

Chem2Bio2RDF: A Linked Open Data Portal for Chemical Biology

Bin Chen¹
David J Wild¹
Qian Zhu¹

Ying Ding²
Xiao Dong¹

Huijun Wang¹
Yuyin Sun²

Madhuvanathi Sankaranarayanan¹

¹School of Informatics and Computing, Indiana University, Bloomington, IN, USA

²School of Library and Information Science, Indiana University, Bloomington, IN, USA
{binchen|dingying|huiwang|djwild|xdong|yuysun|qianzhu|madhsank}@indiana.edu

ABSTRACT

The Chem2Bio2RDF portal is a Linked Open Data (LOD) portal for systems chemical biology aiming for facilitating drug discovery. It converts around 25 different datasets on genes, compounds, drugs, pathways, side effects, diseases, and MEDLINE/PubMed documents into RDF triples and links them to other LOD bubbles, such as Bio2RDF, LODD and DBPedia. The portal is based on D2R server and provides a SPARQL endpoint, but adds on few unique features like RDF faceted browser, user-friendly SPARQL query generator, MEDLINE/PubMed cross validation service, and Cytoscape visualization plugin. Three use cases demonstrate the functionality and usability of this portal.

Categories and Subject Descriptors

[D.2.12](#) [Interoperability]: Data Mapping

General Terms

Experimentation

Keywords

Semantic Web, RDF, SPARQL, Systems Chemical Biology

1. INTRODUCTION

Curing disease and providing better health care require a synthesis of data and understanding across different disciplines, domains, and applications. The information and data are most useful when they are well integrated with each other. Traditional data integration techniques relying on static mapping methodologies are unlikely to be scalable. Recent advances in biology and medicine have led to an explosion of new data sources, such as genes, proteins, genetic variations, chemical compounds, diseases and drugs. Since these data are disconnected, it is very challenging for biomedical scientists to identify related genes through different diseases, to discover new drugs, and to interconnect compounds with pathways. The recent development of chemical and biological sciences has resulted in the emergence of fields like systems biology which adopt a comprehensive approach to the study of biological systems, chemogenomics, and systems chemical biology [20].

Systems chemical biology is a new and developing area, with current application in polypharmacology [5, 12] and adverse drug reaction [22] addressing problems in efficacy and toxicity, which are the two main reasons accounting for the drug failure in the drug discovery. Polypharmacology and adverse drug reaction involve linking heterogeneous chemical and biological data with broad range of scales from small molecules (small compound, drug), to super-molecules (gene, protein), to biological systems (protein complex, pathway), and to phenotypes (disease, side

effects). In addition, many databases cover similar data (called homogeneous data here) but with slightly different focuses. For instance, DrugBank¹ has drug target association, while PharmGKB² has similar information from different perspectives. All the heterogeneous and homogeneous data are scattered around the web and published in diverse formats (i.e., text file, XML and relational database). Data integration is essential to systems chemical biology.

Semantic Web technologies enable data integration and data interlinking on the Web and demonstrate promising potentials in life sciences, healthcare and drug discovery [17, 19, 23]. Integrating heterogeneous data is obviously necessary for advanced network biology and network medicine research [6]. Bio2RDF [3] manages to integrate public bioinformatics databases and convert them into 46 million RDF triples. Linking Open Drug Data (LODD) [13] links various sources of drug data together to answer interesting scientific and business questions. Bio2RDF already covers most of biological data (e.g., protein and pathway) and LODD has collected various sources relating to chemicals (particularly drugs), although some efforts have been made to link both together [7], its practical application is still not fully explored. The integration of chemical and biology data is actually an interaction between small compounds and proteins, usually called chemogenomics data. However, little effort has been allocated currently to integrate them. In this paper, we build the system called Chem2Bio2RDF to address these issues around system chemical biology [28].

The contribution of this paper can be summarized as follows:

- We have aggregated and converted a variety of public chemogenomics data distributed around the Web into RDF formats, which enables linking with other biological Semantic Web information resources such as Bio2RDF and LODD.
- We built up a portal called Chem2Bio2RDF and designed a faceted browser to display DrugBank data, which can be further extended to other LOD bubbles.
- We prototyped an automatic SPARQL query generator so that user can avoid writing complicated SPARQL queries.
- We provided a MEDLINE/PubMed literature cross-validation service so that the SPARQL query results can be cross-validated by related literatures in MEDLINE/PubMed.
- We prototype a Cytoscape plugin for this portal.

This paper is organized as the following: Section 2 surveys the related works; Section 3 describes the Chem2Bio2RDF portal and its unique features; Section 4 discusses three cases to prove the concept and Section 5 evaluates the methods; Section 6 makes the conclusion and points out future research.

¹ <http://drugbank.ca/>

² <http://www.pharmgkb.org/>

- DrugBank [35]: It combines detailed drug (i.e. chemical, pharmacological and pharmaceutical) data with comprehensive drug target (i.e. sequence, structure, and pathway) information. The database contains nearly 4,800 drug entries and more than 2,500 non-redundant protein.
- ChEMBL [36]: ChEMBL focuses on the interactions and functional effects of small molecules binding to their macromolecular targets. It provides 500,000 bioactive compounds, their quantitative properties and bioactivities (binding constants, pharmacology and ADMET, etc).
- PDSP [37]: PDSP Ki database provides information on the abilities of drugs to interact with an expanding number of molecular targets. It provides Ki value (one measure of binding affinity) for its 766,000 interaction.
- PharmGKB [38]: It curates primary genotype and phenotype data, annotate gene variants and gene-drug-disease relationships via literature review.
- Binding MOAD [39]: Mother Of ALL Databases (MOAD) collects all well resolved protein crystal structures with clearly identified biologically relevant ligands annotated with experimentally determined binding data extracted from literature. All the structures are extracted from Protein Data Bank (PDB), containing very high-quality ligand-protein binding.

The detailed data processing was reported in our previous work [28]. We add simple provenance (What, When, Where, Why, and Who) to the RDF resources, which can be viewed at our website (chem2bio2rdf.org/datasets.html). The user can explore the RDF resources using the filter function. For example, the Figure 2 shows all the chemogenomics RDF resources which are extracted from literature.

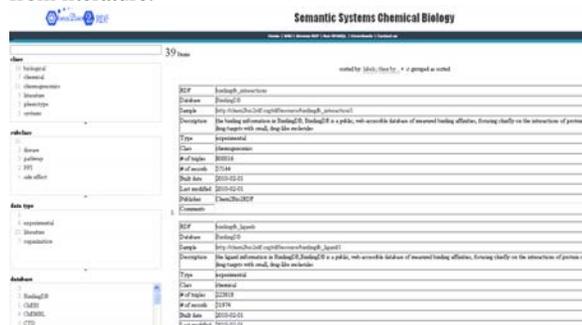


Figure 2. Faceted Browser for the datasets

All the data was processed from its original dataset, however, some of them does not meet the further utilization. For example, BindingDB uses string (i.e., >0.5) instead of number to present its binding affinity, if the user wants to select the data greater than 0.4, this text does not support this kind of search, so we have to manually separated the text into two parts, namely operator and number. In addition to data quality, the missed links between different RDF resources is a serious issue. For instance, many data formats (SMILES, SDF, InChi and CID number) are available to present one chemical, but the machine does not know they actually talk the same subject. Thus we assigned the chemical a CID number (PubChem Compound ID) if available, as PubChem is a hub of the public compounds. Meanwhile, proteins can be presented as GI number, UNIPROT ID, EC number, PDB ID and gene symbol, so we added many separated RDF resource to connect them together. GI2UNIPROT offers the conversion between GI number and UNIPROT ID. All the proteins could be

eventually converted to UNIPROT ID, that allows to link to pathways, protein-protein interaction and diseases.

Chem2Bio2RDF aims to bridge chemical and biological data, we linked our data to LODD and Bio2RDF using owl:sameAs. As LODD and BioRDF have strict namespace definition and dereferenceable URI. For instance, the URI of a drug in Bio2RDF is http://bio2rdf.org/drugbank_drugs; followed by drug id. This allows us simply link our data (i.e. drug, enzyme, protein, gene, pubmed) to them.

3.2 Chem2Bio2RDF Faceted Browser

MIT Simile project toolsets are used to provide user-friendly faceted RDF triple browser. Figure 3 shows the difference between normal SPARQL endpoint RDF browser and the Chem2Bio2RDF faceted RDF browser (<http://chem2bio2rdf.org/exhibit/drugbank.html>). Currently it provides the following views:

- Table view (see Figure 4a): The drug structure and related properties are displayed in a table format and the faceted filters on the right side can narrow down the view. For example, shown here are drugs that been approved by FDA in 2001 with the anti-allergic agent category;
- Timeline view (See Figure 4b): It displays the timeline of drugs based on the year they get approved by FDA;
- Tile view (see Figure 4c): It groups drugs first by year of approval, then by alphabetic order of drug names;
- Thumbnail view (See Figure 4d): It displays drugs based on their drug names with drug structures as thumbnails.

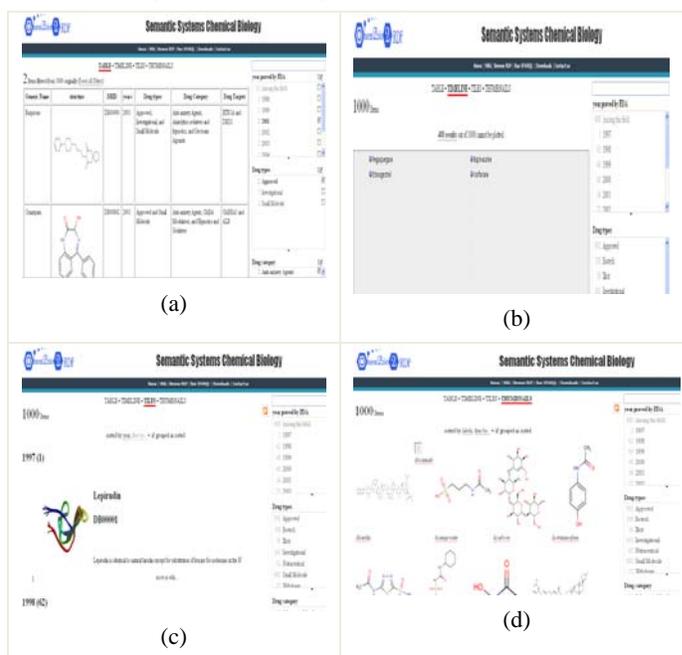


Figure 4. Screenshots of Chem2Bio2RDF Portal

3.3 Automatic SPARQL generator

If the users want to explore the relation between two classes, they do not necessarily need to know the specific dataset and the linkage between them, Link Path Generator (LPG) can automatically explores all possible link paths based on Systems Chemical Biology Data Source Ontology that models the semantic description of the datasets and their linkages. More specifically, LPG contains a graph mining module to enumerate all the possible routes from between two data sources within the

linked graph, which can be translated into concrete SPARQL queries easily. LPG is written as:

$$R(A, B) \xrightarrow{\text{Ontology}} R((A_1, A_1, \dots, A_m), (B_1, B_1, \dots, B_n)) \xrightarrow{\text{Network}} \sum_i^m \sum_j^n \text{path}_{A_i \rightarrow B_j}$$

Where A and B are two classes which can be *Chemical*, *Biological*, *Pathway* and so on. m and n is the number of data sources relating to A and B respectively. The relation between A and B can be derived from the union of the paths between all the sources of A and all the sources of B . We implement a graph mining algorithm to enumerate all the unique, non-redundant (one that doesn't revisit a node multiple times) paths between A_i and B_j . For example, if we are looking for relation between *Pathway* and *Side effect*, 14 paths are generated by LPG as listed below:

1) [sider -> kegg] contains 7 paths:

```
[sider, compound_hub, bindingdb_ligand, bindingdb_protein, uniprot_hub, kegg]
[sider, compound_hub, ctd, gene, gene2uniprot, uniprot_hub, kegg]
[sider, compound_hub, drugbank_drug, drugbank_target, uniprot_hub, kegg]
[sider, compound_hub, matador, uniprot_hub, kegg]
[sider, compound_hub, pubchem_bioassay, gi, gi2uniprot, uniprot_hub, kegg]
[sider, compound_hub, qsar, gene, gene2uniprot, uniprot_hub, kegg]
[sider, compound_hub, ttd_drug, ttd_target, uniprot_hub, kegg]
```

2) [sider -> reactome] contains 7 paths:

```
[sider, compound_hub, bindingdb_ligand, bindingdb_protein, uniprot_hub, reactome]
[sider, compound_hub, ctd, gene, gene2uniprot, uniprot_hub, reactome]
[sider, compound_hub, drugbank_drug, drugbank_target, uniprot_hub, reactome]
[sider, compound_hub, matador, uniprot_hub, reactome]
[sider, compound_hub, pubchem_bioassay, gi, gi2uniprot, uniprot_hub, reactome]
[sider, compound_hub, qsar, gene, gene2uniprot, uniprot_hub, reactome]
[sider, compound_hub, ttd_drug, ttd_target, uniprot_hub, reactome]
```

Based on the link path, SPARQL can be constructed, and then all the results are combined to output.

3.4 MEDLINE/PubMed Literature Cross-validation

MEDLINE/PubMed is the most frequently consulted online scientific medical resource in the world. In our Chem2Bio2RDF portal, we are processing 17,862,546 scientific abstracts from 1966 to present to rdf format. Not only the normal citation information, i.e. authors, journal, publication date, etc. have been processed, but also the Chem/Bio information have also been extracted. Current, our extractions are based on exact dictionary match. We used the PubChem as the dictionary to identify the compounds. The Uniprot human genes are used as the dictionary for the genes. The MESH disease terms and the COSTART side effects are applied to diseases and side effects extraction. As shown in Figure 5, the chemical compounds, genes, diseases and side effects for each abstracts are extracted and then associated with other data sources.

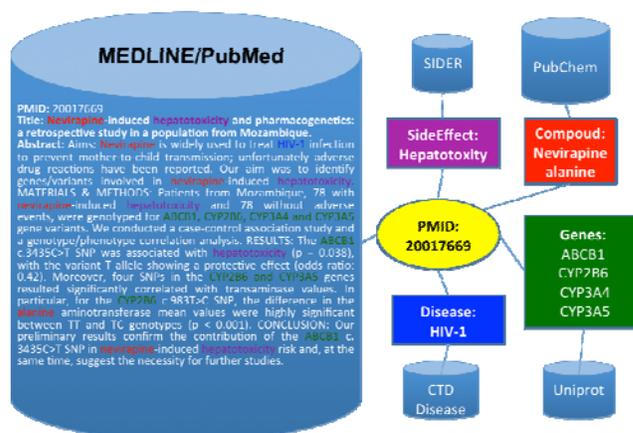


Figure 5, Chem/Bio Extraction for MEDLINE/PubMed

After the Chem/Bio terms extraction, each literature is connected to other notes (chemical, genomic, phenotype, etc.) in the Chem2Bio2RDF graph, which provides a more meaningful way to search and classify the literature based on domain knowledge. For instance, if users want to study the literatures resource for a given drug, they can retrieve the literatures that contain either the given drug or the compound with a similarity value greater than a given threshold. If users are interested in studying certain disease, both the literatures contain the given disease and the literatures contain the compounds/genes that associated with the given disease can be retrieved. Those potential search methods enable a broad way to retrieve documents based on domain knowledge.

Meanwhile, the Chem/Bio information extraction for MEDLINE/PubMed provides a cross-validation for our Chem2Bio2RDF search result. As shown in Figure 6, the left window allows users to specify two search terms, for example, find the relationships for the compound, doxazosin, and the side effect, necrosis. The right window shows the network visualization based on our RDF data, which connect the compound with its targets, targets with pathways, and then associated the pathways with side effects. More details about this network validation are shown in case 2. The bottom window provides the MEDLINE/PubMed results for this given input by retrieving the literatures that contain both inputs and the associated genes/pathways. This MEDLINE/PubMed validation can also used to rank the generated paths based on literature confidence.

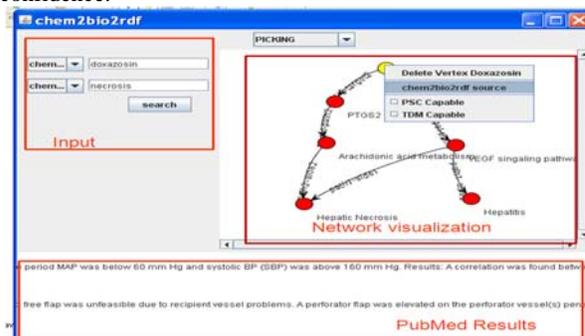


Figure 6. MEDLINE/PubMed cross-validation service

3.5 Cytoscape Plugin

Cytoscape is an open source bioinformatics tool originally developed for visualizing and analyzing biological networks. Over the past few years, Cytoscape has been used for the analysis

of chemgenomics data which explores the relationships between compounds, genes and diseases. Since Cytoscape allows users to develop plugins that are best suited to the needs of their own dataset, we developed a plugin called Chem2Bio2RDFViz that allows users to interact with the Chem2Bio2RDF dataset in the most efficient manner possible (see Figure 8).

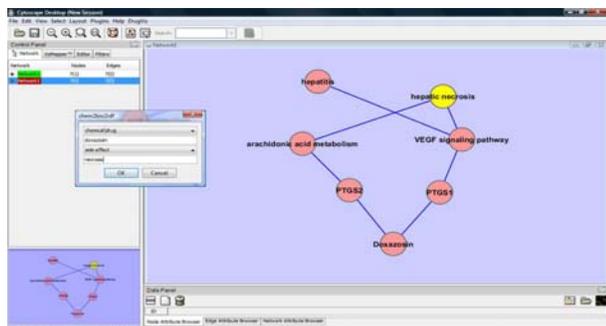


Figure 8. Cytoscape plugin prototype

Chem2Bio2RDFViz allows you to enter the names of the drug/protein/side effect/pathway/disease/gene whose network you would like to visualize and analyze. Once you enter the text, a network will be generated in the Cytoscape desktop which can be further manipulated according to the user's needs. Right clicking on each node leads to a range of choices that give you more information related to the network. The figure below shows the visualization of the network Doxazosin and its side effect Necrosis. Right clicking on the node Doxazosin will give you the MEDLINE/PubMed option that displays the MEDLINE/PubMed articles related to Doxazosin and Necrosis. Also, we can use the myriad inbuilt features of Cytoscape to understand the networks better.

4. Use cases

Here we provide three use cases to demonstrate the different levels of drug discovery.

4.1 Compound targets exploration

Given a drug, its targets are needed to be explored to understand its real mechanism, while for a drug-like compound, finding all its possible targets in the early discovery stage is desired in order to avoid the unexpected side effects in the clinical experiments. It's possible that many efforts have been made to identify the targets of a drug, but the results are scattered in the literature or hidden in some experimental records. This raises the following question:

Question: Given a drug (or chemical), find its possible targets (i.e., Gefitinib).

SPARQL:

```
SELECT ?uniprot WHERE {
  {?compound compound:CID ?compound_cid . FILTER
  (?compound_cid= 123631) .
  ?chemical bindingdb_ligand:cid ?compound .
  ?target bindingdb_interaction:monomerid ?chemical.
  ?target bindingdb_interaction:uniprot ?uniprot.
  ?target bindingdb_interaction:ic50_value ?ic50 . FILTER
  (?ic50<10000) . }
  UNION {?compound compound:CID ?compound_cid . FILTER
  (?compound_cid= 123631) .
  ?drug drugbank_drug:CID ?compound .
  ?drugtarget drugbank_interaction:DBID ?drug.
```

```
?drugtarget drugbank_interaction:human ?human . FILTER
(?human="1") .
```

```
?drugtarget drugbank_interaction:SwissProt_ID ?uniprot. }}
GROUP BY ?uniprot
```

Gefitinib (CID=123631) is widely investigated in non-small cancer of the lung, colon cancer, breast cancer and cancer of the head and neck. We are only using DrugBank and BindingDB dataset to find its targets in this example. DrugBank offers approved drug targets from literature. As it considers non-human target, we used a filter to select only human targets. BindingDB provides drug target experimental interaction results. IC50 is a measurement of binding affinity, we selected ic50<10000nm (nm is the default unit) to guarantee the compound is able to interact the target.

While using only DrugBank, it yields one result P00533 (EGFR: Epidermal growth factor receptor), whose well known inhibitor is Gefitinib. While adding another dataset BindingDB, We got a new possible target P04626 (ERBB2: Receptor tyrosine-protein kinase erbB-2). EGFR (the known Gefitinib target) is part of the ERBB receptor family, which has four closely related members: EGFR (ERBB1), HER2 (ERBB2), HER3 (ERBB3) and HER4 (ERBB4), thus Gefitinib is active against ERBB2 is not surprising, however, DrugBank does not have this information, thus it's necessary to consider all the datasets as much as possible.

4.2 Disease specific chemical discovery

Drug discovery process generally starts with the disease target identification, followed by hits selection (active compounds against the disease target). Finding the potential chemicals is the main task in the early drug discovery stage. As so many public data emerges, it becomes possible to find some intriguing chemicals by running one SPARQL. This is particularly of interest to some lower-profitable disease like malaria that is paid little attention by pharmaceutical companies. But academia and the small companies do not have adequate funding to screen millions of compounds in order to find the hits. Chem2Bio2RDF attempts to collect all this kind of public data and the compounds can be easily obtained by the following SPARQL.

Questions: Find a disease (i.e., Malaria) specific chemicals

SPARQL:

```
SELECT * WHERE {
  ?chemogenomics chemogenomics:CID ?compound_cid .
  ?chemogenomics chemogenomics:GENE ?gene_symbol .
  ?omim omim:gene ?gene_symbol .
  ?omim omim:Disorder_name ?disease . FILTER
  regex(?disease,"Malaria","i") .
}
```

The SPARQL starts to find disease causing genes in OMIM and then select the associated compounds of the genes. 1,003 compounds and 5 disease causing genes are returned from this query, which enable scientists to do further experiments.

4.3 Adverse drug reaction

Adverse reaction in drug usage, as the notion self-suggests, has serious consequence and is often subject to rigorous investigation in pharmaceutical R&D processes. In this case, we integrate into Chem2Bio2RDF another scenario to study the most significant pathways that are associated with a given drug affect. The association between side effect and pathway is made possible using pathway's gene components that are targets of related drugs. More specifically, we consider a gene is related to a certain side effect if at least two drugs targeting this gene incur the side effect. On top of that, if there exists a pathway that contains more

than 2 gene targets that associated with that side effect, an associative relationship between the pathway and side effect can be drawn. Clearly, the more these associative paths can be discovered, the stronger the evidence of such pathway-adverse drug effect association it becomes.

Question: Find the top 5 pathways in the KEGG pathway contain at least two of the efficient target that associated with a given side effect (i.e. *hepatomegaly*). A gene target is consider as efficient if the gene is targeted by at least two drugs that cause the given side effect.

SPARQL:

```
SELECT ?pathway_id (count(?pathway_id) as ?count) WHERE {
  ?sider2compound sider:side_effect ?side_effect . FILTER
  regex(?side_effect,"hepatomegaly","i") .
  ?sider2compound sider:cid ?compound .
```

```
?drug drugbank_drug:CID ?compound .
?drug2target drugbank_interaction:DBID ?drug .
?drug2target drugbank_interaction:SwissProt_ID ?uniprot .
?kegg_pathway kegg_pathway_protein:Uniprot ?uniprot .
?kegg_pathway kegg_pathway_protein:PathwayID ?pathway_id .
} GROUP BY ?pathway_id ORDER BY ?count
```

In this SPARQL, we need to first map side effect to drugs, and then map drug to targets. At the end, the target will be mapped to pathway. All the possible connections between side effect and pathway are counted to get the most reasonable pathways that associated with the given side effect.

Results: Available studies suggest that hepatic toxicity has been the most frequent single cause of safety-related drug withdrawal (e.g., ticrynafen, benoxaprofen, bromfenac, troglitazone, nefazodone). Hepatotoxicity discovered after approval from marketing also has limited clinical use of many drugs, including isoniazid, labetalol, trovafloxacin, tolcapone, and felbarmate. Thus, it is important to systematically review the compounds, gene target and pathways that associated with liver injury.

In our mapping process, we use the hepatic necrosis, hepatitis and hepatomegaly as an example to study the relationship between drugs, targets, pathways and side effect. The graph is shown as following:

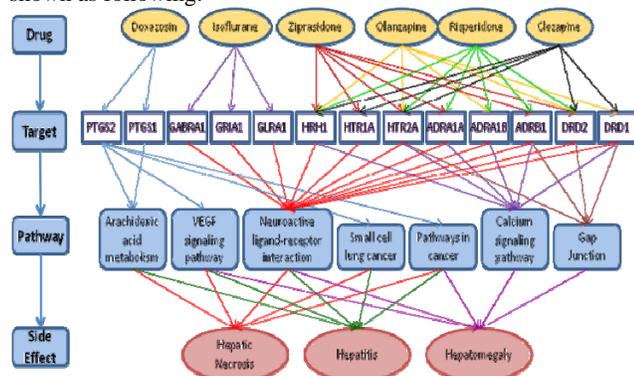


Figure 9: Adverse Drug Reaction

Figure 9 shows the top 5 pathways that associated with the three majority liver injury. It demonstrates that the mechanisms for hepatic necrosis and hepatitis are very close. They share the top 5 pathways: Arachidonic acid metabolism, VEGF signaling pathway, Neuroactive ligand-receptor interaction, small cell lung cancer, and pathways in cancer. The mechanism for

hepatomegaly is a little bit different. The top 5 pathways of hepatomegaly contain the calcium signaling and gap junction pathway, which are not involved in the hepatic necrosis and hepatitis. Literature review [14] shows that those pathways are highly correlated with liver injury. For instance, the increase concentration of calcium in the calcium signaling pathway will cause the damage of hepatic cell. The targets we discovered are also known as the major targets for liver diseases based on literature reviews [11].

5. Evaluation

We evaluate the outcomes from the previous case studies against the confirmed evidence in literature search as well as salient domain knowledge. In addition, we also developed comprehensive assessment to examine the coverage of datasets in their associated domains. Finally we illustrate the improvement gained through Chem2Bio2RDF as opposed to our previous work where no semantic web framework is deployed.

5.1 Study of systems chemical biology

The difficulties of polypharmacology are to explore the combination of targets and then to identify active compounds against the sets of targets. Linking between chemical, biological, systems, and phenotype data is demonstrated to be a promising way to address the problems. For example, linking between bioassay data and market drug data enables to explore the compounds similar to the drug that already shows polypharmacology. Quinacrine, which has been used as an anthelmintic and in the treatment of giardiasis and malignant effusions, shows polypharmacology. One compound loxapine (CID 71399) is found to show similar polypharmacology with quinacrine. Loxapine is active in both BioAssay 859 and BioAssay 377, whose targets are CHRM1 and ABCB1 respectively. As loxapine tends to be hydrophobic molecules, medicinal chemists would not be surprised that it is active in BioAssay 377 which identifies substrates (or inhibitors) for multidrug resistance transporter. It is also reported that loxapine might get metabolized to amoxapine that is a considerably weak antagonist in BioAssay 859 [9]. Other than loxapine, many identified compounds such as oxybutynin and dexamethasone were proved to show polypharmacology by literature reviews.

While we link bioassay data to pathways, we could identify the compounds that inhibit at least two of proteins in a pathway, leading to the pathway dysfunction. For example, compound CID 6,419,769 could interact with proteins HSD11B1 and AKR1C4, which are in the different branches of C21-Steroid hormone metabolism pathways. The blocking of the pathway might be able to partially explain why CID 6,419,769 has side effect [1]. In protein-protein interaction network, two proteins are connected if both are physically interacted. In terms of polypharmacology, the deletion of one protein does not affect the whole network, but if two connected nodes with high degree were deleted, the network would be disturbed. For example, by linking bioassay to PPI, we found that two compound (CID 460,747 and CID 9,549,688) are active against two high degree proteins (PLK1 and TP53) which are associated with cancer. Via linking data sources among different domains, not only promising compounds to be high effective could be identified but also the risk of compounds could be somehow evaluated.

5.2 Dataset and result coverage

There are parallel contributions from different data sources and vendors toward same domain (for example, KEGG and Reactome

are two independent data vendor provides same datum coverage over pathway domain), therefore the approach in which one single data source is made delegating the entire domain could produce incomplete outcome. To avoid this potential pitfall, for each domain in Chem2Bio2RDF, we have collected from a variety of data sources to make the ensemble as complete as possible. In particular, coverage for systems and chemogenomics is under close scrutiny due to their central roles in system chemical biology. Here we list the percentage coverage for PPI, pathway and chemogenomics and demonstrate the significance of integration using semantic web.

In PPI (Table 1), HPRD and DIP have 35,645 and 32,976 unique protein pairs respectively, and the total number of unique pairs in two datasets is 67,769. Each dataset contributes almost half of the pairs, and both share very little number of common pairs. The PPI network would not be complete if either dataset were ignored. Pathway is more complicated than PPI, since each organization could have its own definition of pathway, which makes the whole integration very difficult. For example, the pathway in Reactome is usually composed by a small number of proteins, although the total number of pathways is more than KEGG, the proteins involved in Reactome are far less than KEGG (Table 2). We are not able to judge which one is better, thus we have to consider all pathway datasets together. Figure 10 shows the dataset distribution of chemogenomics data. A chemical protein (gene) interaction is recorded as one entry, and all the unique interactions were derived from 10 datasets. We did not consider another two chemogenomics data sets (KEGG Ligand and PharmGKB), as KEGG Ligand includes only metabolic molecules rather than chemicals designed for drug discovery and many drugs in PharmGKB only provide names from which the chemical identifier is not able to be linked to compound. Many datasets only contribute a small portion of interactions so that it is not able to represent all chemogenomics data. Some of the datasets are quite small but they cannot be ignored. For example, BindingMOAD provides the protein and ligand complex crystal information which is the most accurate binding data. Kidb (called PDSP) is only interested in the receptors rather than all the protein spaces. DrugBank provides the most comprehensive drug and its target interactions.

Table 1: PPI data source distribution

Data source	# of records	percentage
HPRD	35645	52.6%
DIP	32976	48.7%
ALL	67769	

Table 2: Pathway data source distribution

Data source	protein		pathway	
	# of records	percentage	# of records	percentage
KEGG	8172	81.0%	192	34.8%
Reactome	4397	43.6%	360	65.2%
ALL	10091		552	

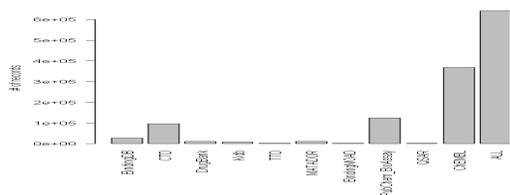


Figure 10: chemogenomics data source distribution

One dataset hardly be able to cover all the records of one domain, thus missing some of the datasets might neglect many link paths, sometimes, the results might even be changed significantly. We designed the following experiments to answer to the case 1 question, where we are searching for disease specific chemicals. In this task, we need to first identify disease genes (OMIM and PharmGKB) and then find chemicals interacting with the gene or gene expression products (BindingDB, DrugBank, ChEMBL and PubChem BioAssay). In our evaluation, we either selected one dataset in each step or selected all the datasets. As the table 4 shows, the number of chemicals retrieved from all datasets is far more than that only one dataset is considered.

Table 4: Results of discovering disease specific chemicals

Dataset used	# of chemical discovered
OMIM, DrugBank	77
OMIM, BindingDB	175
OMIM, PubChem	3577
OMIM, ChEMBL	2036
PharmGKB, DrugBank	292
PharmGKB, BindingDB	1256
PharmGKB, PubChem	28410
PharmGKB, ChEMBL	20513
ALL	45606

5.3 Comparison with previous work

In the previous study, only 4 datasets (PubChem, DrugBank, KEGG and HPRD) were studied [5], no attempt is made to integrate multiple datasets within one domain simultaneously. In Chem2Bio2RDF, efforts have been made to provide coverage over the comprehensive landscape of system chemical biology, 23 datasets so far has been integrated, as a direct result; the expected entity pairs discovered has increased significantly. Moreover, LPG would answer user's queries by automatically locating related datasets.

In addition, the performance has been improved as well as lots of time was required to parse the heterogeneous datasets previously. But now, since data sources have been published as RDF triples, SPARQL queries can be used to directly identify hidden knowledge.

In the work of adverse drug reaction study[22], it relies on the internal data source which includes 1,458,680 unique compounds and 2,190 unique targets, Chem2Bio2RDF currently covers 235,313 compounds, 19,534 genes, 641, 855 chemogenomics which are public accessible. However, we must take care of the quality of data in public domain. We also need to build robust chemogenomics predictive model based on the rich chemogenomics data source, as it is impossible that every compound is tested against all the targets.

6. Conclusion and Future work

This paper discusses our on-going effort to create a user-friendly Linked Open Data portal for chemical biology. It applies semantic web technologies to integrate data and identify hidden association. Comparing with other LOD portals, currently Chem2Bio2RDF portal contains the following features: a RDF faceted browser, an automatic SPARQL query generator, a MEDLINE/PubMed literature cross-validation service, a cytoscape plugin, and a WENDI 2.0 service for disease-drug-compound discovery.

Our future work will focus on: (1) extending cytoscape visualization plugin to enable analytical graph mining and association detection; (2) adding a predictive model for compound gene associations, as not every compound is tested against all the targets, a predictive model is necessary; (3) applying CRF(common random field) name entity method to identify the chem./bio terms to provide an extensible gene, compound, pathway, disease and drug extraction for MEDLINE/PubMed literatures; (4) enriching provenance data to the current Chem2Bio2RDF data to enable provenance-based filtering and visualization; and (5) identifying semantic associations among different types of entities to enable complex drug discovery.

Acknowledgement

This work is funded by NIH VIVO project (UF09179) and Eli Lilly.

7. References

- [1] R. C. Andrews, O. Rooyackers and B. R. Walker. Effects of the 11 Beta-hydroxysteroid Dehydrogenase Inhibitor Carbenoxolone on Insulin Sensitivity in Men with Type 2 Diabetes. *J. Clin. Endocrinol. Metab.*, 88, 285-291, 2003.
- [2] E. Antezana, W. Blondé, M. Egaña, A. Rutherford, R. Stevens, B. De Baets, V. Mironov and M. Kuiper. BioGateway: a semantic systems biology tool for the life sciences. *BMC Bioinformatics*, 10 Suppl 10:S11, 2009.
- [3] F. Belleau, M. A. Nolin, N. Tourigny, P. Rigault and J. Morissette. Bio2RDF: Towards a mashup to build bioinformatics knowledge systems. *Journal of Biomedical Informatics*, 41, 706-716, 2008.
- [4] C. Bizer and R. Cyganiak. D2R Server - Publishing Relational Databases on the Semantic Web. *Poster at the 5th International Semantic Web Conference*, 2006.
- [5] B. Chen, D. J. Wild, and R. Guha. PubChem as a Source of Polypharmacology. *Journal of Chemical Information and Modeling*, 49(9), pp 2044-2055, 2009.
- [6] H. Chen, L. Ding, Z. Wu, T. Yu, L. Dhanapalan and J. Y. Chen. Semantic web for integrated network analysis in biomedicine, *Brief Bioinform.*, 10(2):177-92, 2009.
- [7] K. Cheung, H. R. Frost, M. S. Marshall, E. Prud'hommeaux, M. Samwald, J. Zhao and A. Paschke. A journey to Semantic Web query federation in the life sciences. *BMC Bioinformatics*, 10 (Suppl 10): S10, 2009.
- [8] K. Cheung, K. Yip, Smith A., R. Deknikker, A. Masiar, and M. Gerstein. YeastHub: A semantic web use case for integrating data in the life sciences domain. *Bioinformatics*, 21(Suppl 1): i85-96, 2005.
- [9] J. Coupet, S. K. Fisher, C. E. Rauh, F. Lai and B. Beer. Interaction of Amoxapine with Muscarinic Cholinergic Receptors - an in Vitro Assessment. *Eur. J. Pharmacol.*, 112, 231-235, 1985.
- [10] B. J. Druker, M. Talpaz, D. J. Resta, B. Peng, E. Buchdunger, J. M. Ford, N. B. Lydon, H. Kantarjian, R. Capdeville, S. Ohno-Jones and C. L. Sawyers. Efficacy and Safety of a Specific Inhibitor of the BCR-ABL Tyrosine Kinase in Chronic Myeloid Leukemia. *N. Engl. J. Med.*, 344, 1031- 1037, 2001.
- [11] G. Gong, G. Waris, R. Tanveer and A. Siddiqui. Human hepatitis C virus NS5A protein alters intracellular calcium levels, induces oxidative stress, and activates STAT-3 and NF-kappa B. *Proc Natl Acad Sci U S A.*, 98(17):9599-604, 2001.
- [12] A. L. Hopkins. Network Pharmacology: The Next Paradigm in Drug Discovery. *Nat. Chem. Biol.*, 4, 682- 690, 2008.
- [13] A. Jentzsch, J. Zhao, O. Hassanzadeh, K. Cheung, M. Samwald and B. Andersson. Linking open drug data. Graz, Austria, 2009.
- [14] B. E. Jones and M. J. Czajka. III. Intracellular signaling in response to toxic liver injury. *Am J Physiol.*, 275(5 Pt 1):G874-8, 1998.
- [15] C. T. Keith, A. A. Borisy and B. R. Stockwell. Multicomponent Therapeutics for Networked Systems. *Nat. Rev. Drug. Discovery*, 4, 71-110, 2005.
- [16] E. K. Neumann. A life science semantic web: are we there yet? *Science*, 283:22-5, 2005.
- [17] E. K. Neumann and D. Quan. Biodash: a semantic web dashboard for drug development. *Pac Symp on Biocomput*, 11: 176-187, 2006.
- [18] E. K. Neumann, E. Miller and J. Wilbanks. What the semantic web could do for the life sciences. *Drug Discovery Today:BIOSILICO*, 2:228-34, 2006.
- [19] T. I. Oprea, A. Tropsha, J. Faulon and M. D. Rintoul. Systems chemical biology. *Nat Chem Biol*, 3:447-450, 2007.
- [20] D. J. Wild. Mining large heterogenous datasets in drug discovery. *Expert Opinion on Drug Discovery*, 4(10), pp 995-1004, 2009.
- [21] J. Scheiber, B. Chen, M. Milik, S. C. Sukuru, A. Bender, D. Mikhailov, S. Whitebread, J. Hamon, K. Azzouli, L. Urban, M. Glick, J. W. Davies and J. L. Jenkins. Gaining insight into off-target mediated effects of drug candidates with a comprehensive systems chemical biology analysis. *J Chem Inf Model.*, 49(2):308-17, 2009.
- [22] T. Slater, C. Bouton and E. S. Huang. Beyond data integration. *Drug Discovery Today*, 13(13-14):584-9, 2008.
- [23] A. Smith, K. Cheung, K. Yip, M. Schultz and M. Gerstein. LinkHub: A semantic web system that facilitates cross-database queries and information retrieval in proteomics. *BMC Bioinformatics*, 8(Suppl 3): S5, 2007.
- [24] N. Villanueva-Rosales, K. Osbahr and M. Doumontier. Towards a Semantic Knowledge base for Yeast biologists. *J Biomed Inform.*, 41(5):779-89, 2008.
- [25] J. Wang, J. Y. Zhou and G. S. Wu. ERK-Dependent MKP-1-Mediated Cisplatin Resistance in Human Ovarian Cancer Cells. *Cancer Res.*, 67, 11933-11941, 2007.
- [26] L. Xie, J. Li, L. Xie and P. E. Bourne. Drug discovery using chemical systems biology: identification of the protein-ligand binding network to explain the side effects of CETP inhibitors. *PLoS Comput Biol.*, 5(5), 2009.
- [27] B. Chen, X. Dong, D. Jiao, H. Wang, Q. Zhu, Y. Ding, D. J. Wild. Chem2Bio2RDF: a semantic framework for linking and data mining chemogenomic and systems chemical biology data *BMC Bioinformatics*, forthcoming
- [28] Y. Wang, J. Xiao, T. O. Suzek, J. Zhang, J. Wang and S. H. Bryant. PubChem: a public information system for analyzing bioactivities of small molecules. *Nucl Acids Res*, 2009, 37:W623-W633.
- [29] M. Kanehisa, S. Goto, M. Hattori, K. F. Aoki-Kinoshita, M. Itoh, S. Kawashima, T. Katayama, M. Araki, and M. Hirakawa. From genomics to chemical genomics: new developments in KEGG. *Nucleic Acids Res*. 2006, 34:D354-357.
- [30] C. J. Mattingly, G. T. Colby, J. N. Forrest and J. L. Boyer. The Comparative Toxicogenomics Database (CTD). *Environ Health Perspect*. 2003, 111(6):793-795.
- [31] T. Liu, Y. Lin, X. Wen, R. N. Jorissen and M. K. Gilson. BindingDB: a web-accessible database of experimentally determined protein-ligand binding affinities. *Nucl Acids Res* 2007, 35:D198-D201.
- [32] S. Günther, M. Kuhn, M. Dunkel, M. Campillos, C. Senger, E. Petsalaki, J. Ahmed, E. G. Urdiales, A. Gewiss, L. J. Jensen, R. Schneider, R. Skoblo, R. B. Russell, P. E. Bourne, P. Bork and R. Preissner. SuperTarget and Matador: resources for exploring drug-target relationships. *Nucl Acids Res*. 2008, 36:D919-922
- [33] X. Chen, Z. L. Ji and Y. Z. Chen. TTD: Therapeutic Target Database. *Nucleic Acids Res*. 2002 1;30(1):412-5.
- [34] D. S. Wishart, C. Knox, A. C. Guo, S. Shrivastava, M. Hassanali, P. Stothard, Z. Chang and J. Woolsey. DrugBank: A Comprehensive Resource for in Silico Drug Discovery and Exploration. *Nucleic Acids Res*. 2006, 34,
- [35] www.ebi.ac.uk/chembl/
- [36] http://pdsp.med.unc.edu/pdsp.php
- [37] M. Hewett, D. E. Oliver, D. L. Rubin, K. L. Easton, J. M. Stuart, R. B. Altman and T. E. Klein. PharmGKB: the Pharmacogenetics Knowledge Base. *Nucleic Acids Res*. 2002, 30(1):163-5.
- [38] M. L. Benson, R. D. Smith, N. A. Khazanov, B. Dimcheff, J. Beaver, P. Dresslar, J. Nerothin, and H. A. Carlson. Binding MOAD, a high-quality protein-ligand database. *Nucleic Acids Res*. 2008, 36(Database issue): D674-D678.