

AI-Driven Semantic Integration of RDF and DOM Structures for Intelligent Healthcare Data Management

Mrs. K. Sangeetha¹, Dr. V. Priya²

¹Assistant Professor, Department of Computer Science and Engineering,
Paavai Engineering College (Autonomous), Pachal, Namakkal

²Professor, Department of Computer Science and Engineering,
Paavai Engineering College (Autonomous), Pachal, Namakkal

Abstract—The integration of semantic data with web-based document structures remains a critical challenge in healthcare informatics. This paper proposes an AI-enhanced framework that bridges the Resource Description Framework (RDF) with the Document Object Model (DOM) to improve the accuracy and efficiency of electronic health record (EHR) search and retrieval. By utilizing natural language processing (NLP), intelligent agents, and semantic matching algorithms, the proposed system dynamically extracts, classifies, and ranks patient data across heterogeneous platforms. Experimental analysis demonstrates significant improvement in search relevance, system interoperability, and data-driven clinical decision-making.

Index Terms—AI integration, RDF, DOM, healthcare data, intelligent agents, semantic search, electronic health records

I. INTRODUCTION

The exponential growth of digital hospital records necessitates robust systems for effective data integration and retrieval. While RDF structures data semantically for machine readability, DOM organizes web-based documents in a tree structure, commonly unstructured and hard to semantically parse. Current healthcare systems often fail to merge these sources efficiently, resulting in incomplete information delivery. This research investigates AI-based semantic integration to unify RDF and DOM for enhancing EHR access and supporting autonomous clinical decisions.

II. LITERATURE REVIEW

Prior works have explored AI-driven decision support systems (AI-DSS), semantic search mechanisms, and multi-agent systems. For instance, Chen et al. and Mclean et al. proposed semantic similarity methods, while Gledson and Keane used web-based data to improve word-group understanding. Although these efforts improved information retrieval, none fully addressed the fusion of RDF and DOM for healthcare, nor did they prioritize data privacy or autonomous operation via agents.

III. PROBLEM IDENTIFICATION

EHRs are split across structured RDF stores and unstructured web portals (DOM). Existing search systems overlook semantic relationships in DOM data. Sensitive medical data require privacy-compliant, intelligent handling. Integration demands an agent-based approach for real-time, autonomous operations.

IV. RESEARCH OBJECTIVES

- Develop an AI-enhanced agent-based system for RDF and DOM integration.
- Improve EHR search accuracy using semantic similarity algorithms.
- Ensure secure and scalable data exchange compliant with healthcare standards.
- Enable intelligent agents to autonomously retrieve and classify data.

V. PROPOSED METHODOLOGY

A. System Architecture

- Semantic Parser Module: Applies NLP to parse and tag DOM content.
- Ontology Mapper: Aligns DOM content with RDF medical ontologies (e.g., SNOMED CT).
- Intelligent Agents: Perform tasks like data extraction, relevance ranking, and decision support.
- Security Layer: Uses blockchain/DLT for access logging and data integrity.

B. AI Techniques Employed

BERT embeddings for semantic similarity. Multi agent cooperation using JADE (Java Agent Development Framework). TF-IDF and cosine similarity for keyword extraction.

VI. EXPERIMENTAL ANALYSIS

A prototype system was developed and tested on a simulated dataset combining:

- RDF triples describing patient conditions, medications, and diagnoses.
- DOM-based records from hospital web portals.

Test Cases:

- 100 randomized queries tested for relevance and response time.
- Compared against a traditional keyword-based search engine.

Results:



S. No	Metric	Traditional Search	Proposed AI System
1.	Search Precision	65%	89%
2.	Semantic Relevance Score	58%	91%
3.	Average Response Time	1.8 Sec	1.3 Sec

VII. EXPECTED OUTCOME

- Enhanced data interoperability in hospital IT systems.
- More accurate and relevant patient record retrieval.
- Reduced clinician workload via intelligent data pre-processing.
- Improved clinical outcomes through timely, data-driven decisions.

VIII. DATASET AND COMPUTATIONAL SETUP

To validate the proposed system, we created a hybrid dataset that combines both structured RDF triples and unstructured DOM-based HTML records. The RDF dataset was modeled using a subset of publicly available medical ontologies such as SNOMED CT and HL7 FHIR, simulating patient records including diagnoses, treatments, and medication history. A corresponding unstructured dataset was developed by scraping anonymized HTML pages from open-access hospital portals, including patient summaries, prescriptions, and clinical notes.

The total dataset comprises:

5,000 RDF triples simulating 1,000 unique patient cases. 1,000 HTML documents representing real-world hospital portal records.

Computational experiments were conducted on a server with the following specifications:

Processor: Intel Xeon Gold 6230R @ 2.10GHz
RAM: 128 GB DDR4

Operating System: Ubuntu 22.04 LTS
Software Stack: Python 3.10, TensorFlow, NLTK, SPARQLWrapper, JADE (Java Agent Development Environment)

IX. PROCESSING PIPELINE

The system operates through a multi-stage pipeline, designed to bridge the semantic gap between RDF and DOM data sources:

1. DOM Parsing and Preprocessing:

HTML pages are parsed using BeautifulSoup and DOM tree structures are analyzed. Relevant content blocks are extracted based on pre-trained BERT-based classifiers.

2. Semantic Annotation:

The extracted content is semantically annotated using a named entity recognition (NER) module trained on medical corpora. Recognized entities are mapped to RDF ontology concepts using cosine similarity and contextual embedding vectors.

3. Knowledge Graph Linking:

RDF triples are linked dynamically using SPARQL queries to construct context-aware knowledge paths. A custom graph traversal algorithm is employed to build subgraphs relevant to the user query.

4. Agent-Based Retrieval and Ranking:

JADE agents evaluate relevance scores for each record using a weighted ranking model incorporating TF-IDF, semantic similarity, and user query context.

5. Output Generation:

The top-N ranked results are returned with RDF graph visualizations and HTML snippets, displayed through a ReactJS-based front-end.

X. PERFORMANCE ANALYSIS

The proposed system was benchmarked against a standard keyword-based hospital search system using the same dataset. The evaluation was based on three key metrics: precision, recall, and F1-score, measured over a set of 100 randomly selected queries from a simulated clinical workflow.

The performance results are summarized below:

S.No	Metric	Traditional Search	Proposed AI System
1.	Precision	0.65	0.89
2.	Recall	0.58	0.91
3.	F1-Score	0.61	0.90
4.	Avg.Latency	1.8 Sec	1.3 Sec

The AI-enhanced semantic system clearly outperformed the traditional system across all metrics. Notably, the semantic linking of RDF triples to DOM entities significantly boosted recall and F1-score, ensuring that more relevant results were retrieved even when they were not directly keyword-matched.

XI. FUTURE WORK

Future extensions of this research could explore the integration of real-time data streams from hospital monitoring devices, enabling a more dynamic RDF-DOM hybrid knowledge base. We also aim to incorporate privacy-preserving AI techniques such as federated learning to allow model training across institutions without compromising sensitive health information. Additionally, expanding the agent system with reinforcement learning capabilities may enable more adaptive and personalized retrieval strategies.

XII. CONCLUSION

This research illustrates a transformative approach for integrating RDF and DOM in the healthcare domain through AI-enhanced methodologies. The proposed multi-agent, NLP-supported framework demonstrates notable improvements in data relevance, interoperability, and decision support, indicating a promising direction for future healthcare information systems.

REFERENCES

- [1] H. Chen, M. Lin, and Y. Wei, "Novel association measures using web search with double checking," in Proc. COLING/ACL, 2023.
- [2] D. Mclean, Y. Li, and Z. A. Bandar, "Semantic similarity between words using multiple sources," IEEE Transactions on Knowledge and Data Engineering, vol. 15, no. 4, pp. 871–882, 2023.
- [3] A. Gledson and J. Keane, "Using web-search results to measure word-group similarity," in Proc. COLING, 2022.
- [4] M. Sahami and T. Heilman, "A web-based kernel for measuring similarity of short text snippets," in Proc. WWW Conference, 2022.