

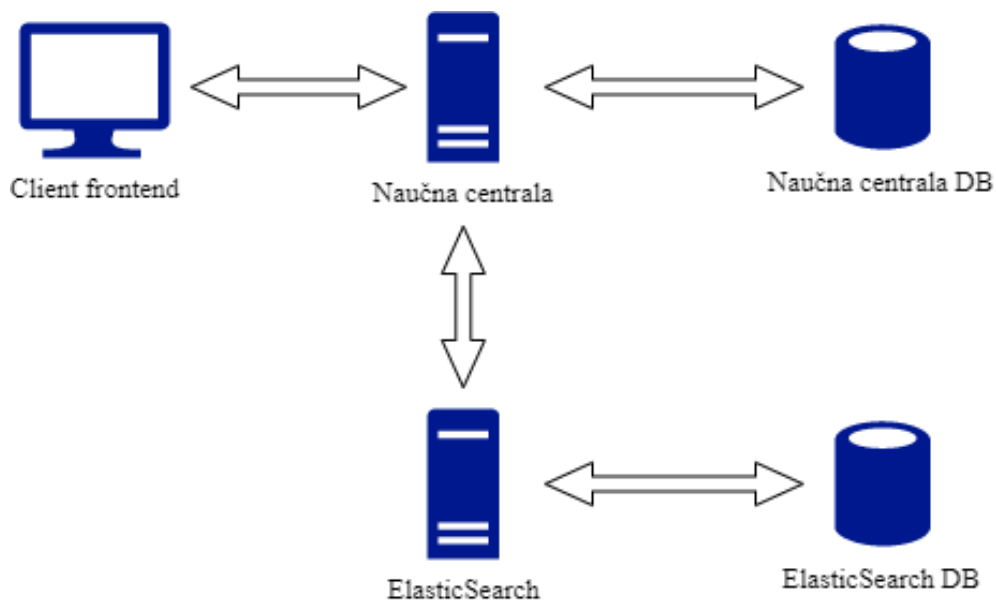
Upravljanje digitalnim dokumentima

Za projekat iz predmeta Upravljanje digitalnim dokumentima je potrebno napraviti aplikaciju koja će omogućiti pretraživanje radova u okviru Naučne centrale. Naučna centrala je sistem koji sadrži veliku količinu radova iz različitih naučnih oblasti. Kako bi korisniku bilo što lakše da pronađe odgovarajući rad, potrebno mu je omogućiti pretragu radova po različitim kriterijumima. Da bi se to postiglo, u okviru ovog projekta će se koristiti Elasticsearch platforma. Ova platforma omogućava pretraživanje kolekcije naučnih radova na osnovu postavljenih upita.

Arhitektura aplikacije

Aplikacija se sastoji iz narednih modula:

1. Client frontend – Angular aplikacija sa korisničkim interfejsom koji omogućava zadavanje različitih upita za pretragu
2. Naučna centrala – SpringBoot aplikacija u kojoj se obavlja celokupna logika i koja prosleđuje zahteve ka Elasticsearch-u
3. Elasticsearch – Platforma koja omogućava pretraživanje digitalnih radova koji se nalaze u njenoj bazi na osnovu postavljenog upita



Komunikacija između Client frontend-a i Naučne centrale se obavlja preko REST API-a.

Naučna centrala će komunicirati sa Elasticsearch-om putem REST API-a. Da bi ovo bilo moguće, potrebno je ubaciti odgovarajuće dependency-e unutar SpringBoot aplikacije za Naučnu centralu. Takođe je potrebno u konfiguracionom fajlu u okviru Naučne centrale napraviti klijenta koji omogućava slanje REST zahteva. Zahtevi će se slati na <http://localhost:9200>, obzirom da Elasticsearch sluša na tom portu ukoliko mu se prisupa putem REST-a.

Skladištenje podataka

U okviru baze podataka Naučne centrale će se čuvati sledeći podaci u vezi sa naučnim radom:

- Nalov rada
- Imena koautora, kao i njihove email i fizičke adrese
- Ključni pojmovi
- Apstrakt
- Naučna oblast rada
- PDF dokument

Digitalna biblioteka prihvaćenih radova će se skladištiti u okviru Elasticsearch-a, kako bi bilo omogućeno pretraživanje podataka. Da bi bila omogućena geoprostorna pretraga, neophodno je na Elasticsearch platformi čuvati *latitude* i *longitude* od svih recenzenata u sistemu.

SerbianAnalyzer

Da bi se obavilo pretprocesiranje tekstova na srpskom jeziku, najpre je instaliran plugin koji to omogućava (Serbian Analyzer). Analyzer koji će se koristiti u okviru ovog projekta je kompatibilan sa Java verzijom 13.0.1, Gradle verzijom 6.0 i Elasticsearch verzijom 7.4.0. i dostupan je na sledećem linku: <https://github.com/markomartonos/udd06/tree/plugin-update>.

Prvi korak jeste preuzimanje projekta sa navedenog linka koji sadrži analyzer za srpski jezik. Kako bi se analyzer mogao ubaciti u Elasticsearch, neophodno je izbuildovati odgovarajući fajl pomoću Gradle-a. Potrebno je otvoriti terminal u root folderu preuzetog projekta i izvršiti naredbu *gradlew clean build*. Na taj način je nastala arhivirana distribucija u folderu *build/distributions* unutar root foldera sa nazivom *serbian-analyzer-1.0-SNAPSHOT.zip*. Sledeći korak je instaliranje dobijenog fajla kao plugin-a unutar Elasticsearch-a. To se ostvaruje pozicioniranjem unutar *bin* foldera u okviru root foldera u koji je preuzet Elasticsearch. Pomoću naredbe *elasticsearch-plugin install file:<absolute path of distribution archive>* se instalira plugin.

ElasticSearch

Kako bi se koristio Elasticsearch, prvo je potrebno pokrenuti server. To se postiže pozicioniranjem unutar *bin* foldera u okviru root foldera u koji je preuzet Elasticsearch i pokretanjem *elasticsearch.bat* fajla. Da bi Elasticsearch sa instaliranim plugin-om imao mogućnost analiziranja teksta na srpskom jeziku, potrebno je još putem PUT zahteva naglasiti da se, pored podrazumevanog analyzera za engleski jezik, koristi i analyzer za srpski jezik.

Na slici ispod je prikazano kako je moguće „obavestiti“ Elasticsearch da prilikom dodavanja novog dokumenta u index *casopis* koristi i analyzer za srpski i analyzer za engleski jezik.

```

1 {
2   "mappings":{
3     "properties":{
4       "content":{
5         "type":"text","fields":{
6           "sr":{
7             "type":"text","analyzer":"serbian"
8           },
9           "en":{
10            "type":"text","analyzer":"english"
11          }
12        }
13      }
14    }
15  }
16 }

```

Kao što je već rečeno, komunikacija sa Elasticsearch-om će se obavljati putem REST API-a. Takođe, koristiće se i ElasticsearchTemplate kao i ElasticsearchRepository koji omogućavaju indeksiranje i pretraživanje. U okviru ovog projekta će se koristiti Elasticsearch verzija 7.4.0.

Indexing unit

Prilikom dodavanja novog rada u izabrani časopis, indexing unit će imati ovakav oblik.

```

1 {
2   "naziv_casopisa": "Casopis1",
3   "naslov_rada": "Rad1",
4   "autor":{
5     "ime": "Petar",
6     "prezime": "Petrovic"
7   },
8   "koautori":[
9     {
10      "ime": "Marko",
11      "prezime": "Markovic"
12    },
13    {
14      "ime": "Nikola",
15      "prezime": "Nikolic"
16    }
17  ],
18   "kljucni_pojmovi": "programiranje",
19   "sadrzaj": "tekst rada",
20   "naucna_oblast": "informatika"
21 }

```

Obzirom da je za geoprostornu pretragu potrebno pronaći recenzente čija je adresa na određenoj udaljenosti od datih adresa, u Elasticsearch će biti potrebno dodavati i recenzente sistema. Za njihovo dodavanje će biti korišćen indexing unit sledećeg oblika:

```

1 {
2   "id_recenzent": "1",
3   "location" : {
4     "lat" : 45.41,
5     "lon" : 19.70
6   }
7 }

```

More like this

U okviru ovog projekta potrebno je omogućiti *More like this* pretragu. Ovo je omogućeno Elasticsearch platformom putem MLT query-a. Ovaj upit pronalazi dokumente koji su poput datog skupa dokumenata. Da bi to učinio, MLT bira skup reprezentativnih termina ovih ulaznih dokumenata, formira upit koristeći te termine, izvršava ga i vraća rezultate. Korisnik je taj koji kontroliše ulazne dokumente, kako se biraju termini i kako se upit formira.

Za projekat Naučne centrale potrebno je pronaći sve recenzente koji su recenzirali radove slične podnetom tekstu rada. To znači da će se u upitu za MLT u okviru parametra *fields* (spisak polja za preuzimanje i analiziranje teksta) nalaziti polje koje sadrži tekst naučnog rada. Time smo ograničili Elasticsearch da pretragu po sličnosti bazira samo na ovom polju. Jedini obavezan parametar *More like this* upita jeste parametar *like*. U ovom parametru korisnik može navesti tekst u slobodnoj formi i/ili jedan ili više dokumenata sa kojima je potrebno naći sličnost. U ovom projektu će parametar *like* imati vrednost teksta podnetog dokumenta.

```
1 {
2   "query": {
3     "more_like_this": {
4       "fields": ["sadrzaj"],
5       "like": "tekst podnetog dokumenta"
6     }
7   }
8 }
```

Geoprostorna pretraga

Za projekat Naučne centrale je takođe potrebno podržati geoprostornu pretragu. Elasticsearch ima podršku za takve upite. Ovakva pretraga u okviru Elasticsearch-a omogućava filtriranje dokumenata, tako da se kao rezultat dobiju samo oni dokumenti koji postoje na određenoj udaljenosti od zadate.

Za potrebe ovog projekta će se koristiti lokacija kao *geo_point* tačka. To znači da podatke o lokaciji možemo indeksirati u različitim formatima, konkretno u ovom slučaju će se koristiti *latitude* i *logitude*. Za svakog recenzenta u sistemu će se na Elasticsearch platformi čuvati podatak o njegovom ID-u, kao i *latitude* i *longitude* njegove adrese. Da bi se izvršila geoprostorna pretraga, iz baze podataka Naučne centrale će se dobiti autor i svi koautori podnetog rada, odnosno njihovi parametri *latitude* i *longitude*. Potom će se vršiti geoprostorna pretraga na osnovu tih parametara, kako bismo dobili recenzente koji se nalaze na traženoj udaljenosti od autora i svih koautora.

```
1 {
2   "query": {
3     "bool": {
4       "must": {
5         "match_all": {}
6       },
7       "must_not": {
8         "geo_distance": {
9           "distance": "100km",
10          "autor.lokacija": {
11            "lat": 43.320904,
12            "lon": 21.895760
13          }
14        }
15      }
16     }
17   }
18 }
```

Na slici iznad je prikazan primer geoprostorne pretrage, gde se kao rezultat vraćaju svi dokumenti čija je udaljenost veća od 100km od date lokacije.