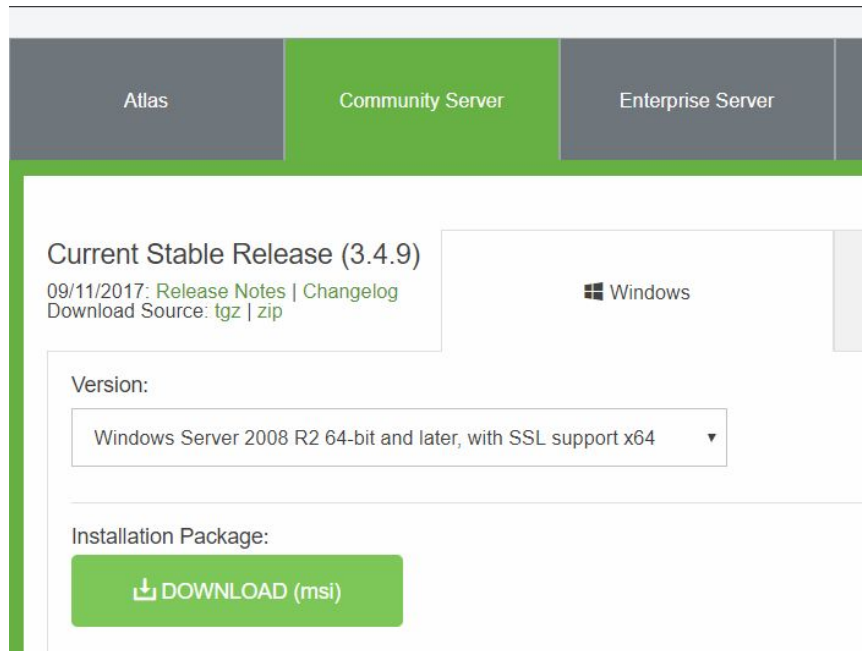# Introduction to MongoDB

26/03/2018

# What is MongoDB

-   It's a free NoSQL document oriented database.

    -   Not only SQL, since SQL queries are supported.

-   It uses JSON-like documents to save data.

    -   Flexible schemas - Allows missing fields.

-   High performance, high availability, and automatic scaling.

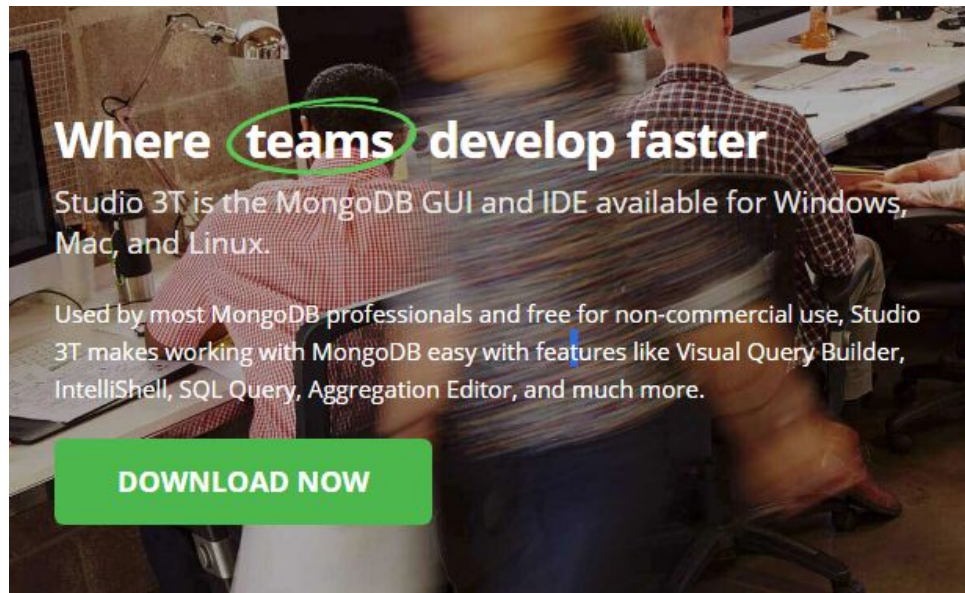    -   Used by large companies like Amazon.

# Installation

## MongoDB community server

## Studio 3T

# Installation

Step 1

Step 2

# Execution

Step 3 ->   Create folder "C:\data\db\"  ->  Open cmd -> run "mongod"

# What is a document

- A sum of key / value pairs

```
{
    name: "sue",              ⟵  field: value
    age: 26,                  ⟵  field: value
    status: "A",              ⟵  field: value
    groups: [ "news", "sports" ]  ⟵  field: value
}
```

# What is a collection

- A grouping of MongoDB **documents**.
- A **collection** is the equivalent of an RDBMS table.
- A **collection** exists within a single database.
- **Collections** do not enforce a schema.
- **Documents** within a collection can have different fields.
- **Collections** are grouped into **databases**.

Example: **db** = tweets, **collections** = greek_tweets, english_tweets, etc.

# Operations

- Run / Connect to mongo

- Import documents

- Insert documents

- Create queries

- Quick data lookups

- Indexing

# Connection

Studio 3T:

- Studio 3T -> Connect -> New connection -> (enter name) -> Save -> Connect
- Right click -> Add database -> (enter name) -> ok

Python:

```python
import pymongo

client = pymongo.MongoClient('localhost', 27017)
db = client['yelp']
```

# Import dataset

Studio 3T:

- Connect -> Select DB -> Import collections -> JSON -> Add sources ->
  Rename target collections -> Rename collections* -> Start import

* Rename collections to "restaurants", "reviews" and "users".

# Insert documents

Python:

```python
        # Find the coordinates for each restaurant and
        # save them to an external collection
        all_restaurants = find_all_restaurants()
        for restaurant in all_restaurants:
            json_obj = {
                'name': restaurant['name'],
                'business_id': restaurant['business_id'],
                'longitude': restaurant['longitude'],
                'latitude': restaurant['latitude']
            }
            insert_to_db(json_obj, 'restaurants_coordinates')
```

# Querying

Studio 3T:

- Right click collection -> Open Intellishell
- Examples:
  - db.restaurants.find( {} )
  - db.restaurants.find( {"neighborhood": "Downtown"} )
  - db.restaurants.find({ **$and**: [ {"neighborhood": "Southeast"}, {"city": "Las Vegas"} ]})

Python:

```
24        db['restaurants'].find({'neighborhood': neighborhood})
25        db['restaurants'].find_one({'business_id': restaurant_id})
26        db['reviews'].find({"$and": [{"business_id": restaurant_id}, {"stars": 5}]})
```

# Data lookups - For quick data checking

Studio 3T:

- Right click collection -> Open Intellishell
    - Ex: db.collection_name.find( { "Search_Field": "value" }, { "Field_to_display": 1 } )
    - db.restaurants.find({ "neighborhood": "Downtown" }, { "name": 1 })
    - db.restaurants.find({ "name": /.*pollos.*/ }, { "text": 1 } )
        - /.* is the regex equivalent for: "any single character, 0 or more times"

Python:

```
27      def find_reviews_that_contain_a_word(word):
28          return db['reviews'].find({'text': {'$regex': '.*' + word + '.*'}})
29
```

# Indexing

- Basic indexing
    - Speeds up queries on specific fields
    - Mostly used fields should be indexed
    - Index responsibly : Too many indices might slow down the database

- Text indexing
    - Faster text search queries on string content
    - Term lookup
    - Exact phrase lookup
    - Automatic relevance sorting

# How to - Text indexing

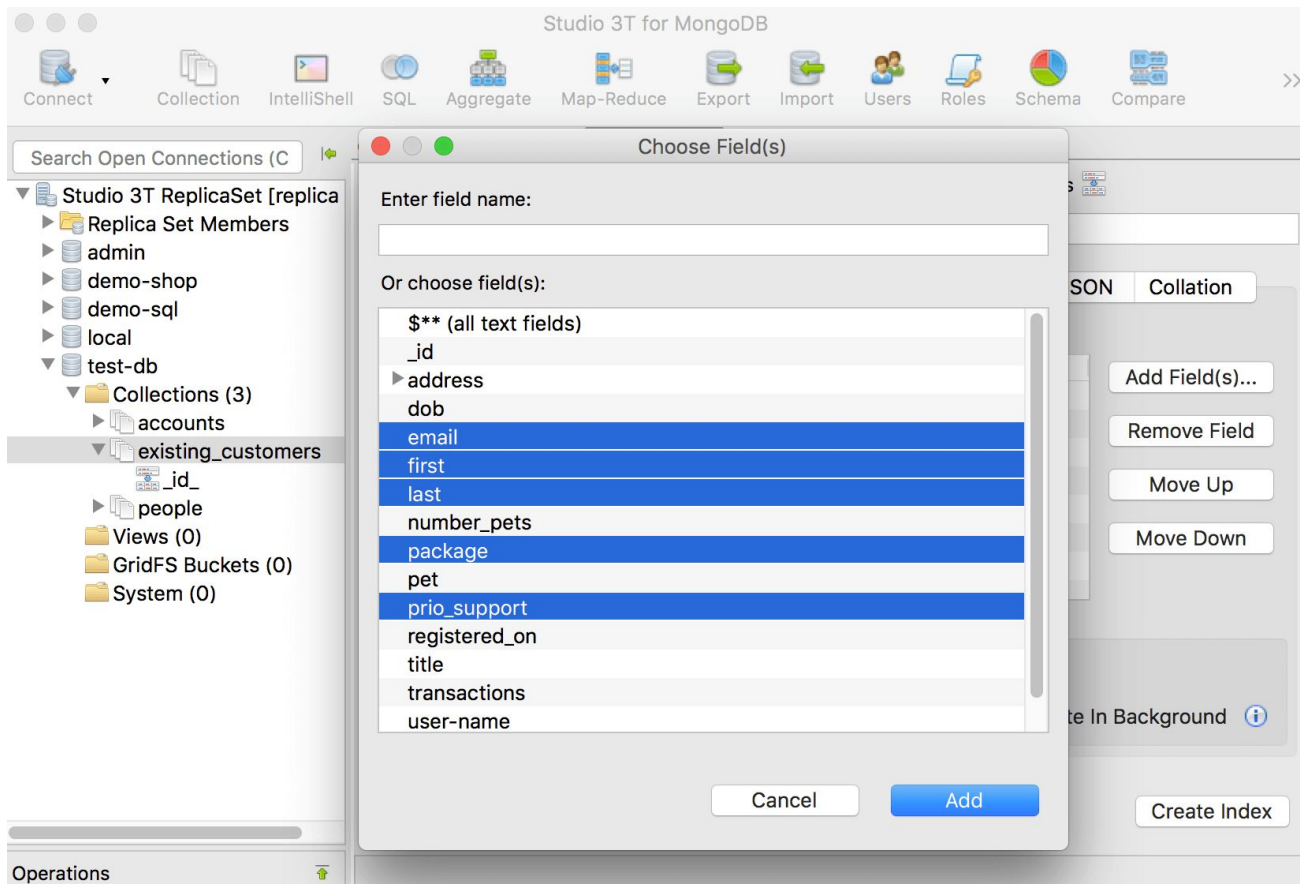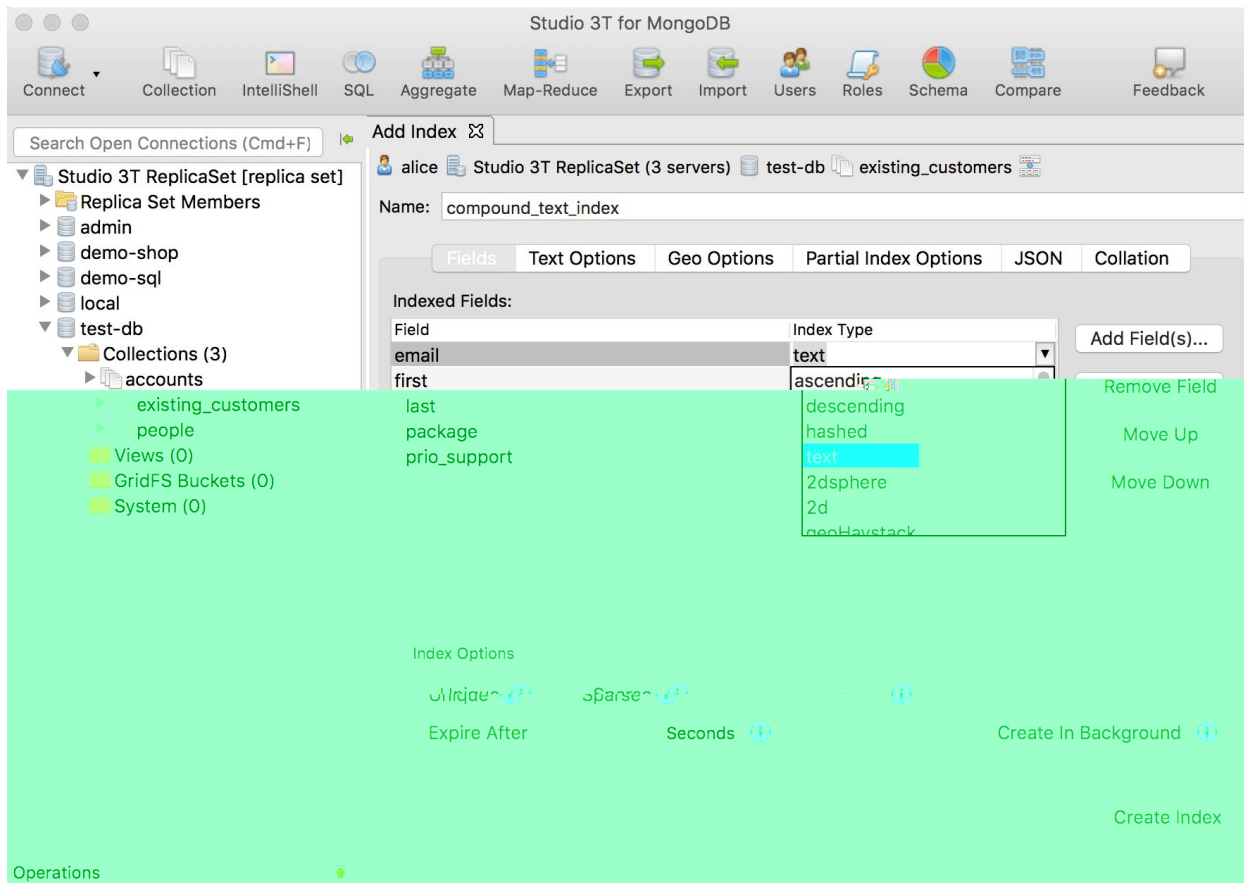# How to - Text indexing

# How to - Text indexing

# How to - Text indexing

# How to - Text indexing

# How to - Text indexing

# Querying with text indices

- db.stores.find( { $text: { $search: "java coffee shop" } } )
    - $text tokenizes the search string and performs a logical OR
    - Will search on all indexed fields
    - Results include a relevance score for each record

- db.stores.find( { $text: { $search: "java \"coffee shop\"" } } )
    - Will match exact phrases "java" OR "coffee shop"

- db.stores.find( { $text: { $search: "java shop -coffee" } } )
    - Will match ("java" OR "shop") AND NOT("coffee")

# Web app development



## Possibility #1

Static HTML pages. All results are hardcoded

## Possibility #2

## Possibility #3

Some backend language

23

A python-based microframework, suitable for small-scale applications.

How to create a simple Flask App:

**app.py**

```python
from flask import Flask, render_template

app = Flask(__name__, template_folder='views')

@app.route("/")
def home():
    uni_name = "Aristotle University of Thessaloniki"
    return render_template('some_html_file.html', uni_name=uni_name)

if __name__ == "__main__":
    app.run(debug=True, host='127.0.0.1', port=5110)
```

24

# some_html_file.html

```html
<!doctype html>
<html>
    <head>
        <meta charset="utf-8">
    </head>
    <body>
        <p>Welcome to your first Flask application!</p>
        <!-- Variable is passed, using the jinja2 templating engine -->
        <h2>{{ uni_name }}</h2>
    </body>
</html>
```
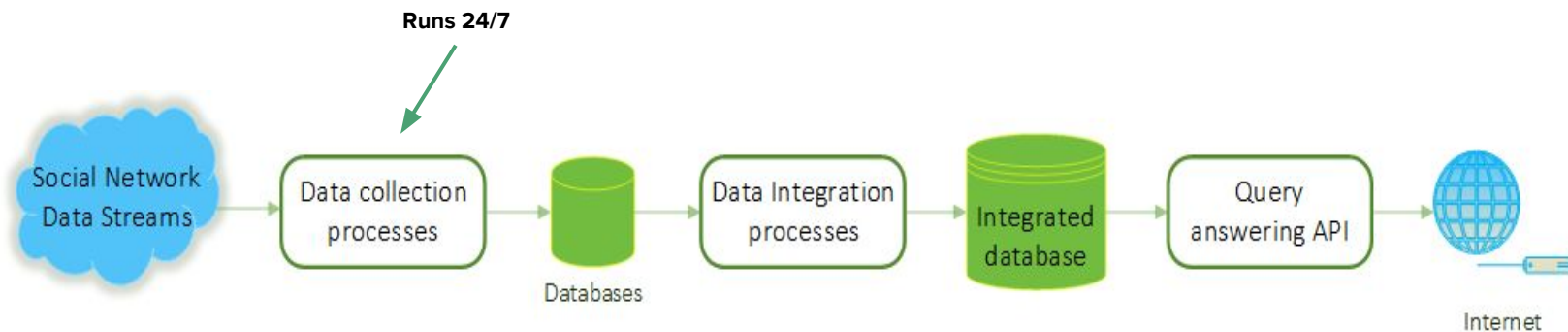
127.0.0.1:5110

← → C ⌂ ⓘ 127.0.0.1:5110

Welcome to your first Flask application!

## Aristotle University of Thessaloniki

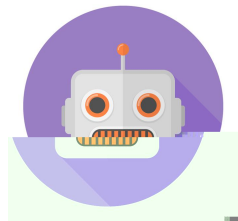# Diligent - A Social Media Data Integration platform

- Social media data explosion - Information Era
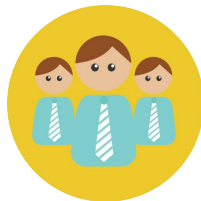    - All data and metadata are saved into databases
    - Including private information (fingerprint, face, political views, etc.)


- Various data sources can be Integrated
    - Extract valuable insights
    - Personalized advertising
    - Sentiment analysis
    - -> Build a complete image - Increase data value $$


- Machine Learning opportunities
    - Integrated data as input

# Diligent - System pipeline

**Runs 24/7**

Social Network Data Streams → Data collection processes → Databases → Data Integration processes → Integrated database → Query answering API → Internet
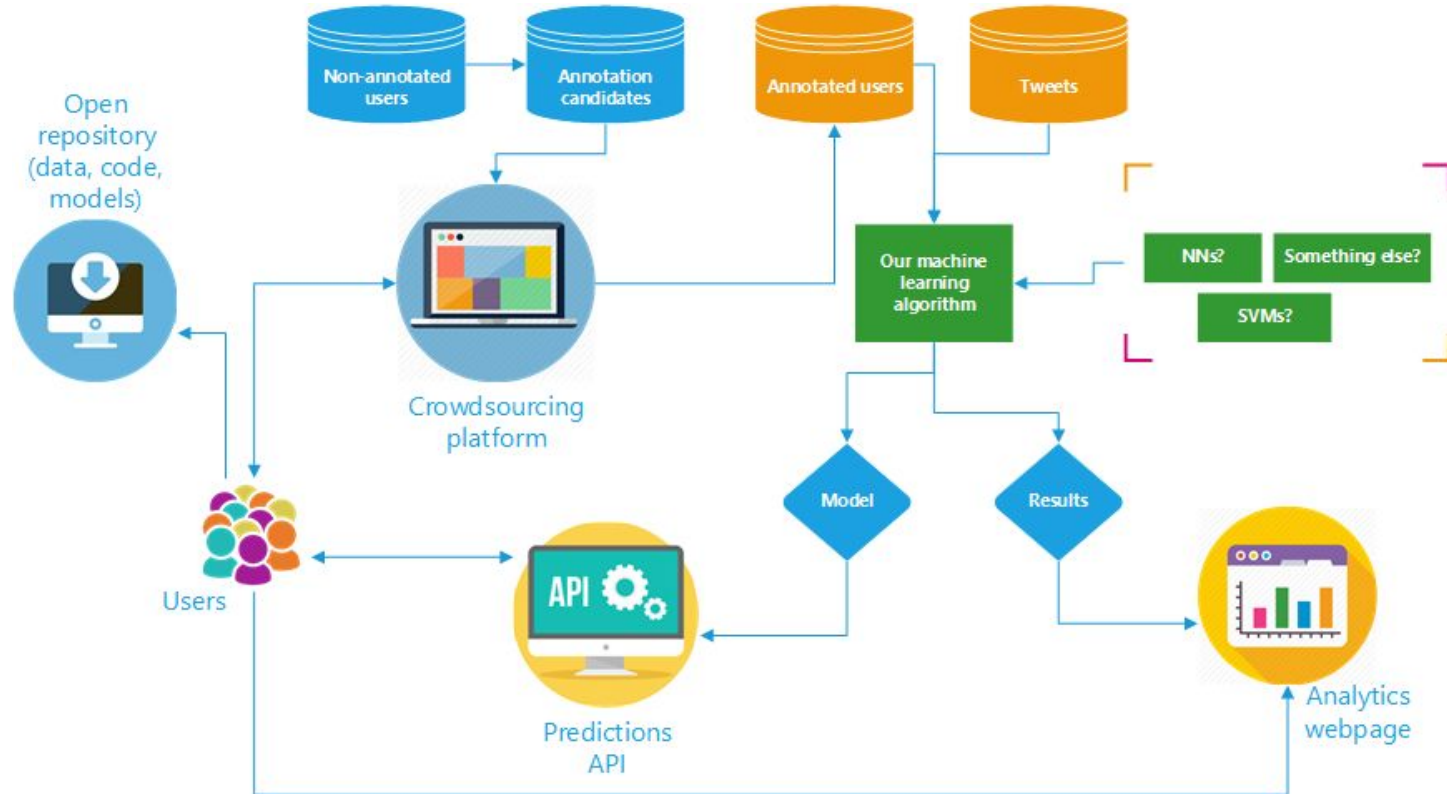
# Twitter bots detection

- Own a set of pre-annotated users as bots or humans
- Extract features from each user (focus on sentiment)
  - E.g. "num of tweets", "tweets entropy", "average sentiment polarity"
- Feed a machine learning algorithm all the features vectors
- Evaluate and extract model

# Full-scale architecture

# Resources

- [Dataset jsons](#)
- [Github link](#)
- [MongoDB](#)
- [Flask](#)
- [Flask-RESTful](#)
- [Python](#)