

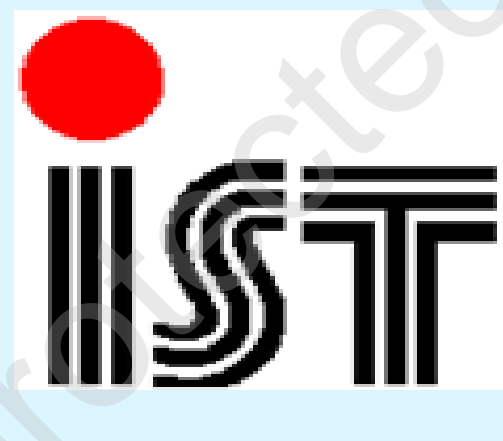
# Towards Linked Open Gene Mutations Data

Achille Zappa<sup>1,2</sup>, Andrea Splendiani<sup>3</sup>, Paolo Romano<sup>1</sup>

<sup>1</sup> Bioinformatics, IRCCS AOU San Martino - IST National Cancer Research Institute, Genoa, Italy

<sup>2</sup> Department of Informatics, Systems and Telematics, University of Genoa, Genoa, Italy

<sup>3</sup> Rothamsted Research, Harpenden, Hertfordshire, United Kingdom



## Background

The vision of the Semantic Web is to evolve the Web into a distributed knowledge base, and this vision relies on the evolution of the Web into a Web of Data. A great number of information resources are being made available in the LOD (Linked Open Data) cloud. Biological information is also being converted and it constitutes a major component of the cloud. The Human Genome Variation Society and Human Variome Project have produced recommendations for nomenclature of variations and for contents of mutation databases. They also outlined conditions for the integration of Locus Specific Data Bases (LSDB) with other biological databases. However, human variation data, despite its relevance for medicine, has not yet been adequately taken into account. This poster is a first attempt in this direction. [Fig. 1]

## Second stage: native Triple Store

The query performance of D2RQ is lower than that which can be achieved by using a devoted RDF triple store. In order to evaluate reliability and performances of an on-the-fly mapping system, such as D2RQ, compared to a native RDF framework, a dump of all triples generated by D2RQ was loaded into a dedicated Jena TDB triple store.

## Fourth stage: RDF enrichment via SPARQL Update

We have then enriched our datasets by adding triples connecting samples to related UMLS concepts. This mapping was implemented by means of a SPARQL Update federated query interconnecting our datasets with the Linked Life Data [LLD] endpoint.

## Fifth stage: Linked Data Interface

Finally, we exposed the content of our TDB triple store as a Linked Data interface using Pubby is used to add a Linked Data interface to our SPARQL endpoint. It handles external requests by connecting to the SPARQL endpoint, issuing a SPARQL "DESCRIBE" query about the requested URI, and showing the result in a HTML or RDF page, supporting Linked Data compliant content resolution and negotiation procedures. [Fig. 2]

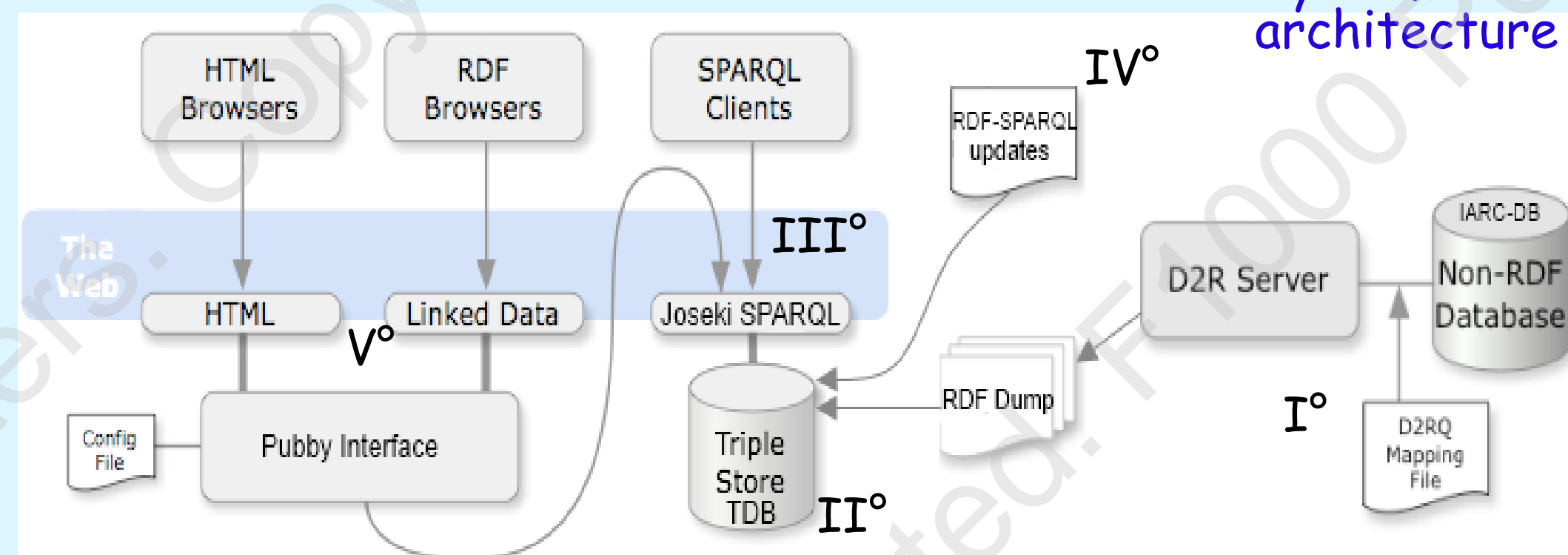
## IARC TP53 Mutations database

The IARC TP53 Mutation Database has been maintained at the International Agency for Research on Cancer in Lyon, France, since 1994. The database compiles all TP53 mutations that have been reported in the published literature since 1989. The database includes annotations on functional impact of mutations, either predicted or experimentally assessed, clinico-pathologic characteristics of tumors and demographic and life-style information on patients. A relational implementation is available at the National Cancer Research Institute of Genoa.

## First stage: RDB to RDF

In our case, automatic RDB to RDF mappings were first created by using D2RQ. The mapping file was then manually refined to improve its commitment to a shared representation and, thus, to encode in the mapping some semantics that is not expressed in the RDB schema. By customizing predicates we have been able to better represent the semantics of our data, according to shared ontologies.

Figure 1: System architecture



## Third stage: SPARQL Endpoint

Once the RDF database was created, it was made accessible as a SPARQL endpoint using Joseki, a Jena tool that provides support for SPARQL queries through an HTTP engine. Joseki was configured to connect to the TDB database to create an implementation of SPARQLer, a user friendly interface to a SPARQL server.

Figure 2: the Pubby interface of the prototype. You can browse among entities in different datasets that are linked each others.

The figure shows three screenshots of the Pubby interface. The first screenshot shows the 'individual/847' page with a table of properties and values. The second screenshot shows the 'samples/KIK92-MY279' page with a table of properties and values. The third screenshot shows the 'somatic\_mutations/863' page with a table of properties and values. Each screenshot includes a 'Property' column and a 'Value' column, with various biological and clinical data points.

## A prototype implementation

We have therefore implemented a semantic infrastructure for TP53 variation data as a prototype for studying issues related to the publication of mutation data on the LOD cloud. It includes data on somatic mutations and related bibliographic references, bio-samples, and patients demography. We also published summary gene variations data. Specific biological ontologies on gene variation does not exist. Only a limited set of external ontologies and terminologies have been taken into account. These include the NCI Thesaurus (NCIT) for medical terminology, the Bibliographic Ontology (BIBO) and the BibTeX definition in Web Ontology Language (OWL) for bibliographic references, the Disease ontology and MIO. Moreover, external links were set to LOD implementations of DBpedia, a system including all structured information which is present in Wikipedia pages, PubMed, the Human Genome Nomenclature Committee (HGNC) database, the On-line Mendelian Inheritance in Man (OMIM) system, UniProt, and the Unified Medical Language System (UMLS). All links were defined to the Bio2RDF entry points of these databases, expressed by using the unified Bio2RDF URI style, with the exception of DBpedia, and UMLS, that was connected through LinkedLifeData.

Figure 3: SPARQL fed-query sample

```
SELECT DISTINCT ?variation_label ?neoplasm ?clinical_trial
WHERE {
  ?SERVICE <http://linkedlifedata.com/sparql> {
    ?clinical_trial rdfs:label ?label .
    ?variation_label rdfs:label ?variation_label .
    ?neoplasm rdfs:label ?neoplasm .
    ?variation_label rdfs:label ?variation_label .
    ?neoplasm rdfs:label ?neoplasm .
  }
  ORDER BY ?variation_label ?neoplasm
```

?variation_label	?neoplasm	?clinical_trial
NM_000546.1:c.507G>A	Adenocarcinoma, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001332>
NM_000546.1:c.838A>G	Adenocarcinoma, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001332>
NM_000546.1:c.507G>A	Adenocarcinoma, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001428>
NM_000546.1:c.428G>A	Adenocarcinoma, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001428>
NM_000546.1:c.482C>A	Dysplasia, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001932>
NM_000546.1:c.482C>A	Dysplasia, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00003076>
NM_000546.1:c.482C>A	Dysplasia, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00003094>
NM_000546.1:c.482C>A	Dysplasia, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00003223>
NM_000546.1:c.482C>A	Dysplasia, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00003239>
NM_000546.1:c.469G>T	Hyperplasia, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001378>
NM_000546.1:c.469G>T	Hyperplasia, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001780>
NM_000546.1:c.469G>T	Hyperplasia, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00003641>
NM_000546.1:c.451C>G	Squamous cell carcinoma, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00006929>
NM_000546.1:c.451C>G	Squamous cell carcinoma, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00006929>
NM_000546.1:c.47A_+76ins1	Squamous cell carcinoma, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001450>
NM_000546.1:c.488A>G	Squamous cell carcinoma, NOS	<http://data.linkedlifedata.com/resource/trials/NCT00001450>

## Availability

Presently, we offer access to our database via two distinct modalities. The first interface is a SPARQL endpoint available at <http://bioinformatics.istge.it/logvdsparql/sparql>. It is implemented via Joseki on TDB. Interfaces to validators for SPARQL queries and for RDF data are also available. The second interface is a Linked Data representation available at <http://bioinformatics.istge.it/logvd/>. It is based on the Pubby frontend.

## Some perspectives

In the short time, the exploitation of SPARQL searches on mutation data and other biological databases may support some interesting and useful data retrieval presently not possible. [Fig. 3] In a longer time frame, reasoning on integrated variation data may also support discoveries towards personalized medicine. For this to happen, however, a long way is still needed, due to the relative isolation of variation data sources and lack of standardization both in terminologies and data schema.

## Acknowledgements

This work has been partially supported by the Italian Ministry of Health (project RNBIO -Rete Nazionale di Bioinformatica Oncologica) and by the Liguria region (project Liguria eScience).