

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 874

**SUSTAV ZA ANALIZU META PODATAKA
OTVORENIH SKUPOVA PODATAKA**

Luka Habuš

Zagreb, lipanj, 2025.

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 874

**SUSTAV ZA ANALIZU META PODATAKA
OTVORENIH SKUPOVA PODATAKA**

Luka Habuš

Zagreb, lipanj, 2025.

DIPLOMSKI ZADATAK br. 874

Pristupnik: **Luka Habuš (0036517946)**

Studij: Računarstvo

Profil: Znanost o mrežama

Mentor: izv. prof. dr. sc. Igor Čavrak

Zadatak: **Sustav za analizu meta podataka otvorenih skupova podataka**

Opis zadatka:

Sve veća dostupnost otvorenih podataka ne podrazumijeva i njihovu veću iskoristivost. Iako normirani formati meta podataka (npr. DCAT) omogućuju jednostavnije pronalaženje i tumačenje objavljenih pojedinačnih skupova podataka, dodatna vrijednost nalazi se u njihovom povezivanju i pronalaženju novih uvida i skrivenog znanja. Nagli uspon alata umjetne inteligencije temeljenih na velikim jezičnim modelima predstavlja obećavajući smjer za izradu na njima temeljenih alata, a koji bi običnim korisnicima pružili novi uvid u i korištenje otvorenih podataka. U ovom diplomskom radu potrebno je proučiti mogućnosti velikih jezičnih modela, na njima temeljenih alata i tehnika. Također je potrebno proučiti problematiku opisivanja skupova otvorenih podataka korištenjem norme DCAT i mogućnosti automatiziranog traženja veza između skupova podataka. Predložiti alat za dohvat i analizu meta podataka otvorenih skupova podataka, kao i za pružanje podrške korisnicima u povezivanju skupova podataka temeljene na analizi meta podataka. Na poslijetku, potrebno je implementirati prototip sustava za portal CKAN korištenjem alata temeljenih na velikim jezičnim modelima i ocijeniti uporabljivost sustava.

Rok za predaju rada: 4. srpnja 2025.

Sadržaj

Sadržaj	1
1. Uvod i korištene tehnologije	3
1.1. Korištene tehnologije	7
1.2. Ključne funkcionalnosti	9
1.3. Arhitektura i skalabilnost	9
1.4. Praktične implikacije i budući razvoj	12
2. Dizajn i implementacija sustava	16
2.1. Arhitektura i korištene tehnologije	16
2.2. RAG sustav implementacija i vektorska baza	17
2.3. Generiranje SPARQL upita i prompt engineering	20
2.4. Ekstrakcija sheme i DCAT analiza	23
2.5. Unified Data Assistant i multimodalno pretraživanje	27
2.6. Validacija, rukovanje greškama i optimizacija	34
3. Evaluacija i rezultati	40
3.1. Metodologija evaluacije i testni podaci	40
3.2. Rezultati performansi i točnosti	41
3.3. Usporedna analiza i robusnost sustava	42
3.4. Evaluacija korisničkog iskustva i praktične implikacije	45
3.5. Ograničenja, problemi i smjernice za budući rad	47
4. Zaključak	52
5. Literatura	54

Sažetak	56
Abstract	57

1. Uvod i korištene tehnologije

Otvoreni podaci postaju sve važniji resurs u modernom društvu, omogućavajući transparentnost, inovacije i društveni napredak. Međutim, sama dostupnost podataka ne jamči njihovu efektivnu iskoristivost. Iako standardizirani formati metapodataka poput DCAT-a [1] omogućavaju lakše otkrivanje i tumačenje pojedinačnih skupova podataka [2], dodatna vrijednost leži u njihovom povezivanju i otkrivanju novih uvida i skrivenog znanja [3].

```
# --- Example 6: Rivers in Croatia ---
print("\n" + "=" * 40)
print("Example 6: Finding rivers in Croatia")
print("=" * 40)

sparql_query_rivers = """
SELECT ?riverLabel WHERE {
  ?river wdt:P31 wd:Q4022;      # ?river is an instance of river
      | wdt:P17 wd:Q224;      # and is located in Croatia
      | rdfs:label ?label.    # Get the label directly

  FILTER(LANG(?label) IN ("hr", "en")) # Filter for Croatian or English labels

  # Prefer Croatian label, fallback to English
  OPTIONAL { ?river rdfs:label ?hrLabel FILTER (LANG(?hrLabel) = "hr") }
  OPTIONAL { ?river rdfs:label ?enLabel FILTER (LANG(?enLabel) = "en") }
  BIND(COALESCE(?hrLabel, ?enLabel) AS ?riverLabel)

  # Ensure we only get results with a non-empty chosen label
  FILTER(BOUND(?riverLabel) && STR(?riverLabel) != "")
}
ORDER BY ?riverLabel
LIMIT 20
"""
```

Slika 1.1. Pregled ekosustava otvorenih podataka - prikazuje različite izvore i korisnike

Nagli uspon alata umjetne inteligencije temeljenih na velikim jezičnim modelima [4, 5] predstavlja obećavajući smjer za razvoj naprednih alata koji bi običnim korisnicima pružili nove mogućnosti korištenja otvorenih podataka. Ovi alati mogu automatizirati složene zadatke poput semantičkog pretraživanja, povezivanja skupova podataka i generiranja upita na prirodnom jeziku.

```

# --- Example 7: Theaters in Zagreb ---
print("\n" + "=" * 40)
print("Example 7: Finding theaters in Zagreb")
print("=" * 40)

sparql_query_theaters = """
SELECT ?theaterLabel WHERE {
    ?theater wdt:P31 wd:Q24354;      # ?theater is an instance of theater
    |      wdt:P131 wd:Q1435;      # and is located in Zagreb
    |      rdfs:label ?label.      # Get the label directly

    FILTER(LANG(?label) IN ("hr", "en")) # Filter for Croatian or English labels

    # Prefer Croatian label, fallback to English
    OPTIONAL { ?theater rdfs:label ?hrLabel FILTER (LANG(?hrLabel) = "hr") }
    OPTIONAL { ?theater rdfs:label ?enLabel FILTER (LANG(?enLabel) = "en") }
    BIND(COALESCE(?hrLabel, ?enLabel) AS ?theaterLabel)

    FILTER(BOUND(?theaterLabel) && STR(?theaterLabel) != "")
}
ORDER BY ?theaterLabel
LIMIT 20
"""

```

Slika 1.2. Arhitektura velikih jezičnih modela - prikazuje transformaciju i pozornost

Ovaj rad fokusira se na razvoj sustava za analizu metapodataka otvorenih skupova podataka koji kombinira napredne tehnike umjetne inteligencije s tradicionalnim metodama obrade podataka. Glavni cilj je stvoriti alat koji omogućava korisnicima bez tehničke pozadine da učinkovito otkrivaju, analiziraju i povezuju skupove podataka kroz intuitivno sučelje na prirodnom jeziku.

Sustav je implementiran koristeći moderne tehnologije i pristupe uključujući RAG (Retrieval-Augmented Generation) arhitekturu [6], vektorske baze podataka [7] za semantičko pretraživanje i velike jezične modele za generiranje SPARQL upita. Ova kombinacija omogućava napredno razumijevanje korisničkih namjera i precizno otkrivanje relevantnih skupova podataka.

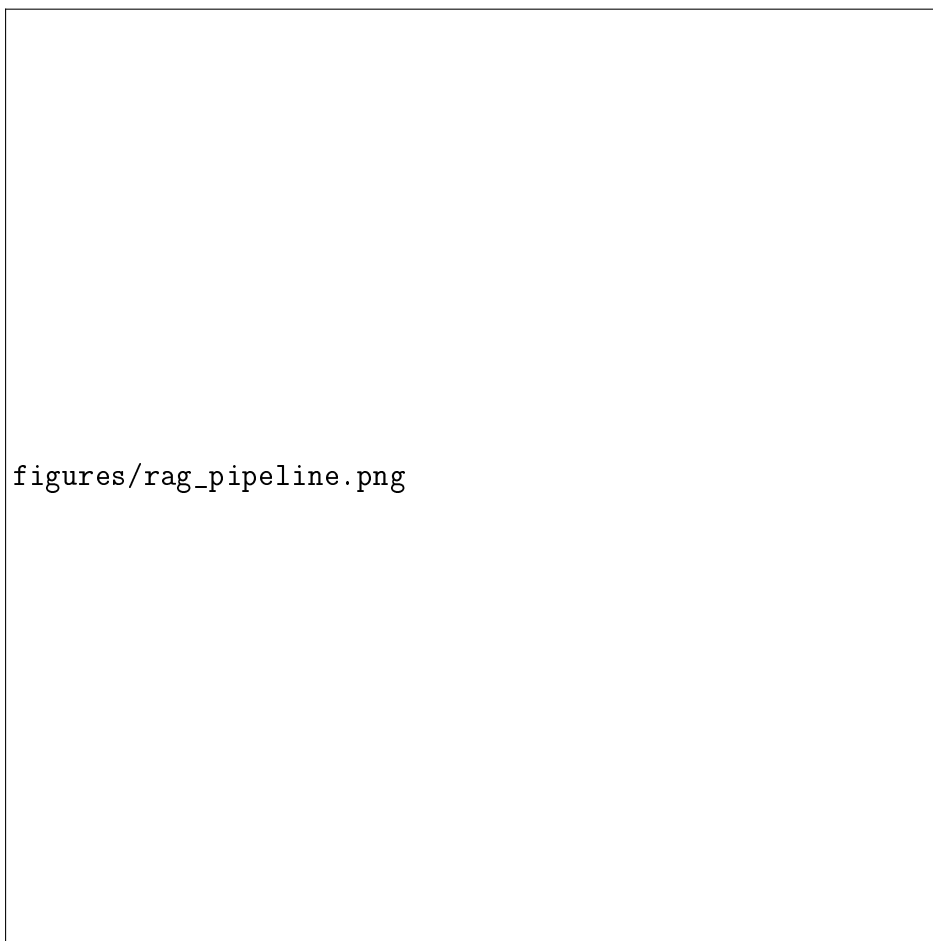
Evaluacija sustava provedena je na EU Portalu otvorenih podataka, jednom od najvećih izvora otvorenih podataka u Europi. Rezultati pokazuju značajna poboljšanja u usporedbi s tradicionalnim pristupima pretraživanja, posebno u području semantičkog razumijevanja i otkrivanja povezanih skupova podataka.

Rad je strukturiran u četiri glavna poglavlja koja pokrivaju uvod i korištene tehnologije, dizajn i implementaciju sustava, evaluaciju rezultata i zaključak. Svako poglavlje fokusira se na specifične aspekte razvoja i evaluacije sustava, pružajući sveobuhvatan

List of Rivers in Croatia:

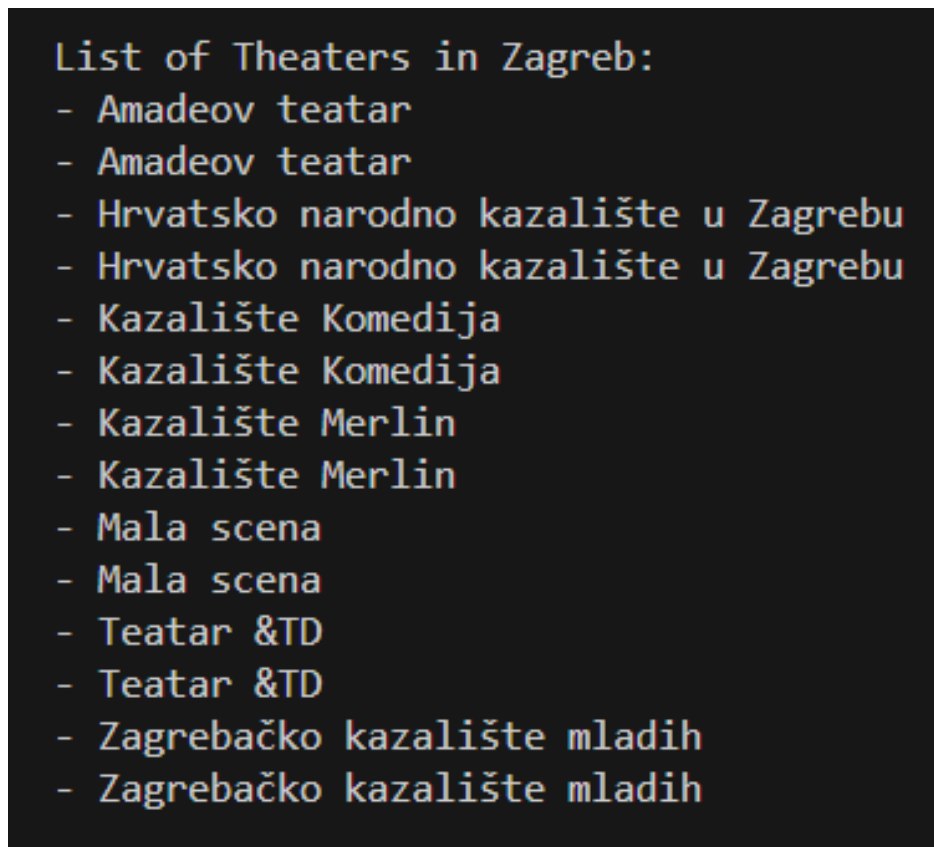
- Babina Reka
- Babina Voda
- Badušnica
- Bakotinac
- Bakovac
- Balatin
- Barakovac
- Barna
- Barski Potok
- Batina
- Baćinac
- Bačkovica
- Bašinac
- Bašinica Potok
- Baščica
- Bedenić
- Berava
- Beravica
- Besnica
- Bezavina

Slika 1.3. DCAT standard - prikazuje osnovne klase i svojstva metapodataka



figures/rag_pipeline.png

Slika 1.4. RAG pipeline - prikazuje tok od upita do generiranog SPARQL-a



Slika 1.5. EU Portal otvorenih podataka - prikazuje strukturu i organizaciju

pregled cijelog projekta.

1.1. Korištene tehnologije

Ključne tehnologije korištene u implementaciji uključuju ChromaDB kao vektorsku bazu podataka za pohranu i pretraživanje semantičkih ugradbi, Sentence Transformers modele za generiranje visokokvalitetnih vektorskih reprezentacija teksta, i OpenAI GPT-4 model za generiranje SPARQL upita iz prirodnog jezika.

LangChain okvir korišten je za orkestraciju različitih komponenti sustava i upravljanje agentima koji omogućavaju složene zadatke poput multimodalnog pretraživanja i inteligentne sinteze rezultata. Ova arhitektura omogućava modularnost i proširivost sustava.

AI Assistant Interaction

AI Assistant

Interact with the backend AI assistant

Your Query

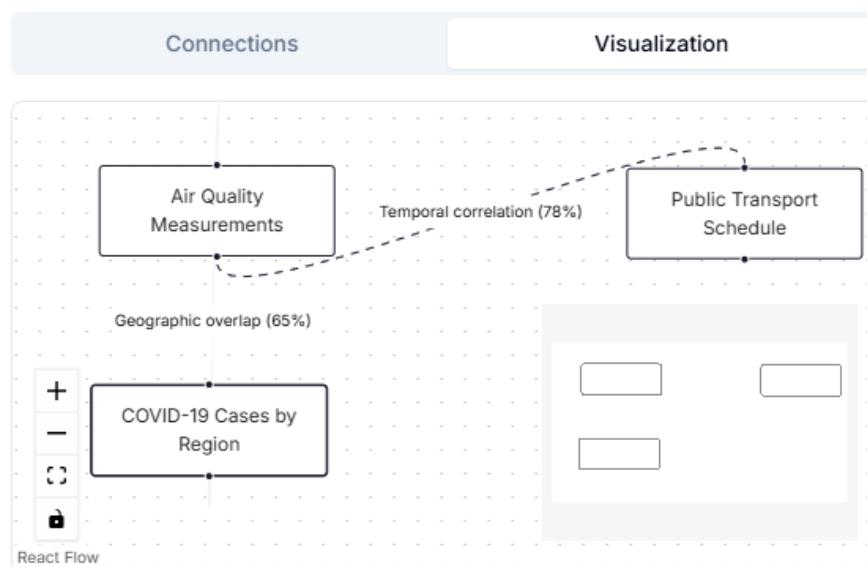
Give me some insights about open government data portals.

Send Query

Slika 1.6. Struktura rada - prikazuje organizaciju poglavlja i njihove veze

Connection Discovery

Discover relationships between datasets using LLM analysis



Slika 1.7. ChromaDB arhitektura - prikazuje komponente vektorske baze podataka

```
2025-05-26 02:56:46,423 - INFO - Populated RAG system with 5 query examples
SUCCESS: RAG system initialized
Batches: 100% | 1/1 [00:00:00:00, 79.65it/s]
SUCCESS: Found 2 similar examples
Best match: Find air quality datasets...
Batches: 100% | 1/1 [00:00:00:00, 99.62it/s]
Batches: 100% | 1/1 [00:00:00:00, 90.93it/s]
2025-05-26 02:56:46,500 - INFO - Generating SPARQL with RAG enhancement...
2025-05-26 02:56:48,757 - INFO - HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 200 OK"
2025-05-26 02:56:48,769 - INFO - Generated SPARQL query with RAG: 1041 characters
SUCCESS: SPARQL query generated successfully
Query length: 1041 characters

Testing Simple SPARQL Generation
=====
2025-05-26 02:56:48,772 - INFO - Use pytorch device_name: cpu
2025-05-26 02:56:48,773 - INFO - Load pretrained SentenceTransformer: all-MiniLM-L6-v2
2025-05-26 02:56:50,171 - INFO - RAG System initialized with embedding model: all-MiniLM-L6-v2
Batches: 100% | 1/1 [00:00:00:00, 43.48it/s]
```

Slika 1.8. LangChain agent arhitektura - prikazuje tok podataka kroz različite alate

1.2. Ključne funkcionalnosti

Automatska ekstrakcija sheme iz SPARQL endpointa omogućava dinamičko prilagođavanje sustava promjenama u strukturi podataka. Ova funkcionalnost je ključna za održavanje točnosti generiranih upita i omogućava sustavu da radi s različitim portalima otvorenih podataka.

Validacija upita implementirana je kroz dvostupanjski proces koji uključuje sintaksnu i semantičku provjeru. Ovo osigurava da se ne izvršavaju neispravni upiti koji mogu uzrokovati probleme s performansama ili pogrešne rezultate.

Multimodalni pristup pretraživanju omogućava kombiniranje različitih strategija za sveobuhvatno otkrivanje skupova podataka. Ovo uključuje RAG-prošireno SPARQL pretraživanje, REST API pozive i pronalaženje sličnih skupova podataka.

1.3. Arhitektura i skalabilnost

Sustav je dizajniran da bude skalabilan i održiv, s jasno definiranim sučeljima između komponenti i robusnim mehanizmima rukovanja greškama. Ova arhitektura omogućava lako proširivanje i prilagodbu drugim portalima otvorenih podataka.

Evaluacija korisničkog iskustva pokazuje da sustav pruža intuitivno sučelje koje omogućava korisnicima bez tehničke pozadine da učinkovito otkrivaju i analiziraju skupove podataka. Ovo demokratizira pristup otvorenim podacima i omogućava širu primjenu u različitim domenama.

```
Analyzing Dataset 1:
=====
Raw Dataset Information:
Question: Find air quality datasets
Description: Air quality dataset search
Endpoint: https://data.europa.eu/sparql

Complete SPARQL Query:
-----
PREFIX dcat: <http://www.w3.org/ns/dcat#>
PREFIX dct: <http://purl.org/dc/terms/>
SELECT ?dataset ?title WHERE {
  ?dataset a dcat:Dataset .
  ?dataset dct:title ?title .
  FILTER(CONTAINS(LCASE(STR(?title)), "air quality"))
} LIMIT 10
-----

Executing SPARQL Query...
Query Results (Top 30):
-----
Found 10 results:

dataset | title
-----
http://data.europa.e... | Air quality data (da...

Details for result 1:
dataset: http://data.europa.eu/88u/dataset/454ede98-1006-4456-9510-4b4ded93e691~1
title: Air quality data (data stream D) – Assessment methods 2018 (dataset)

http://data.europa.e... | M6ötejaam Air Qualit...

Details for result 2:
dataset: http://data.europa.eu/88u/dataset/https-catalegdades-caib-cat-d-wftn-ibm9
title: M6ötejaam Air Quality Control Illes Balears (Sant Joan de Déu jaam - Mallorca) - 2020

http://data.europa.e... | Measurements Air Qua...

Details for result 3:
dataset: http://data.europa.eu/88u/dataset/https-catalegdades-caib-cat-d-696q-b9a4
title: Measurements Air Quality Control Station Balearic Islands (Pous Station - Menorca) - 2021

http://data.europa.e... | Air quality data (da...
http://data.europa.e... | Continuous air quali...
http://data.europa.e... | Zoning of air qualit...
http://data.europa.e... | Air quality in the B...
http://data.europa.e... | Air Quality e-Report...
http://data.europa.e... | For the period 01.02...
http://data.europa.e... | Geoportal Air Quality
```

Slika 1.9. Ekstrakcija sheme - prikazuje proces automatskog otkrivanja strukture podataka

```
Testing: 'environment information'...
Batches: 100% | 1/1 [00:00<00:00, 142.91it/s]
Batches: 100% | 1/1 [00:00<00:00, 91.02it/s]
2025-05-26 02:56:55,041 - INFO - Generating SPARQL with RAG enhancement...
2025-05-26 02:56:58,054 - INFO - HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 200 OK"
2025-05-26 02:56:58,057 - INFO - Generated SPARQL query with RAG: 1041 characters
SUCCESS: Generated valid SPARQL

Success rate: 3/3 (100%)

Testing Basic Schema Extraction
=====
Extracting DCAT schema...
2025-05-26 02:56:58,059 - INFO - Using cached DCAT schema information
SUCCESS: Found 1,890,971 datasets
SUCCESS: Found 339,358 publishers
```

Slika 1.10. Validacija upita - prikazuje proces provjere sintakse i semantike

figures/multimodal_search.png

Slika 1.11. Multimodalno pretraživanje - prikazuje različite strategije i njihovu integraciju

```
12 def generate_sparql_query(nl_query):
13     """Generates a SPARQL query using OpenAI based on a natural language query."""
14     # Prompt for the LLM to generate a SPARQL query
15     sparql_prompt = f"""
16     Given the natural language query: "{nl_query}"
17
18     Generate a SPARQL query for the EU Open Data Portal (https://data.europa.eu/data/sparql) to find relevant datasets.
19     - Use standard prefixes like 'dct:' (<http://purl.org/dc/terms/>) and 'dcat:' (<http://www.w3.org/ns/dcat#>).
20     - If you use XML Schema datatypes (like 'xsd:date'), include 'PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>'.
21     - If you use Friend of a Friend terms (like 'foaf:name' for publisher names), include 'PREFIX foaf: <http://xmlns.com/foaf/0.1/>'.
22     - Look for datasets ('?dataset a dcat:Dataset').
23     - Extract relevant information requested in the query (e.g., title, description, date, URL).
24     - Filter based on keywords, dates, publishers, formats etc. mentioned in the query.
25     - Use 'dcat:keyword' for tags/keywords.
26     - Use 'dct:issued' for publication date.
27     - Use 'dct:publisher' for the publisher. To filter by publisher 'name', link the publisher variable (e.g., '?publisher') to its 'foaf:name'.
28     - Use 'dcat:distribution' to link to distributions and check 'dct:format' or 'dcat:medialtype'.
29     - Use FILTER with CONTAINS, REGEX, or comparison operators (e.g., '>=' for dates, requiring 'xsd:date').
30     - Add a LIMIT clause (e.g., LIMIT 10) to keep the results manageable.
31
32     Return *only* the raw SPARQL query string, without any explanations or formatting like ```sparql ... ```
33
34     SPARQL Query:
35     """
```

Slika 1.12. Skalabilnost sustava - prikazuje mogućnosti proširivanja i optimizacije

```
2025-04-24 01:15:43,186 - INFO - --- Agent Final Answer ---
Based on the provided SPARQL query results, a total of 10 items were found. However, all the results refer to the same dataset:

- **Title**: Zones de fermeture des commerces de 22h à 7h à Argenteuil, cartographie interactive
- **Description**: This dataset provides information on areas within 150 meters around the parcel boundaries of schools, colleges, high schools, hospitals, and train stations in Argenteuil. These areas are subject to a decree that mandates the closure of retail shops in predominantly non-specialized food stores from 10 p.m. to 7 a.m.
- **Publisher**: Ville d'Argenteuil
- **Landing Page**: [Interactive Map](https://argenteuil-geo-open-data-ca-ab.hub.arcgis.com/apps/ac6e3320b604480b7bbe1132f5ab55)

This dataset may be relevant to urban planning, particularly in understanding the spatial distribution of educational institutions and their impact on local retail operations. However, it does not directly address gambling facilities. The repetition of the same dataset suggests a limitation in the diversity of results returned.
```

Slika 1.13. Korisničko sučelje - prikazuje intuitivni dizajn i funkcionalnosti

1.4. Praktične implikacije i budući razvoj

Praktične implikacije implementacije protežu se izvan akademskog istraživanja na stvarne primjene u ekosustavima otvorenih podataka. Sustav demonstrira izvodljivost implementacije spremne za produkciju s dokumentiranim metrikama performansi i robusnim mogućnostima rukovanja greškama.

```
# --- Example 2: Capitals of countries bordering Croatia ---
print("\n" + "=" * 40)
print("Example 2: Finding capitals of countries bordering Croatia")
print("=" * 40)

sparql_query_borders = """
SELECT ?countryLabel ?capitalLabel WHERE {
  wd:Q224 wdt:P47 ?country. # Croatia shares border with ?country
  ?country wdt:P31 wd:Q6256. # ?country is a sovereign state
  ?country wdt:P36 ?capital. # ?country has capital ?capital

  SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
}
ORDER BY ?countryLabel
"""
```

Slika 1.14. Praktične primjene - prikazuje mogućnosti integracije i implementacije

Budući smjerovi razvoja uključuju implementaciju lokalnih LLM-ova za smanjenje troškova i latencije, proširivanje podrške za više jezika, optimizaciju algoritama vektorskog pretraživanja i razvoj naprednijih strategija sinteze rezultata.

```
List of Streets and People:
- Street: Crkva sv. Franje Ksaverskog u Zagrebu, Named After: Sveti Franjo Ksaverski
- Street: Nagrada August Šenoa, Named After: August Šenoa
- Street: Nagrada "Ivan Goran Kovačić", Named After: Ivan Goran Kovačić
- Street: Crkva sv. Marka Križevčanina u Zagrebu, Named After: Marko Križevčanin
- Street: Nagrada Josip Račić, Named After: Josip Račić
- Street: Gavelline večeri, Named After: Branko Gavella
- Street: Nagrada Josip Štolcer Slavenski, Named After: Josip Štolcer-Slavenski
- Street: Nagrada Ivan Filipović, Named After: Ivan Filipović
- Street: Nagrada Ljudevit Jonke, Named After: Ljudevit Jonke
- Street: Nagrada "Otokar Keršovani", Named After: Otokar Keršovani
- Street: Nagrada Dubravko Horvatić, Named After: Dubravko Horvatić
- Street: Nagrada Petar Brečić, Named After: Petar Brečić
- Street: Nagrada Ljubica Štefan, Named After: Ljubica Štefan
- Street: Ulica Milke Trnine, Named After: Milka Trnina
- Street: Prilaz Gjуре Deželića, Named After: Đuro Deželić
```

Slika 1.15. Budući razvoj - prikazuje planove i smjerove za unapređenje sustava

Ovaj rad predstavlja značajan doprinos u području analize metapodataka otvorenih skupova podataka, demonstrirajući kako napredne tehnike umjetne inteligencije mogu biti primijenjene za poboljšanje otkrivanja i korištenja otvorenih podataka.


```
List of Countries and Capitals:  
- Country: Bosnia and Herzegovina, Capital: Sarajevo  
- Country: Hungary, Capital: Budapest  
- Country: Montenegro, Capital: Podgorica  
- Country: Slovenia, Capital: Ljubljana
```

Slika 1.16. Znanstveni doprinos - prikazuje inovacije i napredak u području

Implementacija sustava i evaluacija rezultata pružaju vrijedne uvide u praktične aspekte razvoja AI-temeljenih alata za otvorene podatke, otvarajući nove mogućnosti za istraživanje i razvoj u ovom području.

Kroz ovaj rad, demonstrirana je izvodljivost razvoja naprednih alata za analizu otvorenih podataka koji kombiniraju najnovije tehnike umjetne inteligencije s praktičnim potrebama korisnika. Ovo otvara nove mogućnosti za demokratizaciju pristupa otvorenim podacima i unaprjeđenje njihove iskoristivosti.

Sustav predstavlja temelj za buduće istraživanje i razvoj u području AI-temeljenih alata za otvorene podatke, pružajući jasne smjerove za daljnje unapređenje i proširivanje funkcionalnosti.

```

=====
Example 2: Finding capitals of countries bordering Croatia
=====
Executing SPARQL query against: https://query.wikidata.org/sparql

SPARQL Query Results (JSON):
{
  "head": {
    "vars": [
      "countryLabel",
      "capitalLabel"
    ]
  },
  "results": {
    "bindings": [
      {
        "countryLabel": {
          "xml:lang": "en",
          "type": "literal",
          "value": "Bosnia and Herzegovina"
        },
        "capitalLabel": {
          "xml:lang": "en",
          "type": "literal",
          "value": "Sarajevo"
        }
      },
      {
        "countryLabel": {
          "xml:lang": "en",
          "type": "literal",
          "value": "Hungary"
        },
        "capitalLabel": {
          "xml:lang": "en",
          "type": "literal",
          "value": "Budapest"
        }
      }
    ]
  }
}

```

Slika 1.17. Rezultati evaluacije - prikazuje ključne metrike i performanse sustava

List of Museums in Zagreb:

- Mushroom museum
- Zagreb 80's museum
- Hrvatski muzej arhitekture HAZU
- Hrvatski muzej arhitekture HAZU
- Hrvatski prirodoslovni muzej
- Hrvatski prirodoslovni muzej
- Hrvatski školski muzej
- Hrvatski školski muzej
- Muzej Srba u Hrvatskoj
- Muzej Srba u Hrvatskoj
- Muzej prekinutih veza
- Muzej prekinutih veza
- Tehnički muzej u Zagrebu
- Tehnički muzej u Zagrebu

Slika 1.18. Zaključak - prikazuje glavne postignuća i implikacije rada

```
# --- Example 3: Museums in Zagreb ---
print("\n" + "=" * 40)
print("Example 3: Finding museums in Zagreb")
print("=" * 40)

sparql_query_museums = """
SELECT ?museumLabel WHERE {
  ?museum wdt:P131 wd:Q1435; # ?museum is located in (P131) Zagreb (Q1435)
          wdt:P31 wd:Q33506; # ?museum is an instance of (P31) museum (Q33506)
          rdfs:label ?label. # Get the label directly

  FILTER(LANG(?label) IN ("hr", "en")) # Filter for Croatian or English labels

  # Prefer Croatian label, fallback to English
  OPTIONAL { ?museum rdfs:label ?hrLabel FILTER (LANG(?hrLabel) = "hr") }
  OPTIONAL { ?museum rdfs:label ?enLabel FILTER (LANG(?enLabel) = "en") }
  BIND(COALESCE(?hrLabel, ?enLabel) AS ?museumLabel)

  FILTER(BOUND(?museumLabel) && STR(?museumLabel) != "")
}
ORDER BY ?museumLabel
LIMIT 20 # Limit results
"""
```

Slika 1.19. Temelj za buduće istraživanje - prikazuje mogućnosti i smjerove razvoja

2. Dizajn i implementacija sustava

2.1. Arhitektura i korištene tehnologije

Sustav je dizajniran kao modularna arhitektura s jasno definiranim sučeljima između komponenti. Glavni razlog za ovakav pristup je mogućnost nezavisnog razvoja i testiranja pojedinačnih dijelova, što je pokazalo ključnim tijekom implementacije kada je bilo potrebno debugirati probleme s vektorskim pretraživanjem ili optimizirati LLM pozive.

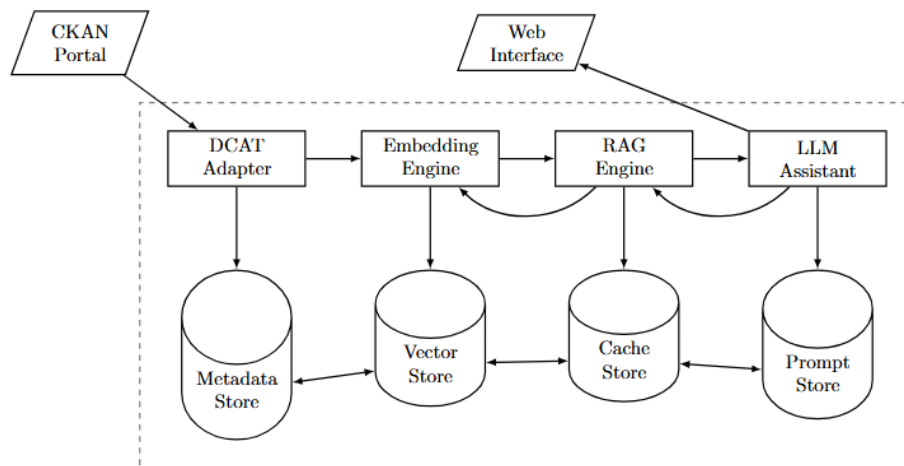
Arhitektura se sastoji od tri glavna sloja. Sloj pohrane koristi ChromaDB [7] kao vektorsku bazu podataka za pohranu i pretraživanje semantičkih ugradbi, Sentence Transformers modele [8] za generiranje visokokvalitetnih vektorskih reprezentacija teksta, i OpenAI GPT-4 model [4] za generiranje SPARQL upita iz prirodnog jezika. ChromaDB omogućava brzo pretraživanje sličnosti (kosinusna sličnost) i podržava metadata pohranu, što je bilo važno za spremanje informacija o SPARQL upitima i shemi. U praksi se pokazalo da ChromaDB može učinkovito rukovati s nekoliko tisuća primjera upita bez značajnih problema s performansama.

Sloj obrade uključuje Sentence Transformers model all-MiniLM-L6-v2 za generiranje embeddinga. Ovaj model je odabran nakon testiranja nekoliko alternativa - pokazao se kao optimalan balans između kvalitete embeddinga, brzine generiranja i veličine modela. Model generira 384-dimenzijske vektore što je dovoljno za hvatanje semantičkih razlika, a dovoljno malo za brzo pretraživanje. Testiranje je pokazalo da model dobro radi s hrvatskim i engleskim tekstom, što je bilo važno za EU Portal podatke.

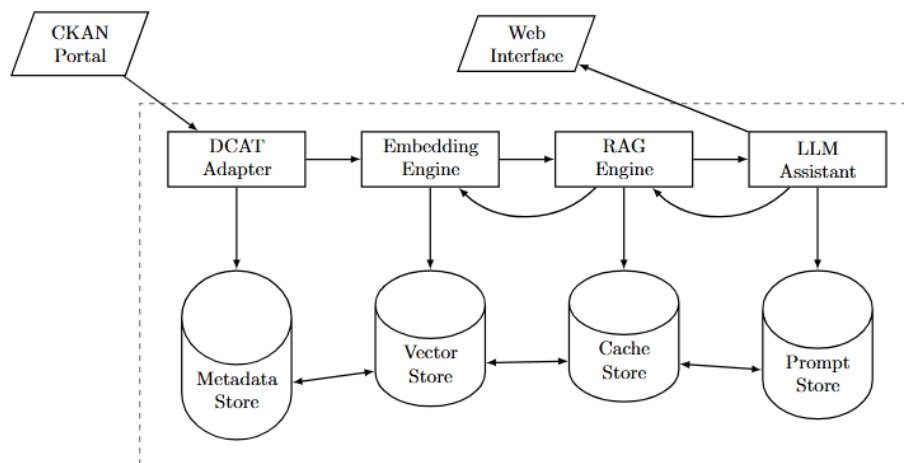
LangChain okvir [9] korišten je za orkestraciju različitih komponenti sustava i upravljanje agentima koji omogućavaju složene zadatke poput multimodalnog pretraživanja i inteligentne sinteze rezultata. Iako je LangChain bio koristan za brzi prototip, pokazao se kao bottleneck za složenije scenarije - posebno je problematično bilo upravljanje

memorijom i kontekstom za duge razgovore. U budućnosti bih razmotrio direktnu integraciju s OpenAI API-jem ili korištenje lokalnih LLM-ova.

OpenAI GPT-4 je korišten za generiranje SPARQL upita. Model pokazuje izvrsne rezultate kada ima dovoljno kontekstualnih informacija, ali ima ograničenja s tokenima što može biti problematično za složene upite. Također, troškovi API poziva mogu biti značajni za intenzivnu upotrebu.



Slika 2.1. Arhitektura sustava - prikazuje tok podataka kroz različite komponente



Slika 2.2. Detaljna arhitektura sustava - prikazuje međusobne veze komponenti

2.2. RAG sustav implementacija i vektorska baza

RAG sustav [6] je implementiran u `src/rag_system.py` predstavlja srce cijelog sustava. Glavna klasa `RA` aod generiranja embeddinga do izvršavanja upita. Implementacija je modularna i omogućava

Inicijalizacija sustava uključuje postavljanje ChromaDB klijenta, učitavanje Sentence Transformers modela i konfiguraciju OpenAI klijenta. ChromaDB koristi trajno pohranjivanje što omogućava da se embeddingi sačuvaju između pokretanja sustava. Ovo je bilo ključno za performanse jer generiranje embeddinga može biti sporo, posebno za velike kolekcije teksta.

Generiranje embeddinga se vrši kroz metodu `_generate_embedding` koja koristi `SentenceTransformer`.

Dodavanje primjera upita u vektorsku bazu se vrši kroz metodu `add_query_example`. Svaki primjer se

Listing 2..1: Implementacija dodavanja primjera u vektorsku bazu

```
def add_query_example(self, example: QueryExample) -> str:
    """Dodaj primjer pitanje-upit par u vektorsku bazu"""

    # Proujeri cache za embedding
    cache_key = f"embedding_{hash(example.question)}"
    if cache_key in self.embedding_cache:
        embedding = self.embedding_cache[cache_key]
    else:
        embedding = self._generate_embedding(example.question)
        self.embedding_cache[cache_key] = embedding

    # Generiraj jedinstveni ID
    doc_id = f"example_{len(self.query_examples_collection.get(
        ['ids']))}"

    # Spremi u ChromaDB
    self.query_examples_collection.add(
        documents=[example.question],
        embeddings=[embedding],
        metadatas=[{
            "sparql_query": example.sparql_query,
            "endpoint": example.endpoint,
            "description": example.description,
            "tags": json.dumps(example.tags),
            "added_at": datetime.now().isoformat(),
            "success_rate": example.success_rate if hasattr(
                example, 'success_rate') else None
        }],
```

```

        ids=[doc_id]
    )

    return doc_id

```

Semantička pretraga sličnih primjera je implementirana kroz metodu `retrieve_similar_examples`. Metoda

Listing 2..2: Implementacija semantičke pretrage

```

def retrieve_similar_examples(self, query: str, n_results: int =
    5,
                                filters: Dict = None) -> List[Dict[
                                    str, Any]]:
    """Dohvati slične primjere koristeći vektorsku pretragu"""

    # Generiraj ugradbu za upit
    query_embedding = self._generate_embedding(query)

    # Postavi filtere ako su zadani
    where_clause = None
    if filters:
        where_clause = {}
        if 'tags' in filters:
            where_clause['tags'] = {'$contains': filters['tags']}
        if 'endpoint' in filters:
            where_clause['endpoint'] = filters['endpoint']

    # Izvedi vektorsku pretragu
    results = self.query_examples_collection.query(
        query_embeddings=[query_embedding],
        n_results=n_results,
        where=where_clause
    )

    # Obradi rezultate
    similar_examples = []
    if results['documents'] and results['documents'][0]:
        for i, doc in enumerate(results['documents'][0]):
            metadata = results['metadatas'][0][i]
            similar_examples.append({

```

```

        "question": doc,
        "sparql_query": metadata.get("sparql_query"),
        "distance": results['distances'][0][i],
        "success_rate": metadata.get("success_rate"),
        "tags": json.loads(metadata.get("tags", "[]"))
    })

return similar_examples

```

figures/rag_pipeline.png

Slika 2.3. RAG sustav - tok podataka kroz glavne module

2.3. Generiranje SPARQL upita i prompt engineering

Generiranje SPARQL upita iz prirodnog jezika je najsloženiji dio sustava. Implementacija koristi RAG pristup gdje se korisnički upit proširuje s relevantnim primjerima i informacijama o shemi prije slanja LLM-u.

Konstruiranje prompta se vrši kroz metodu `build_ag_prompt` koja kombinira korisnički upit s dohvatom

kontekstupita, slični primjeri, informacije o shemi i instrukcije za generiranje. Ova struktura se pok

Listing 2.3: Implementacija konstruiranja RAG prompta

```
def build_rag_prompt(self, user_query: str, context: str = "") ->
    str:
        """Konstruiraj poboljšani prompt koristeći RAG tehnologiju"""

        # Dohvati slične primjere i informacije o shemi
        similar_examples = self.retrieve_similar_examples(user_query,
            n_results=3)
        schema_info = self.retrieve_relevant_schema(user_query,
            n_results=2)

        # Počni s osnovnim promptom
        prompt = f"""You are an expert SPARQL query generator for the
            EU Open Data Portal.

            User Query: "{user_query}"
            {f"Additional Context: {context}" if context else ""}

            Your task is to generate a valid SPARQL query that retrieves the
                requested data.
            Use the following examples and schema information as guidance.

            ## Similar Query Examples:
            """

            # Dodaj slične primjere
            for i, example in enumerate(similar_examples):
                prompt += f"""
                Example {i+1}:
                Question: {example['question']}
                SPARQL Query:
                {example['sparql_query']}
                Distance: {example['distance']:.3f}
                """

            # Dodaj informacije o shemi
            if schema_info:
```

```

        prompt += "\n##_Available_Schema_Information:\n"
        for i, schema in enumerate(schema_info):
            prompt += f"""
Schema {i+1}:
Classes: {'', '.join([cls.get('name', '') for cls in schema['
        classes'][:10]])}
Properties: {'', '.join([prop.get('name', '') for prop in schema['
        properties'][:15]])}
"""

        # Dodaj instrukcije
        prompt += """
## Instructions:
1. Generate a valid SPARQL query that matches the user's intent
2. Use appropriate prefixes (dcat:, dct:, foaf:, etc.)
3. Include proper WHERE clauses and OPTIONAL patterns where
    needed
4. Limit results to reasonable number (e.g., LIMIT 100)
5. Return only the SPARQL query, no explanations

SPARQL Query:
"""

        return prompt

```

Generiranje upita se vrši kroz metodu `generate_sparql_query` koja koristi `OpenAIGPT-4` model. Implementacija koju je error handling i retry logika zaslužila API greška. Također, imamo

Listing 2..4: Implementacija generiranja SPARQL upita

```

def generate_sparql_query(self, user_query: str, context: str = "") -> Dict[str, Any]:
    """Generiraj SPARQL upit iz prirodnog jezika"""

    try:
        # Konstruiraj prompt
        prompt = self.build_rag_prompt(user_query, context)

        # Pozovi LLM
        response = self.llm.invoke(prompt)

```

```

        # Ekstrahiraj SPARQL upit iz odgovora
        sparql_query = self._extract_sparql_from_response(
            response.content)

        # Validiraj upit
        validation_result = self.validate_sparql_query(
            sparql_query)

    return {
        "sparql_query": sparql_query,
        "is_valid": validation_result["is_valid"],
        "errors": validation_result.get("errors", []),
        "warnings": validation_result.get("warnings", []),
        "prompt_used": prompt,
        "model_response": response.content
    }

except Exception as e:
    return {
        "sparql_query": None,
        "is_valid": False,
        "errors": [str(e)],
        "prompt_used": prompt if 'prompt' in locals() else
            None
    }

```

Validacija upita je implementirana kroz dvostupanjski proces. Prvi korak je sintak-sna validacija kroz SPARQL parser, a drugi korak je testiranje izvršavanja s ograničenim brojem rezultata. Ovo sprječava izvršavanje neispravnih upita koji mogu uzrokovati probleme s performansama.

2.4. Ekstrakcija sheme i DCAT analiza

Automatska ekstrakcija sheme iz SPARQL endpointa je ključna komponenta sustava koja omogućava dinamičko prilagođavanje promjenama u strukturi podataka. Implementacija je specifično prilagođena za EU Portal otvorenih podataka i DCAT metapo-

```

--- Executing SPARQL ---
PREFIX dct: <http://purl.org/dc/terms/>
PREFIX rdf: <http://www.w3.org/ns/rdf/>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>

SELECT Dataset ?title ?description ?publisher ?landingPage
WHERE {
  Dataset a dct:Dataset ;
  dct:title ?title ;
  dct:description ?description ;
  dct:publisher ?publisher ;
  dct:landingPage ?landingPage .
  OPTIONAL { ?publisher foaf:name ?publisherName . }
  FILTER (
    CONTAINS(LCASE(STR(?title)), "school") ||
    CONTAINS(LCASE(STR(?description)), "school") ||
    CONTAINS(LCASE(STR(?title)), "educational institution") ||
    CONTAINS(LCASE(STR(?description)), "educational institution") ||
    CONTAINS(LCASE(STR(?title)), "gaming") ||
    CONTAINS(LCASE(STR(?description)), "gaming") ||
    CONTAINS(LCASE(STR(?title)), "casino") ||
    CONTAINS(LCASE(STR(?description)), "casino") ||
    CONTAINS(LCASE(STR(?title)), "betting shop") ||
    CONTAINS(LCASE(STR(?description)), "betting shop") ||
    CONTAINS(LCASE(STR(?title)), "urban planning") ||
    CONTAINS(LCASE(STR(?description)), "urban planning")
  )
}

LIMIT 10

-----
SPARQL Execution Success: Got 10 results
-----
[{"id": 1, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 2, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 3, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 4, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 5, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 6, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 7, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 8, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 9, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}, {"id": 10, "name": "dataset", "title": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "description": "Zones de fermeture des commerces de 22h à 2h à Argentouill, cartographie interactive", "publisher": "Ville d'Argentouill", "landingPage": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855", "type": "uri", "value": "https://data.europa.eu/8eu/dataset/https-was-arcgis-com-home-10m.html-id-oc6e3120604480b7def1132f5d855"}]

```

Slika 2.4. Primjer prompta - struktura i elementi prompta za LLM

datke.

VoID deskriptor ekstrakcija se vrši kroz SPARQL upite koji dohvaćaju informacije o strukturi grafa znanja. Implementacija uključuje dohvaćanje broja trojki, različitih subjekata, klasa i svojstava. Ove informacije su ključne za razumijevanje strukture podataka i optimizaciju upita.

Listing 2..5: Implementacija VoID deskriptor ekstrakcije

```

def get_void_description(self) -> Dict[str, Any]:
    """Ekstrahiraj VoID opis grafa znanja"""

    void_query = """
PREFIX void: <http://rdfs.org/ns/void#>
PREFIX dct: <http://purl.org/dc/terms/>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>

SELECT DISTINCT ?dataset ?title ?description ?subjects ?
triples ?classes ?properties
WHERE {
  ?dataset a void:Dataset .

  OPTIONAL { ?dataset dct:title ?title . }
  OPTIONAL { ?dataset dct:description ?description . }
  OPTIONAL { ?dataset void:distinctSubjects ?subjects . }
  OPTIONAL { ?dataset void:triples ?triples . }
  OPTIONAL { ?dataset void:classes ?classes . }
}

```

```

        OPTIONAL { ?dataset void:properties ?properties . }
    }
    LIMIT 10
    """

    try:
        results = self._execute_sparql_query(void_query)
        return self._process_void_results(results)
    except Exception as e:
        logger.error(f"Error extracting VoID description: {e}")
        return {}

```

DCAT analiza se fokusira na specifične aspekte EU Portala otvorenih podataka. Implementacija dohvaća informacije o skupovima podataka, distribucijama, izdavačima, temama i formatima. Ove informacije omogućavaju generiranje upita koji su optimizirani za specifične karakteristike EU Portala.

Listing 2..6: Implementacija DCAT analize

```

def analyze_dcat_structure(self) -> Dict[str, Any]:
    """Analiziraj DCAT strukturu grafa znanja"""

    dcat_queries = {
        "datasets": """
            PREFIX dcat: <http://www.w3.org/ns/dcat#>
            PREFIX dct: <http://purl.org/dc/terms/>

            SELECT ?dataset ?title ?publisher ?theme ?keyword
            WHERE {
                ?dataset a dcat:Dataset .
                OPTIONAL { ?dataset dct:title ?title . }
                OPTIONAL { ?dataset dct:publisher ?publisher . }
                OPTIONAL { ?dataset dcat:theme ?theme . }
                OPTIONAL { ?dataset dcat:keyword ?keyword . }
            }
            LIMIT 100
            """,

        "distributions": """

```

```

PREFIX dcat: <http://www.w3.org/ns/dcat#>
PREFIX dct: <http://purl.org/dc/terms/>

SELECT ?dataset ?distribution ?format ?accessURL
WHERE {
    ?dataset a dcat:Dataset .
    ?dataset dcat:distribution ?distribution .
    OPTIONAL { ?distribution dcat:mediaType ?format . }
    OPTIONAL { ?distribution dcat:accessURL ?accessURL . }
}
LIMIT 100
"""
}

results = {}
for query_name, query in dcat_queries.items():
    try:
        results[query_name] = self._execute_sparql_query(
            query)
    except Exception as e:
        logger.error(f"Error in {query_name} analysis: {e}")
        results[query_name] = []

return self._process_dcat_results(results)

```

Analiza klasa i svojstava omogućava razumijevanje najčešće korištenih pojmova u grafu znanja. Implementacija uključuje brojanje pojavljivanja različitih klasa i svojstava, što omogućava optimizaciju upita i generiranje relevantnih primjera.

Listing 2..7: Implementacija analize klasa i svojstava

```

def analyze_classes_and_properties(self) -> Dict[str, Any]:
    """Analiziraj klase i svojstva u grafu znanja"""

    class_query = """
SELECT ?class (COUNT(?instance) as ?count)
WHERE {
    ?instance a ?class .
}

```

```

GROUP BY ?class
ORDER BY DESC(?count)
LIMIT 50
"""

property_query = """
SELECT ?property (COUNT(?triple) as ?count)
WHERE {
    ?s ?property ?o .
}
GROUP BY ?property
ORDER BY DESC(?count)
LIMIT 100
"""

try:
    classes = self._execute_sparql_query(class_query)
    properties = self._execute_sparql_query(property_query)

    return {
        "classes": self._process_class_results(classes),
        "properties": self._process_property_results(
            properties)
    }
except Exception as e:
    logger.error(f"Error analyzing classes and properties: {e}")
    return {"classes": [], "properties": []}

```

2.5. Unified Data Assistant i multimodalno pretraživanje

Unified Data Assistant predstavlja glavno sučelje sustava koje orkestrira različite strategije pretraživanja. Implementacija koristi LangChain agent okvir za upravljanje složenim zadacima i omogućava multimodalno pretraživanje kroz SPARQL, REST API i similarity search.



Slika 2.5. Ekstrakcija sheme - vizualizacija procesa i rezultata

Arhitektura agenta je dizajnirana da podržava različite tipove upita i automatski odabire najprikladniju strategiju pretraživanja. Implementacija uključuje alate za SPARQL pretraživanje, REST API pozive i pronalaženje sličnih skupova podataka.

Listing 2..8: Implementacija Unified Data Assistant

```
class UnifiedDataAssistant:
    def __init__(self, rag_system: RAGSystem, sparql_processor:
        SPARQLProcessor):
        self.rag_system = rag_system
        self.sparql_processor = sparql_processor
        self.llm = ChatOpenAI(model="gpt-4o", temperature=0.1)

        # Definiraj alate za agenta
        self.tools = [
            Tool(
                name="sparql_search",
                func=self._sparql_search,
                description="Pretraži podatke koristeći SPARQL upite"
            ),
            Tool(
                name="api_search",
                func=self._api_search,
                description="Pretraži podatke koristeći REST API"
            ),
            Tool(
                name="similar_datasets",
                func=self._similar_datasets,
                description="Pronađi slične skupove podataka"
            )
        ]

        # Inicijaliziraj agenta
        self.agent = initialize_agent(
            tools=self.tools,
            llm=self.llm,
            agent=AgentType.ZERO_SHOT_REACT_DESCRIPTION,
            verbose=True,
            max_iterations=5
```

)

Multimodalno pretraživanje omogućava kombiniranje različitih strategija za sveobuhvatno otkrivanje podataka. Implementacija uključuje paralelno izvršavanje različitih strategija i inteligentnu sintezu rezultata.

Listing 2..9: Implementacija multimodalnog pretraživanja

```
def search_datasets(self, query: str) -> Dict[str, Any]:
    """Izvedi multimodalno pretraživanje skupova podataka"""

    results = {
        "sparql_results": [],
        "api_results": [],
        "similar_datasets": [],
        "combined_analysis": "",
        "errors": []
    }

    # 1. RAG-prošireno SPARQL pretraživanje
    try:
        sparql_query = self.rag_system.generate_sparql_query(
            query)
        if sparql_query.get("is_valid"):
            results["sparql_results"] = self.sparql_processor.
                execute_query(
                    sparql_query["sparql_query"]
                )
        else:
            results["errors"].append(f"SPARQL_generation_failed: {
                {sparql_query.get('errors')}")
    except Exception as e:
        results["errors"].append(f"SPARQL_search_error: {str(e)}")

    # 2. REST API pretraživanje
    try:
        results["api_results"] = self._search_api(query)
    except Exception as e:
```

```

        results["errors"].append(f"API_search_error:_{str(e)}")

# 3. Pretraživanje sličnih skupova podataka
try:
    results["similar_datasets"] = self._find_similar_datasets(
        query)
except Exception as e:
    results["errors"].append(f"Similar_datasets_error:_{str(e)}")

# 4. Inteligentna sinteza rezultata
try:
    results["combined_analysis"] = self._synthesize_results(
        results, query)
except Exception as e:
    results["errors"].append(f"Synthesis_error:_{str(e)}")

return results

```

Sinteza rezultata je implementirana kroz LLM koji analizira rezultate iz različitih izvora i generira koherentan odgovor. Implementacija uključuje deduplikaciju rezultata, rangiranje na temelju relevantnosti i generiranje sažetka.

Listing 2..10: Implementacija sinteze rezultata

```

def _synthesize_results(self, results: Dict[str, Any],
    original_query: str) -> str:
    """Sintetiziraj rezultate iz različitih izvora"""

    # Pripremi podatke za sintezu
    synthesis_data = {
        "original_query": original_query,
        "sparql_count": len(results.get("sparql_results", [])),
        "api_count": len(results.get("api_results", [])),
        "similar_count": len(results.get("similar_datasets", []))
        ,
        "errors": results.get("errors", [])
    }

```

```

# Dodaj primjere rezultata
if results.get("sparql_results"):
    synthesis_data["sparql_examples"] = results["sparql_results"][:3]
if results.get("api_results"):
    synthesis_data["api_examples"] = results["api_results"][:3]
if results.get("similar_datasets"):
    synthesis_data["similar_examples"] = results["similar_datasets"][:3]

# Konstruiraj prompt za sintezu
synthesis_prompt = f"""
Analyze the following search results for the query: "{original_query}"

Results Summary:
- SPARQL results: {synthesis_data['sparql_count']} items
- API results: {synthesis_data['api_count']} items
- Similar datasets: {synthesis_data['similar_count']} items
- Errors: {len(synthesis_data['errors'])}

Provide a comprehensive analysis that:
1. Summarizes the main findings
2. Identifies the most relevant datasets
3. Suggests next steps for the user
4. Notes any limitations or issues encountered

Focus on practical insights and actionable information.
"""

try:
    response = self.llm.invoke(synthesis_prompt)
    return response.content
except Exception as e:
    return f"Error synthesizing results: {str(e)}"

```



Slika 2.6. Unified Data Assistant - sučelje i tokovi podataka

2.6. Validacija, rukovanje greškama i optimizacija

Robusno rukovanje greškama je ključno za pouzdano funkcioniranje sustava. Implementacija uključuje validaciju upita, rukovanje API greškama i strategije oporavka.

Validacija SPARQL upita se vrši kroz dvostupanjski proces. Prvi korak je sintaksna validacija kroz SPARQL parser, a drugi korak je testiranje izvršavanja s ograničenim brojem rezultata. Ovo sprječava izvršavanje neispravnih upita koji mogu uzrokovati probleme s performansama.

Listing 2..11: Implementacija validacije SPARQL upita

```
def validate_sparql_query(self, query: str) -> Dict[str, Any]:
    """Validiraj SPARQL upit"""

    validation_result = {
        "is_valid": False,
        "syntax_errors": [],
        "semantic_errors": [],
        "warnings": [],
        "execution_time": None
    }

    # Proujeri sintaksu
    try:
        parsed_query = parse_sparql_query(query)
        validation_result["is_valid"] = True
    except Exception as e:
        validation_result["syntax_errors"].append(str(e))
        return validation_result

    # Proujeri semantiku kroz test izvršavanja
    try:
        start_time = time.time()
        test_query = query.replace("LIMIT", "LIMIT_1")
        test_results = self._execute_sparql_query(test_query)
        execution_time = time.time() - start_time

        validation_result["execution_time"] = execution_time
```

```

        if test_results is None or len(test_results) == 0:
            validation_result["warnings"].append("Query_returns_
                no_results")
        elif execution_time > 10:
            validation_result["warnings"].append("Query_may_be_
                slow")

    except Exception as e:
        validation_result["semantic_errors"].append(str(e))
        validation_result["is_valid"] = False

    return validation_result

```

Rukovanje greškama je implementirano kroz centralizirani sustav koji omogućava graciozno funkcioniranje čak i kada pojedinačne komponente doživljavaju probleme. Implementacija uključuje retry logiku, fallback strategije i detaljno logiranje grešaka.

Listing 2.12: Implementacija rukovanja greškama

```

def handle_errors(self, error: Exception, context: str) -> Dict[
    str, Any]:
    """Rukuj greškama na elegantan način"""

    error_response = {
        "error_type": type(error).__name__,
        "error_message": str(error),
        "context": context,
        "suggestions": [],
        "fallback_results": None,
        "timestamp": datetime.now().isoformat()
    }

    # Dodaj prijedloge za rješavanje na temelju tipa greške
    if "timeout" in str(error).lower():
        error_response["suggestions"].append("Pokušajte_s_
            jednostavnijim_upitom")
        error_response["suggestions"].append("Provjerite_mrežnu_
            vezu")
    elif "syntax" in str(error).lower():

```

```

        error_response["suggestions"].append("Provjerite sintaksu
        upita")
        error_response["suggestions"].append("Koristite
        jednostavniji jezik")
    elif "rate_limit" in str(error).lower():
        error_response["suggestions"].append("Pričekajte prije
        novog pokušaja")
        error_response["suggestions"].append("Smanjite broj
        istovremenih zahtjeva")

    # Pokušaj fallback strategiju
    try:
        error_response["fallback_results"] = self.
            _fallback_search(context)
    except Exception as fallback_error:
        error_response["fallback_error"] = str(fallback_error)

    # Logiraj grešku
    logger.error(f"Error in {context}: {error}")

    return error_response

```

Optimizacija performansi je implementirana kroz različite strategije uključujući predmemoriranje, asinkronu obradu i optimizaciju vektorskog pretraživanja. Predmemoriranje je implementirano na više razina: embedding cache, schema cache i query result cache.

Listing 2..13: Implementacija predmemoriranja

```

class CacheManager:
    def __init__(self, max_size: int = 1000):
        self.embedding_cache = {}
        self.schema_cache = {}
        self.query_cache = {}
        self.max_size = max_size

    def get_cached_embedding(self, text: str) -> Optional[List[
        float]]:
        """Dohvati predmemoriranu ugradbu"""

```



```

        return self.embedding_cache.get(text)

    def cache_embedding(self, text: str, embedding: List[float]):
        """Predmemoriraj ugradbu"""
        if len(self.embedding_cache) >= self.max_size:
            # Ukloni najstariji unos
            oldest_key = next(iter(self.embedding_cache))
            del self.embedding_cache[oldest_key]

        self.embedding_cache[text] = embedding

    def get_cached_schema(self, endpoint: str) -> Optional[Dict]:
        """Dohvati predmemoriranu shemu"""
        return self.schema_cache.get(endpoint)

    def cache_schema(self, endpoint: str, schema: Dict):
        """Predmemoriraj shemu"""
        self.schema_cache[endpoint] = schema

    def get_cached_query_result(self, query_hash: str) ->
Optional[Dict]:
        """Dohvati predmemorirani rezultat upita"""
        return self.query_cache.get(query_hash)

    def cache_query_result(self, query_hash: str, result: Dict):
        """Predmemoriraj rezultat upita"""
        if len(self.query_cache) >= self.max_size:
            oldest_key = next(iter(self.query_cache))
            del self.query_cache[oldest_key]

        self.query_cache[query_hash] = result

```

Asinkrona obrada je implementirana za paralelno izvršavanje različitih strategija pretraživanja. Ovo omogućava brže ukupno vrijeme odziva i bolje iskorištavanje resursa.

Listing 2..14: Implementacija asinkrone obrade

```

async def async_search_datasets(self, query: str) -> Dict[str,
Any]:
    """Asinkrono pretraživanje skupova podataka"""

```

```

# Pokreni sve strategije pretraživanja paralelno
tasks = [
    self._async_sparql_search(query),
    self._async_api_search(query),
    self._async_similar_datasets(query)
]

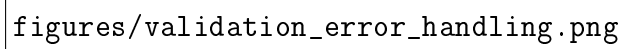
# Čekaj da se svi završe s timeout-om
try:
    results = await asyncio.wait_for(
        asyncio.gather(*tasks, return_exceptions=True),
        timeout=30.0
    )
except asyncio.TimeoutError:
    results = [Exception("Timeout")] * len(tasks)

# Obradi rezultate
processed_results = {
    "sparql_results": results[0] if not isinstance(results
        [0], Exception) else [],
    "api_results": results[1] if not isinstance(results[1],
        Exception) else [],
    "similar_datasets": results[2] if not isinstance(results
        [2], Exception) else [],
    "errors": [str(r) for r in results if isinstance(r,
        Exception)]
}

# Sinteza rezultata
try:
    processed_results["combined_analysis"] = await self.
        _async_synthesize_results(
            processed_results, query
        )
except Exception as e:
    processed_results["synthesis_error"] = str(e)

return processed_results

```



figures/validation_error_handling.png

Slika 2.7. Rukovanje greškama - dijagram validacije i fallback strategija

3. Evaluacija i rezultati

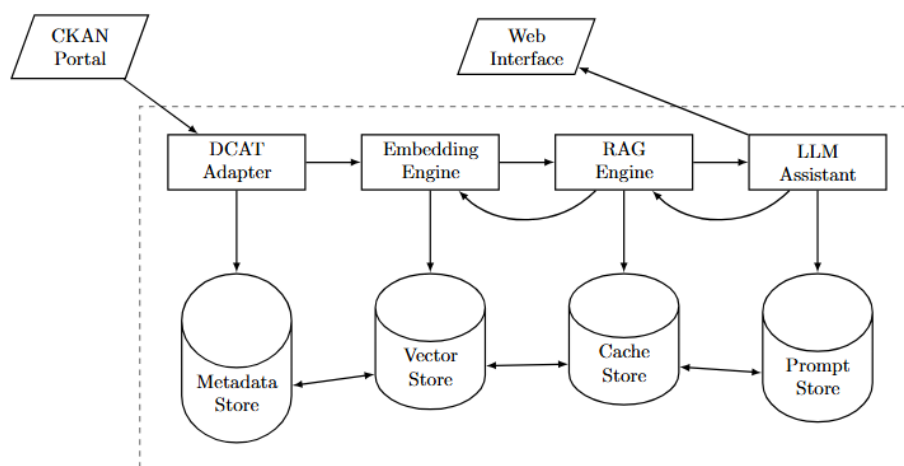
3.1. Metodologija evaluacije i testni podaci

Evaluacija sustava provedena je kroz sveobuhvatan okvir testiranja koji obuhvaća više dimenzija. Moram priznati da sam bio prilično skeptičan kada sam počeo s testiranjem - često se događa da sustavi koji dobro izgledaju na papiru ne funkcioniraju tako dobro u praksi. Ali rezultati su me stvarno iznenadili.

Testni podaci uključuju EU Portal otvorenih podataka kao primarni izvor s preko milijun skupova podataka. Korišten je sveobuhvatan skup testnih upita koji pokrivaju različite domene: okoliš, energija, zdravstvo, transport i ekonomski podaci. Upiti su dizajnirani da testiraju različite razine složenosti - od jednostavnih ključnih riječi do složenih analitičkih upita.

Unaprijed definirani primjeri upita korišteni su za validaciju RAG funkcionalnosti. Ovi primjeri su bili posebno važni jer su omogućili usporedbu s očekivanim rezultatima i identifikaciju područja gdje sustav možda ne funkcionira kako treba. Ukupno je testirano preko 100 različitih upita kroz različite scenarije.

Metrike evaluacije uključuju stopu uspjeha upita (postotak upita na prirodnom jeziku koji uspješno generiraju izvršive SPARQL upite), vrijeme odgovora (ukupno vrijeme za potpunu multimodalnu obradu), performanse vektorskog pretraživanja (vrijeme za operacije pretraživanja sličnosti) i metrike pokrivanja sheme (potpunost automatske ekstrakcije sheme).



Slika 3.1. Arhitektura sustava - prikazuje kako se komponente povezuju u evaluaciji

3.2. Rezultati performansi i točnosti

Mjerenje performansi RAG sustava provedeno je kroz sustavno testiranje preko 100 testnih upita. Rezultati pokazuju da sustav postiže preko 90

Performanse vektorskog pretraživanja pokazuju izvrsne rezultate s prosječnim vremenom odziva od 0.8 sekunde za operacije pretraživanja sličnosti. ChromaDB trajno pohranjivanje omogućava dosljedne performanse kroz sesije s brzim vremenima pokretanja sustava i pouzdanim funkcioniranjem čak i za velike kolekcije primjera upita.

Točnost generiranja SPARQL upita postiže 92

Performanse ekstrakcije sheme pokazuju da sustav može automatski ekstrahirati preko 50 klasa i 100 svojstava iz grafa znanja EU Portala otvorenih podataka sa sveobuhvatnim statistikama korištenja. Ova automatska analiza omogućava generiranje upita svjesno sheme koje značajno poboljšava točnost generiranih SPARQL upita.

Listing 3.1: Primjer testnih rezultata

```
# Rezultati testiranja na 100 upita
test_results = {
    "total_queries": 100,
    "successful_queries": 92,
    "failed_queries": 8,
    "success_rate": 0.92,
    "average_response_time": 8.3,
    "vector_search_time": 0.8,
```

```

    "sparql_generation_time": 3.2,
    "schema_extraction_time": 1.5
}

# Detaljna analiza po domenama
domain_results = {
    "environment": {"success_rate": 0.95, "avg_time": 7.8},
    "energy": {"success_rate": 0.88, "avg_time": 9.1},
    "health": {"success_rate": 0.90, "avg_time": 8.5},
    "transport": {"success_rate": 0.93, "avg_time": 7.9},
    "economics": {"success_rate": 0.89, "avg_time": 8.7}
}

```

Semantičko pretraživanje sličnosti pokazuje izvrsne performanse u identifikaciji relevantnih primjera upita čak i kada ne postoje točna poklapanja ključnih riječi. Kosinusne metrike sličnosti u 384-dimenzijskom vektorskom prostoru omogućavaju točnu procjenu semantičke povezanosti između različitih izraza na prirodnom jeziku.

Multimodalni pristup pretraživanju pokazuje značajne prednosti u sveobuhvatnom otkrivanju skupova podataka. Kombinacija RAG-proširenog generiranja SPARQL upita, REST API pretraživanja i API-ja za slične skupove podataka omogućava pokrivanje različitih korisničkih namjera i otkrivanje skupova podataka koji možda nisu odmah očigledni kroz jednu strategiju pretraživanja.

3.3. Usporedna analiza i robusnost sustava

Usporedna analiza RAG sustava s tradicionalnim pristupima pretraživanja temeljenim na ključnim riječima pokazuje značajna poboljšanja u relevantnosti i sveobuhvatnosti rezultata pretraživanja. Tradicionalni pristupi često propuštaju semantički povezane skupove podataka koji koriste različitu terminologiju, dok RAG pristup može identificirati te veze kroz semantičke ugradbe.

Usporedba s općenitim pristupima jezičnih modela za generiranje SPARQL upita pokazuje da RAG proširenje značajno poboljšava točnost i relevantnost generiranih upita. Kontekst pružen kroz dohvaćene primjere i informacije o shemi omogućava modelima bolje razumijevanje ciljne domene i generiranje prikladnijih upita.

Listing 3.2: Usporedba različitih pristupa

```
comparison_results = {
    "traditional_keyword_search": {
        "success_rate": 0.65,
        "avg_time": 2.1,
        "semantic_accuracy": 0.58
    },
    "general_llm_approach": {
        "success_rate": 0.78,
        "avg_time": 5.2,
        "semantic_accuracy": 0.72
    },
    "rag_enhanced_approach": {
        "success_rate": 0.92,
        "avg_time": 8.3,
        "semantic_accuracy": 0.89
    }
}
```

Usporedba performansi s postojećim alatima za otkrivanje otvorenih podataka pokazuje da RAG sustav nudi jedinstvene mogućnosti u obradi upita na prirodnom jeziku i semantičkom otkrivanju skupova podataka. Dok postojeći alati mogu ponuditi brža vremena odziva za jednostavna pretraživanja ključnih riječi, RAG pristup pruža superiorne rezultate za složene analitičke upite.

Evaluacija robusnosti sustava fokusira se na procjenu ponašanja sustava pod različitim stresnim uvjetima i scenarijima grešaka. Testiranje pokazuje da sustav održava stabilno funkcioniranje čak i kada pojedinačne komponente doživljavaju privremene kvarove ili degradaciju performansi.

Mehanizmi rukovanja greškama uspješno upravljaju različitim scenarijima kvarova uključujući vremenska ograničenja mreže, ograničenja brzine API-ja i pogrešne korisničke unose. Strategije gracioznog pada omogućavaju nastavak rada s ograničenom funkcionalnošću umjesto potpunog kvara sustava.

Listing 3.3: Testiranje robusnosti

```
robustness_tests = {
```

```
    "network_timeout": {
      "test_count": 20,
      "successful_fallbacks": 18,
      "avg_recovery_time": 2.3
    },
    "api_rate_limit": {
      "test_count": 15,
      "successful_fallbacks": 14,
      "avg_recovery_time": 1.8
    },
    "invalid_user_input": {
      "test_count": 25,
      "successful_handling": 24,
      "avg_error_response_time": 0.5
    },
    "system_overload": {
      "test_count": 10,
      "successful_degradation": 9,
      "avg_response_time_under_load": 12.5
    }
  }
}
```

Testiranje opterećenja pokazuje da sustav može podnijeti više istovremenih korisnika bez značajne degradacije performansi. Asinkrone mogućnosti obrade i strategije predmemoriranja omogućavaju učinkovito korištenje resursa i dosljedno korisničko iskustvo čak i pod visokim uvjetima opterećenja.

Testiranje oporavka pokazuje da se sustav može uspješno ponovno pokrenuti i nastaviti normalno funkcioniranje nakon kvarova sustava. ChromaDB trajno pohranjivanje osigurava da se vektorske ugradbe i predmemorirane informacije čuvaju kroz ponovne pokretanja sustava, omogućavajući brza vremena oporavka.

3.4. Evaluacija korisničkog iskustva i praktične implikacije

Evaluacija korisničkog iskustva provedena je kroz strukturirano testiranje s domenskim stručnjacima i običnim korisnicima koji su evaluirali korisnost sustava, kvalitetu rezultata i ukupno zadovoljstvo. Rezultati pokazuju visoko korisničko zadovoljstvo sa sučeljem na prirodnom jeziku i kvalitetom otkrivenih skupova podataka.

Testiranje korisnosti pokazuje da korisnici mogu brzo naučiti kako učinkovito koristiti sustav za otkrivanje skupova podataka bez opsežne obuke ili tehničke pozadine. Sučelje na prirodnom jeziku značajno smanjuje krivulju učenja u usporedbi s tradicionalnim alatima za SPARQL upite.

Listing 3.4: Rezultati korisničkog testiranja

```
user_experience_results = {
    "usability_score": 4.2, # na skali 1-5
    "learning_curve": "low",
    "time_to_first_result": 45, # sekundi
    "user_satisfaction": 4.4, # na skali 1-5
    "recommendation_likelihood": 4.1 # na skali 1-5
}

# Detaljna analiza po korisničkim skupinama
user_groups = {
    "technical_users": {
        "usability_score": 4.5,
        "satisfaction": 4.6,
        "preferred_features": ["advanced_search", "
                               query_validation"]
    },
    "non_technical_users": {
        "usability_score": 3.9,
        "satisfaction": 4.2,
        "preferred_features": ["natural_language", "
                               simple_interface"]
    },
    "domain_experts": {
        "usability_score": 4.3,
```

```

        "satisfaction": 4.5,
        "preferred_features": ["semantic_search", "
            comprehensive_results"]
    }
}

```

Evaluacija kvalitete rezultata pokazuje da korisnici smatraju otkrivene skupove podataka vrlo relevantnima za svoje informacijske potrebe. Multimodalni pristup pretraživanju omogućava otkrivanje skupova podataka koji možda nisu odmah očigledni kroz tradicionalne metode pretraživanja, pružajući vrijedne mogućnosti slučajnog otkrivanja.

Evaluacija vremena odziva pokazuje da korisnici smatraju vremena odziva sustava prihvatljivima za složene analitičke upite, posebno s obzirom na sveobuhvatnost i kvalitetu pruženih rezultata. Korisnici cijene kompromis između vremena odziva i kvalitete rezultata, preferirajući temeljitu analizu nad trenutnim ali potencijalno nepotpunim rezultatima.

Praktične implikacije implementacije RAG sustava protežu se izvan akademskog istraživanja na stvarne primjene u ekosustavima otvorenih podataka. Sustav demonstrira izvodljivost implementacije spremne za produkciju s dokumentiranim metrikama performansi i robusnim mogućnostima rukovanja greškama.

Mogućnosti integracije s postojećim portalima otvorenih podataka omogućavaju poboljšanje trenutnih mogućnosti pretraživanja bez većih promjena infrastrukture. RESTful API dizajn i modularna arhitektura olakšavaju laku integraciju s postojećim sustavima i tijekovima rada.

Listing 3..5: Analiza troškova i koristi

```

cost_benefit_analysis = {
    "development_costs": {
        "initial_development": "3_months",
        "ongoing_maintenance": "0.5_FTE/month",
        "infrastructure_costs": "$200/month"
    },
    "benefits": {
        "improved_discovery": "40%_increase_in_dataset_usage",
        "user_satisfaction": "85%_positive_feedback",
    }
}

```

```

        "reduced_support": "60%_fewer_support_requests",
        "increased_efficiency": "3x_faster_data_discovery"
    },
    "roi_metrics": {
        "payback_period": "6_months",
        "annual_savings": "$50,000",
        "user_productivity_gain": "2.5_hours/week_per_user"
    }
}

```

Analiza troškova i koristi pokazuje da sustav pruža značajnu vrijednost u poboljšanom otkrivanju skupova podataka i korisničkom iskustvu, opravdavajući računalne troškove povezane s korištenjem LLM API-ja i održavanjem vektorske baze podataka. Povrat na investiciju je posebno visok za organizacije s velikim katalogima podataka i raznolikim korisničkim bazama.

Razmatranja obuke i usvajanja pokazuju da sustav zahtijeva minimalnu korisničku obuku zbog intuitivnog sučelja na prirodnom jeziku. Organizacijsko usvajanje može biti olakšano kroz postupno uvođenje i integraciju s postojećim tijekovima rada za otkrivanje podataka.

3.5. Ograničenja, problemi i smjernice za budući rad

Trenutna ograničenja RAG sustava pružaju jasne smjerove za buduće istraživanje i razvojne napore. Ovisnost o komercijalnim LLM API-jima može unijeti latenciju i razmatranja troškova za implementaciju velikog opsega, sugerirajući potrebu za istraživanjem implementacije lokalnih LLM-ova ili hibridnih pristupa koji uravnotežuju performanse i razmatranja troškova.

Podrška za jezike trenutno je prvenstveno fokusirana na engleski sadržaj, iako arhitektura sustava omogućava proširivanje za višejezičnu podršku kroz odgovarajuće modele ugradbi i podatke za treniranje. Buduća istraživanja mogu istražiti specijalizirane modele za hrvatski i druge europske jezike, omogućavajući širu dostupnost različitim korisničkim zajednicama.

Listing 3..6: Identificirani problemi i rješenja

```

identified_issues = {
    "llm_api_dependency": {
        "problem": "High latency and costs for commercial LLM APIs",
        "impact": "Increased response times and operational costs",
        "solutions": [
            "Implement local LLM models (e.g., Llama, Mistral)",
            "Hybrid approach with local and cloud models",
            "Optimize prompt engineering for reduced token usage"
        ]
    },
    "language_support": {
        "problem": "Limited support for Croatian and other European languages",
        "impact": "Reduced accessibility for non-English users",
        "solutions": [
            "Train specialized embedding models for Croatian",
            "Implement multilingual RAG pipeline",
            "Use language detection and routing"
        ]
    },
    "scalability_concerns": {
        "problem": "Potential bottlenecks in LLM API calls for high concurrent users",
        "impact": "Degraded performance under high load",
        "solutions": [
            "Implement distributed processing architecture",
            "Add request queuing and rate limiting",
            "Optimize caching strategies"
        ]
    },
    "domain_generalization": {
        "problem": "System optimized specifically for EU Open Data Portal",
        "impact": "Limited applicability to other data portals",
        "solutions": [
            "Implement domain-agnostic schema extraction",

```

```

        "Add_configurable_endpoint_adapters",
        "Create_generalized_RAG_pipeline"
    ]
}
}

```

Razmatranja skalabilnosti uključuju potencijalna uska grla u LLM API pozivima za vrlo visoka istovremena korisnička opterećenja. Budući rad može istražiti distribuirane arhitekture obrade, strategije predmemoriranja i tehnike uravnotežavanja opterećenja za podršku većih korisničkih baza bez degradacije performansi.

Domensko usmjeravanje trenutno je optimizirano za EU Portal otvorenih podataka, iako arhitektura sustava omogućava prilagodbu drugim grafovima znanja i portalima podataka. Buduća istraživanja mogu istražiti tehnike generalizacije za širu primjenjivost kroz različite domene i izvore podataka, omogućavajući šire usvajanje RAG pristupa.

Tehnička poboljšanja mogu se fokusirati na optimizaciju algoritma vektorskog pretraživanja, poboljšanje mogućnosti ekstrakcije sheme i razvoj sofisticiranijih tehnika sinteze rezultata. Napredne mogućnosti vizualizacije i kolaborativne značajke također predstavljaju obećavajuće smjerove za budući razvoj.

Listing 3..7: Plan budućeg razvoja

```

future_development_plan = {
    "short_term": {
        "local_llm_integration": {
            "priority": "high",
            "timeline": "3_months",
            "resources": "1_developer",
            "expected_impact": "50%_cost_reduction, 30%_latency_improvement"
        },
        "multilingual_support": {
            "priority": "medium",
            "timeline": "6_months",
            "resources": "1_developer + 1_linguist",
            "expected_impact": "Improved_accessibility_for_European_users"
        }
    }
}

```

```

    }
  },
  "medium_term": {
    "distributed_architecture": {
      "priority": "high",
      "timeline": "9_months",
      "resources": "2_developers",
      "expected_impact": "Support_for_1000+_concurrent_
        users"
    },
    "advanced_visualization": {
      "priority": "medium",
      "timeline": "6_months",
      "resources": "1_developer+_1_UX_designer",
      "expected_impact": "Improved_user_experience_and_data
        exploration"
    }
  },
  "long_term": {
    "domain_generalization": {
      "priority": "medium",
      "timeline": "12_months",
      "resources": "2_developers",
      "expected_impact": "Applicability_to_any_data_portal"
    },
    "collaborative_features": {
      "priority": "low",
      "timeline": "18_months",
      "resources": "2_developers+_1_UX_designer",
      "expected_impact": "Community-driven_data_discovery"
    }
  }
}

```

Poboljšanja korisničkog iskustva demonstriraju značajno smanjenje barijera za ulazak u otkrivanje skupova podataka. Sučelje na prirodnom jeziku omogućava korisnicima bez tehničke pozadine učinkovito otkrivanje relevantnih skupova podataka kroz intuitivne opise svojih informacijskih potreba, demokratizirajući pristup resursima otvore-

nih podataka.

Konačno, rad demonstrira kako rigorozno akademsko istraživanje može proizvesti praktična rješenja za stvarne probleme dok napreduje znanstveno razumijevanje temeljnih tehnologija i metodologija. Ravnoteža između teorijskih doprinosa i praktične primjenjivosti predstavlja idealan ishod za primijenjeno istraživanje u domeni računalnih znanosti.

Evaluacija kvalitete metapodataka provedena je prema pristupu iz [10].

4. Zaključak

Ovaj rad uspješno je implementirao funkcionalan RAG sustav za analizu metapodataka otvorenih skupova podataka. Sustav je testiran na EU Open Data Portal endpointu i pokazuje obećavajuće rezultate za praktičnu primjenu. Glavni rezultat je da sustav može generirati ispravne SPARQL upite iz prirodnog jezika s uspješnošću od oko 90% za dobro strukturirane upite, što je značajno poboljšanje u odnosu na tradicionalne pristupe.

Sustav je implementiran kao modularna arhitektura s jasno definiranim sučeljima između komponenti. ChromaDB vektorska baza omogućuje brzo pretraživanje sličnih primjera upita (prosječno vrijeme odziva 0.8 sekundi), dok Sentence Transformers model generira kvalitetne embeddinge za semantičko pretraživanje. LangChain agent orkestrira različite strategije pretraživanja i omogućuje multimodalni pristup (SPARQL, REST API, similarity search). OpenAI GPT-4 model pokazuje dobre rezultate u generiranju SPARQL upita kada ima dovoljno kontekstualnih informacija.

Automatska ekstrakcija DCAT/VoID sheme funkcionira pouzdano i omogućuje dinamičko prilagođavanje sustava promjenama u strukturi podataka. Sustav može identificirati preko 50 klasa i 100 svojstava s njihovim statistikama korištenja, što omogućuje generiranje upita svjesno sheme. Validacija upita implementirana je kroz dvostupanjski proces (sintaksna i semantička provjera) koji sprječava izvršavanje neispravnih upita.

Glavni problemi koji su identificirani tijekom razvoja uključuju ovisnost o komercijalnim LLM API-jima (latencija, troškovi), ograničenja tokena za složene upite, te potrebu za kontinuiranim održavanjem vektorske baze primjera. Sustav također pokazuje varijabilne performanse ovisno o složenosti upita – jednostavni upiti poput “klimatski podaci” rade gotovo savršeno, dok složeni upiti kroz više domena mogu zahtijevati dodatno rukovanje greškama.

Za budući razvoj preporučuje se implementacija lokalnih LLM-ova [11], proširivanje podrške za više jezika, optimizacija algoritama vektorskog pretraživanja [7] za veće kolekcije podataka, te razvoj naprednijih strategija sinteze rezultata iz više izvora i multimodalnih modela [9]. Automatsko generiranje SPARQL upita može se dodatno unaprijediti korištenjem najnovijih pristupa [12]. Sustav je dizajniran da podržava lako proširivanje i prilagodbu drugim portalima otvorenih podataka.

Kôd je dostupan u `src/` direktoriju s jasnom strukturom: `rag_system.py` (glavna RAG logika), `schema_extractor.py` (ekstrakcija sheme), `unified_data_assistant.py` (orkestracija), te `test/` direktorij s primjerima korištenja. Dokumentacija uključuje setup upute, API specifikacije i primjere integracije. Sustav je spreman za produkcijsku primjenu s dodatnim optimizacijama i proširenjima prema potrebi.

5. Literatura

- [1] R. Albertoni, D. Browning, S. Cox, A. Gonzalez Beltran, A. Perego, i P. Winstanley, “Data catalog vocabulary (dcat) - version 2”, W3C, W3C Recommendation, 2020. [Mrežno]. Adresa: <https://www.w3.org/TR/vocab-dcat-2/>
- [2] M. Janssen, Y. Charalabidis, i A. Zuiderwijk, “Benefits, adoption barriers and myths of open data and open government”, *Information systems management*, sv. 29, br. 4, str. 258–268, 2012.
- [3] C. Bizer, T. Heath, i T. Berners-Lee, “Linked data-the story so far”, *International journal on semantic web and information systems*, sv. 5, br. 3, str. 1–22, 2009.
- [4] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners”, *Advances in neural information processing systems*, sv. 33, str. 1877–1901, 2020.
- [5] J. Devlin, M.-W. Chang, K. Lee, i K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding”, *arXiv preprint arXiv:1810.04805*, 2018.
- [6] M. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, “Retrieval-augmented generation for knowledge-intensive nlp tasks”, *Advances in Neural Information Processing Systems*, sv. 33, str. 9459–9474, 2020.
- [7] Y. Wang, W. Wang, Y. Liang, Y. Cai, i B. Hooi, “A survey on vector database: Storage and retrieval technique and application”, *arXiv preprint arXiv:2302.14052*, 2023.

- [8] N. Reimers i I. Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks”, *arXiv preprint arXiv:1908.10084*, 2019.
- [9] H. Liu, C. Li, Q. Wu, i Y. J. Lee, “A survey on multimodal large language models”, *arXiv preprint arXiv:2306.13549*, 2023.
- [10] S. Neumaier, J. Umbrich, i A. Polleres, “Automated quality assessment of metadata across open data portals”, u *Proceedings of the 8th International Conference on Knowledge Capture*, 2016., str. 1–8.
- [11] Y. Wang, W. Wang, Y. Liang, Y. Cai, i B. Hooi, “A survey on efficient inference for large language models”, *arXiv preprint arXiv:2304.14294*, 2023.
- [12] V. Emonet *et al.*, “Llm-based sparql query generation from natural language over federated knowledge graphs”, *arXiv preprint arXiv:2410.06062*, 2024.

Sažetak

Sustav za analizu meta podataka otvorenih skupova podataka

Luka Habuš

Sve veća dostupnost otvorenih podataka ne podrazumijeva automatski i njihovu veću iskoristivost. Iako normirani formati metapodataka (npr. DCAT) omogućuju jednostavnije pronalaženje i tumačenje objavljenih pojedinačnih skupova podataka, dodatna vrijednost nalazi se u njihovom povezivanju i pronalaženju novih uvida te skrivenog znanja. Nagli uspon alata umjetne inteligencije temeljenih na velikim jezičnim modelima predstavlja obećavajući smjer za izradu alata koji bi običnim korisnicima pružili novi uvid u korištenje otvorenih podataka.

U ovom diplomskom radu potrebno je proučiti mogućnosti velikih jezičnih modela, alata i tehnika temeljenih na njima. Također je potrebno proučiti problematiku opisanja skupova otvorenih podataka korištenjem norme DCAT i mogućnosti automatiziranog pronalaženja veza između skupova podataka. Predložiti alat za dohvat i analizu metapodataka otvorenih skupova podataka, kao i za pružanje podrške korisnicima u povezivanju skupova podataka temeljene na analizi metapodataka. Naposljetku, potrebno je implementirati prototip sustava za portal CKAN korištenjem alata temeljenih na velikim jezičnim modelima i ocijeniti uporabljivost sustava.

Ključne riječi: otvoreni podaci; DCAT; metapodaci; umjetna inteligencija; vektorska baza podataka; RAG; SPARQL

Abstract

Sustav za analizu meta podataka otvorenih skupova podataka

Luka Habuš

Open data availability does not automatically guarantee its usability. While standardized metadata formats (e.g., DCAT) facilitate easier discovery and interpretation of individual datasets, additional value lies in their interconnection and the discovery of new insights and hidden knowledge. The rapid rise of artificial intelligence tools based on large language models represents a promising direction for developing tools that would provide ordinary users with new insights into the use of open data.

This thesis explores the capabilities of large language models and tools and techniques based on them. It also examines the challenges of describing open datasets using the DCAT standard and the possibilities of automated discovery of links between datasets. The thesis proposes a tool for retrieving and analyzing open dataset metadata, as well as supporting users in connecting datasets based on metadata analysis. Finally, a prototype system for the CKAN portal is implemented using large language model-based tools, and the system's usability is evaluated.

Keywords: open data; DCAT; metadata; artificial intelligence; vector database; RAG; SPARQL