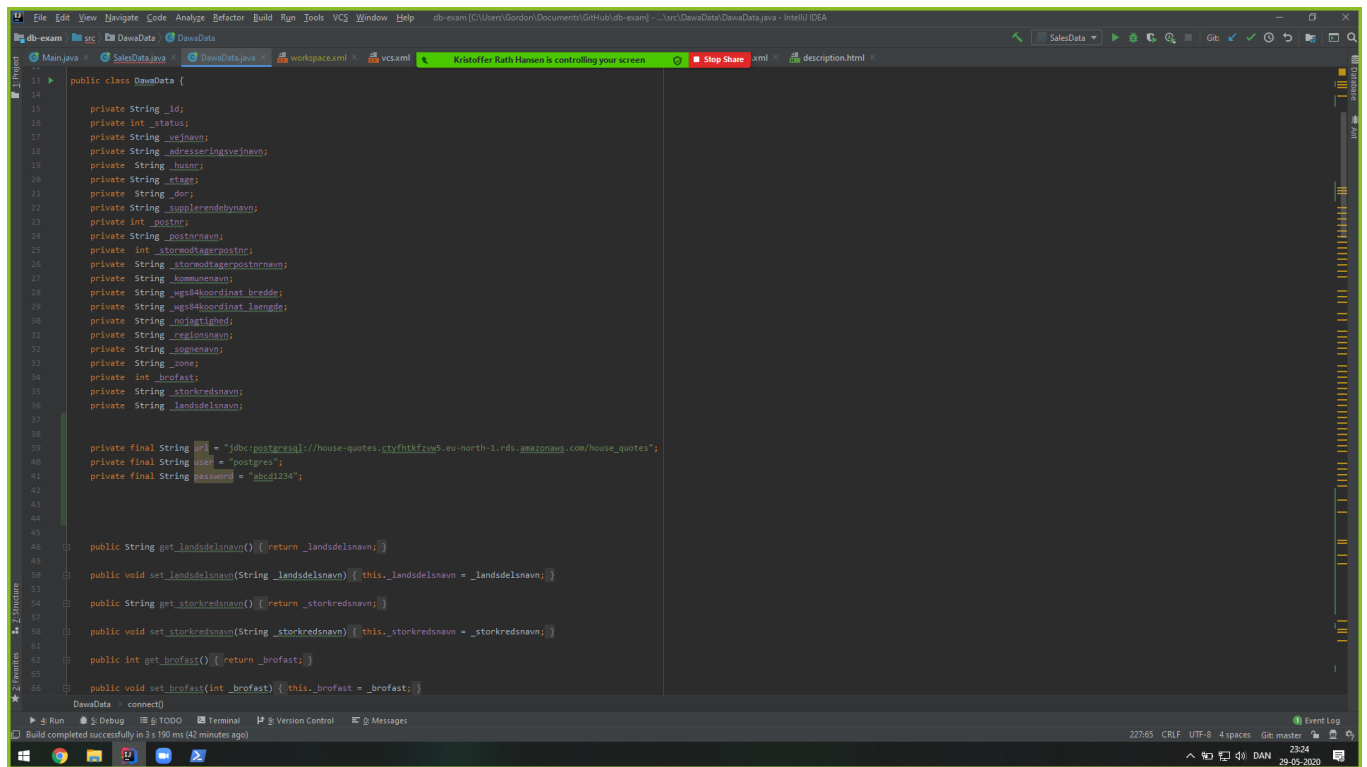


```
1 package SalesData;
2 import com.mongodb.BasicDBObject;
3 import com.mongodb.DBCursor;
4 import com.mongodb.client.MongoDatabase;
5 import com.mongodb.MongoClient;
6 import com.mongodb.MongoCredential;
7
8 import java.sql.Date;
9
10 public class SalesData {
11
12     private String _address;
13     private int __zip_code;
14     private double _price;
15     private Date _sell_date;
16     private String _sell_type;
17     private int _price_per_sqm;
18     private int __no_rooms;
19     private String _housing_type;
20     private int _size_in_sqm;
21     private int _year_of_construction;
22     private int _price_change_in_pct;
23
24     public String get_address() { return _address; }
25
26     public void set_address(String _address) {
27         this._address = _address;
28     }
29
30     public int get __zip_code() { return __zip_code; }
31
32     public void set_zip_code(int __zip_code) { this.__zip_code = __zip_code; }
33
34     public double get_price() { return _price; }
35
36     public void set_price(double _price) { this._price = _price; }
37
38     public Date get_sell_date() { return _sell_date; }
39
40     public void set_sell_date(Date _sell_date) { this._sell_date = _sell_date; }
41
42     public String get_sell_type() { return _sell_type; }
43
44     public void set_sell_type(String _sell_type) { this._sell_type = _sell_type; }
45
46     public void getAveragePricedInZip()
47
48     }
49 }
```



```
114 public String get_id() { return _id; }
115
116 public void set_id(String _id) { this._id = _id; }
117
118
119
120
121
122
123
124
125
126 public Connection connect() throws SQLException {
127     return DriverManager.getConnection(url, user, password);
128 }
129
130
131 private void displayQueryData(ResultSet rs) throws SQLException {
132     while (rs.next()) {
133         System.out.println(rs.getInt(1) + "count");
134     }
135 }
136
137 public boolean getDavaDataRecords(String vejnavn, String husnr, Integer postnr, String postnavn){
138     return getDavaDataRecords(vejnavn, husnr, etage, null, denr, null, postnr, postnavn);
139 }
140
141 public boolean getDavaDataRecords(String vejnavn, String husnr, String etage, String denr, Integer postnr, String postnavn) {
142     String SQL = "";
143
144     if(etage == null) {
145         SQL = "SELECT COUNT(*) FROM house_quotes where vejnavn = '"+vejnavn+"' and husnr = '"+husnr+"' and postnr = '"+postnr+"' and postnavn = '"+postnavn+"'";
146     } else {
147         SQL = "SELECT COUNT(*) FROM house_quotes where vejnavn = '"+vejnavn+"' and husnr = '"+husnr+"' and etage = '"+etage+"' and dr = '"+denr+"' and postnr = '"+postnr+"' and postnavn = '"+postnavn+"'";
148     }
149
150     int count = 0;
151
152     try {
153         Connection conn = connect();
154         Statement stat = conn.createStatement();
155         ResultSet rs = stat.executeQuery(SQL);
156         // display actor information
157         rs.next();
158         count = rs.getInt("count");
159     } catch (SQLException ex) {
160         System.out.println(ex.getMessage());
161     }
162 }
```


Github link: <https://github.com/lakridserne/db-exam>



```
13 public class DmsData {
14
15     private String _id;
16     private int _status;
17     private String _velum;
18     private String _adresseringsnavn;
19     private String _husnr;
20     private String _stage;
21     private String _dor;
22     private String _supplerendebnavn;
23     private int _postnr;
24     private String _postnavn;
25     private int _storkredsegerpostnr;
26     private String _storkredsegerpostnavn;
27     private String _kommunenavn;
28     private String _wgs84koordinat_bredde;
29     private String _wgs84koordinat_laengde;
30     private String _nojagtighed;
31     private String _regionsnavn;
32     private String _sogneavn;
33     private String _zone;
34     private int _brofast;
35     private String _storkredseger;
36     private String _landsdelsnavn;
37
38
39     private final String url = "jdbc:postgresql://house-quotes-gtyhkfz9y5.eu-north-1-rds.amazonaws.com/house_quotes";
40     private final String user = "postgres";
41     private final String password = "abcd1234";
42
43
44
45
46     public String get_landsdelsnavn() {return _landsdelsnavn;}
47
48     public void set_landsdelsnavn(String _landsdelsnavn) {this._landsdelsnavn = _landsdelsnavn;}
49
50
51     public String get_storkredseger() {return _storkredseger;}
52
53     public void set_storkredseger(String _storkredseger) {this._storkredseger = _storkredseger;}
54
55     public int get_brofast() {return _brofast;}
56
57     public void set_brofast(int _brofast) {this._brofast = _brofast;}
58
59     public void connect() {
60         // ...
61     }
62 }
```