### DBCLSにおけるデータベース RDF化への取組みと今後の展望

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 山口 敦子

### Web of Data (1)

インターネットの出現



ケーブルを意識せずに

コンピュータをつなぐ世界の出現

各コンピュータに蓄積される文書を得るためには、 文書の在処を何らかの手段で知り、

アクセスするしかなかった

ftp telnet gopher archie

...

### Web of Data (2)

### WWWの出現



文書がどこにあるかを意識せずに ししく 利用できる世界(Web of Documents)の出現

各文書に含まれるデータを得るためには, 文書ごとに解析してデータを 抽出していくしかなかった. HTMLパーザ 自然言語処理

### Web of Data (3)

ライフサイエンスDBにおいて,

現状: DB数增大&大規模化

課題: データを合わせることで新しい知識が生まれる

そのためにはデータ共有と統合的利用の仕組みが必須

データの置き場や形式を意識せずに

データをつないで利用できる世界(Web of Data)



RDF 技術の利用

### **RDF**

データ記述・交換のための標準的枠組み(W3C仕様)

主語 ● 目的語 の三つ組で もの(ex. 遺伝子, タンパク質, 生物種, 化合物…)の関係を表現 →シンプルな形式

主語, 述語は IRI(ex. URL) 目的語は IRI かリテラル(具体的な値)

→ものに対してグローバルにユニークなID付け

三つ組の集合=データ

→ひとつひとつのものにURLなどのグローバルなIDをつけ, グラフとしてつなぐことで, ひとつのグローバルなデータ空間(Linked Data)を実現

### 世界的にRDF化が加速

UniProt (2008-)

Bio2RDF (2008-)



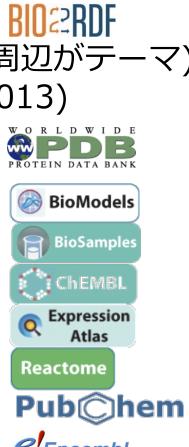
NBDC 基盤技術開発プログラム(2011-2013)

wwPDB (2011-)

EBI RDF platform (2013-)

- BioModels
- BioSamples
- ChEMBL
- Expression Atlas
- Reactome

PubChem (2014-) Ensembl (TBA)

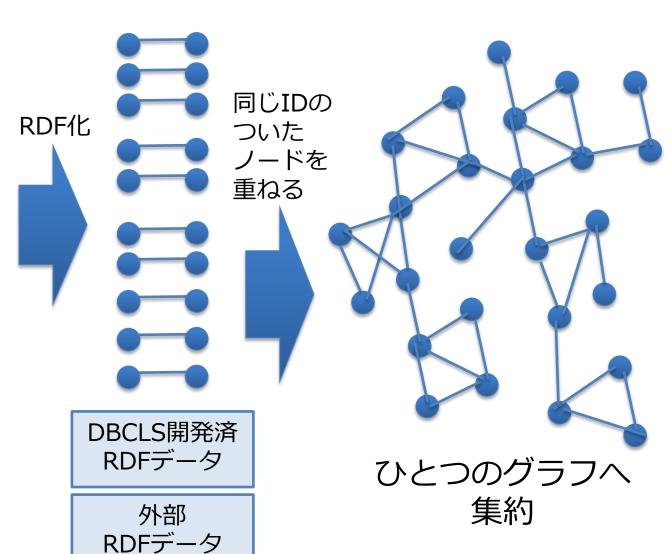


RDF化されていな い有用DB

文献

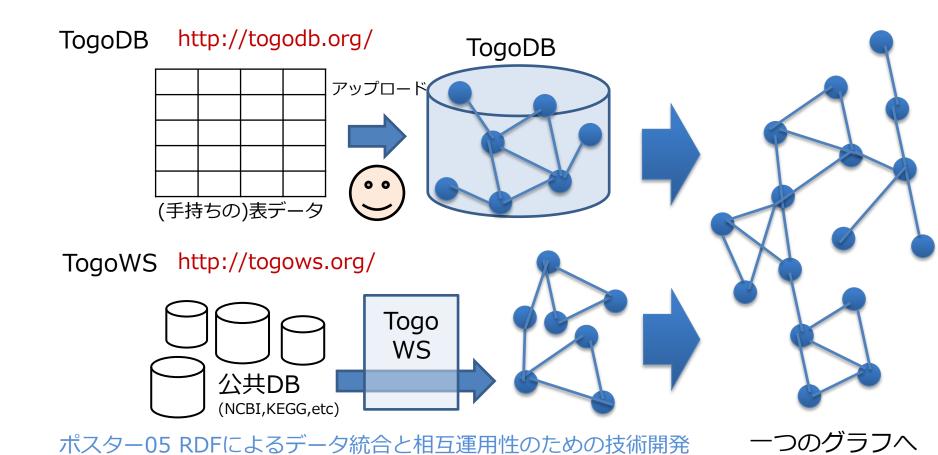


統合化推進DB



RDF化されていな い有用DB

RDF化ツールの開発, データのキュレーション

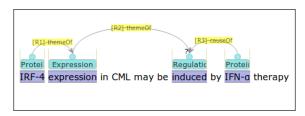


#### 文献

文献内の知識をRDF化するツールの開発

PubAnnotation http://www.pubannotation.org/

PubMedの文献に埋もれた知識のRDF化





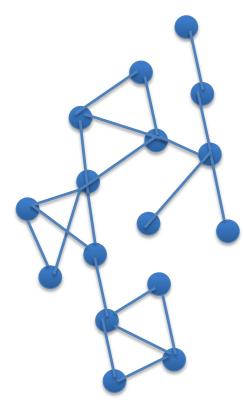
新着論文レビューのナビゲーションシステム開発

http://navi.first.lifesciencedb.jp/stanza/top

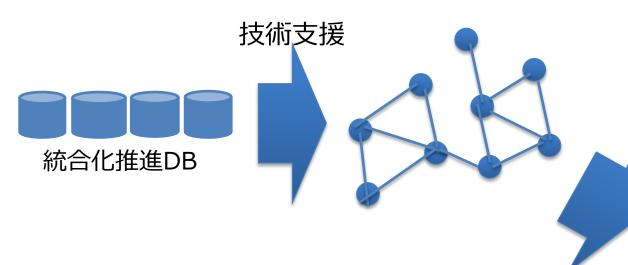
DBCLSで作成している日本語文献のRDF化



ポスター07 日本語コンテンツに対する セマンティックウェブ技術の適用



一つのグラフへ

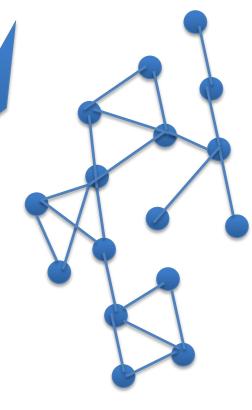


### SPARQLthon

一か月に一回二日間DBCLS 主催で行う 技術交流/開発会合

http://wiki.lifesciencedb.jp/mw/SPARQLthon

次回は10/27,28



### めざす世界

#### データの置き場や形式を意識せずにデータをつないで利用できる世界

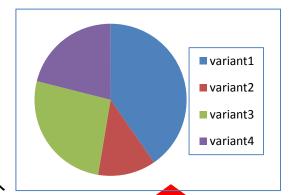
日本人のBRAFのSNPs の統計情報とそれに関連する病気を知りたい.





RDF検索言語(SPARQL)へ

変換



| variant  | disease                                |
|----------|--|
| variant1 | colorectal cancer                      |
| variant2 | cardiofacioc<br>utaneous<br>syndrome 1 |
| variant3 | lung cancer                            |
| variant4 | Noonan<br>syndrome 7                   |

### 可視化

| person  | variant  | disease                            |
|---------|----------|------------------------------------|
| person1 | variant1 | colorectal cancer                  |
| person2 | variant3 | lung cancer                        |
| person3 | variant2 | cardiofaciocutaneous<br>syndrome 1 |
| person4 | variant1 | colorectal cancer                  |
| person5 | variant3 | lung cancer                        |

検索



結果

ひとつのグラフ (LinkedData)

### H28年度達成目標

#### DBのRDF化による分野を超えた統合化の実現

RDFデータを揃えるのみならず, 実験系生物研究者のニーズに沿ったサービスを開発

実験系生物研究者 向けDBサービス

質問応答システム

引き続き開発する RDFサービス (TogoGenome等) 日本語 コンテンツ

NGS データ 利用環境

情報系生物研究者向けDBサービス

RDFデータ

SPARQL エンドポイント (Web API)

NGSメタデータ

# 利用可能な技術/サービス (検索言語変換): LODQA

http://lodqa.org/

質問応答システム

# **LOD**QA

Querying linked open data (LOD) using natural language. Sounds great? Let's realize it!

Front Motivation Participants References github

LODQA (Linked Open Data Question Answering) is an open source project aiming at developing a system to generate SPARQL queries from natural language queries.

#### News

- (02/17-21/2014) OKBQA 2014 Hackathon is held in Jeju.
- (01/10/2014) Relation detection is improved to find the right pairs of terms in relation.
- (12/20/2013) LODQA is launched as an open source project.

#### Prototype Demo (targeting OMIM)

what anatomical abnormality is associated with kabuki syndrome?

submit

- what genes are associated with kabuki syndrome?
- what sign is associated with kabuki syndrome?
- what cellular dysfunction is associated with kabuki syndrome?

#### Example

- what pathologic function is associated with kabuki syndrome?
- · what neoplastic process is associated with kabuki syndrome?
- what anatomical abnormality is associated with kabuki syndrome?
- what genes are associated with alzheimer disease?

### 利用可能な技術/サービス (検索言語変換): LODQA

入力:"what anatomical abnormality is associated with kabuki syndrome?"

http://lodqa.org/



### 利用可能な技術/サービス(検索): TogoGenome

TOGO GENOME http://togogenome.org/

RDFで集積したゲノム情報を検索するシステム

遺伝子機能,生物種,表現型,環境等の条件の組み合わせによる ゲノム情報の絞り込み検索

**SPARQL** 

検索

絞り込み条件軸 (オントロジー)

検索結果を可視化した 一覧性の高い レポート表示

環境

GO細胞局在

GO分子機能

GO生命プロセス

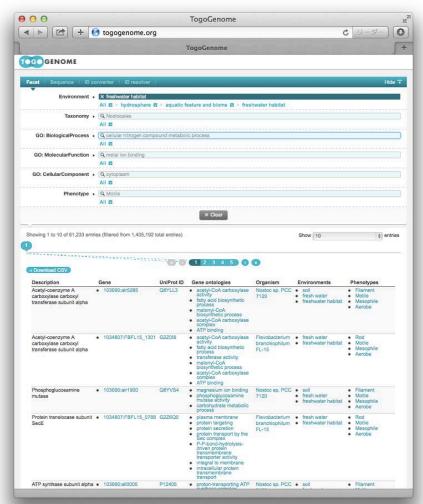
微生物表現型

生物分類

遺伝子情報

生物種情報

環境情報

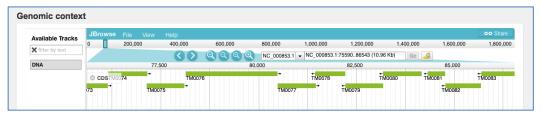


## 利用可能な技術/サービス(可視化): TogoStanza



SPARQL検索した結果を可視化するためのシステム

TogoGenome, MicrobeDB.jp, CyanoBase等での表示に利用されている



Genomic context スタンザ

Protein orthologs

D9TMZ1, D9TEN6, E1QW86, K0J5R1, B1ZRT3, D9T8Z3, Q987M4, C6D060, A6WFR1, Q9FC27, Q9WXT2, I0I5K4, Q7D1J0, E4N533, E4N6L6, E8NDX 3, G8S8S3, F8ECV6, C7QZA6, I3YCG4, D9T1B2, E8NDX2, A4XM78, F4L1I5, F6FSR0, Q251A2, F0SWU9, D6Y1S1, D2B4Z9, D2QWF7, D1BAG5, A9A YF2, P94388, D1BIQ2, A6WDI8, D1BDY6, I2ETZ3, D1CDW7, C6W605, C0QVP2, D1BXU6, Q8AAS2, C7RYP7, A5URF7, D9TNY9, D5UKL6, Q9WYR6, C6W7C0, C7PZN0, G0J7C6, B5YCD5, A5CMG3, I4BAP7, C8W7B9, Q8F524, E1QW87, A0LAM5, F4FAX0, E8NCP5, Q47PW9, K0J4P6, C7RZ94, D2P T82, C5C0D3, C6CYB6, G0LBL3, C6WE33, Q022X6, D6ZHW0, C5C5T9, F4C6N3, Q835Y2, C7MFS9, D3QAK3, G0A485, A0K059, D2QEJ8, A0LWK7, E8MZH0, C7QDS2, C6WFJ7, C7QDC4, B1ZQH5

Protein orthologs スタンザ

Nucleotide sequence スタンザ

現在、TogoGenome上では約40種のスタンザを提供

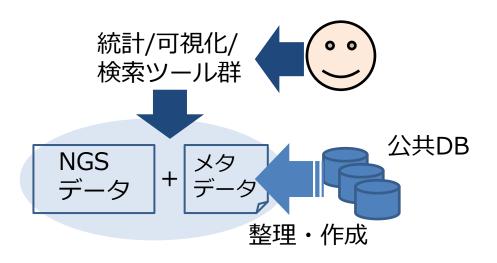
ポスター06 TogoGenome / TogoStanza の新規機能開発

# 実験系生物研究者ニーズに沿った サービス開発

NGS データ 利用環境

NGS実験データを簡単に再利用/検索できる環境の構築

- ・メタデータの整理
- ・DDBJ との連携
- ・検索・解析ツールの開発



ポスター27 NGSデータの利用を促進する統合環境の構築とサービスの提供 28 公共NGSデータの活用を促進する検索システムの構築

### 利用可能な技術/サービス

NGS データ 利用環境

遺伝子発現データのリファレンス RefEx http://refex.dbcls.jp/



異なる実験手法の遺伝子発現データを統合し, 参照データとして整理

検索技術 GGRNA, GGGenome



超絶高速ゲノム配列検索

http://GGRNA.dbcls.jp/ http://GGGenome.dbcls.jp/

ポスター29 遺伝子発現リファレンスデータセット RefEx 30 統合遺伝子検索GGRNAと高速塩基配列検索GGGenome: 塩基配列データベースをすばやく検索するウェブサーバ

# 実験系生物研究者ニーズに沿った サービス開発

日本語コンテンツ

#### 日本語コンテンツの拡充

・ライフサイエンス分野のDBやツールの動画教材TogoTV の整備 http://togotv.dbcls.jp/

・ライフサイエンスを解説する日本語テキストの整備



ポスター03 新しい日本語Webコンテンツ、
「新着論文レビュー」と「領域融合レビュー」
07 日本語コンテンツに対するセマンティックウェブ技術の適用

### まとめ

有用な統合データベースの実現には統合化推進プログラムのデータを含め、網羅的に幅広い種類のデータをRDF化することが望ましい(Web of Data).

#### DBCLSはその実現へ向けて

- ・DBのRDF化の取り組み
- ・RDFでアプリケーション開発
- ・実験系生物研究者ニーズに沿ったサービス開発これまでの技術の蓄積を生かして行う