

生命科学系研究データ循環基盤としてのアーカイブと横断検索

○大波純一¹、八塚茂¹、杉崎太一朗²、友田史緒里²、牧口大旭²、加藤健弘³、井上圭介⁴、大久保克彦³、川本祥子^{5,6}、畠中秀樹¹

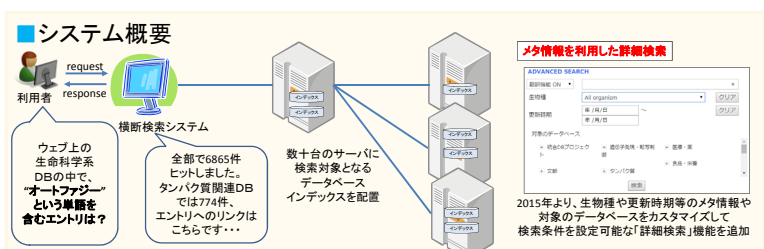
1. 国立研究開発法人科学技術振興機構バイオサイエンスデータベースセンター(NBDC)、2. 三井情報株式会社、3. 株式会社日立製作所、4. 株式会社日立公共システム、5. 大学共同利用法人情報・システム研究機構ライフサイエンス統合データベースセンター(DBCLS)、6. 情報・システム研究機構 国立遺伝学研究所

要旨

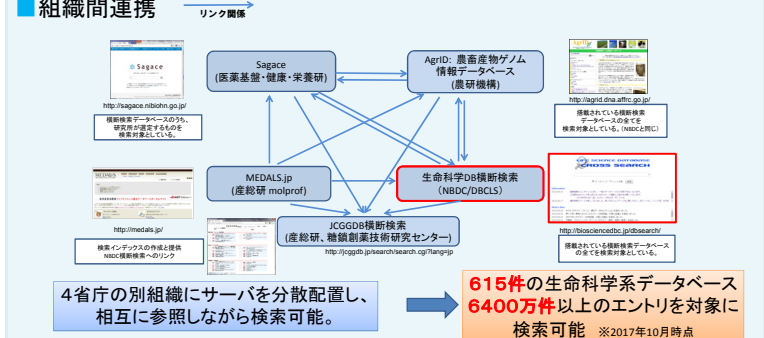
近年のオープンサイエンス推進の潮流や情報通信技術の発展に伴い、生命科学系研究データの規模は増加の一途を辿っている。バイオサイエンスデータベースセンター(NBDC)は、この流れの中で、生命科学データベース横断検索と生命科学系データベースアーカイブ、及びIntegbio データベースカタログを公開し、データの利用者だけでなく開発者の課題に応える活動を行っている。横断検索では今年、データベース横断検索機構の開発と、検索クエリログの分析に基づく検索結果スコアリングアルゴリズムの改定を実施し、検索対象の拡大と検索結果の適正化を達成した。アーカイブでは寄託件数が130件を突破し、Scientific Dataの推奨レポジトリに認められ、コミュニティにおける基盤の1つとして認知されつつある。

国際社会では研究データの適切な循環のため、FAIR 原則 (Findable, Accessible, Interoperable, Reusable) に従うことが標準となりつつある。横断検索とアーカイブはこの条件をサポートし、日本の生命科学界の研究データ循環のメインパイプとなるよう、引き続き改善と品質向上に努めていく。

生命科学データベース横断検索とは



キーワードに対し、データベースの一括検索を行い、ヒット件数と、検索結果(対象へのリンク、ヒット箇所)を返す**Web検索システム**



検索利用者の意図抽出アルゴリズムの構築

横断検索のページランクスコア算出方法の改善

横断検索の検索結果の表示順を改善し、利用者が素早く目的の情報にアクセス可能とするため、検索エンジンアプリのElasticsearchにおけるアルゴリズム "TF/IDF practical scoring function of Lucene" のパラメータ設定に検討した。

$$score(q, d) = queryNorm(q) \cdot coord(q, d) \cdot \sum_{t \in d} (tf(t \text{ in } d) \cdot idf(t))^2 \cdot t.getBoost() \cdot norm(t, d)$$

The relevance score of document d for query q

The query normalization factor $1 / \sqrt{sumOfSquaredWeights}$

The coordination factor (add the rate of matching query words, relates to partial match)

The term frequency for term t in d

The inverse frequency for t

The boost that has been the field-length norm, combined with the index-time field-length boost, for each t in q or d

The sum of the weights applied to the query

https://www.elastic.co/guide/en/elasticsearch/guide/2.x/practical-scoring-function.html

動的に設定できるboostパラメータを検索クエリによって変化するように修正を実施した。

複数語クエリのパターンニングによる検索意図の解釈

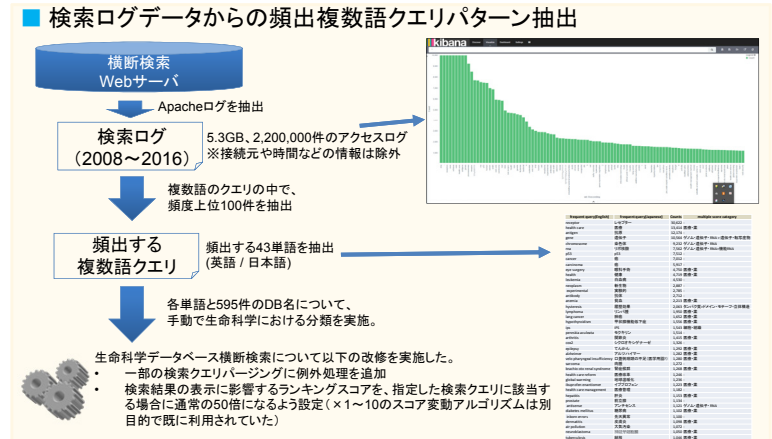
Chemical-chemical context pattern
#C and #C
#C versus #C
#C and #C interaction
#C and #C combination
#C plus #C
#C with #C
comparison of #C and #C
interaction between #C and #C
#C oxidase #C
#C dehydrogenase #C
combine #C and #C
#C transporter #C

Chemical-disease context pattern
#D and #C
#C induce #D
#D treatment #C
treatment of #D with #C
#C #D review
#D with #C
#D child #C
#D induce by #C
#D due to #C
#D treatment with #C
if any
#C metabolism and #D

NCBIでは、2015年にPubMedの検索クエリを解析し、Chemical (#C)とDisease (#D)を含む複数語検索のペアを抽出した。このペアを分類し、左に示すような頻出パターンに分類した。この情報を元にPubMedの検索結果表示方式への反映を行った。同様に2単語以上の複数語検索では、GoogleのVertical search (画像検索・ニュース検索等)に対応するような、特徴的な「サブクエリ」(「画像」や「ニュース」が含まれている)は、意図の推定が行いやすいことが知られている。[This concept is inspired by Yamamoto et al., "Overview of the NTCIR-12 Mine-2 Task", NTCIR 2016]

検索ログを元に複数語クエリのカテゴリと解釈を実施し、横断検索のboostパラメータの変動アルゴリズムを構築した。

複数検索クエリを利用した意図抽出の例



修正前後の検索結果順の比較

search query	Before the tuning			After the tuning			Difference
	1st result	2nd result	3rd result	1st result	2nd result	3rd result	
toxic データベース	医学-漢字検索全文データベース 毒物学データベース 毒物学データベース (2) : データベースの検索	医学-漢字検索全文データベース 毒物学データベース 毒物学データベース (2) : データベースの検索	医学-漢字検索全文データベース 毒物学データベース 毒物学データベース (2) : データベースの検索	Integbio database catalog "Open TG-GATEs"	Integbio database catalog "Open TG-GATEs"	Integbio database catalog "Open TG-GATEs"	"データベース" multiplied the score of database catalog score.
結核 症状	I-STAGE "病原体菌種群の1種 A case of atypical tuberculosis"	I-STAGE "病原体菌種群の3種 Three cases of atypical tuberculosis"	"医学-漢字検索全文データベース" "病原体菌種群の3種 Three cases of atypical tuberculosis"	"UMIN-CTD臨床試験データベース" UMIN00005454 結核関連薬に於ける代謝経路分類 (結核) Mycobacterium avium complex(結核菌)の薬物相互作用に関する研究論文 検索結果を明らかにす 研究論文	"UMIN-CTD臨床試験データベース" UMIN00005454 結核関連薬に於ける代謝経路分類 (結核) Mycobacterium avium complex(結核菌)の薬物相互作用に関する研究論文 検索結果を明らかにす 研究論文	"UMIN-CTD臨床試験データベース" UMIN00005454 結核関連薬に於ける代謝経路分類 (結核) Mycobacterium avium complex(結核菌)の薬物相互作用に関する研究論文 検索結果を明らかにす 研究論文	"gene" multiplied the score of "gene/branch" category. But score of documents are higher than the multiplied scores.
gene enolase	"医学-漢字検索全文データベース" "酵素名変異体と:myc" タンパク質と変異するエンザイム	"医学-漢字検索全文データベース" "酵素名変異体と:myc" タンパク質と変異するエンザイム	"医学-漢字検索全文データベース" "酵素名変異体と:myc" タンパク質と変異するエンザイム	"医学-漢字検索全文データベース" "酵素名変異体と:myc" タンパク質と変異するエンザイム	"医学-漢字検索全文データベース" "酵素名変異体と:myc" タンパク質と変異するエンザイム	"医学-漢字検索全文データベース" "酵素名変異体と:myc" タンパク質と変異するエンザイム	"gene" multiplied the score of "gene/branch" category. But score of documents are higher than the multiplied scores.

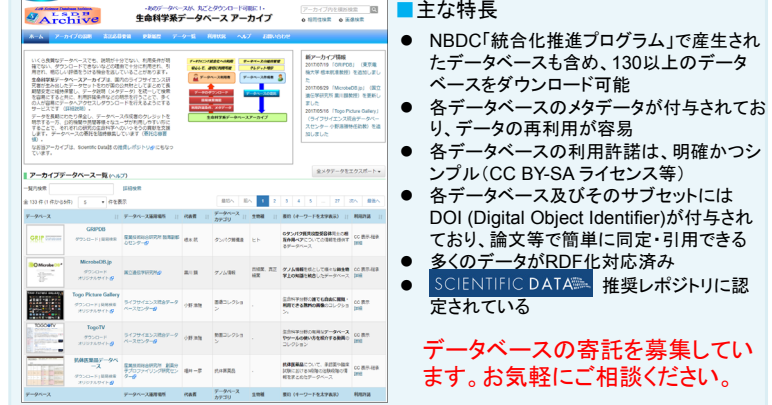
検索語によって、明らかに改善されたと判断できるケースと、変化の無かったケースが観察された

2017年5月、本番環境に反映

生命科学系データベースアーカイブとは

サービス概要

生命科学系データベースアーカイブは、国内のライフサイエンス研究者が生み出したデータセットをわが国の公共財としてまとめて長期安定に維持保管し、データ説明(メタデータ)を統一して検索を容易にすると共に、利用許諾条件などの明示を行うことで、多くの人が容易にデータへアクセスしダウンロードを行えるようにしている。データを長期にわたり保全し、データベース作成者のクレジットを明示する一方、公的機関や民間等様々なユーザが利用しやすい形にすることで、それぞれの研究の生命科学へのいっそうの貢献を支援する。当アーカイブは2017年度、Nature publishingのScientific Data誌の推奨レポジトリに認定された。



主な特長

- NBDC「統合化推進プログラム」で産されたデータベースも含め、130以上のデータベースをダウンロード可能
- 各データベースのメタデータが付与されており、データの再利用が容易
- 各データベースの利用許諾は、明確かつシンプル(CC BY-SA ライセンス等)
- 各データベース及びそのサブセットにはDOI (Digital Object Identifier)が付与されており、論文等で簡単に同定・引用できる
- 多くのデータがRDF化対応済み
- SCIENTIFIC DATA 推奨レポジトリに認定されている

データベースの寄託を募集しています。お気軽にご相談ください。

まとめ・今後の方針

- 生命科学データベース横断検索では、検索ログを利用した意図抽出アルゴリズムを設定し、より利用者がデータを発見し易い基盤となるよう、改善を継続している。
- 生命科学系データベースアーカイブでは、寄託データを長期にわたり保持する活動をおこない、Scientific Data誌の推奨レポジトリに認定され、利用者の継続的なアクセスに貢献している。
- 国際社会では研究データの適切な循環のため、FAIR 原則 (Findable, Accessible, Interoperable, Reusable) に従うことが標準となりつつある。横断検索とアーカイブがこの条件をサポートし、日本の生命科学界の研究データ循環のメインパイプとなるよう、引き続き改善と品質向上に努めていく。