

SPARQL Builder: 生物学研究者が SPARQL を使いこなすための補助ツール

SPARQL Builder: A tool for biologists to build a SPARQL query

山口敦子¹ 古崎晃司² 呉紅艶¹ 小林紀郎³

Atsuko Yamaguchi¹, Kouji Kozaki², Hongyan Wu¹ and Norio Kobayashi³

¹ ライフサイエンス統合データベースセンター

¹ Database Center for Life Science

² 大阪大学 産業科学研究所

² The Institute of Scientific and Industrial Research, Osaka University

³ 理化学研究所 情報基盤センター

³ Advanced Center for Computing and Communication, RIKEN

概要: 生命科学で扱う多様なデータを有機的に扱い統合するために、セマンティックウェブ技術を活用したデータベース構築が広く行われるようになってきた。RDF 化され、SPARQL エンドポイントを公開する生命科学系データベースの数は年々増加している。そこで、生物学研究者の多種多様なユースケースに対応したデータ処理を行うためには、SPARQL エンドポイントを有効に活用したアプリケーション開発が必要となる。しかし、SPARQL クエリは RDF データのグラフ構造を熟知していなければ記述することができず、生物研究者が欲しいデータの取得やさらに取得したデータを効果的に処理するアプリケーションの構築は一般に困難である。そこで、あらかじめ RDF データ構造を知らなくても生物学研究者が欲するデータを取得する SPARQL クエリを生成できるツール SPARQL Builder の試作を行った。本発表では、システム設計と、試作で得られた問題点や知見について述べる。

1. はじめに

1.1 背景

生命科学研究においては、分野の細分化や計測機器の発展に伴い、多種多様かつ膨大なデータが産出される。これらのデータを集約し公開するデータベースは、生物学研究者に必要な不可欠なインフラとして利用されている。しかしながら、従来、プロジェクトや目的ごとなど個別にデータベースが作成、公開されてきたため、データベースの所在や利用方法がわかりにくいなど多くの問題を抱えていた。

近年では、生命科学の多種多様かつ膨大なデータを統合的に扱うために、セマンティックウェブ技術に基づいたデータの記述や公開が推進されてきた。タンパク質配列データベース UniProt [1,2] では、2008 年頃から大量のデータベース間のリンク情報を扱うために RDF データモデルを採用している。2013 年 10 月にはヨーロッパ最大のバイオインフォ

マティクスの拠点である EMBL-EBI が、UniProt に加えて、ChEMBL, Expression Atlas など生物学でよく使われるデータベースを RDF 化し、その RDF データにアクセスできる SPARQL エンドポイントを公開した[3]。

このような状況の下、RDF 化され SPARQL エンドポイントで公開されるデータを一般の生物学研究者が有効活用できるようにするためには、SPARQL に準拠した研究者の要求に沿ったアプリケーション開発が必要である。そのために、アプリケーションで必要と思われる SPARQL クエリを予め用意しておく方法が現在は主流である。例えば、BioGateway では、データセットと共に SPARQL クエリのセットが用意されている[4]。しかしながら、生物学研究者の要求するデータは非常に多種多様であり、RDF のグラフ構造もデータ毎に異なり、データ統合によってグラフ構造も順次拡張される。このため、予めすべての SPARQL クエリをアプリケーション側で用意することは不可能である。また、SPARQL クエリ

の構築に当たっては、RDFデータの構造や仕様を熟知し、自分の欲しいデータ構造を精緻に記述することが求められるが、一般の生物学研究者にとってこの作業は技術的に難しい。

我々は、RDFに準拠したアプリケーションTogoTableを例に取り、TogoTableへのSPARQLクエリの半自動的追加を目的として、一般の生物学研究者が自分の欲しいデータをSPARQLで記述できる手法を確立する研究を行っている。本発表ではこれまでに検討した手法を述べ、それをSPARQL Builderと呼ぶツールとして試作した結果、実装に伴う問題点、さらその解決方法について述べる。

1.2 表データ注釈付加ツール「TogoTable」

TogoTableはユーザがアップロードした表形式のデータに対して、公開されているSPARQLエンドポイントから関連するデータを付加できるウェブツール(図1)である[3]。

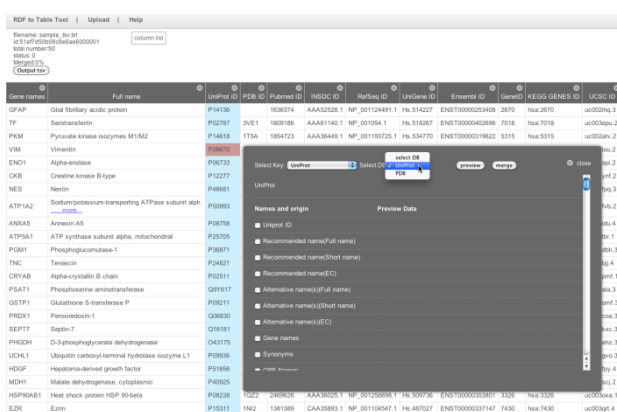


図 1: TogoTable

典型的な使い方としては、(1) 自分の実験データを表形式でアップロードする、(2) アップロードした表の中からデータベースのIDをデータとして含むカラムを1つ選び、データを取得したいデータベースをリストから選ぶ、(4) 取得可能なデータの種類の表示されるので、取得したいデータの種類を選ぶ、(5) ユーザが選んだカラム(= IDの集合)とデータの種類を利用して、予め作成されたSPARQLクエリリストからクエリを選び、ユーザが選んだデータベースのSPARQLエンドポイントから取得して、表の左端に新たなカラムとして表示する、である。

データベースごとに取得可能なデータの種類はそれを取得可能なSPARQLクエリとペアで予めアプリケーション内に登録されている。このSPARQLク

エリのリストは、段階(4)での取得可能なデータの種類の表示、および、段階(5)でのSPARQLエンドポイントからのデータ取得に利用されている。

2. SPARQL Builder

2.1 開発要件

SPARQLを全く書けないユーザに任意のSPARQLエンドポイントから任意のSPARQLクエリが書けるような補助をするのは様々な面で困難である。例えば、エンドポイントごとにRDFデータをブラウズさせ、入力データと目的のデータをユーザが示すことでSPARQLクエリを自動記述する方法を考える。しかしながら、予めユーザがデータベースの全体像を知らなければ、目的のデータがそのデータベースにあるかどうかさえ不明であり、あるとしても、どちらへどうブラウズすれば目的のデータにたどり着けるか不明である。

そこで、本発表では、目的を「ユーザがTogoTableのSPARQLクエリリストを半自動的に拡充可能とするツール」と状況を限定することで、必要とされるSPARQLクエリの形式の限定を試みた。

表形式データに対して、各カラムは何らかのクラスに対応すると仮定する。この場合、ユーザがTogoTableインタフェースを通じて指定するのは、SPARQLエンドポイント、入力クラス(前述したTogoTableの段階(2)に相当)および出力クラス(TogoTableの段階(4)に相当)である。しかしながら、あるRDFグラフにおいて、入力クラスや出力クラスを`rdfs:domain`, `rdfs:range`にもつことができるPropertyは一般に複数あり、また、途中で複数のクラスを挟んで関連をもつこともありえる。そして、入力クラス内のインスタンスと出力クラスのインスタンスがグラフ上でどの経路を通じてつながっているかによって、インスタンス間の関係性は異なる。したがって、ユーザが要求する関連データを出力させるには、SPARQLエンドポイント、入力クラス、出力クラス、入力クラスと出力クラスの間の経路、の4つを指定させる必要がある。

本発表では、この4つをユーザに指定させ、それをもとにSPARQLクエリを作成するシステムSPARQL Builderについて考える。

2.2 システム概要

前述した開発要件を満たすため、SPARQL Builderは、(1)ユーザがSPARQLエンドポイント、入力クラス、出力クラス、クラス間経路を指定する画面、

(2) 経路計算モジュール, (3) 経路から SPARQL クエリを作成するモジュールの3つの部分から構成される。

図2は SPARQL Builder のフローを表したものである。まず、画面から SPARQL エンドポイント, 入力クラス, 出力クラスをユーザより取得する。取得した情報は経路計算モジュールに引き渡される。経路計算モジュールは, 入力クラスの隣接クラスをあるステップ数まで SPARQL エンドポイントに問い合わせ取得し, その情報から入力クラスを始点とし出力クラスを終点とする経路をすべて計算する。計算された経路群は, 画面に表示される(図3)。ユーザが表示された経路群の中から一つの経路を選ぶと, その経路から SPARQL クエリが作成され, 画面に表示される。また, 作成されたクエリは SPARQL エンドポイントに投げられ, 得られた結果も画面に表示される(図4)。

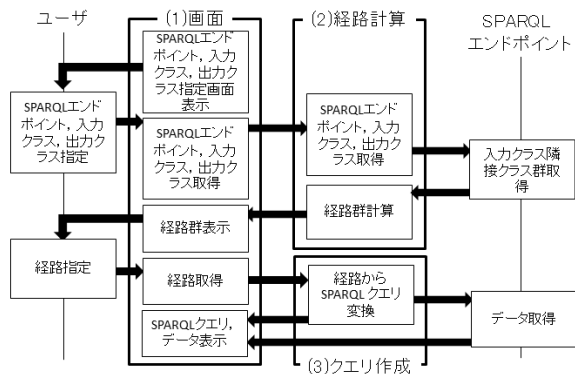


図 2: SPARQL Builder フロー図

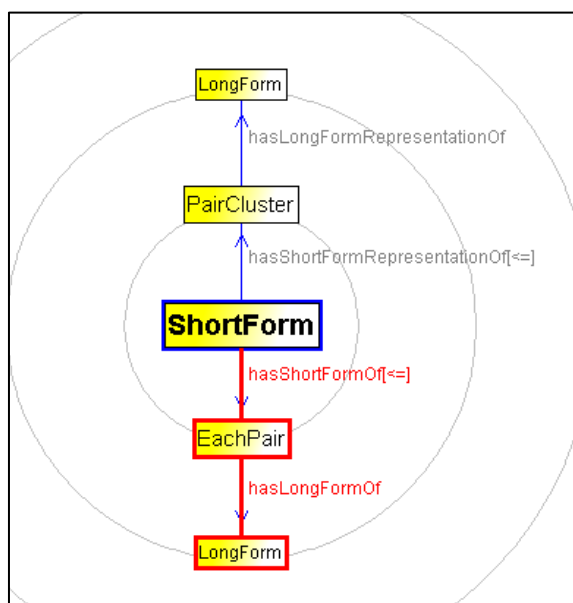


図 3: 経路選択画面

c1	c2	c3	i2	i3	i1
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	BMP-decapentaplegic@en	lpp@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C-1-3a-DE, followed by C@en	def@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D-4flurophenyl@en	dff@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C-1-deoxyinosine-5-phosphate@en	des-c@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C-4-carboxylate permease@en	dca@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C-4-carboxylate transport@en	dct@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C-4-carboxylate transport mutants@en	dctm@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C-4-carboxylate transporter@en	dcta@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C-4-carboxylate transporter system@en	dctrs@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C-4-dicarboxylic acid transport genes@en	dctrd@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	C8-H-studeralad da110 2@en	dt@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D binding protein@en	dtp@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D mutant@en	dpm@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	The expansion work of coil swelling RA@en	delta.s@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D-1-thr-2-amino-1-g-methylamino-1,3-propane	denf@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D-1-deoxyribose 5-phosphate synthase gene@en	des@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D-3-acetoxys-cis-2-thydro-5-(2-dimethylamino-e	dmsam@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D-4-deficient L14 mutant by inactivating the first	dmb@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D-4a-D-4 ligase@en	dml@en	
http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	http://purl.org/alle/ids/...	D, 4a, D, 4a, D, 4a, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100	del70@en	

図 4: 検索結果表示画面

2.2 経路計算モジュール

経路計算モジュールでは, ユーザが指定する任意の SPARQL エンドポイントにおいて, 入力クラスを始点とし, 出力クラスを終点とし, クラス間のインスタンスをつなぐプロパティを辺とするグラフ内のクラス間経路を計算する。

任意の SPARQL エンドポイントをターゲットとするため, 対象の SPARQL エンドポイントから経路計算に必要なデータを取得する必要がある。しかしながら, 通常の経路探索問題と異なり, 同じクラスにたどり着いた別の経路をカットすることができないため, データによっては経路数が爆発的に増えてしまう。そのため, 現時点では, 経路の長さの最大数を指定し(デフォルトは3), 経路数の爆発的増大を抑え, その範囲内で出力クラスを終点とする経路を探索する。

2.2 経路表示画面

経路表示画面では, システムが提示した経路群からユーザが必要とするデータを表す一つの経路を選ぶような表示をする必要がある。その一方で, 全ての経路をリスト化すると重複した情報が多く, 経路をユーザが細かく見比べて違いを見ながら選ぶ必要がある。そのため, 現在は, 経路群を樹状に表し, 葉(出力クラスに対応)をユーザが選ぶことで, 経路を指定できるようにした。樹状に表すことで, 重複した情報はまとめられ, かつ, 根から各葉への経路は一つであるため, 経路を一意に指定することができる。

3. 考察

SPARQL Builder の試作に伴い、いくつかの問題点が明らかとなった。これらの問題点は、おおまかにシステム上の問題点とデータベース側の問題点に分けることができる。

システム上の最も大きな問題は経路数の爆発により、経路群表示までの時間がかかりがちであり、かつユーザが選ぶ際にも候補が多すぎて選びにくいという点である。たとえば、全てのクラスは `rdfs:Class` につながり、かつ、`rdfs:Class` はそれ自身が `rdfs:Class` のインスタンスであることより、クラス間グラフ上で `rdfs:Class` は非常に巨大な次数を持つ頂点となっている。そのため、素朴にクラス間経路を探索すると、一度 `rdfs:Class` を経由する経路が大量に出現することとなる。また、次数が巨大な頂点については、経路数爆発の問題がおこるだけではない。巨大な次数を持つ頂点を通過した経路は、情報がその頂点を経由する際に失われるため、クエリとしては意味がないことが多い。そのため、適切な方法でこのような経路をユーザに提示する前に削除していく必要がある。この削除の基準をどう決めていくかは現在も検討中である。

データベース側の問題点としては、入力クラスの隣接クラス群取得の際、クラスが定義されていないインスタンスは探索できない、また、プロパティに `rdfs:domain`, `rdfs:range` が定義されていない場合、インスタンスを辿る必要があるため、隣接クラスの取得に時間がかかる、などがある。後者の問題点は今後の計算機やトリプルストアの性能向上、あるいは隣接クラス取得クエリの改良によって、改善の可能性があるものの、前者については、システム側で解決するのは非常に困難である。したがって、クラスが定義されていないインスタンスに対しては、現在は探索対象から外すしかなく、データはあるのに取得できない、あるいは経路はあるのに、途中で分断されてしまうことが起きてしまう状況になっている。この問題を解決するため、データベース探索法の研究を目下進めている。

4. まとめ

生物学研究者で SPARQL に習熟していないデータベースユーザを対象に、任意の SPARQL エンドポイントに対して、ユーザが求めるデータを取得できる SPARQL を記述する補助ツール SPARQL Builder を、TogoTable 上での利用に限定し、SPARQL の形式を制限したうえで試作した。

今後の課題としては、前節で述べた問題点、特に

経路探索の問題点について改良をすすめていきたいと考えている。経路の枝刈りの基準として、情報量などの統計的基準と、実データを見てその傾向から示される基準の両方から検討をすすめる予定である。

また、クラス指定や経路指定の画面の設計について、SPARQL に不慣れなユーザにより使いやすいように変更を加えていきたいと考えている。例えば、入力クラスと出力クラスの指定の際、何らかの基準でランキングを行って候補クラスを表示するなど、よりスムーズにユーザが欲しいデータの指定が行われるような画面システム設計が必要であると考えている。

参考文献：

- [1] UniProt, <http://www.uniprot.org/>
- [2] Redaschi, Nicole and Consortium UniProt: UniProt in RDF: Tackling Data Integration and Distributed Annotation with the Semantic Web. Nature Precedings, <<http://dx.doi.org/10.1038/npre.2009.3193.1>> (2009)
- [3] RDF Platform <EMBL-EBI, <http://www.ebi.ac.uk/rdf/>
- [4] Erick Antezana, Ward Blondé, Mikel Egaña, Alistair Rutherford, Robert Stevens, Bernard De Baets, Vladimir Mironov and Martin Kuiper: BioGateway: a semantic systems biology tool for the life sciences, BMC Bioinformatics, 10(Suppl 10):S11, (2009)
- [5] TogoTable <http://togotable.dbcls.jp/>