

# From Web 1.0 → 3.0: Is RDF access to RDB enough?

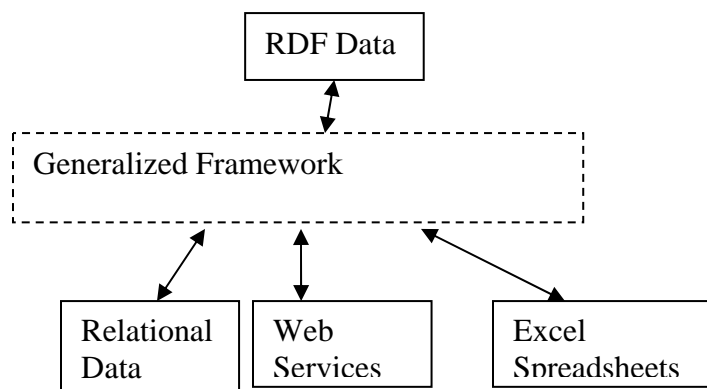
Vipul Kashyap, Senior Medical Informatician, Partners Healthcare System, [vkashyap1@partners.org](mailto:vkashyap1@partners.org)  
Martin Flanagan, CTO, InSilico Discovery, [mflanagan@insilicodiscovery.com](mailto:mflanagan@insilicodiscovery.com)

## Abstract:

There have been a host of initiatives led by the W3C to create standards and infrastructure that enable the transition from the current web to the semantic web. In this paper, we present a position that argues for a generalized approach for RDF access to; (a) Relational Databases; (b) Web Services; and (c) Tabular data sources such as Excel Spreadsheets. A Life Sciences use case is presented and the use of a common framework is demonstrated through illustrations from a working system. We believe that the effort in evolving a common framework is incremental, and extends the reach and value of Semantic Web technologies in the Healthcare and Life Sciences.

## 1. Introduction

The wide variety of life science data using different types of data storage and access schemes are a barrier to understanding of biological knowledge and its application to clinical research and practice. For example how does software discover the biological reasons explaining the statistical correlations observed between multiple genes in different types of studies? A large amount of biological data is currently available through web-based data repositories and web services [1,2]. The value of this information would increase dramatically if they were accessible to systems which might not directly support the native data storage format or access mechanism, but which do support RDF [3]. There is a need for a generalized framework for providing RDF access to various types of data storage and access schemes.



We believe that the effort for evolving a generalized framework is incremental and can be integrated with the efforts to create RDF access to Relational Databases. We now discuss a use case from the healthcare and life sciences and demonstrate through illustrations of a working system, how this common framework is achieved today. A live demonstration of the system will be presented at the workshop, if selected for the workshop.

## 2. Use Scenario: Biological Explanations for Statistical Correlations

This use case scenario is based on a typical Life Science Researcher, Professor Genomus who performs various types of gene correlation studies and experiments. One of the key motivations of Professor Genomus is to determine the possible biological explanations of the correlations observed between two biological artifacts, for e.g., a Gene and a SNP Variant; or understand the clinical implications, for e.g., whether a Gene or a Variant is implicated in some disease. Some typical questions that a professor might answer along with different types of data sources that might contain a potential answer are:

- What is the location of a given Gene, e.g., CPNE1 on the Human Genome?  
**Data Repository:** NCBI Entrez, **Access Mechanism:** Web Services
- For what gene(s) is a given SNP, e.g., rs6060535 in the upstream regulatory region?  
**Data Repository:** RDBMS containing dbSNP and regulatory region data,  
**Access Mechanism:** JDBC/SQL
- What genes have been found to be "coexpressed" with CPNE1 and in what study?  
**Data Repository:** Excel Spreadsheet containing the co-expression patterns of various genes in various studies.  
**Access Mechanism:** .NET API, MS Office API

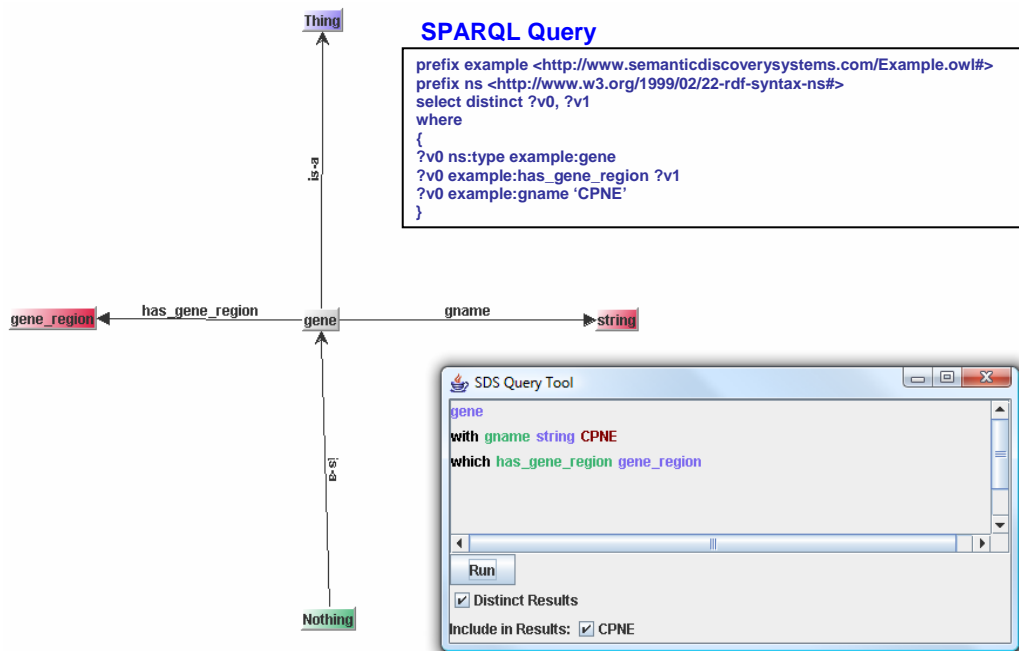
We now demonstrate a solution approach and mapping framework. Given the theme of the workshop, we will focus on a simple example that retrieves information about gene regions and use it to illustrate the mapping framework.

## 3. Solution Approach

The various steps in the solution approach are presented below.

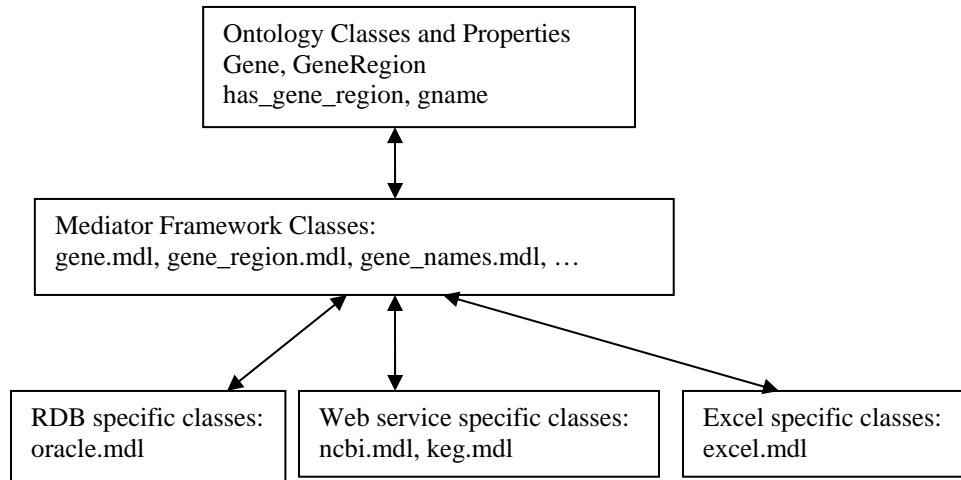
### 3.1 Ontology-based Query Specification

Professor Genomus invokes the Semantic Discovery System (SDS) [4] and as illustrated below, specifies a query based on the ontology of his interest.



## 3.2 Mapping Framework

Professor Genomus then gets in touch with Mr. Information Architect and asks him to include the relevant data sources into the system. Mr. Information Architect taps into the generalized framework [5.6] supported by the SDS as illustrated below. A detailed step through of the mapping framework is available at [7].



### 3.2.1 Mapping to Relational Databases

Mr. Information Architect maps the mediator class to a relational database (Oracle).

```
import orasample.HR.GENENAMES;
import orasample.HR.ORACLE_GENE_NAMES;

define GeneNames as select distinct x.name from excelj("GeneNames.xls") x
where x.name != "";
materialize GeneNames;
*/
/* While this next statement does exactly the same from an Oracle database */
define GeneNames as select * from ORACLE_GENE_NAMES;
materialize GeneNames;

// This extract is to show the start of what would typically be a very large corporate Oracle gene database
// We show 20+ records only to get the 'shape' of the data across for the example only.

orasample("drop table GeneNames");
orasample("create table GeneNames( Name char(20), Polymer char(20), Starting number(9), Ending number(9))");
orasample("INSERT INTO GeneNames VALUES( 'A1BG', '19q, 13, 4 ) ");
orasample("INSERT INTO GeneNames VALUES( 'A2M', '12p13.3-p12.3', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'A2M', '14q, 32, 33 ) ");
orasample("INSERT INTO GeneNames VALUES( 'A2ML1', '12p13.31', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'A2MP', '12p13.3-p12.3', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'A3GALT2', '1p35.1', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'A4GALT', '22q11.2-q13.2', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AAAS', '12q13', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AACP', '8p22', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AADAC', '3q, 25, 11.3 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AADAC', '3q21.3-q25.2', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AADACL1', '3q, 26, 31 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AADACL2', '3q, 25, 1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AADAT', '4q33', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AAK1', '2q35', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AAMP', '2p14', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AANAT', '17q25', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AARS', '16q22', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AARS2', '6p21.1', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AARSD1', '17q, 21, 31 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AASHPPT', '11q22', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AASS', '7q, 31, 3 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AARS2', '6p21.1', -1, -1 ) ");
orasample("INSERT INTO GeneNames VALUES( 'AATF', '17q11.2-q12', -1, -1 ) ");
```

### 3.2.2 Mapping to Web Services

Mr. Information Architect maps the mediator class to a web service.

The screenshot shows the Query Server Explorer on the left and the ncbi.mdl mediator class on the right. The Explorer shows a tree view with 'ncbi.run\_eFetch' selected, showing its input and output structures. The mediator class code on the right includes:

```
import ncbi.run_eSearch as eSearch;
import ncbi.run_eFetch;

define HUMAN as human;

/*
 * Given a term, look its IDs up in the taxonomy database.
 */
define getTaxIds(term) as
select y from
eSearch(struct(eSearchRequest:(struct(db:"taxonomy", sort:
null, retType: null, retMax: null, email: null, retStart: null,
tool: null, term: term + "[*]", datatype: null, maxdate: null,
mindate: null, reldate: null, field: null, usehistory: null,
queryKey: null, webEnv: null)))) x,
x.eSearchResult.idList y;

/*
 * Given a gene name and a species name, get the genbank ids for the gene
 */
define getGeneId(name, species) as
select y from eSearch(struct(eSearchRequest:(struct(db:"gene", sort: null,
retType: null, retMax: null, email: null, retStart: null, tool: null,
term: name || "[*]" AND toid + species + "[Organism]", datatype: null,
maxdate: null, mindate: null, reldate: null, field: null, usehistory: null,
queryKey: null, webEnv: null)))) x,
x.eSearchResult.idList y;

define getGeneIds(name, species) as
Flatten(select getGeneId(name, s) from getTaxIds(species) s);

/*
 * Get the locations in the given species genome of a gene name. Use geneTrack status
value to remove the duplicated report, and only get the primary report.
 */
define extractGenomeLocation(s) as
if (s like "[0-9]+q[0-9]+.[0-9]+")
then
let
{
t1 as tokenize(s, "q");
t2 as tokenize(t1[1], ".");
}
in
struct(polymer:t1[0]+"q", start:stringToInt(t2[0]), end:stringToInt(t2[1]))
else
struct(polymer:s, start:-1, end:-1);

define getGenomeLocations(geneName, species) as
select extractGenomeLocation(z.maps_displayStr) from
getGeneIds(geneName, species) gid,
run_eFetch(struct(eFetchRequest: struct(report: null, complexity: null,
seq_stop: null, seq_start: null, strand: null, retmax: null, retstart: null,
query_key: null, email: null, tool: null, rathmax: null, id: null, db: "ncbi"))
```

Annotations with arrows point from the Explorer's 'ncbi.run\_eFetch' structure to the 'eSearch' and 'run\_eFetch' calls in the code, and from the 'extractGenomeLocation' function to the 'extractGenomeLocations' function.

### 3.2.3 Mapping to Excel Spreadsheets

Mr. Information Architect then maps the mediator class to an excel spreadsheet.

The screenshot shows the Query Server Explorer on the left and the gene\_names.mdl mediator class on the right. The Explorer shows a tree view with 'excelj' selected, showing its input and output structures. The mediator class code on the right includes:

```
define GeneNames as select distinct x.name from excelj("GeneNames.xls") x
where x.name != "";

materialize GeneNames;

select * from GeneNames;
```

Annotations with arrows point from the Explorer's 'excelj' structure to the 'excelj' call in the code, and from the 'select \* from GeneNames;' line to the 'GeneNames' table reference.

### 3.3 SPARQL Translation to Access Mechanisms and Execution

Once the mappings from the mediator classes (and hence to ontological elements) to the underlying data sources have been established, the SPARQL query is automatically translated into the local query languages and access mechanism supported by the data sources. For instance in the example discussed above:

- The Oracle database executes a SQL query of the form:  

```
SELECT * FROM GENE_NAMES
```
- The NCBI web service, `run_efetch` is invoked
- The MS Office API is invoked to retrieve the spreadsheet

After the data is retrieved from the underlying data sources, the generalized mediator framework is responsible for merging the retrieved data. It also optimizes the query processing to return results in a reasonable amount of time.

## 4. Conclusions

In this paper, we have proposed that the mapping between RDF data to RDB storage/access be viewed as a special case of a generalized mapping mechanism which seeks to map ontological concepts to data retrieval and access via RDBs, Web Services and Structured Tabular data as exemplified by Spreadsheets. A system that supports this mapping framework is also presented. The central vision is one of enabling point and click mapping using a GUI to map ontological classes and properties to different types of data sources. Some of these interfaces enable functionality for a subject matter expert such as Professor Genomus, whereas other interfaces would enable functionality for an architect such as Mr. Information Architect. The key to building out such user interfaces and web-enabling them is the design of a generic mapping framework to handle different types of data resources on the web. We believe RDBs, web services and structured tabular data sources such as Excel spreadsheets are easily incorporated into a generalizable framework as demonstrated by the system discussed in this paper. The longer term goal is to incorporate frameworks for semi-structured and unstructured content such as GRDDL [8] and GATE [9], a framework for NLP functionality.

## References

- [1] [http://www.ncbi.nlm.nih.gov/entrez/query/static/esoap\\_help.html](http://www.ncbi.nlm.nih.gov/entrez/query/static/esoap_help.html)
- [2] <http://www.genome.jp/kegg/soap/>
- [3] <http://www.w3.org/RDF/>
- [4] <http://www.insilcodiscovery.com>
- [5] A Guide and Reference Manual to the SDS SQL and MDL languages.  
[http://www.insilcodiscovery.com/v2/index.php?option=com\\_kb&Itemid=149&page=articles&articleid=44](http://www.insilcodiscovery.com/v2/index.php?option=com_kb&Itemid=149&page=articles&articleid=44)
- [6] A Guide to the operating principles and administration of the SDS Server,  
[http://www.insilcodiscovery.com/v2/index.php?option=com\\_kb&Itemid=149&page=articles&articleid=46](http://www.insilcodiscovery.com/v2/index.php?option=com_kb&Itemid=149&page=articles&articleid=46)
- [7] [http://www.insilcodiscovery.com/v2/index.php?option=com\\_content&task=view&id=82&Itemid=151](http://www.insilcodiscovery.com/v2/index.php?option=com_content&task=view&id=82&Itemid=151)
- [8] <http://www.w3.org/TR/grddl/>
- [9] GATE: A General Architecture for Text Engineering, <http://gate.ac.uk>