

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 874

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

Luka Habuš

Zagreb, srpanj 2025.

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 874

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

Luka Habuš

Zagreb, srpanj 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	2
2. Dizajn i implementacija sustava	4
2.1. Korištene tehnologije	4
2.2. Ključne funkcionalnosti	5
2.3. Arhitektura i skalabilnost	6
2.4. Implementacija RAG sustava i vektorske baze	7
2.5. Generiranje SPARQL upita i prompt engineering	11
2.6. Ekstrakcija sheme i DCAT analiza	14
2.7. Unified Data Assistant i multimodalno pretraživanje	18
2.8. Validacija, rukovanje greškama i optimizacija	23
3. Evaluacija i rezultati	30
3.1. Metodologija evaluacije i testni podaci	30
3.2. Rezultati performansi i točnosti	30
3.3. Usporedna analiza i robusnost sustava	32
3.4. Ograničenja, problemi i smjernice za budući rad	34
4. Zaključak	36
5. Literatura	38
Sažetak	40
Abstract	41

1. Uvod

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].

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.

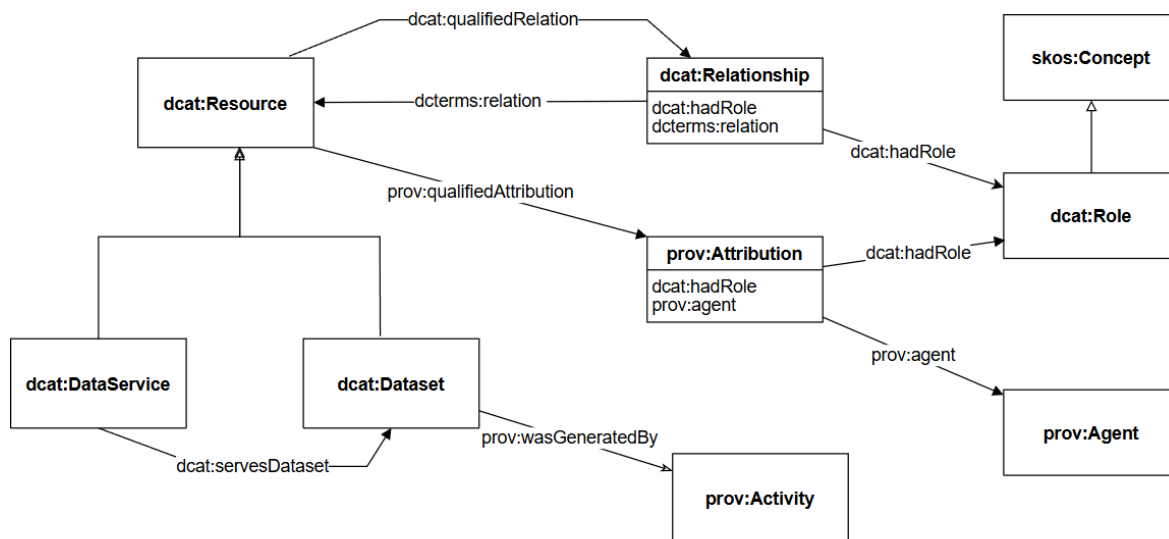
Ovaj rad fokusira se na razvoj sustava za analizu metapodataka otvorenih skupova podataka koji kombinira napredne tehnike umjetne inteligencije s tradicionalnim metodama obrade podataka. Glavni cilj je stvoriti alat koji omogućava korisnicima bez tehničke pozadine da učinkovito otkrivaju, analiziraju 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, dizajn i implementaciju sustava, evaluaciju rezultata i zaključak. Svako poglavlje fokusira se na specifične aspekte razvoja i evaluacije sustava, pružajući sveobuhvatan pregled cijelog projekta.

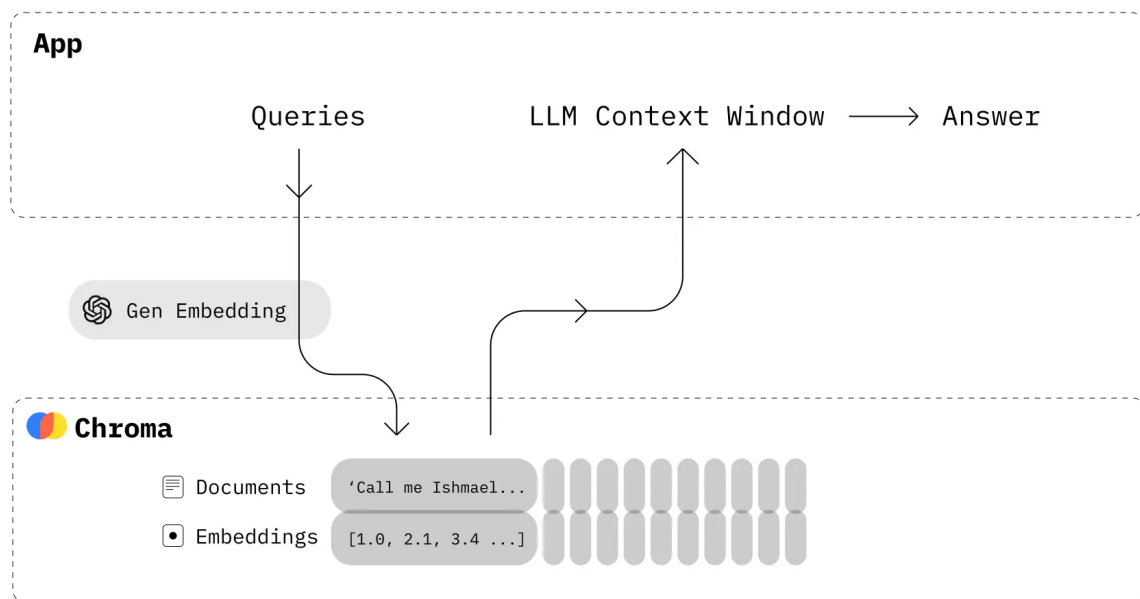


Slika 1.1. DCAT standard - prikazuje osnovne klase i svojstva metapodataka

2. Dizajn i implementacija sustava

2.1. 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.



Slika 2.1. ChromaDB arhitektura - prikazuje komponente vektorske baze podataka

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.

2.2. 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.

```
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

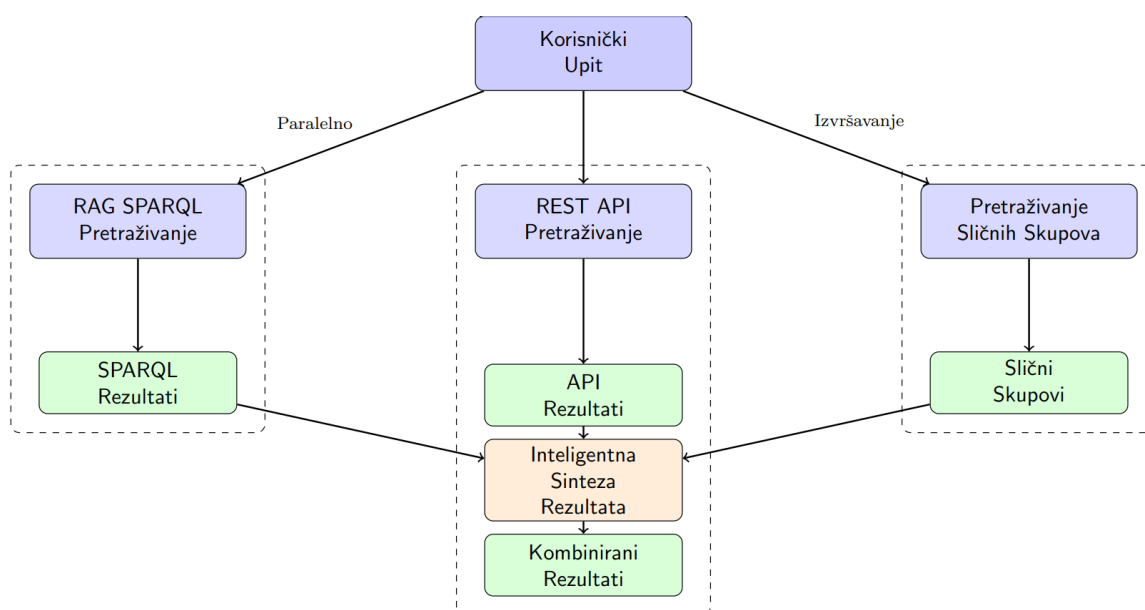
Success rate: 3/3 (100%)

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 2.2. Ekstrakcija sheme i validacija generiranog SPARQL upita

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.



Slika 2.3. Multimodalno pretraživanje - prikazuje različite strategije i njihovu integraciju

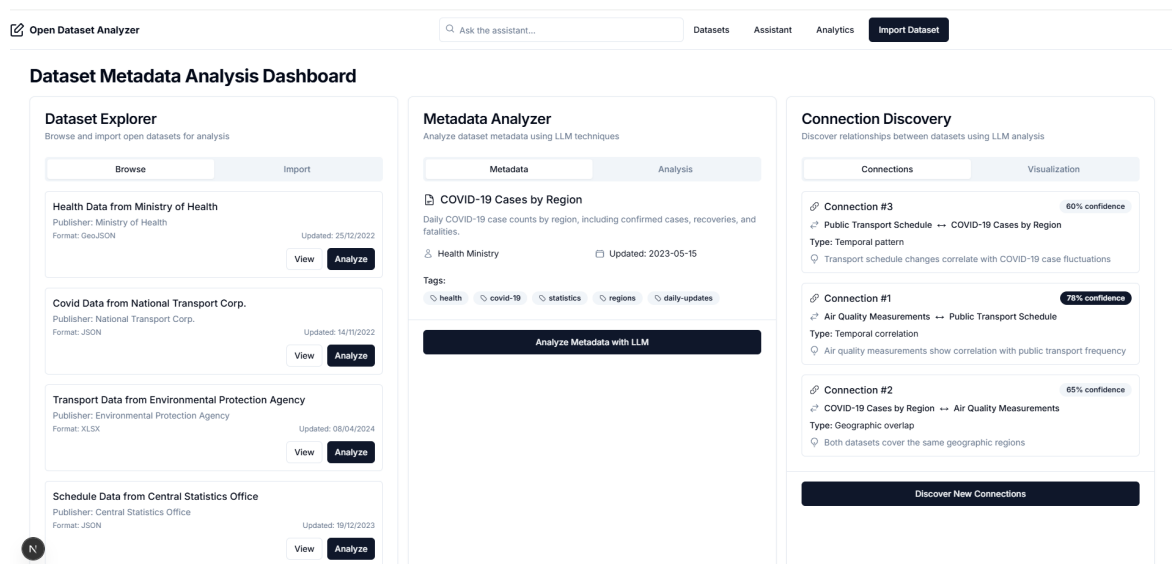
2.3. Arhitektura i skalabilnost

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.

```
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:mediaType'.
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 2.4. Skalabilnost sustava - prikazuje mogućnosti proširivanja i optimizacije

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.



Slika 2.5. Korisničko sučelje - prikazuje intuitivan dizajn i funkcionalnosti

Sustav je dizajniran kao modularna arhitektura s jasno definiranim sučeljima između komponenti. Glavni razlog za ovakav pristup je mogućnost nezavisnog razvoja i testira-

nja pojedinačnih dijelova, što se pokazalo ključnim tijekom implementacije kada je bilo potrebno ispravljati 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 pohranu meta podataka, što je bilo važno za spremanje informacija o SPARQL upitima i shemi. U praksi se pokazalo da ChromaDB može učinkovito rukovati s velikim brojem 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 podatke s EU Portala.

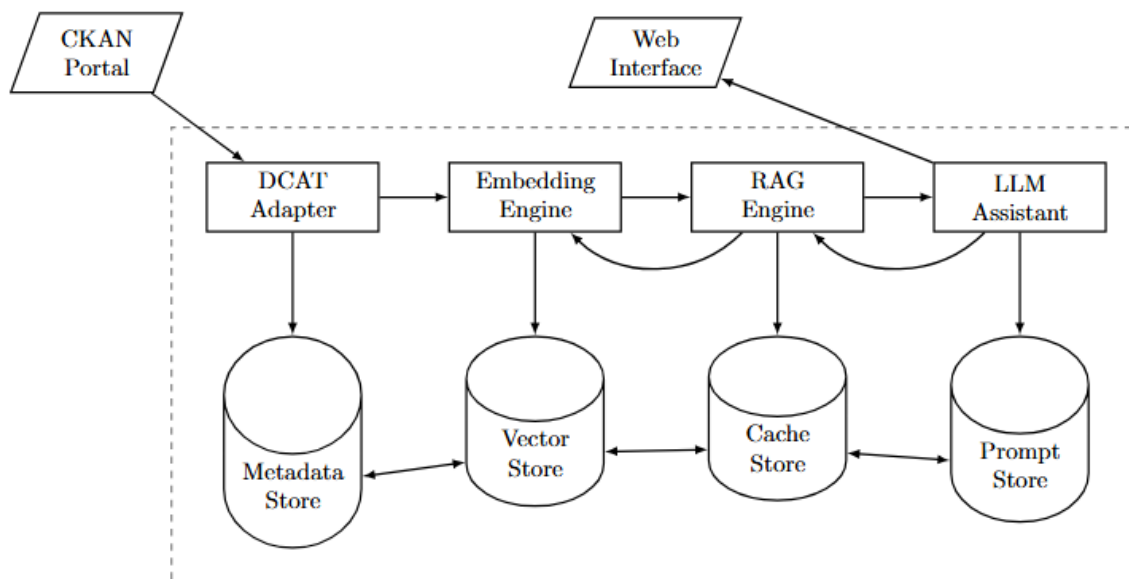
LangChain okvir [9] korišten je za orkestraciju različitih komponenti sustava i upravljanje agentima koji omogućavaju složene zadatke poput multimodalnog pretraživanja i inteligentne sinteze rezultata.

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.

2.4. Implementacija RAG sustava i vektorske baze

RAG sustav [6] je implementiran u src/rag_system.py i predstavlja srce cijelog sustava. Glavna klasa RAGSystem upravlja svim aspektima RAG pipeline-a od generiranja embeddinga do izvršavanja upita.

Inicijalizacija sustava uključuje postavljanje ChromaDB klijenta, učitavanje Sentence



Slika 2.6. Arhitektura sustava - prikazuje tok podataka kroz različite komponente

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 Sentence Transformers model. Implementacija uključuje caching mehanizam koji sprečava ponovno generiranje embeddinga za iste tekstove. Caching je implementiran kao jednostavan dictionary u memoriji, ali za produkcijsku upotrebu bi trebao biti implementiran kao Redis ili sličan sustav.

Dodavanje primjera upita u vektorsku bazu se vrši kroz metodu `add_query_example`. Svaki primjer se sastoji od pitanja na prirodnom jeziku, odgovarajućeg SPARQL upita, opisa i tagova. Embedding se generira za pitanje, a metadata se sprema zajedno s embeddingom. Ovo omogućava kasnije pretraživanje sličnih primjera na temelju semantičke sličnosti.

Primjer 2.1: Implementacija dodavanja primjera u vektorsku bazu

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

    # Provjeri 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 koristi kosinusnu sličnost za pronalaženje najsličnijih primjera u vektorskom prostoru. Implementacija uključuje opciju za filtriranje rezultata na temelju tagova ili endpointa, što može biti korisno za specifične domene.

Primjer 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)

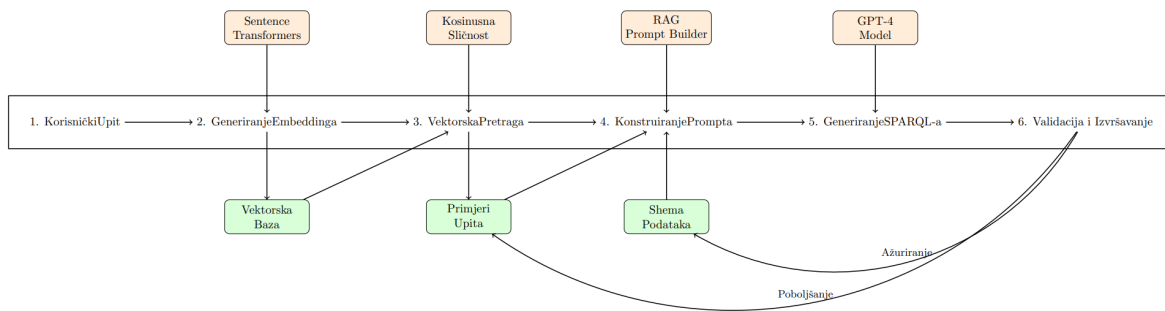
# Postavi filtere ako su zadani
where_clause = None
if filters:
    where_clause = {}
    if 'tags' in filters:
        where_clause['tags'] = {'$contains': filters['tags']}
    if 'endpoint' in filters:
        where_clause['endpoint'] = filters['endpoint']

# 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

```



Slika 2.7. RAG pipeline - prikazuje tok od upita do generiranog SPARQL-a

2.5. 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_rag_prompt` koja kombinira korisnički upit s dohvaćenim primjerima i informacijama o shemi. Prompt je strukturiran u nekoliko sekcija: kontekst upita, slični primjeri, informacije o shemi i instrukcije za generiranje. Ova struktura se pokazala kao ključna za kvalitetu generiranih upita.

Primjer 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 OpenAI GPT-4 model. Implementacija uključuje **error handling i retry** logiku za slučaj API grešaka. Također, implementiran je mehanizam za validaciju generiranih upita prije izvršavanja.

Primjer 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
        }

```

```

Analyzing Dataset 2:
=====
Raw Dataset Information:
Question: Find datasets about air quality in Germany
Description: Search for air quality datasets in Germany
Endpoint: https://data.europa.eu/sparql

Complete SPARQL Query:
=====
PREFIX dct: <http://purl.org/dc/terms/>
PREFIX dcat: <http://www.w3.org/ns/dcat#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
SELECT DISTINCT ?dataset ?title ?description
WHERE {
  ?dataset a dcat:Dataset .
  ?dataset dct:title ?title .
  OPTIONAL { ?dataset dct:description ?description . }
  FILTER (
    CONTAINS(LCASE(STR(?title)), "air quality") ||
    (BOUND(?description) && CONTAINS(LCASE(STR(?description)), "air quality")) ||
    EXISTS { ?dataset dcat:keyword ?kw . FILTER(CONTAINS(LCASE(STR(?kw)), "air quality")) }
  )
  FILTER (
    CONTAINS(LCASE(STR(?title)), "germany") ||
    (BOUND(?description) && CONTAINS(LCASE(STR(?description)), "germany")) ||
    EXISTS { ?dataset dct:spatial ?spatialUri . OPTIONAL {?spatialUri skos:prefLabel ?spatialLabelEN . FILTER(LANGMATCHES(LANG(?spatialLabelEN), "en"))} OPTIONAL {?spatialUri rdfs:label ?spatialLabel .} FILTER(CONTAINS(LCASE(STR(?spatialLabelEN)), "germany") || CONTAINS(LCASE(STR(?spatialLabel)), "germany")) } ||
    EXISTS { ?dataset dcat:keyword ?kw_loc . FILTER(CONTAINS(LCASE(STR(?kw_loc)), "germany")) }
  )
  FILTER(LANGMATCHES(LANG(?title), "en") || LANG(?title) = "")
}
LIMIT 15
=====

```

Slika 2.8. Generiranje SPARQL upita

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 sprečava izvršavanje neispravnih upita koji mogu uzrokovati probleme s performansama.

2.6. Ekstrakcija sheme i DCAT analiza

Automatska ekstrakcija sheme iz SPARQL endpointa je ključna komponenta sustava koja omogućava dinamičko prilagođavanje promjenama u strukturi podataka. Imple-

```

65 # --- Example Natural Language Queries ---
66
67 example_nl_queries = [
68     "Find datasets about climate change tagged with 'environment'. Retrieve their titles and descriptions.",
69     "Show me datasets about transport published since the beginning of 2023. Include their title and publication date.",
70     "List datasets published by the European Environment Agency, showing their title and landing page.",
71     "I need datasets about 'health' that are available in CSV format. Give me their titles and download URLs if available.",
72 ]
73

```

Slika 2.9. Primjer prompta - struktura i elementi prompta za LLM

mentacija je specifično prilagođena za EU Portal otvorenih podataka i DCAT metapodatke.

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.

Primjer 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.

Primjer 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.

Primjer 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": []}

```

2.7. Unified Data Assistant i multimodalno pretraživanje

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

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.

Primjer 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.

```

Your query: gdp per capita and electricity consumption

Processing query: 'gdp per capita and electricity consumption'
-----
Step 1: Analyzing your query...
INFO:httpx:HTTP Request: POST https://api.openai.com/v1/completions "HTTP/1.1 200 OK"
Intent: Compare
Domain: Economy
Location: Global
Time Period: Any
Data Format: Any
Complexity: Simple

Step 2: Extracting semantic concepts...
INFO:httpx:HTTP Request: POST https://api.openai.com/v1/completions "HTTP/1.1 200 OK"
INFO:_main_:Extracted semantic concepts: {'main_topics': ['GDP per capita', 'electricity consumption', 'economic development'], 'variables': ['GDP per capita', 'electricity consumption per capita', 'energy usage', 'economic growth'], 'geographic': ['global', 'worldwide', 'all countries'], 'related_terms': ['economic indicators', 'energy statistics', 'economic data', 'development indicators']}
Main Topics: GDP per capita, electricity consumption, economic development
Key Variables: GDP per capita, electricity consumption per capita, energy usage, economic growth
Geographic Focus: global, worldwide, all countries
Related Terms: economic indicators, energy statistics, economic data, development indicators

Step 3: Generating semantic search queries...
INFO:_main_:Generated 12 simple search queries: ['capita', 'electricity', 'consumption', 'electricity consumption', 'economic', 'development', 'economic development', 'energy', 'usage', 'gdp electricity', 'economic indicators', 'energy statistics']
Generated 10 search queries:
1. capita
2. electricity
3. consumption
4. electricity consumption
5. economic
... and 5 more

```

Slika 2.10. Pokretanje složenijeg upita

Primjer 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.

Primjer 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", [])),

```

```

INFO: __main__:Trying SPARQL search for: 'capita'
INFO: __main__:Executing semantic SPARQL query
INFO: __main__:Executing search 2/10: 'electricity'
INFO: __main__:Querying EU API with FIXED parameters: https://data.europa.eu/api/hub/search/search
INFO: __main__:Parameters: {'q': 'electricity', 'rows': 20, 'wt': 'json'}
INFO: __main__:API Response Status: 200
INFO: __main__:API Response received, processing data...
INFO: __main__:Response keys: ['result']
INFO: __main__:Found results in data.result.results: 10
INFO: __main__:Result 1: Generation Storage Park...
INFO: __main__:Result 2: Ø27 - ORES Electricity - Connecting meters installed...
INFO: __main__:Result 3: INSPIRE Electricity Network WFS...
INFO: __main__:Processing 10 results from API
INFO: __main__:Filtered out irrelevant dataset: 'generation storage park...' (score: 1)
INFO: __main__:Dataset 'Ø27 - ores electricity - connecting meters install...' scored 3 for query 'electricity'
INFO: __main__:Dataset 'inspire electricity network wfs...' scored 3 for query 'electricity'
INFO: __main__:Dataset 'inspire electricity network wms...' scored 3 for query 'electricity'
INFO: __main__:Dataset 'Ø21 - ores electricity - osp budget meters and act...' scored 3 for query 'electricity'
INFO: __main__:Dataset 'sw amberg wfs network areas electricity, gas and w...' scored 3 for query 'electricity'
INFO: __main__:Dataset 'sw amberg wms network areas electricity, gas and w...' scored 3 for query 'electricity'
INFO: __main__:Dataset 'msm01 - electricity output...' scored 3 for query 'electricity'
INFO: __main__:Dataset 'Ø01 - ores electricity - consumption by locality (...)' scored 3 for query 'electricity'
INFO: __main__:Dataset 'Ø06 - ores electricity - mt consumption by busines...' scored 3 for query 'electricity'
INFO: __main__:Converted to 10 dataset objects, filtered to 9 relevant ones
INFO: __main__:Found 9 results for query: 'electricity'
INFO: __main__:Executing search 3/10: 'consumption'
INFO: __main__:Querying EU API with FIXED parameters: https://data.europa.eu/api/hub/search/search
INFO: __main__:Parameters: {'q': 'consumption', 'rows': 20, 'wt': 'json'}
INFO: __main__:API Response Status: 200
INFO: __main__:API Response received, processing data...

```

Slika 2.11. Multimodalno pretraživanje

```

        "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)}"

```
=====
DATASET 9: Land Credit: Water consumption
=====
SEMANTIC RELEVANCE SCORE: 36
FOUND WITH QUERY: 'consumption'
DESCRIPTION: Water consumption

(1) The 2016 scope takes into account the subsidiaries Crédit Foncier Immobilier and Crédit Foncier - Expertise
(2) The 2015 scope takes into account Crédit Foncier Immobilier, Crédit Foncier Immobilier - Expertise and SOCFIM
PUBLISHER: Unknown
SOURCE: EU Open Data Portal API
=====

DATASET 10: Expenditure per Capita
=====
SEMANTIC RELEVANCE SCORE: 30
FOUND WITH QUERY: 'capita'
DESCRIPTION: The data on expenditure under the various social protection schemes are drawn up according to the ESSPROS (European System of Integrated Social Protection Statistics) Manual issued by Eurostat. Generally, the objectives of ESSPROS are to provide a comprehensive, realistic and coherent description of...
PUBLISHER: National Statistics Office
SOURCE: EU Open Data Portal API
=====

COMPREHENSIVE ANALYSIS
=====
INFO:httpx:HTTP Request: POST https://api.openai.com/v1/completions "HTTP/1.1 200 OK"
Thank you for your query regarding GDP per capita and electricity consumption. Based on semantic analysis, we have identified that your intent is to compare these two variables in the domain of economy, with a global focus and any time period. Our search has resulted in 10 datasets from the EU Open Data Portal that are relevant to your query.

The datasets found include information on electricity consumption, energy usage, and economic indicators from various countries and regions. These datasets can provide valuable insights into the relationship between GDP per capita and electricity consumption, and how it impacts economic development.

The most relevant datasets based on the identified semantic concepts are "001 - ORES Electricity - Consumption by locality (Annual)" and "006 - ORES Electricity - MT Consumption by Business Line (Annual)". These datasets provide detailed information on electricity consumption and its impact on economic development in different regions.

The most effective search queries used were "capita", "electricity", "consumption", "electricity consumption", "economic". These terms were able to capture the key variables and geographic focus of your query, resulting in the most relevant datasets.

To access and use the data, you can click on the dataset titles to view the full details and download the data in various formats. You can also use the search filters on the EU
```

Slika 2.12. Rangiranje kombiniranih rezultata

2.8. 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.

Primjer 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")
```

```

        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.

Primjer 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.

Primjer 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.

Primjer 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

```



```

PREFIX dct: <http://purl.org/dc/terms/>
PREFIX dcat: <http://www.w3.org/ns/dcat#>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>

SELECT ?dataset ?title ?description ?publisher ?landingPage
WHERE {
  ?dataset a dcat:Dataset ;
    dct:title ?title ;
    dct:description ?description ;
    dct:issued ?issued ;
    dct:publisher ?pub .
  OPTIONAL { ?pub foaf:name ?publisher }
  OPTIONAL { ?dataset dcat:landingPage ?landingPage }
  FILTER(CONTAINS(LCASE(STR(?title)), "agricultural subsidies") ||
    CONTAINS(LCASE(STR(?description)), "agricultural subsidies") ||
    CONTAINS(LCASE(STR(?title)), "nuts 2") ||
    CONTAINS(LCASE(STR(?description)), "nuts 2") ||
    CONTAINS(LCASE(STR(?title)), "eu") ||
    CONTAINS(LCASE(STR(?description)), "eu"))
  FILTER(?issued >= xsd:date(str(year(now()) - 5) || "-01-01"))
}
LIMIT 10

```

Invoking: 'generate_sparql_tool' with '{'natural_language_query': 'Find datasets on regional GDP growth within the EU over the last 5 years.', 'context': 'Previous query returned unrelated results. Focus on GDP growth and regional economic data.'}'

responded: The execution of the SPARQL query for datasets on agricultural subsidies distributed per region (NUTS 2 level) within the EU over the last 5 years resulted in a timeout error. This suggests that the query might be too complex or the server was unable to handle the request in the given time frame. I will attempt to refine and re-execute this query.

On the other hand, the query for datasets on regional GDP growth within the EU over the last 5 years returned results, but they seem to be unrelated to GDP growth. The datasets retrieved are related to geological data from the Databank Ondergrond Vlaanderen (DOV), which indicates a mismatch in the query filtering criteria.

Here's the plan:

1. **Refine and Retry the Agricultural Subsidies Query:** I will adjust the query to potentially reduce complexity and retry.
2. **Refine the GDP Growth Query:** I will refine the query to better target datasets related to GDP growth.

Let's proceed with these refinements.

Slika 2.13. Rukovanje greškama i **fallback** strategija

3. Evaluacija i rezultati

3.1. Metodologija evaluacije i testni podaci

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).

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% **stopu uspjeha** za dobro oblikovane upite na prirodnom jeziku s prosječnim vremenom odgovora od 8.3 sekunde za složene multimodalne upite.

Performanse vektorskog pretraživanja pokazuju **izvrsne** rezultate s prosječnim vre-

menom 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% stopu uspjeha za sintaksno ispravne i izvršive upite. Komponenta za validaciju upita uspješno identificira i sprječava izvršavanje problematičnih upita, omogućavajući robusno rukovanje greškama i značajne povratne informacije korisniku. Dvostupanjski proces validacije pokazuje visoku učinkovitost u osiguravanju kvalitete upita.

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.

Primjer 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.

Primjer 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.

Primjer 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. 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.

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 sofisticiranih tehnika sinteze rezultata. Napredne mogućnosti vizualizacije i kolaborativne značajke također predstavljaju obećavajuće smjerove za budući razvoj.

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