

○平川 英樹¹、藤澤 貴智²、長崎 英樹¹、森 宙史²、福島 敦史³、ジェルフィ アンドレア¹、市原 寿子¹、
中村 保一²、金谷 重彦⁴、有田 正規²、黒川 顕²、磯部 祥子¹、田畑 哲之¹

1. かずさDNA研究所
2. 国立遺伝学研究所
3. 理化学研究所環境資源科学研究センター
4. 奈良先端科学技術大学院大学先端科学技術研究科

要旨

統合化推進プログラムにおいて、植物ではPlant GARDEN (<https://plantgarden.jp>)、微生物ではMicrobeDB.jp (<https://microbedb.jp>)、メタボロームではMetaboBank (<https://www.ddbj.nig.ac.jp/metabobank/>) が構築されており、それぞれが扱う以下の主要なデータセットについてRDF化が行なわれている。

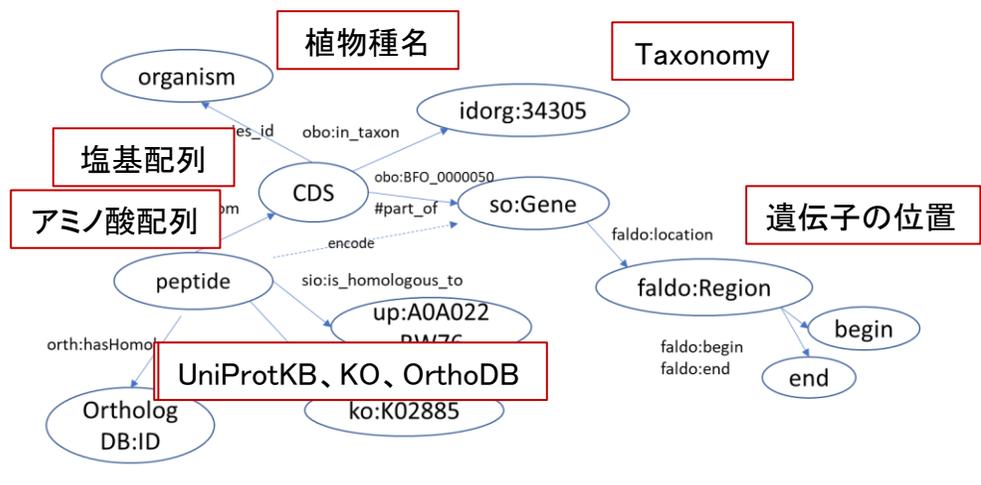
- ・Plant GARDEN: 植物ゲノム、遺伝子、アノテーション、DNAマーカー、QTL
- ・MicrobeDB.jp: メタ16Sデータ、メタゲノム解析データ、オントロジーアノテーションされたサンプル属性情報
- ・MetaboBank: 実験で得られた代謝産物に関する生データや実験に関するメタデータ

植物ゲノム統合では、PubTatorにより各植物種に関する病気や遺伝子、化合物などの情報について文献キュレーションを行い、DBCLSのTogoDBを用いて得られたデータに対してRDF化を行っている。現在、SPARQLエンドポイントを構築し、各統合化データベースの連携によって、更なるデータ統合を図っている。SPARQLを用いることで、植物種名や遺伝子、化合物(代謝産物)などでクエリ検索し、新たな知識発見ができるようになることを試みている。

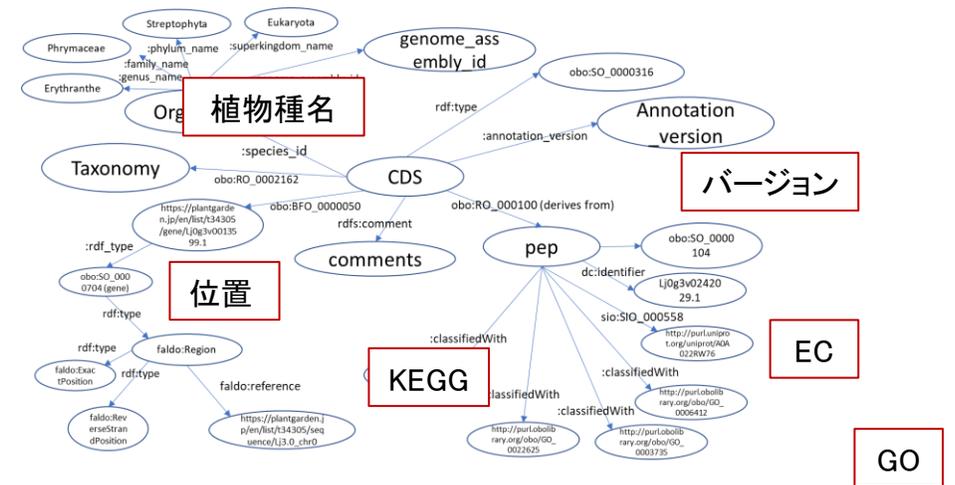
植物ゲノム関連データのRDF化

データ連携を行うため、Plant GARDENで公開している遺伝子、アノテーション、DNAマーカー、QTL情報のRDF化を行った。国立遺伝学研究所とDBCLSからのご協力を頂いた。

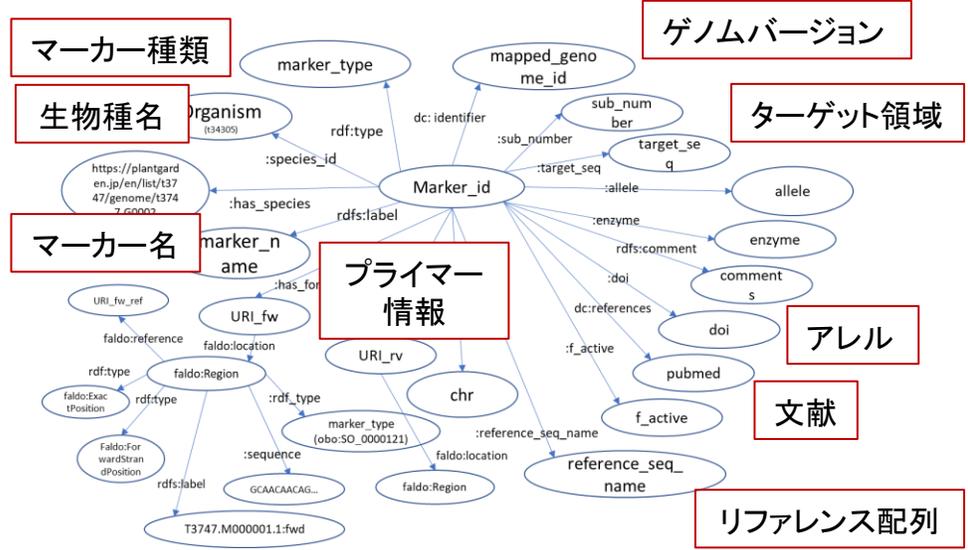
① 遺伝子情報(gff)



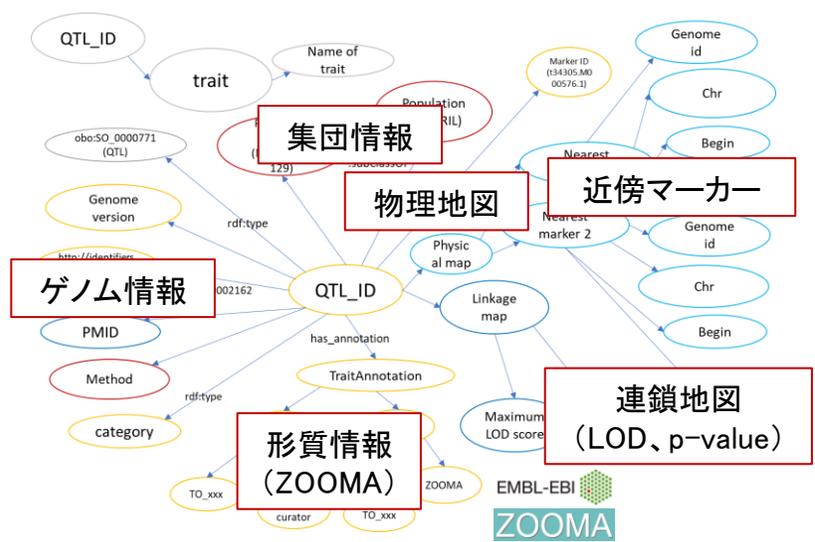
② アノテーション情報 (Hayai-Annotation)



③ DNAマーカー情報



④ QTL情報



RDFを用いたデータ連携

植物は微生物と相互作用しており、様々な二次代謝産物を産生する。RDFを用いて植物、微生物、メタボローム統合データベースを繋ぐことを試みている。

Plant GARDEN

- ・ゲノム情報(134種)
ゲノム配列、CDS、PEP、アノテーション
- ・DNAマーカー(38種)
- ・QTL(27種)

PubTatorを用いた文献キュレーションのRDF

- ・生物種、病気、遺伝子、化合物など
- ・植物と微生物の相互作用
Plant Growth Promoting Bacteria (PGPB; 共生、成長促進)、病原性細菌(細菌、カビ、ウイルス)
- ・代謝産物(MESH ID)

連携



MetaboBank

- ・実験で得られた代謝産物に関する生データや実験に関するメタデータ
- ・データ数
- ・KNASAcKのRDF (DBCLSにて実施)



連携



- ・系統・遺伝子・環境の3つの軸に沿って整理・統合されたフルRDFのデータベース
- ・メタゲノム解析データ(メタ16S rRNA、遺伝子)
- ・オントロジーアノテーションされたサンプル属性情報
- ・単細胞の真菌類、藻類のゲノムデータ
- ・解析プロトコルの標準化
- ・解析パイプラインの開発

外部データベースのRDF



化合物(11,886個)
酵素反応(13,673個)

Medical Subject Headings RDF



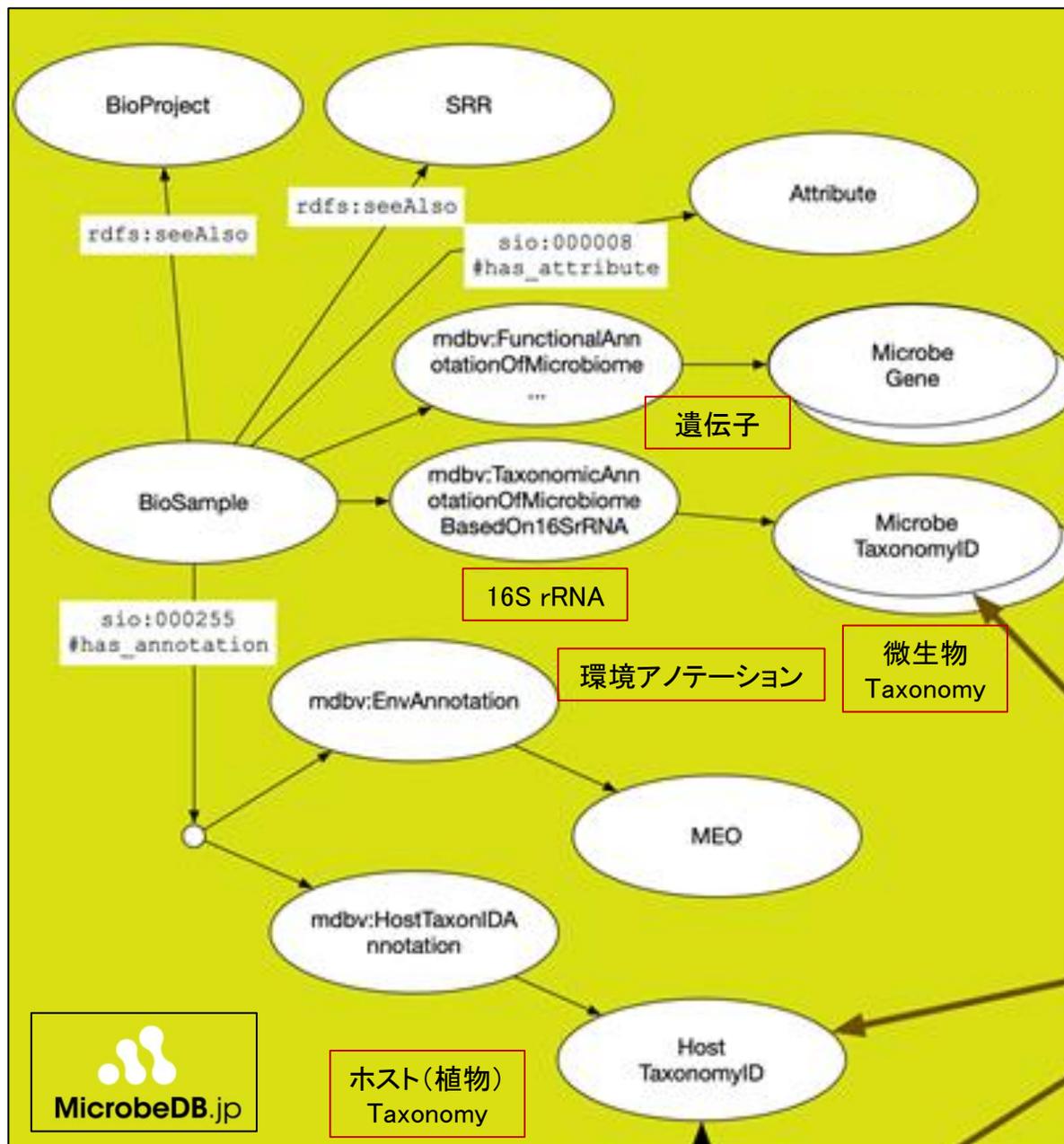
National Bioscience
Database Center

統合化推進プログラムの
他のDB

MicrobeDB.jpにおけるRDFの関係図

以下の項目についてRDF化を行った。

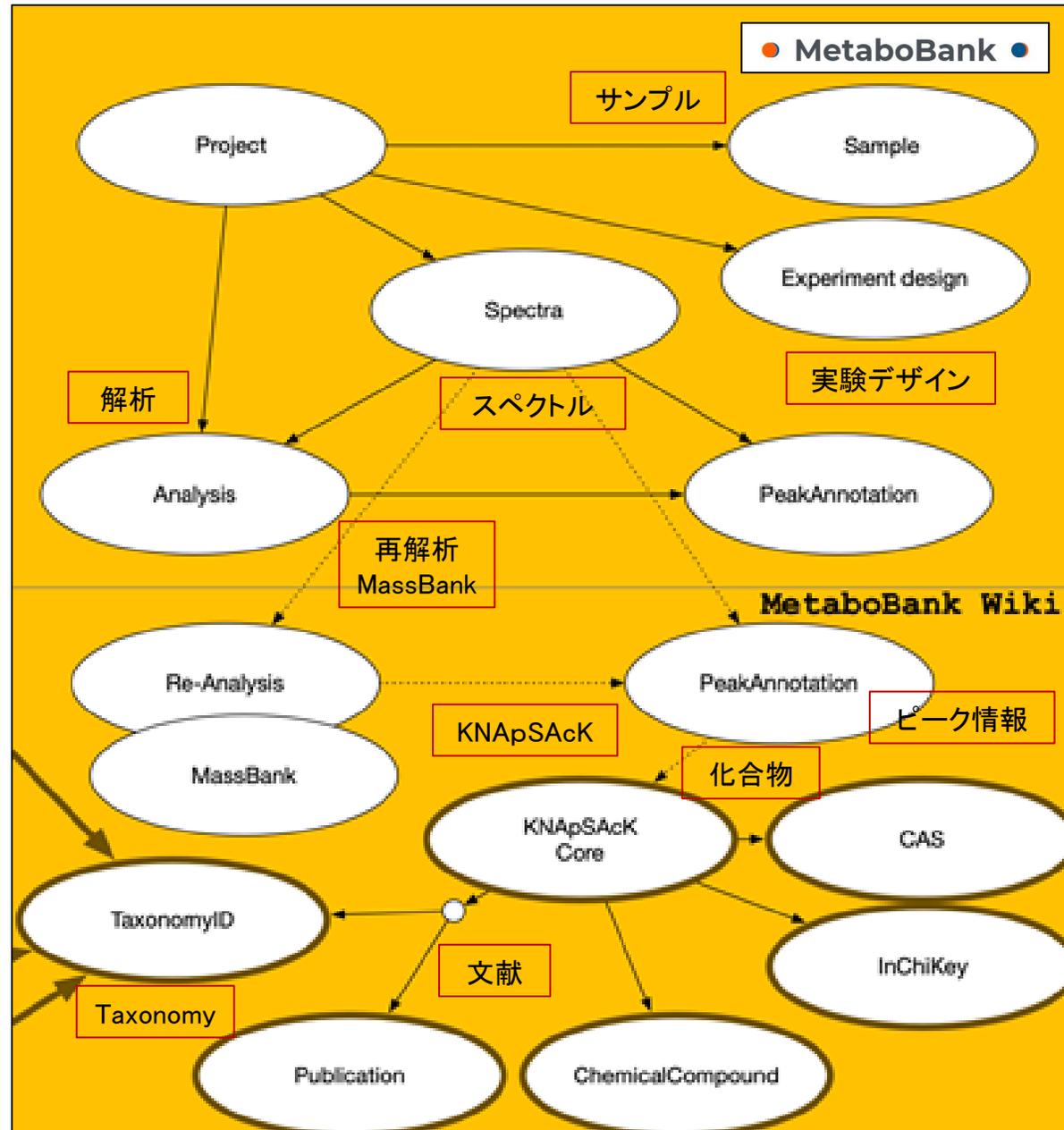
- ・ホスト(植物種)名
Taxonomy ID
- ・環境アノテーション
MEO
- ・BioSample
遺伝子
16S rRNA
微生物名、Taxonomy ID
- ・BioProject
- ・SRR



MetaboBankにおけるRDFの関係図

以下の項目についてRDF化を行った。

- ・生物種名
Taxonomy ID
- ・Project ID
サンプル情報
実験デザイン
スペクトル
解析内容
ピークアノテーション
再解析 (MassBank)
再解析されたpeakアノテーション
- ・KNApSack Core
文献情報
化合物
CAS
InChiKey



エンドポイントに格納したトリプル

各統合データベースにおいてエンドポイントを作成した。以下の情報についてのトリプルデータを格納した。

MetaboBank(メタボローム)
<https://mbs1.ddbj.nig.ac.jp/sparql>

Plant GARDEN(植物)
<http://plantgarden.jp/sparql>

MicrobeDB.jp(微生物)
<https://microbedb.jp/sparql>

Abbreviate URIs	#triples
http://plantgarden.jp/resource/gene	323,678,344
http://metadb.riken.jp/db/plantMetabolomics	33,893,537
http://ddbj.nig.ac.jp/ontologies/taxonomy	33,503,477
http://plantgarden.jp/resource/pg_marker	12,435,112
http://mb-wiki.nig.ac.jp/resource	3,489,616
http://plantgarden.jp/resource/rgaugury	2,337,589
http://plantgarden.jp/resource/pubtator	447,403
http://plantgarden.jp/resource/genome	4,396
http://localhost:8890/DAV/	2,959
http://www.openlinksw.com/schemas/virttrdf#	2,479
http://plantgarden.jp/resource/species	1,792
http://plantgarden.jp/resource/subspecies	1,639
http://plantgarden.jp/resource/trait	1,378
http://www.w3.org/2002/07/owl#	160
http://localhost:8890/sparql	14
http://www.w3.org/ns/ldp#	3

遺伝子

MetaboWikiエントリー情報

Taxonomy

DNAマーカー

リソース

病害抵抗性遺伝子

センテンス

ゲノム

植物種

植物種(subspecies)

形質

格納されたデータの確認 (Endpoint browser)

DBCLSから提供されているEndpoint browserを用いてデータが繋がっていることを確認した。

The image shows a screenshot of the DBCLS Endpoint browser interface. The top navigation bar includes 'DBCLS', 'Research', 'Services', 'Contact', and 'About'. The main header area contains the 'Endpoint browser' title and the DBCLS logo. Below the header, there are two input fields: 'Endpoint' with the URL 'http://plantgarden.jp/sparql' and 'Start node' with the URL 'http://plantgarden.jp/resource/t1000413.G001/gene#t1000413.G001.011346'. Red arrows point from Japanese text to these fields. Below the input fields, there are several toggle switches for 'Control' (Property, RDF-config, Layer arrangement, Force sim., Scroll zoom) and an 'Edge length' slider. The main content area displays a graph visualization of the data. The graph starts with a node labeled '起点 (イチゴの遺伝子)' (node_0, URI: p0:t3747.G001.000001). This node is connected to several other nodes: 'SO (Sequence Ontology)' (obo:SO_0000704), '遺伝子の位置' (faldo:Region), and '遺伝子領域' (faldo:Region). The '遺伝子の位置' node is further connected to 'faldo:ExactPosition' nodes, which are then connected to 'faldo:ForwardStrandP...' nodes. The '遺伝子領域' node is connected to 'faldo:ExactPosition' nodes, which are then connected to 'faldo:ForwardStrandP...' nodes. The 'faldo:ForwardStrandP...' nodes are connected to 'Integer' nodes (node_10: 266, node_14: 97) and 'URI' nodes (node_11: p2:). The graph is annotated with Japanese text: '起点 (イチゴの遺伝子)', 'SO (Sequence Ontology)', '遺伝子の位置', '遺伝子領域', '読み方向', and '長さ'. A link 'about Endpoint browser |' is visible at the bottom left of the interface.

Endpoint
http://plantgarden.jp/sparql

Start node
http://plantgarden.jp/resource/t1000413.G001/gene#t1000413.G001.011346

<Option>

Control: Property RDF-config Layer arrangement Force sim. Scroll zoom

Edge length:

about Endpoint browser |

起点 (イチゴの遺伝子)

SO (Sequence Ontology)

遺伝子の位置

遺伝子領域

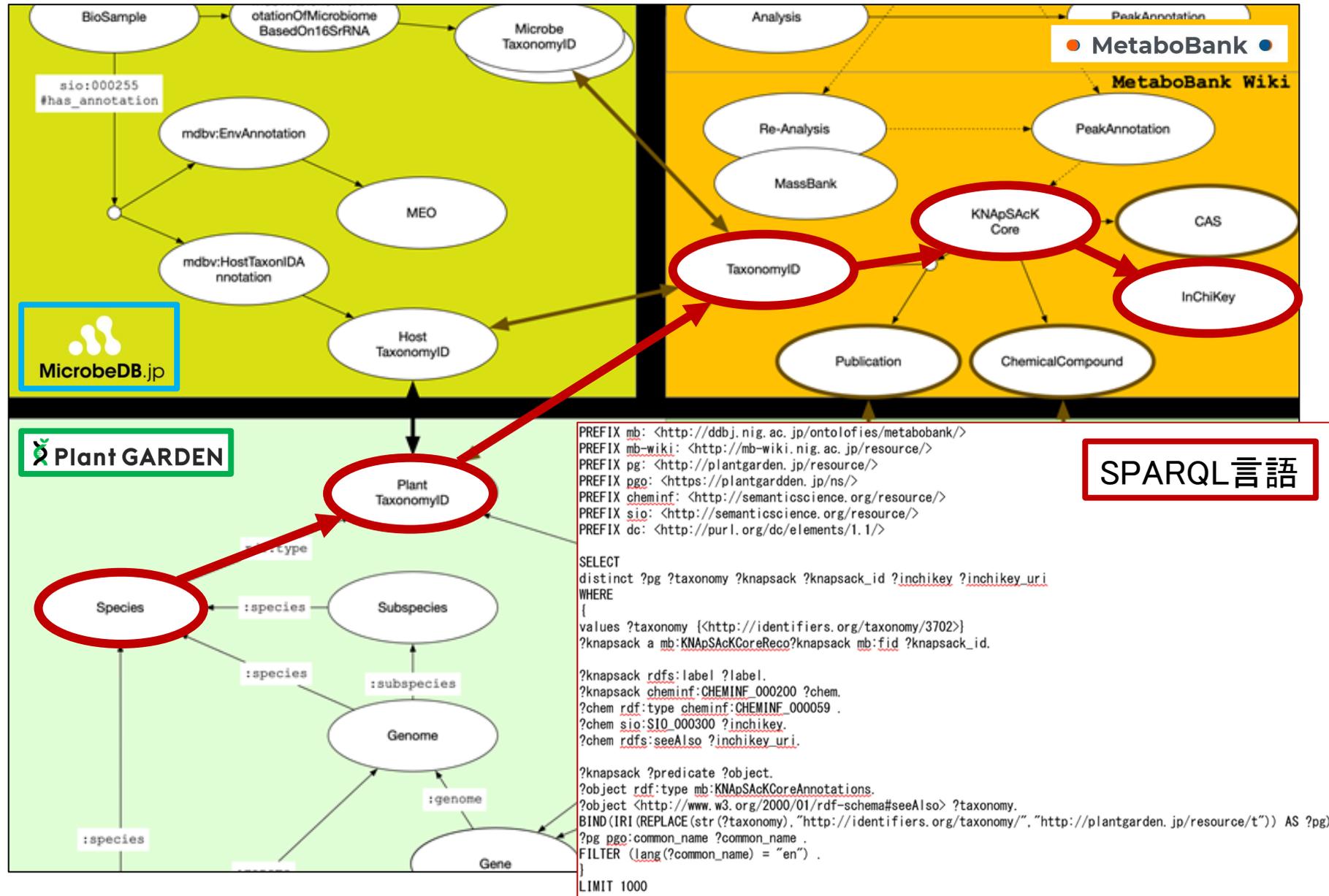
読み方向

長さ

<https://sparql-support.dbcls.jp/endpoint-browser.html>

3者間におけるRDFの繋がり

データが繋がっていることを確認するため、植物種名を起点として、KNApSackの化合物情報を検索することを試みた。



統合化データベース連携検索システムの構築

SPARQL検索には言語に関する知識が必要であり、一般ユーザが容易に検索できるようにする必要性がある。

Plant GARDEN

- ・ゲノム情報(134種)
ゲノム配列、CDS、PEP、アノテーション
- ・DNAマーカー(38種)
- ・QTL(27種)

PubTatorを用いた文献キュレーションのRDF

- ・生物種、病気、遺伝子、化合物など
- ・植物と微生物の相互作用
Plant Growth Promoting Bacteria (PGPB; 共生、成長促進)、病原性細菌(細菌、カビ、ウイルス)
- ・代謝産物(MESH ID)

連携



MetaboBank

- ・実験で得られた代謝産物に関する生データや実験に関するメタデータ
- ・データ数
- ・KNApSAcKのRDF
(DBCLSにて実施)



連携



- ・系統・遺伝子・環境の3つの軸に沿って整理・統合されたフルRDFのデータベース
- ・菌叢メタゲノムデータ(メタ16S rRNA、遺伝子)
- ・ゲノム・ドラフトゲノムデータ
- ・単細胞の真菌類、藻類のゲノムデータ

外部データベースのRDF

Rhea SIB 化合物(11,886個)
酵素反応(13,673個)

Medical Subject Headings RDF NIH



NBDC National Bioscience Database Center

統合化推進プログラムの他のDB

統合化データベース連携検索システムの構築

基礎研究から育種などの応用研究に至る幅広い分野での新たな知識発見を目指す

今後

3者の統合データベースのエンドポイントにトリプルデータを格納し、横断検索できることを確認した。
今後、以下の内容を実施し、データ統合化により有益な情報が得られるようになることを目指す。

RDF化

- ・有益な情報を得ることができる検索の経路、ステップを検討する
- ・スキーマの改変（必要があれば）
- ・センテンスキュレーション、病害抵抗性遺伝子との繋がり

連携(統合化データベース連携検索システムの構築)

- ・一般ユーザにとって使い易いシステムを構築する
- 植物、微生物、メタボローム、糖鎖やフェノタイプなど他の統合化データベース、UniProtなどRDFに対応した外部データベースに対する横断検索を可能にする。

- 検索の例) 作物についてのDNAマーカー・形質に関わる遺伝子 → 産生される代謝産物の検索
→ 菌叢において見られる微生物の検索 → 有益性または病原性の検索
→ 作物の生育に適した土壌(環境)の推定に役立てる

3者間の統合データベースや外部データベースとデータを繋げることによって、様々なデータベースをボーダレスに活用することが可能となり、幅広い研究開発の現場でこれまでにない新たな知識発見が行なえるようになることを目指す。

謝辞: 本研究の一部は、戦略的イノベーション創造プログラム(SIP)「スマートバイオ産業・農業基盤技術」(管理法人:生研支援センター)によって実施されました。