

## Pokročilé databázové technológie

Zadanie 6 – MongoDB

Marek Adamovič

## Obsah

1. Dátový model.....	3
2. Import .....	5
3. Dopyty .....	8

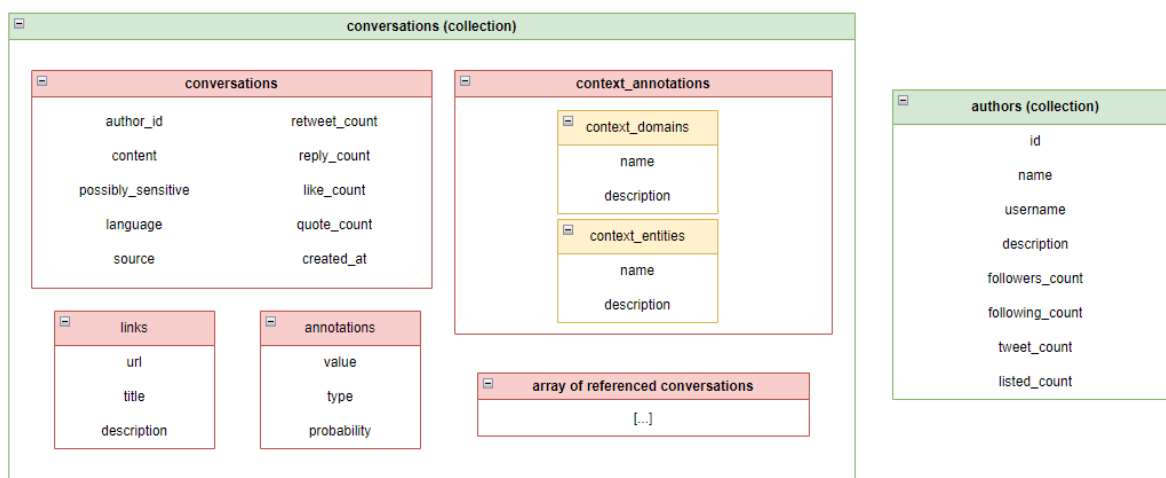
# 1. Dátový model

## Otázka:

Navrhnete dátový model (kolekcie a formát dokumentov) v MongoDB pre dataset tweetov, ktorý bude využívaný mobilnou aplikáciou, ktorá bude:

- Zobrazovať tweety jednotlivých používateľov vo forme feedov
- Zobrazovať jednotlivé tweety a ich retweets

## Odpoveď:



Obrázok 1 Dátový model Mongo

Pri tvorbe dátového modelu (obrázok č.1) sme zvolili prístup, v ktorom väčšina dát bude obsiahnutých v jednej collection (embedded prístup), a to v conversations. To hlavne z dôvodu, že tieto dáta sa v našej pomyslenej aplikácii nebudú (často) meniť a bez samotných tweetov by nemali veľký zmysel. To sa však netýka autorov-užívateľov, ktorých ukladáme (referenced prístup) do vlastnej collection authors. Je to najmä kvôli tomu, aby sme pri zmene údajov autora (čo je prípad, ktorý nastáva veľmi často) nemuseli meniť tieto údaje v každom jednom tweete, ktorý tento autor posol. Taktiež majú dáta autorov zmysel aj bez dát o tweetoch, čo je ďalší dôvod, prečo je dobré mať ich vo vlastnej collection. Čo sa týka referencovaných tweetov, tiež sme zvolili embedded prístup, keďže MongoDB neumožňuje referenced prístup v rámci jednej collection. Teda ak by sme chceli využiť referenced prístup, museli by sme mať vlastnú kolekciu pre referencované tweety. To by veľmi nedávalo zmysel, mať dve kolekcie na tweety, kvôli čomu sme zvolili embedded prístup (s tým, že referencované tweety nebudú mať len základné informácie o nich, aby nám napríklad nevznikali nekonečný embedding v rámci jednej collection). Na obrázku č.2,3 si vieme pozrieť príklad štruktúry importovaných dát.

```
_id: 224
name: "Dave Pell"
username: "davepell"
description: "Managing Editor, Internet."
followers_count: 60064
following_count: 1054
tweet_count: 54618
listed_count: 2016
```

Obrázok 2 Príklad štruktúry dát z collection authors

```
_id: 1496779955958603778
content: "RT @dlepeska: #Ukraina"
possibly_sensitive: false
language: "qht"
source: "Twitter for Android"
retweet_count: 2
reply_count: 0
like_count: 0
quote_count: 0
created_at: "2022-02-24T10:32:04+01:00"
author_id: 182981546
  links: Array
  annotations: Array
  context_domains: Array
    0: Object
      id: 123
      name: "Ongoing News Story"
      description: "Ongoing News Stories like 'Brexit'"
    1: Object
      id: 156
      name: "Cities"
      description: "Cities around the world"
  context_entities: Array
  conversation_references: Array
    0: Object
      id: 96560
      type: "retweeted"
      parent_conversation_author_id: 19939989
      parent_conversation_author_name: "David Lepeska"
      parent_conversation_author_username: "dlepeska"
      parent_conversation_id: 1496779491594874881
      parent_conversation_content: "#Ukraina https://t.co/5SwNvpqIhQ"
      parent_conversation_created_at: "2022-02-24T10:30:14+01:00"
```

Obrázok 3 Príklad štruktúry dát z collection conversations

## 2. Import

### Otázka:

Nainštalujte alebo využite online inštanciu MongoDB servera, do ktorého importujte všetky tweets (a s nimi spojené data – anotácie, referencie, odkazy a informácie o kontexte) zo dňa 24.02.2022.

### Odpoveď:

Na túto úlohu sme využili online inštanciu MongoDB servera. Najskôr sme exportovali dáta v požadovanej štruktúre pomocou query (obrázok č.4 -> authors collection, obrázok č.5 -> conversations collection) z postgres databázy.

```
1 COPY(
2   EXPLAIN ANALYZE SELECT json_build_object(
3     '_id', author.id,
4     'name', author.name,
5     'username', author.username,
6     'description', author.description,
7     'followers_count', author.followers_count,
8     'following_count', author.following_count,
9     'tweet_count', author.tweet_count,
10    'listed_count', author.listed_count
11  )
12  FROM conversations c
13  JOIN authors author ON author.id = c.author_id
14  WHERE c.created_at >= '2022-02-24 00:00:00' AND c.created_at < '2022-02-25 00:00:00'
15  GROUP BY author.id
16  --LIMIT 50
17 ) TO 'D:\skola2022_2023\PDT\zadanie6\authors.json' WITH (FORMAT CSV, QUOTE ' ');
18
19 -- Create date index
20 --CREATE INDEX conversation_created ON conversations USING btree(created_at);
```

Obrázok 4 Export dát potrebných pre author collection

```
1 COPY(
2   SELECT json_build_object(
3     '_id', c.id,
4     'content', c.content,
5     'possibly_sensitive', c.possibly_sensitive,
6     'language', c.language,
7     'source', c.source,
8     'retweet_count', c.retweet_count,
9     'reply_count', c.reply_count,
10    'like_count', c.like_count,
11    'quote_count', c.quote_count,
12    'created_at', c.created_at,
13    'author_id', c.author_id,
14    'links', l_l,
15    'annotations', a_l,
16    'context_domains', cd_l, -- sometimes null
17    'context_entities', ce_l, -- sometimes null
18    'conversation_references', cr_l
19  )
20 FROM conversations c
21 LEFT JOIN LATERAL(
22   SELECT COALESCE(json_agg(json_build_object('id', l.id, 'url', l.url, 'title', l.title, 'description', l.description)), '[]'::json) l_l
23   FROM links l
24   WHERE l.conversation_id = c.id
25 ) AS l_l ON true
26 LEFT JOIN LATERAL(
27   SELECT COALESCE(json_agg(json_build_object('id', a.id, 'value', a.value, 'type', a.type, 'probability', a.probability)), '[]'::json) a_l
28   FROM annotations a
29   WHERE a.conversation_id = c.id
30 ) AS a_l ON true
31 LEFT JOIN LATERAL(
32   SELECT COALESCE(json_agg(json_build_object('id', cd.id, 'name', cd.name, 'description', cd.description)), '[]'::json) cd_l
33   FROM context_domains cd
34   JOIN context_annotations ca ON ca.context_domain_id = cd.id AND ca.conversation_id = c.id
35 ) AS cd_l ON true
36 LEFT JOIN LATERAL(
37   SELECT COALESCE(json_agg(json_build_object('id', ce.id, 'name', ce.name, 'description', ce.description)), '[]'::json) ce_l
38   FROM context_entities ce
39   JOIN context_annotations ca ON ca.context_domain_id = ce.id AND ca.conversation_id = c.id
40 ) AS ce_l ON true
41 LEFT JOIN LATERAL(
42   SELECT COALESCE(json_agg(json_build_object(
43     'id', cr.id,
44     'type', cr.type,
45     'parent_conversation_author_id', author2.id,
46     'parent_conversation_author_name', author2.name,
47     'parent_conversation_author_username', author2.username,
48     'parent_conversation_id', c2.id,
49     'parent_conversation_content', c2.content,
50     'parent_conversation_created_at', c2.created_at
51   )
52   ), '[]'::json) cr_l
53   FROM conversation_references cr
54   JOIN conversations c2 ON c2.id = cr.parent_id
55   JOIN authors author2 ON author2.id = c2.author_id
56   WHERE (c2.created_at >= '2022-02-24 00:00:00' AND c2.created_at < '2022-02-25 00:00:00') AND cr.conversation_id = c.id
57 ) AS cr_l ON true
58 WHERE c.created_at >= '2022-02-24 00:00:00' AND c.created_at < '2022-02-25 00:00:00'
59 ) TO 'D:\skola2022_2023\PD\zadanie6\conversations.json' WITH (FORMAT CSV, QUOTE '');
```

Obrázok 5 Export dát potrebných pre conversations collection

Ked' sme dokončili export dát z postgresu, pomocou python scriptu (obrázok č.6) sme nahrali dáta do online MongoDB clusteru. Stav databázy po importe si vieme pozrieť na obrázku č.7.

# Slovenská Technická Univerzita v Bratislave

## Fakulta Informatiky a Informačných Technológií

```
mongo.py X
mongo.py > ...
1  from pymongo import MongoClient
2  import json
3  import time
4
5  BATCHSIZE = 5000
6  client = MongoClient('mongodb+srv://admin:admin@cluster0.fqbgwci.mongodb.net/test')
7
8  db = client['zadanie6']
9
10
11 def import_conversations():
12     collection = db['conversations']
13     with open('conversations.json', 'r', encoding='utf8') as f:
14         x = 0
15         start = time.time()
16         documents = []
17         for line in f:
18             x += 1
19             documents.append(json.loads(line))
20
21             if len(documents) == BATCHSIZE:
22                 collection.insert_many(documents)
23                 print(x)
24                 documents = []
25
26     #send final data
27     collection.insert_many(documents)
28     print('Final time:')
29     print(time.time() - start)
30
31
32 def import_authors():
33     collection = db['authors']
34     with open('authors.json', 'r', encoding='utf8') as f:
35         x = 0
36         start = time.time()
37         documents = []
38         for line in f:
39             x += 1
40             documents.append(json.loads(line))
41
42             if len(documents) == BATCHSIZE:
43                 collection.insert_many(documents)
44                 print(x)
45                 documents = []
46
47     #send final data
48     collection.insert_many(documents)
49     print(x)
50     print('Final time:')
51     print(time.time() - start)
52     return
53
54 import_conversations()
55 import_authors()
```

Obrázok 6 Python script pre import dát na MongoDB cluster

zadanie6							
LOGICAL DATA SIZE: 2,45GB		STORAGE SIZE: 917,47MB	INDEX SIZE: 1191MB	TOTAL COLLECTIONS: 2		CREATE COLLECTION	
Collection Name	Documents	Logical Data Size	Avg Document Size	Storage Size	Indexes	Index Size	Avg Index Size
authors	981080	228,19MB	244B	151,81MB	1	34,56MB	34,56MB
conversations	1840128	2,23GB	1,27KB	765,66MB	1	84,54MB	84,54MB

Obrázok 7 Stav databázy po importe

### 3. Dopyty

#### Otázka:

Napište dotaz, ktorý nad importovanou databázou:

- Vypíše posledných 10 tweetov pre autora, ktorý má username Newnews\_eu
- Vypíše posledných 10 retweetov pre tweet, ktorý má id 1496830803736731649.

#### Odpoveď:

Pre úlohu a) sme vytvorili dopyt (obrázok č.8), ktorý najskôr spojí collection conversation s collection authors na základe autorovho id (podobne ako join pri relačných databázach). Následne pomocou unwind nahradí pole autorov (v tomto prípade je autor vždy len jeden, takže je to skôr kozmetická úprava, vďaka ktorej vieme pracovať priamo so záznamom a nie s poľom jedného prvku) prvkom autora z tohto poľa. Pomocou match ponecháme záznamy daného autora (Newnews\_eu) a v závere usporiadame výsledky podľa dátumu created\_at a ponecháme 10 výsledkov.

```
first_query = '''[
    {
        "$lookup": {
            "from": "authors",
            "localField": "author_id",
            "foreignField": "_id",
            "as": "author"
        }
    },
    {"$unwind": "$author"},
    {"$match": {"author.username": "Newnews_eu"}},
    {"$sort": {"created_at": -1}},
    {"$limit": 10}
]'''
results1 = collection.aggregate(json.loads(first_query))
for r in results1:
    print(r)
```

Obrázok 8 Query pre vypísanie posledných 10 tweetov autora Newnews\_eu

Pre úlohu b) sme vytvorili dopyt (obrázok č.9), ktorý rozbalí pole s referenciami pomocou unwind a pozrie sa, či v týchto referenciách je retweet konverzácie s požadovaným id. Následne tieto výsledky zoradí podľa času created\_at a vráti prvých 10.



```
second_query = '''[
    {"$unwind": "$conversation_references"},
    {"$match": {
        "$and": [
            {"conversation_references.parent_conversation_id": 1496830803736731649},
            {"conversation_references.type": "retweeted"}
        ]
    }},
    {"$sort": {"created_at": -1}},
    {"$limit": 10}
]'''
results2 = collection.aggregate(json.loads(second_query))
for r in results2:
    print(r)
```

Obrázok 9 Posledných 10 retweetov tweetu s ID 1496830803736731649