

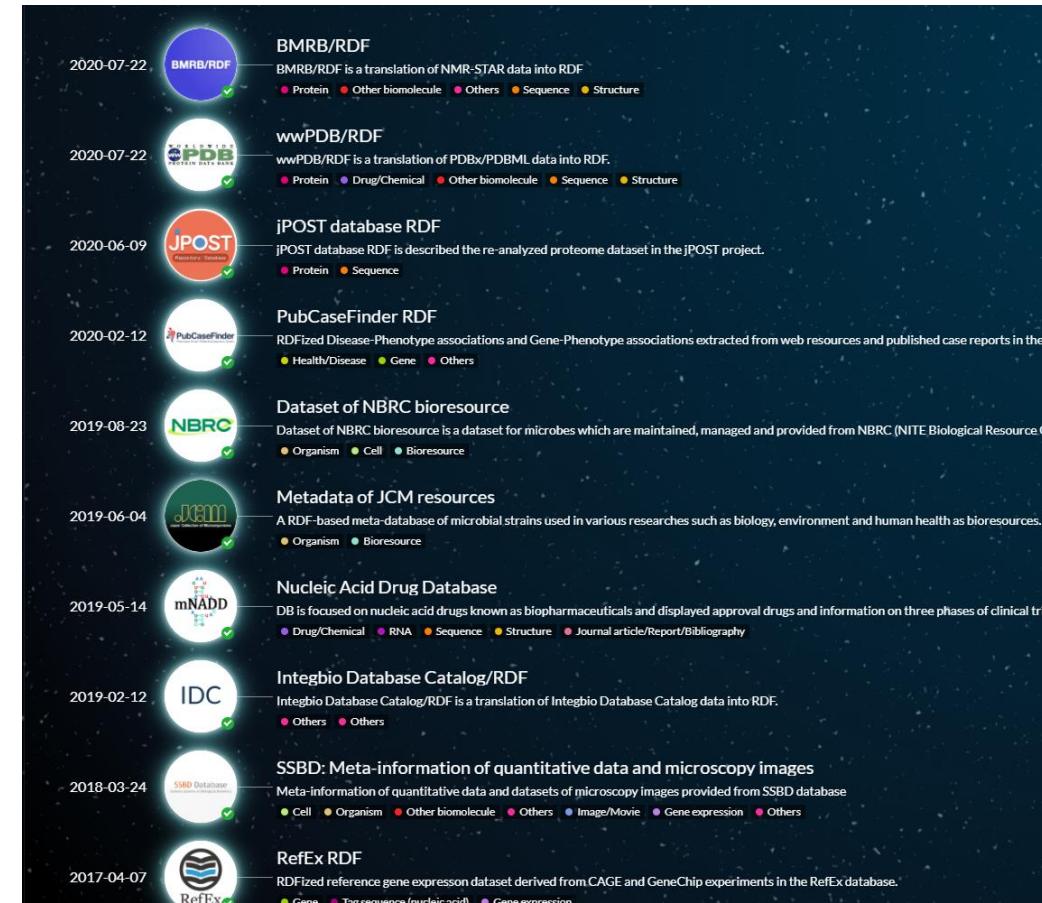
○八塚 茂<sup>1</sup>、片山 俊明<sup>2</sup>、川島 秀一<sup>2</sup>、畠中 秀樹<sup>2</sup>

1. 国立研究開発法人科学技術振興機構バイオサイエンスデータベースセンター(NBDC)、2. 大学共同利用法人情報・システム研究機構ライフサイエンス統合データベースセンター(DBCLS)

## ■ システムおよびRDFデータセットの現状

2020年9月初旬現在で、27のRDFデータセットが公開されており、総データサイズは1000億トリップルである。2015年11月の運用開始時からは約17倍のデータサイズとなっている。

SPARQLエンドポイントとしては、Virtuosoを利用している。すべてのRDFを同じVirtuosoインスタンスにロードしてある方が検索の利便性は高いが、一方でデータサイズが巨大になると運用上に支障をきたすことがあり、現状では、DDBJ、KERO、UniProtなどの大きなデータセットについて個別のインスタンスをたてることで対処している。個々のURIは、  
<http://integbio.jp/rdf/{データセット名}/sparql> のように設定している。



## ■ ガイドラインの設定とレビューの実施

NBDC RDFポータルでは、登録依頼のあったRDFデータセットに対して独自のガイドラインに基づくレビューを実施している。これにより、RDFデータの質とデータ間の相互運用性を高めることができた。

The screenshot shows a GitHub repository interface. At the top, there's a commit from skwsm titled "Correct URI for purl.jp" made on May 10, 2018, with 22 commits. Below the commit details, there are three pull requests listed:

- RDF-portal-guidelines-en.md: Add about HCO in the RDF portal guidelines (2 years ago)
- README.md: Edit the README.md (2 years ago)
- dbcls-rdfizing-db-guidelines-ja.md: Correct URI for purl.jp (4 months ago)

Below the pull requests, there's a file named README.md with the title "DBCLS guidelines for RDFizing databases". The content of this file is as follows:

```
DBCLS guidelines for RDFizing databases

The DBCLS guidelines for RDFizing databases are a collection of useful practices for developers who expose life science databases as RDF. DBCLS have been compiled the guidelines based on experience and knowledge gathered from several hackathon events organized by DBCLS and NBDC such as the DBCLS/NBDC BioHackathons, the SPARQLthons and the RDF summits.
```

## ■ SPARQLエンドポイントのフロントエンドにSPARQL-proxyを利用

SPARQLエンドポイントのフロントエンドとして、DBCLSで開発しているSPARQL-proxyを利用している。これにより、クエリ実行のジョブスケジューリング、キャッシュによる応答性能向上、実行結果の複数フォーマットによるダウンロードを実現することができた。

The screenshot shows the SPARQL Proxy interface. At the top, there's a dark header bar with the title "SPARQL Proxy". Below it, there's a query editor window containing the following SPARQL code:

```
1 # Retrieve a full list of IDs with Japanese and English labels.
2
3 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
4 PREFIX dcterms: <http://purl.org/dc/terms/>
5 PREFIX owl: <http://www.w3.org/2002/07/owl#>
6
7 SELECT STR(?id) AS ?ID STR(?nameJ) AS ?Label_Ja STR(?nameE) AS ?Label_En
8 FROM <http://nanbyodata.jp/ontology/nando>
9 WHERE
10 {
11   ?s a owl:Class ;
12   dcterms:identifier ?id
13   OPTIONAL {
14     ?s rdfs:label ?nameJ
15     FILTER (lang(?nameJ) = "ja")
16   }
}
```

Below the query editor, there's a "Run Query (Ctrl+Enter)" button. To the right, the results are displayed in a table format:

ID	Label_Ja	Label_En
"0000001"	"難病"	"Intractable disease"
"0000002"	"指定難病"	"obsolete class"
"1000001"	"神経・筋疾患"	"designated intractable disease"
"1100001"	"代謝系疾患"	"Neuromuscular disease"
"1100002"	"皮膚・結合組織疾患"	"Metabolic disease"
"1100003"		"Skin and connective tissue disease"

On the far right, there's a blue callout bubble with the text "ダウンロード フォーマット を選択" (Select download format). A dropdown menu next to it shows "Download" selected, with other options like "JSON", "CSV", and "TSV" available.

## ■ 登録対象データセットの拡張

2019年度からは、登録対象について大きく方針転換し、これまで登録依頼のあったRDFに限っていたものを、UniProtやPubChemといった一般に公開されているRDF（この場合、SPARQLエンドポイントのみ提供）や、NBDC等が独自に開発したRDFにまで対象を広げることにした。

### ➤ 既公開

- UniProt
- PubChem

### ➤ 今後公開予定

- BioSamples
- Ensembl
- Rhea

など

データの登録依頼や  
その他のお問い合わせは、  
rdf-portal@integbio.jp  
までお気軽に

The screenshot shows the UniProt RDF dataset page. At the top, it says "UniProt" and "Original site". Below that is a brief description: "The Universal Protein Resource (UniProt) is a comprehensive resource for protein sequence and annotation data. For the detailed information, please visit the original site." A "Specification" section follows, listing the following details:

Tags	<span style="color: pink;">●</span> Protein <span style="color: orange;">●</span> Sequence
Data provider	European Bioinformatics Institute, Swiss Institute of Bioinformatics, and Protein Information Resource
Creators	UniProt Consortium
Version	2019-12-19
Issued	2019-12-19
License	Creative Commons Attribution 4.0 International
Status	<span style="color: green;">✓</span> Unreviewed <span style="color: grey;">●</span> 3rd party dataset
Download file	0 bytes

Below this is a "Linked datasets" section with two entries:

- Reactome: 282,912 links
- NCBI Gene: 15,648 links

## ■ インフラの増強（予定）

データ登録対象の拡張に伴い、2021年度からはインフラを増強して、安定したSPARQL検索の提供と、よりタイムリーなデータ追加・更新を目指す。

