

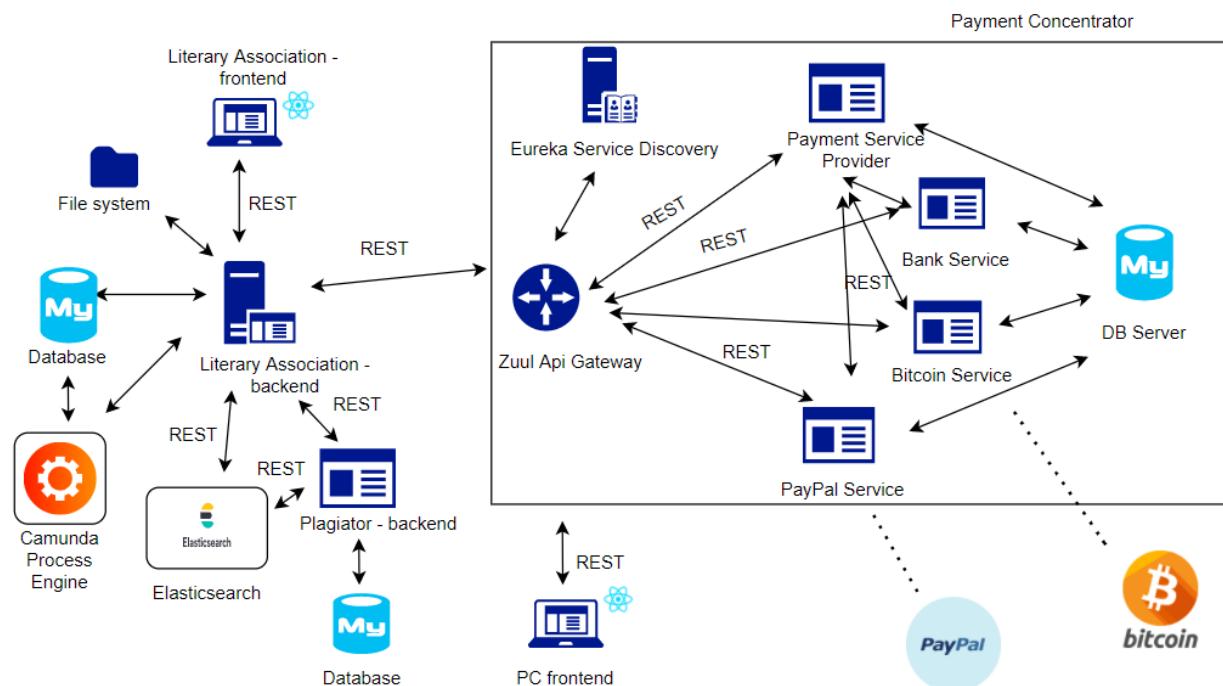
# Upravljanje digitalnim dokumentima

## Kontrolna tačka – Marina Bartulov E2 81/2020

### 1. Arhitektura aplikacija i komunikacija

Literarno udruženje je realizovano kao monolitna aplikacija koja se sastoji od frontend i backend aplikacije koji komuniciraju putem REST-a. Backend aplikacija je klasična Spring Boot aplikacija koja se sastoji od 3 sloja: sloj kontrolera, servisni sloj i sloj podataka koji ostvaruje konekciju ka bazi. Literarno udruženje komunicira sa drugim sistemom koji predstavlja koncentrator plaćanja. On je realizovan kao mikroservisna aplikacija i vrši usluge plaćanja za klijente kao što je Literarno udruženje. Koncentrator plaćanja i Literarno udruženje komuniciraju takođe putem REST-a. Backend obe aplikacije je Spring Boot, a frontend aplikacije su rađene u React-u. Ceo ovaj sistem prikazan je na Slici 1.

Knjige koje se mogu pretraživati, tj. PDF fajlovi će biti čuvani u lokalnom fajl sistemu. Podaci o entitetima sistema će se čuvati u relacionoj bazi, koja će biti MySQL. Određeni podaci o entitetima koji se pretražuju biće ili čuvani samo u bazi, ili samo u indeksnoj strukturi, ili na oba mesta. Sloj podataka kojeg čine repozitorijumi će ostvarivati komunikaciju sa bazom. Literarno udruženja će putem REST-a da komunicira sa Elasticsearch-om prilikom indeksiranja i pretraživanja dokumenata, a takođe će da komunicira i sa sistemom za detekciju potencijalnih plagijarizama koji predstavlja posebnu web aplikaciju, i ta komunikacija će takođe da se ostvaruje putem REST-a.



Slika 1. Arhitektura sistema

## 2. Pokretanje i konfiguracija ElasticSearch-a

Elasticsearch možemo da pokrenemo tako što skinemo zip fajl sa zvaničnog sajta i otpakujemo, i zatim odатle pokrenemo. Drugi način je da ga pokrećemo zajedno sa Spring Boot aplikacijom. Zbog kasnije jednostavnosti instaliranja plugin-a ja sam se odlučila na prvi način. Skinula sam verziju 7.4.0 jer je ona kompatibilna sa pluginom za srpski jezik. Kako je Elasticsearch RESTful search engine komunikacija između backend aplikacije i njega će se odvijati takođe preko REST-a.

Da bi se pokrenuo Elasticsearch potrebno je pozicionirati se u bin folder otpakovanog zip fajla i ukucati komandu `elasticsearch.bat`. Elasticsearch se pokrene na portu 9200. Da bi se proverilo da li je Elasticsearch pokrenut može se poslati GET zahtev na <http://localhost:9200> i dobiti odgovor kao na Slici 2.

```
1  {
2      "name": "MARINA-PC",
3      "cluster_name": "elasticsearch",
4      "cluster_uuid": "KC_88jLKSZWNjKXBocP-4w",
5      "version": {
6          "number": "7.4.0",
7          "build_flavor": "default",
8          "build_type": "zip",
9          "build_hash": "22e1767283e61a198cb4db791ea66e3f11ab9910",
10         "build_date": "2019-09-27T08:36:48.569419Z",
11         "build_snapshot": false,
12         "lucene_version": "8.2.0",
13         "minimum_wire_compatibility_version": "6.8.0",
14         "minimum_index_compatibility_version": "6.0.0-beta1"
15     },
16     "tagline": "You Know, for Search"
17 }
```

Slika 2. Provera da li je Elasticsearch pokrenut

Elasticsearch ima dobre deafult opcije i zahteva jako malo konfiguracije. Što se tiče konfiguracije Elasticsearch-a ona se može obavljati u `elasticsearch.yml` fajlu koji se nalazi u config folderu otpakovanog zip fajla. Tu se mogu konfigurisati neke naprednije opcije, a za naše potrebe je za sada dovoljno da ostane default. Tu se mogu npr. specificirati putanje gde Elasticsearch upisuje podatke i logove. Polja kojima se to konfiguriše jesu path.data i path.logs i preporuka je da u produkciji ona budu podešena na lokacije izvan Elasticsearch foldera.

U Spring Boot projektu za komunikaciju sa Elasticsearch-om ću koristiti `ElasticsearchRepository` za indeksiranje, a za upite `ElasticsearchTemplate` jer on omogućava dinamički sažetak koji nam je potreban u projektu. Dependecy za to je prikazan na Slici 3.

```
<dependency>
    <groupId>org.springframework.boot</groupId>
    <artifactId>spring-boot-starter-data-elasticsearch</artifactId>
</dependency>
```

Slika 3. Dependency za komunikaciju sa Elasticsearch-om

### 3. Serbian Analyzer plugin

Serbian Analyzer plugin se koristi za analiziranje i preprocesiranje unetog teksta na srpskom jeziku. Za njegovo instaliranje potrebno je da imamo instaliran build alat Gradle. Potrebno je da se pozicioniramo u root direktorijum od Serbian Analyzer plugin projekta i izvršimo komandu:

```
./gradlew clean build
```

Ova komanda će da kreira zip fajl i da ga smesti u build/distributions folder. Nakon toga da bismo podesili da ovaj plugin radi sa Elasticsearch-om, potrebno je da se pozicioniramo u njegov /bin direktorijum i da izvršimo komandu:

```
./elasticsearch-plugin install file:<absolute path of distribution archive>
```

Na različite načine se može reći Elasticsearch-u koji analyzer da gleda prilikom indeksiranja i pretraživanja. Za indeksiranje moguće je specificirati analyzer za polje ili default analyzer za indeks. Za pretraživanje se može još specificirati i za query. Najčešće specificiranje različitog analyzer-a za pretraživanje nije potrebno. U većini slučajeva jednostavan pristup je najbolji, a to je da se specificira analyzer za svako tekst polje posebno. Moj plan je da koristim taj pristup. Za svako tekst polje se kaže da je "analyzer" = "serbian" i "search\_analyzer" = "serbian" (Slika 4). Mada nije čak ni potrebno navoditi za search, jer ukoliko se on ne navede podrazumeva se da je isti kao i za indeksiranje.

The screenshot shows a Postman request configuration for a PUT request to `http://localhost:9200/my_index`. The 'Body' tab is selected, showing the following JSON payload:

```
1 {  
2   "mappings": {  
3     "properties": {  
4       "title": {  
5         "type": "text",  
6         "analyzer": "serbian",  
7         "search_analyzer": "serbian"  
8       }  
9     }  
10   }  
11 }
```

Slika 4. Način specificiranja analyzer-a

## 4. Book index unit

Što se tiče knjiga neki podaci će biti čuvani samo u bazi, jer neće biti korišćeni prilikom pretraživanja, niti će biti bitni da se prikažu u sklopu rezultata. Neki od tih podataka su ISBN, mesto i godina izdavanja i broj stranica. Dok će se u indeksu strukturi pored podataka po kojima će se pretraživati knjige, nalaziti i oni podaci koji će biti bitni da se prikažu i na neki način koriste u rezultatu pretrage, a to su podatak da li je knjiga u open access mode-u, kao i putanja do pdf fajla kako bi se ta knjiga mogla download ukoliko jeste u open access mode-u, zatim sinopsis i cena. Takođe će se čuvati i sadržaj pdf fajla, kojeg neće biti u bazi. Na Slici 5 sam prikazala kako bi izgledalo kreiranje indeksa za knjige.



The screenshot shows a Postman request for creating a book index. The method is PUT, and the URL is `http://localhost:9200/my_books_index`. The 'Body' tab is selected, showing a JSON mapping configuration. The JSON code is as follows:

```
1  {
2     "mappings": {
3         "properties": {
4             "title": {
5                 "type": "text",
6                 "analyzer": "serbian"
7             },
8             "writers": {
9                 "type": "text",
10                "analyzer": "serbian"
11            },
12            "content": {
13                "type": "text",
14                "analyzer": "serbian"
15            },
16            "genres": {
17                "type": "text",
18                "analyzer": "serbian"
19            },
20            "openAccess": {
21                "type": "boolean",
22                "index": false
23            },
24            "synopsis": {
25                "type": "text",
26                "index": false
27            },
28            "pdfPath": {
29                "type": "text",
30                "index": false
31            },
32            "price": {
33                "type": "double",
34                "index": false
35            }
36        }
37    }
38 }
```

Slika 5. Kreiranje indeksa za knjige

Na sledećoj slici (Slika 6) je prikazano dodavanje jedne knjige u indeks i odgovor koji se dobija.

The screenshot shows a POST request to `http://localhost:9200/my_books_index/_doc/1`. The request body contains a JSON object representing a book:

```
1 [ {  
2   "title": "Moja prva knjiga",  
3   "writers": "Marina Bartulov",  
4   "content": "Ovo je sadržaj moje prve knjige...",  
5   "genres": "psiologija,filozofija",  
6   "openAccess": true,  
7   "pdfPath": "Neka putanja...",  
8   "price": 700.00  
9 } ]
```

The response body shows the result of the indexing operation:

```
1 {  
2   "_index": "my_books_index",  
3   "_type": "_doc",  
4   "_id": "1",  
5   "_version": 1,  
6   "result": "created",  
7   "_shards": {  
8     "total": 2,  
9     "successful": 1,  
10    "failed": 0  
11  },  
12  "_seq_no": 0,  
13  "_primary_term": 1  
14 }
```

Slika 6. Dodavanje nove knjige u indeks

Na sledećoj slici (Slika 7) je prikazan jedan primer zadavanja upita nad ovom indeksom strukturom.

The screenshot shows a GET request to `http://localhost:9200/my_books_index/_search`. The request body contains a search query in JSON:

```
1 {  
2   "query": {  
3     "bool": {  
4       "must": [  
5         {  
6           "match": {  
7             "content": "sadrzaj"  
8           }  
9         },  
10        {  
11          "match": {  
12            "writers": "marina"  
13          }  
14        },  
15        {  
16          "match": {  
17            "genres": "filozofija"  
18          }  
19        }  
20      ]  
21    }  
22  }
```

The response body shows the search results, including the total count, shards information, and the hit details:

```
1 {  
2   "took": 118,  
3   "timed_out": false,  
4   "_shards": {  
5     "total": 1,  
6     "successful": 1,  
7     "skipped": 0,  
8     "failed": 0  
9   },  
10  "hits": {  
11    "total": {  
12      "value": 1,  
13      "relation": "eq"  
14    },  
15    "max_score": 0.8630463,  
16    "hits": [  
17      {  
18        "_index": "my_books_index",  
19        "_type": "_doc",  
20        "_id": "1",  
21        "_score": 0.8630463,  
22        "_source": {  
23          "title": "Moja prva knjiga",  
24          "writers": "Marina Bartulov",  
25          "content": "Ovo je sadržaj moje prve knjige...",  
26          "genres": "psiologija,filozofija",  
27          "openAccess": true,  
28          "pdfPath": "Neka putanja...",  
29          "price": 700.00  
30        }  
31      }  
32    }  
33  }
```

Slika 7. Primer upita za knjige

## 5. Geoprostorna pretraga

Geoprostorna pretraga u ElasticSearch-u se vrši pomoću geo-distance query koji u rezultat pretrage uključuje samo ono što postoji unutar određene udaljenosti od zadate geo tačke. Prilikom indeksiranja kao tip podatka se koristi geo\_point. Primer kako bih kreirala indeks za beta-čitaoca koji treba da se geoprostorno pretražuju je prikazan na Slici 8. Id od beta-čitaoca i ime se neće indeksirati, nego će se oni koristiti u prikazu rezultata. Što se tiče žanrova za sada će to biti string koji će sadržati id-jeve žanrova koje je beta-čitalac izabrao, ali će razmisliti još da li je ovo najbolje rešenje.

Podaci za beta-čitaoce kao što su username, password, grad, država, email i kazneni poeni će biti čuvani samo u bazi jer oni nisu važni za pretragu i prikaz rezultata pretrage.

Na Slici 9 sam prikazala kako bih dodala novog beta-čitaoca u postojeći indeks.

Na Slici 10 kako bih napravila upit, a na Slici 11 rezultat upita. U upitu govorim da se pronađu svi oni koji ne nalaze okolini od 100 km od zadatke tačke.

The screenshot shows a Postman request configuration. The method is set to PUT, and the URL is `http://localhost:9200/my_beta_readers`. The 'Headers' tab is selected, showing 9 headers. The 'Body' tab is also selected, indicating raw JSON data. The body contains the following JSON mapping:

```
1 {
2   "mappings": {
3     "properties": {
4       "location": {
5         "type": "geo_point"
6       },
7       "betaReaderId": {
8         "type": "long",
9         "index": false
10      },
11      "name": {
12        "type": "text",
13        "index": false
14      },
15      "genres": {
16        "type": "text"
17      }
18    }
19  }
20 }
```

Slika 8. Indeksiranje beta-čitala

PUT http://localhost:9200/my\_beta\_readers/\_doc/1

Params Authorization Headers (9) **Body** Pre-request Script Tests Settings

none form-data x-www-form-urlencoded raw binary GraphQL **JSON**

```

1  {
2   "location": {
3     "lat": 40.12,
4     "lon": -71.34
5   },
6   "genres": "1,2,3",
7   "name": "Marina Bartulov",
8   "betaReaderId": 1
9 }
--
```

Slika 9. Dodavanje novog beta-čitaoca

GET http://localhost:9200/my\_beta\_readers/\_search

Params Authorization Headers (9) **Body** Pre-reqes

none form-data x-www-form-urlencoded raw

```

1  {
2   "query": {
3     "bool": {
4       "must": {
5         "match_all": {}
6       },
7       "filter": {
8         "bool": {
9           "must_not": {
10            "geo_distance": {
11              "distance": "100km",
12              "location": {
13                "lat": 50,
14                "lon": -70
15              }
16            }
17          }
18        }
19      }
20    }
21  }
22 }
```

GET http://localhost:9200/my\_beta\_readers/\_search

Params Authorization Headers (9) Body Pre-request Scr

Body Cookies Headers (3) Test Results

Pretty Raw Preview Visualize JSON

```

10  "hits": {
11   "total": {
12     "value": 1,
13     "relation": "eq"
14   },
15   "max_score": 1.0,
16   "hits": [
17     {
18       "_index": "my_beta_readers",
19       "_type": "_doc",
20       "_id": "1",
21       "_score": 1.0,
22       "_source": {
23         "location": {
24           "lat": 40.12,
25           "lon": -71.34
26         },
27         "genres": "1,2,3",
28         "name": "Marina Bartulov",
29         "betaReaderId": 1
30       }
31     }
32   ]
33 }
```

Slika 11. Upit za geoprostornu pretragu

Slika 10. Rezultat pretrage

## 6. Sistem za detekciju potencijalnih plagijarizama

Sistem za detekciju potencijalnih plagijarizama predstavlja posebnu web aplikaciju, koja se sastoji od Spring Boot backend aplikacije i Angular frontend aplikacije. Za potrebe projekta jaću koristiti samo backend aplikaciju sa kojom će backend aplikacija Literarnog udruženja da komunicira preko REST-a. Pre pokretanja Plagiatora neophodno je podesiti konekciju sa MySql bazom i promeniti mejl za slanje mejlova u application.properties fajlu.

Da bi se koristio Plagiator neophodno se prvo registrovati i aktivirati nalog (Slika 12 i Slika 13)

The screenshot shows the Postman interface for a POST request to `http://localhost:8080/api/signup`. The 'Body' tab is selected, showing a JSON payload:

```
1 {
2     "name": "Marina",
3     "lastName": "Bartulov",
4     "password": "marinabartulov123",
5     "repeatedPassword": "marinabartulov123",
6     "email": "bartulovmarina@gmail.com",
7     "phoneNumber": "0637253077"
8 }
```

Below the body, the status is shown as `201 Created`.

Slika 12. Registracija na Plagiator

The screenshot shows the Postman interface for a GET request to `http://localhost:8080/api/registration/activate/3`. The 'Query Params' tab is selected, showing a table:

KEY	VALUE	DESCRIPTION
Key	Value	Description

Below the table, the status is shown as `200 OK`.

Slika 13. Aktivacija naloga

Nakon toga da bi se koristio sistem neophodno je da se korisnik uloguje. Pošto će backend Literarnog udruženja da komunicira sa ovom aplikacijom i šalje zahteve, prepostavljam da će morati da registrujem moju aplikaciju kao jednog od korisnika i prilikom slanja poziva moraće prvo da se šalje zahtev za logovanje na sistem kako bi se dobio jwt token (Slika 14).

The screenshot shows the Postman interface with a successful API call. The URL is `http://localhost:8080/api/login`. The 'Body' tab is selected, containing the following JSON payload:

```
1 {  
2     "password": "marinabartulov123",  
3     "email": "bartulovmarinal@gmail.com"  
4 }  
5  
6
```

The response details are as follows:

- Status: 200 OK
- Time: 1408 ms
- Size: 508 B
- Raw response content: `"eyJhbGciOiJIUzI1NiJ9.  
eyJzdWIiOiJYJ0dlxvdm1hcmluYTFAZ21halwuY29tIiwiYXV0aCI6W3siYXV0aG9yaXR5IjoiUk9MRV9MT0dHRUQifV0sIm1hdCI6MTYwODU4NzczlISwiZXhwIjoxIjA4NjczMTM1fQ.  
MxhfWVgxpe460pd61wNUecU2xB5FL_eyhjNouSyTU"`

*Slika 14. Logovanje na sistem*

Moja aplikacija će da koristi endpoint za upload fajlova i dobijanje odgovora koji sadrži listu sličnih radova. Ovaj Plagiator ne govori da li je nešto plagijat ili ne, nego čovek to mora da odluči sam na osnovu dobijenih sličnih radova. Ukoliko želim da upload-ujem fajl na sistem, šaljem POST zahtev na url <http://localhost:8080/api/file/upload/new> i podatak koji je neophodno poslati je file (Multipart), koji predstavlja knjigu. Primer jednog takvog zahteva prikazan je na Slici 15.

The screenshot shows a Postman request to `http://localhost:8080/api/file/upload/new` using the `POST` method. The `Body` tab is selected, showing a `form-data` entry named `file` containing the file `KT1.pdf`. The response body is displayed in JSON format, showing a document structure with fields like `similarPapers`, `uploadedPaper`, and `user`. A red arrow points to the `similarPapers` field.

```
21 "similarPapers": [],
22 "uploadedPaper": {
23     "id": 1,
24     "title": "KT1.pdf",
25     "file": null,
26     "pathForPDF": "\\\opt\\\\plagiator\\\\files\\\\2020_12_21_23_51_24-KT1.pdf",
27     "searchHits": 0.0,
28     "similarProcent": 0.0,
29     "user": {
30         "id": 3,
31         "name": "Marina",
32         "lastName": "Bartulov",
33         "email": "bartulovmarina@gmail.com",
34         "phoneNumber": "0637253077",
35         "active": true,
36         "role": {
37             "id": 1,
38             "userType": "ROLE_LOGGED"
39         }
40     }
41 }
```

Slika 15. Upload fajla

Na prethodnoj slici se vidi da kada upload-ujem prvi fajl, da nema sličnih fajlova u odgovoru. Međutim kada ponovo pošaljem takav zahtev sa drugim fajлом dobijam odgovor kao što je na Slici 16.

POST http://localhost:8080/api/file/upload/new

Params Authorization Headers (10) Body Pre-request Script Tests Settings

none form-data x-www-form-urlencoded raw binary GraphQL

KEY	VALUE	DESCRIPTION	***
<input checked="" type="checkbox"/> file	KT2.pdf X		

Body Cookies Headers (9) Test Results

Pretty Raw Preview Visualize JSON ↻

```
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
```

„similarPapers“: [ { „id“: 1, „title“: "KT1.pdf", „file“: null, „pathForPDF“: null, „searchHits“: 7.580690383911133, „similarProcent“: 0.07867288524551223, „user“: { „id“: 3, „name“: "Marina", „lastName“: "Bartulov", „email“: "bartulovmarina@gmail.com", „phoneNumber“: "0637253077", „active“: true, „...“: ... } }

Status: 200 OK Time: 2.04 s Size: 1.34 KB Save

Slika 16. Upload fajla

Sada se u similarPapers nalazi prethodno upload-ovan fajl i nalazi se podatak o tome koliko su slični (similarProcent).