

# RDF resources for biomedical research

Mayumi Kamada<sup>1,\*</sup>, Shuichi Kawashima<sup>2</sup> and Toshiaki Katayama<sup>2</sup>

<sup>1</sup>School of Frontier Engineering, Kitasato University, 1-15-1 Kitazato, Minami, Sagamihara, Kanagawa, 252-0373, Japan

<sup>2</sup>Database Center for Life Science, 178-4-4 Wakashiba, Kashiwa, Chiba 277-0871, Japan

## Abstract

The integration of biomedical databases is essential for drug discovery and understanding disease mechanisms. Med2RDF addresses this by proposing a unified data model and providing tools to convert key biomedical datasets into Resource Description Framework (RDF) format. Target databases include those for disease variants, molecular interactions, cancer omics, and genomic variations. To enhance usability, Med2RDF offers Med2RDF-ontology for organizing shared concepts and the RDF-config for generating SPARQL queries and schema diagrams from YAML-based models. Converted RDF data is available through SPARQL endpoints on the DBCLS RDF Portal, enabling seamless integration with other life science data. This standardized and semantic approach facilitates efficient data utilization, advancing biomedical research.

## Keywords

biomedical, integrated data utilization, structural variation, RDF

## 1. Introduction

Integrating life science and biomedical databases is essential in the early stages of drug discovery and in elucidating disease mechanisms in medical research. However, major biomedical databases are provided in diverse formats using various database systems. Furthermore, utilizing these databases requires a thorough understanding of domain-specific data contexts and the relationships between data items.

Med2RDF (<http://med2rdf.org/>) proposes a data model to integrate and semantically enrich (add semantics to) key biomedical databases. Additionally, it develops and provides software tools to convert these datasets into semantic data representations in the Resource Description Framework (RDF) format (<https://github.com/med2rdf>).

## 2. Target Databases

Med2RDF focuses on major databases in the biomedical field, particularly those relevant to genomic medicine. It proposes data models and provides conversion programs<sup>2</sup> for the RDF format. Table 1 lists the databases that developed the conversion program as of Dec. 2024.

ClinVar and MGenD are databases of disease-associated variants. COSMIC is a database for somatic mutations in cancer. HINT (High-quality INteractomes) and CPDB (ConsensusPathDB) are databases providing molecular interactions information. ICGC (International Cancer Genome Consortium) and TCGA (The Cancer Genome Atlas) are large-scale consortia collecting multi-omics data for various cancers, providing omics datasets. CCLE (Cancer Cell Line Encyclopedia) and GDSC (Genomics of Drug Sensitivity in Cancer) provide profiles of cancer cell lines, including genomic, gene expression, and drug sensitivity data. dbSNP and dbVar are widely used catalogs of genomic variants. dbNSFP and dbSNV provide predictive scores of variant effects from bioinformatics tools. HGNC (HUGO Gene Nomenclature Committee) records all human gene names and symbols approved by the committee.

*SWAT4HCLS 2025: 16th International Conference on Semantic Web Applications and Tools for Health Care and Life Sciences 2025*

\*Corresponding author.

✉ [kamada.mayumi@kitasato-u.ac.jp](mailto:kamada.mayumi@kitasato-u.ac.jp) (M. Kamada)

ORCID [0000-0002-2555-7345](https://orcid.org/0000-0002-2555-7345) (M. Kamada); [0000-0001-7883-3756](https://orcid.org/0000-0001-7883-3756) (. S. Kawashima); [0000-0003-2391-0384](https://orcid.org/0000-0003-2391-0384) (. T. Katayama)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

**Table 1**

The converted databases of Med2RDF project (2024)

Database	URL
ClinVar	<a href="https://www.ncbi.nlm.nih.gov/clinvar/">https://www.ncbi.nlm.nih.gov/clinvar/</a>
MGeND	<a href="https://mgend.ncgm.go.jp/">https://mgend.ncgm.go.jp/</a>
COSMIC	<a href="https://cancer.sanger.ac.uk/cosmic">https://cancer.sanger.ac.uk/cosmic</a>
ICGC	<a href="https://dcc.icgc.org/">https://dcc.icgc.org/</a>
TCGA	<a href="https://portal.gdc.cancer.gov/">https://portal.gdc.cancer.gov/</a>
CCLC	<a href="https://sites.broadinstitute.org/cclc/">https://sites.broadinstitute.org/cclc/</a>
GDSC	<a href="https://www.cancerrxgene.org/">https://www.cancerrxgene.org/</a>
dbSNP	<a href="https://www.ncbi.nlm.nih.gov/snp/">https://www.ncbi.nlm.nih.gov/snp/</a>
dbVar	<a href="https://www.ncbi.nlm.nih.gov/dbvar/">https://www.ncbi.nlm.nih.gov/dbvar/</a>
HGNC	<a href="https://www.genenames.org/">https://www.genenames.org/</a>
HINT	<a href="https://hint.yulab.org/">https://hint.yulab.org/</a>
CPDB	<a href="http://cpdb.molgen.mpg.de/CPDB">http://cpdb.molgen.mpg.de/CPDB</a>
dbNSFP	<a href="https://sites.google.com/site/jpopgen/dbNSFP">https://sites.google.com/site/jpopgen/dbNSFP</a>
dbscSNV	<a href="http://www.liulab.science/dbscsnv.html">http://www.liulab.science/dbscsnv.html</a>

### 3. Ontology

Med2RDF introduces Med2RDF-ontology, which organizes common concepts across the aforementioned biomedical databases (Figure 1). Furthermore, to facilitate connectivity with other public databases, it has developed ontologies such as HCO (Human Chromosome Ontology), MCO (Mouse Chromosome Ontology), and GVO (Genome Variation Ontology), organizing shared concepts.

HCO provides simple and stable URIs for human reference genome versions to identify human chromosomes semantically. In medical research, linking human and mouse data is often necessary. MCO similarly provides URIs for mouse versions. GVO is an ontology that systematically describes various genomic variations, including complex structural variations in genomes.

### 4. Conversion Programs

The RDF conversion programs are publicly available on the Med2RDF GitHub repository (<https://github.com/med2rdf>). The programs are implemented in multiple languages, including Ruby, Python, and Rust. Schema diagrams are provided for each database to facilitate the use of converted RDF data. However, understanding a data model from scratch requires domain knowledge, posing a time cost for users other than the original developers. To address this, The Database Center for Life Science (DBCLS) has developed a tool called “RDF-config” (<https://github.com/dbcls/rdf-config>), which provides practical features for using knowledge graphs. RDF-config can automatically generate SPARQL queries and schema diagrams from data models described in YAML format. Figure 2 shows an example of a schema generated by RDF-config. The YAML definition file includes the RDF graph structure, key elements and their attributes with cardinalities, variable names for retrieved values, and example values for inferencing datatypes. These features help new users understand data models more quickly.

The benefits of RDF-config extend beyond data model comprehension. It standardizes RDF conversion methods and reduces dependency on programming languages and toolkits, thereby enhancing maintainability. We developed RDF-config to create RDF and JSON-LD data for databases including TCGA, GDSC, dbNSFP, and dbscSNV. Going forward, Med2RDF plans to unify the conversion process for future databases using RDF-config.

Figure 1: The Schema of Med2RDF-ontology

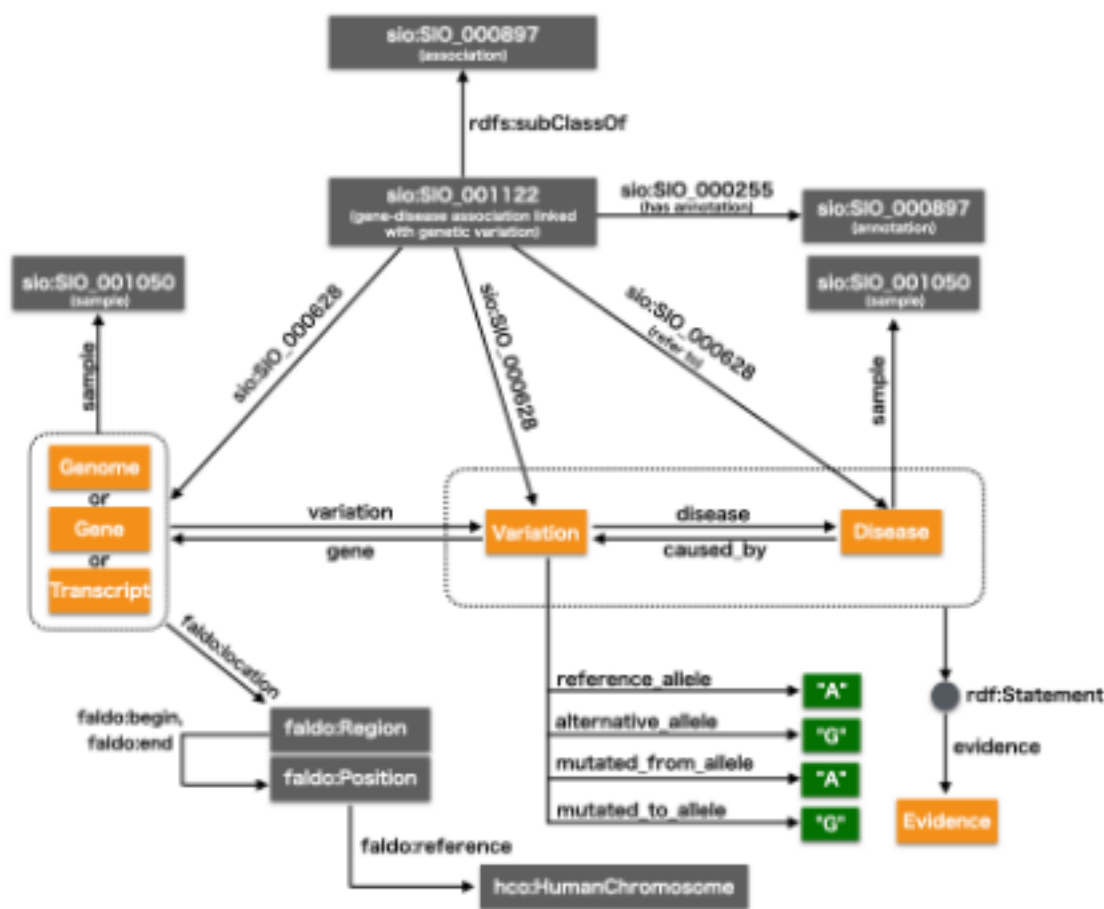
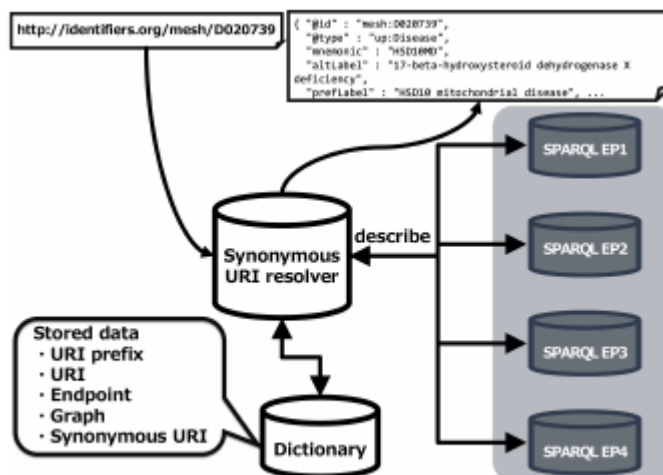


Figure 2: An example of a schema of dbscSNV generated by RDF-config



## 5. Endpoint

The RDF data for each database converted by Med2RDF programs is made available through SPARQL endpoints on DBCLS RDF Portal (<https://rdfportal.org/>) [1]. The RDF Portal regularly updates the endpoints to accommodate new databases and changes to source data. It provides an environment for

SPARQL querying and data download. Users can leverage this endpoint to combine data with other life science databases for integrated analyses.

## **Acknowledgments**

This research is supported by a Grant-in-Aid for Transformative Research Areas (A) “Latent Chemical Space” [JP23H04880 and JP23H04886] for MK from the Ministry of Education, Culture, Sports, Science and Technology, Japan

## **Declaration on Generative AI**

The authors have not employed any Generative AI tools.

## **References**

- [1] S. Kawashima, T. Katayama, H. Hatanaka, T. Kushida, T. Takagi, Ndbc rdf portal: a comprehensive repository for semantic data in life sciences, Database 2018 (2018) bay123. doi:10.1093/database/bay123.