

82. 生命科学研究データ再利用性担保プラットフォーム構築

水口 賢司

大阪大学 蛋白質研究所 計算生物学研究室

Key words : FAIR 原則, データ共有, データ再利用性担保, データウェアハウス, 機械学習

緒言

生命科学分野では、研究成果であるデータを適切にオープンにすることが求められている。このような動きは、オープンサイエンス推進における国内外の議論でも取り上げられ、FAIR 原則である「Findable (見つけられる)、Accessible (アクセスできる)、Interoperable (相互運用できる)、Reusable (再利用できる)」がしばしば言及される [1]。

単一の実験結果であれば、データを公共リポジトリに登録することが有効である。タンパク質立体構造 (worldwide Protein Data Bank) や塩基配列 (DDBJ、NCBI、ENA/EBI) など、古くから確立された国際協調によるデータベース構築が、最近ではプロテオームデータなど他の分野にも広がりつつある。しかし、近年の生命科学研究においては、遺伝子機能とタンパク質間相互作用など、異なった複数種類のデータを組み合わせる事例が増えている。FAIR 原則に則って、研究データの再利用性を担保し共有・活用を促進するためには、異種の生物学データを統合して、共有・活用できる形に整理する必要がある。その作業には、一般に多大な労力が求められるという問題があった。

一方、こうして整理されたデータから、新たな解析目的に応じて柔軟に部分データを再構築できれば、研究データの再利用性を超えた利活用が期待できる。様々な研究分野で人工知能 (artificial intelligence : AI) の活用が進んでいるが、生命科学分野における AI・機械学習の応用においては、学習 (正解) データ作成のコストの大きさが深刻な課題と認識されている。先に述べた複数種類の研究成果データが統合・公開されたシステムから、コンピュータ解析に適した形式に必要な部分データを自由に生成できれば、機械学習モデル構築のボトルネックが解消され、生命科学分野のデータ駆動型研究に変革がもたらされると考えられる。

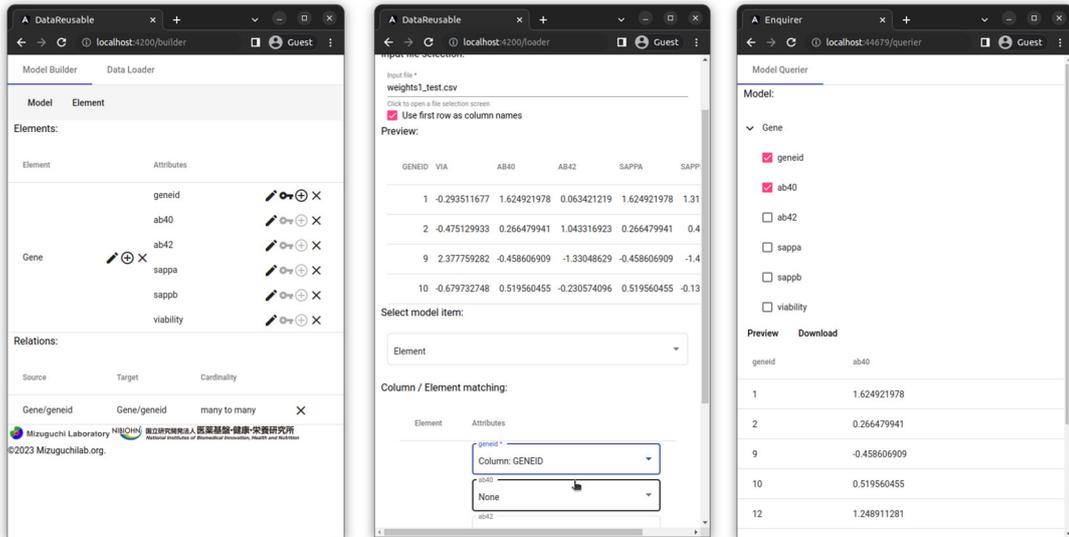
私たちは 10 年以上にわたり、創薬支援データウェアハウス TargetMine の研究開発を続けてきた [2, 3]。TargetMine は、30 以上の公共データベースから得られる遺伝子、タンパク質、化合物、疾患などのデータを統合し、各項目間の関連性を効率的に探索する機能を提供する。現在では、創薬支援だけでなく、基礎生物学分野のオミックスデータ解析などにも広く利用されている。本研究では、TargetMine の開発経験を生かして、生命科学分野の研究データ再利用性担保及び学習データ自動生成プラットフォームを構築することを目的とした。異種データの組み合わせを取得する検索機能を実装し、カスタマイズした多種類のデータの統合とオンライン共有を実現するシステムを開発した。

具体例として、アミロイド前駆体タンパク質 (amyloid-beta precursor protein : APP) の産生メカニズム解明を目的とした網羅的遺伝子ノックダウンとタンパク質間相互作用データを組み合わせたデータセット [4] を用いて、上記の実現可能性を証明した。また、一般のユーザーが自分自身のデータを使用して同様のデータウェアハウスを構築するための詳細なプロトコルの公開を予定している。

方法および結果

1. プラットフォームの定義

従来、データの共有は、一連のファイルまたはデータベースの単純な交換であると考えられてきた。しかし、生命科学研究で FAIR 原則に適切に従うためには、基礎となるデータにアクセスして再利用可能にするツールや方法の共有も考慮する必要がある。



a) データ提供者インターフェイス : モデル保存
 b) データ提供者インターフェイス : データアップロード
 c) データ照会者インターフェイス

図 2. データ提供者と照会者用アプリケーションのスクリーンショット

- a) 提供者アプリケーションに含まれるモデル構築用ビュー。
- b) 提供者アプリケーションに含まれるデータアップロード用ビュー。
- c) 照会者アプリケーションに含まれるデータ検索のためのビュー。

3. 実装

以下では、本プラットフォームの技術的な実装の詳細について説明する。

データ提供者と照会者のインターフェイスは、Angular フレームワークを用いた Web アプリケーションとして実装した。プラットフォームが使用する統合データベースの管理には PostgreSQL を使用した。最後に、接続インターフェース (API) は、Node.js フレームワークを使用して実装した。

図 3 は、すべてのモジュールが連携している様子を示している。



図 3. データ利用性担保プラットフォームアーキテクチャ

データ提供者であるユーザーとデータ照会者であるユーザーに特化した異なるアプリケーションやデータベースサービスを組み合わせることで、統合的な共有機能を提供する。

考 察

これまで、データを共有する場合、データ提供者は、データをプレーンテキストや表計算ソフト (Microsoft Excel など) で提供することが多かった。しかし、共有されるデータが複数の種類を含み、かつ、それらの種類間に何らかの関係がある場合、ユーザーは内容を再構成する必要があり、データ再利用のハードルが高くなる。データをリレーショナルデータベースに変換することで、データ共有を統合的な形にすることができるが、データベースを操作するためには専門的な知識を必要とする。Microsoft Access や Claris FileMaker などの市販ソフトウェアを使えば、少ない労力でデータベースを作成することができるが、これらのソフトウェアはアプリケーションを開発するためのもので、データを共有するためのものではない。この問題を解決するために、本研究では、生物学データをオンラインで共有する汎用的な方法を提案し、実証した。

今回の実装では、技術的な理由から、より網羅的なデータ統合や高度なデータクエリへの対応が可能な InterMine [5] フレームワークを組み込むことはできなかったが、今回提案したデータ共有のコンセプトに InterMine を導入することは可能と考えている。データ提供者がフロントエンドアプリケーションでデータモデルを定義すれば、その情報を用いて InterMine フレームワークがバックエンドと照会者アプリケーションを生成するという枠組みが想定される。さらに、各種データの分布を可視化する機能についても実装を計画している。

将来的には、本プラットフォームを通してデータ提供者が Docker [6] イメージを作成し、公開できる機能の付加を予定している。これにより、共有データに興味を持つユーザーは、Docker をインストールするだけで、簡単にデータの全部または一部を利用することができる。さらに、プラットフォームを構成する各モジュールのソースコードは、GitHub 上で MIT License のもとの公開を予定している。

共同研究者・謝辞

本研究の共同研究者は、医薬基盤・健康・栄養研究所、AI 健康・医薬研究センターの陳怡安及びロドルフォ・アジェンデスである。

文 献

- 1) Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, Blomberg N, Boiten JW, da Silva Santos LB, Bourne PE, Bouwman J, Brookes AJ, Clark T, Crosas M, Dillo I, Dumon O, Edmunds S, Evelo CT, Finkers R, Gonzalez-Beltran A, Gray AJ, Groth P, Goble C, Grethe JS, Heringa J, 't Hoen PA, Hooft R, Kuhn T, Kok R, Kok J, Lusher SJ, Martone ME, Mons A, Packer AL, Persson B, Rocca-Serra P, Roos M, van Schaik R, Sansone SA, Schultes E, Sengstag T, Slater T, Strawn G, Swertz MA, Thompson M, van der Lei J, van Mulligen E, Velterop J, Waagmeester A, Wittenburg P, Wolstencroft K, Zhao J, Mons B. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data*. 2016 Mar 15;3:160018. doi: 10.1038/sdata.2016.18. Erratum in: *Sci Data*. 2019 Mar 19;6(1):6. PMID: 26978244; PMCID: PMC4792175.
- 2) Chen YA, Tripathi LP, Mizuguchi K. TargetMine, an integrated data warehouse for candidate gene prioritisation and target discovery. *PLoS One*. 2011 Mar 8;6(3):e17844. doi: 10.1371/journal.pone.0017844. PMID: 21408081; PMCID: PMC3050930.
- 3) Chen YA, Allendes Osorio RS, Mizuguchi K. TargetMine 2022: a new vision into drug target analysis. *Bioinformatics*. 2022 Sep 15;38(18):4454-4456. doi: 10.1093/bioinformatics/btac507. PMID: 35894632; PMCID: PMC9477527.

- 4) Camargo LM, Zhang XD, Loerch P, Caceres RM, Marine SD, Uva P, Ferrer M, de Rinaldis E, Stone DJ, Majercak J, Ray WJ, Chen YA, Shearman MS, Mizuguchi K. Pathway-based analysis of genome-wide siRNA screens reveals the regulatory landscape of APP processing. *PLoS One*. 2015 Feb 27;10(2):e0115369. doi: 10.1371/journal.pone.0115369. Erratum in: *PLoS One*. 2015;10(6):e0129641. PMID: 25723573; PMCID: PMC4344212.
- 5) Smith RN, Aleksic J, Butano D, Carr A, Contrino S, Hu F, Lyne M, Lyne R, Kalderimis A, Rutherford K, Stepan R, Sullivan J, Wakeling M, Watkins X, Micklem G. InterMine: a flexible data warehouse system for the integration and analysis of heterogeneous biological data. *Bioinformatics*. 2012 Dec 1;28(23):3163-5. doi: 10.1093/bioinformatics/bts577. Epub 2012 Sep 27. PMID: 23023984; PMCID: PMC3516146.
- 6) <https://www.docker.com/>