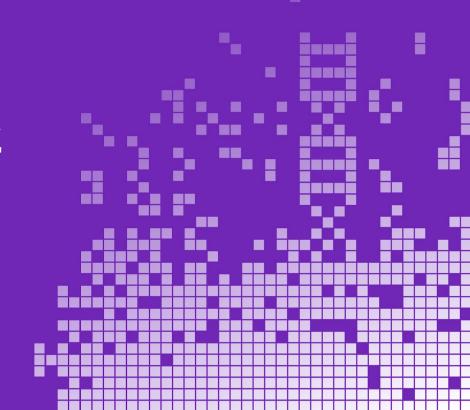


トーゴーの日シンポジウム2025 マルチモーダルデータ X AI

DBCLSでのデータ統合と AIに向けた取り組み

片山俊明 (ROIS・BSI/DBCLS)

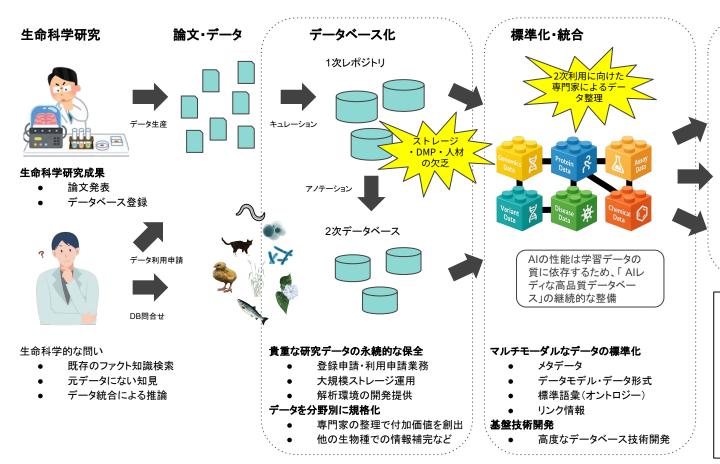
2025-10-20



生命科学×AIの根幹となるデータベース運用

BSI

生命科学の研究成果を集積し、適切にキュレーションすることで、研究者にもAIにも利用可能なデータを創出



研究利用·AI利用

データ利用の基盤を運用・開発中

科学データの取得・再利用

- 生データへのアクセス
- 統合データの機械学習
- 生物・AI由来データ選別

質問応答システム

- 自然言語問合せ
- QAタスクのデータ生成

AIエージェント

- MCPサーバ整備
- SPARQL自動生成

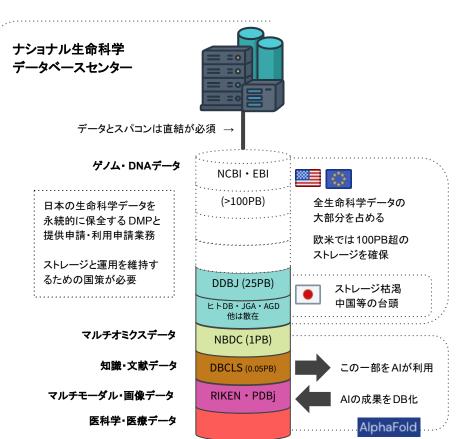
ナショナル DBに求められること

- ・国家予算で科学研究費を出す以上、 得られた研究データを保全する恒久的 なストレージと運用も国が責任を持つ べき。
- ・膨大なデータは移動が困難なため計算環境がデータのそばに必要。
- ・生物のドメイン知識と情報技術をもつ 専門家によるキュレーションの維持が 重要。

マルチモーダルな生命科学データベースの種類

ナショナルセンターによるデータ統合の意義と、そもそも研究者によって「データベース」のイメージが異なる問題

実験由来データとAI由来データの識別



DBの管理機関や省庁が 分散し、データが散在して いることと、永続的な運用 が課題。予算人員に比し て対応分野は広い。

欧米はナショナルセンターで国家レベルのデータマネージメン トを実施、主要なデータベースが集約され相乗効果で付加価 値創出や、製薬などの産業応用にも繋がっている。 アジアでも、インド IBDC・中国 CNCB・韓国 KOBICなどナショナ ルセンターの設立が進む。



DDBJ+DBCLS 15億円・50人



DDBJ. DRA. JGA



BioProject, BioSample



NCBI

Taxonomy

GenBank, SRA, dbGaP

BioProject, BioSample

300億円・350人



EMBL-EBI 180億円 • 650人



ENA. EGA BioProject, BioSam























(EVA

GlyCosmos Microbiome Pub Chem fanta.bio





AGD, CANNDs















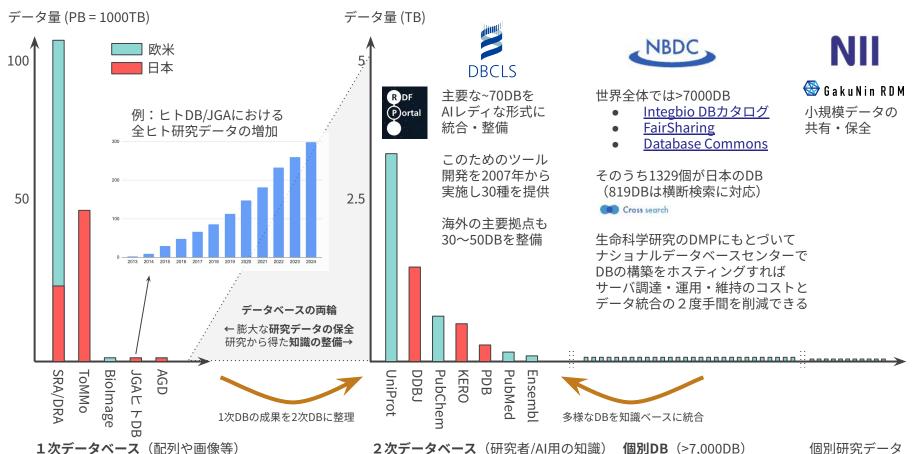
MedGen

ClinVar

Pub Med ed

マルチモーダルな生命科学データベースの量

1次データ(実験の生データ)のレポジトリ運用と、2次データ(研究成果を整理した知識)の整備



個別研究データ

マルチモーダルな生命科学データベースの統合

知識グラフによるデータ統合を技術開発と国際連携による標準化で実現

>2500 DBのカタログ化



- <u>DBカタログ</u> 2582DBs
- DB横断検索 819DBs
- DBアーカイブ 157DBs

DB間のリンク関係を整備



TogoID 114DBs

主要DBを知識グラフで標準化

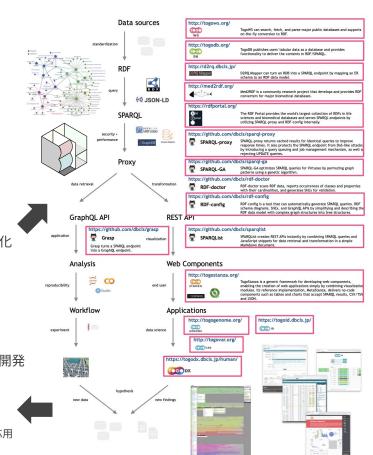


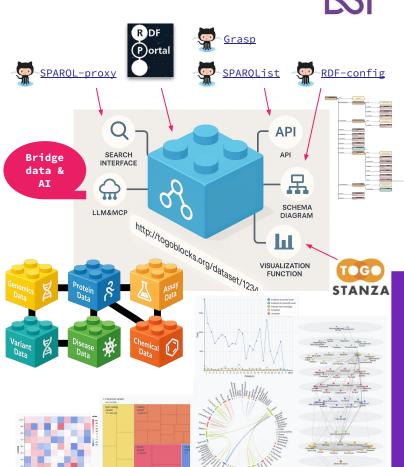
RDFポータル 70DBs + オントロジー開発

統合DBのアプリケーション開発



- TogoVar 約20DB統合
- TogoDX 約20DB統合
- <u>Tabulae</u> 統合DBの創薬応用





補足資料

マルチモーダルな生命科学データベースとAI

AIで用いる科学的知識データの提供と、データベースの品質向上や運用を効率化するAIの利用



1. 次世代AIを開発するための「学習・検証データ」となる情報基盤の 構築

- 正解となる典拠の提供と、良質な Alを育てるための構造化されたデータ整備
 → Alは学習元の質に依存、高品質なデータがハルシネーション抑制の鍵
- DBCLSの知識グラフは AI活用に最も適した形式で整備
- GBCやDICPなどの国内外の公的 DBを収載して拡張・カバレッジ強化

- 2. 登録が容易で、AIにも再利用しやすいメタデータ付きレポジトリの整備

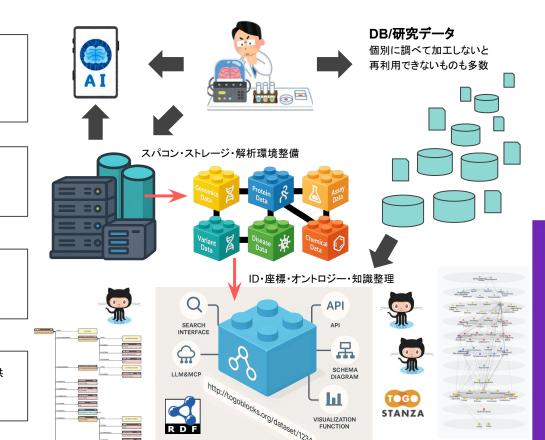
- 共通DMPの策定により、日本のすべての生命科学データを集約可能に
- データ登録を AI・UIで支援し入力負担を軽減、品質の高いメタデータを自動取得
- 集積データは永続的に管理され、他のデータと連結・検索・再利用可能に

3. 研究者が「使える」データベースの提供

- AIエージェントや LLM、LMM、MCPからのアクセスを想定した API・UI整備
- 情報科学の知識がなくても検索・再利用・発見ができる「ユニバーサルな設計」
- 産学連携に対応し、医療・創薬・基礎研究すべてのニーズに応える基盤

- 4. スパコンとストレージの拡充によるデータ利活用インフラの整備

- 大学共同利用機関法人で大学院生~研究者が公平に使える共用計算環境を提供
- ゲノムグラフによる構造多型解析や AI創薬など先端研究に直結
- 大規模データを散逸させずに利活用するための永続的なストレージ戦略を整備



知識グラフ(RDF)で統合された主要な生命医科学DB

BS

RDF: Resource Description Framework

- 塩基配列とアノテーション
 - INSDC (DDBJ/DBCLS)
- ゲノム情報
 - Ensembl (EBI)
 - RefSeq (TogoGenome)
- ▼ アミノ酸配列とアノテーション
 - UniProt (SIB)
 - タンパク質立体構造
 - o PDB (PDBj)
 - O BMRB (PDBj)
 - FAMSBASE (Chuo U)
- 化合物
 - PubChem (NCBI)
 - ChEMBL (EBI)
 - O Nikkaji (JST)
- 遺伝子発現
 - RefEx, GTEx (DBCLS)
 - ExpressionAtlas (EBI)
- サンプル
 - BioSamples (EBI/DDBJ)
 - JCM (RIKEN)

•



























- 医科学 (Med2RDF)
 - TCGA, ICGC, COSMIC, CIViC
 - DGIdb, OpenTG-Gates
 - ClinVar, dbSNP, dbVar
 - ExAC, gnomAD
- \circ HiNT, INstruct
- 糖鎖
 - GlyTouCan, GlycoEpitope, WURCS, GGDonto, PAConto
- ・ プロテオーム ○ iPOST
 - o The Human Protein Atlas
- パスウェイ
 - Reactome (EBI)
- その他
 - PubMed/MeSH (NCBI)
 - BioModels (EBI)
 - MBGD (NIBB/DBCLS)
 - Quanto (DBCLS)
 - 0:

https://rdfportal.org/



国際的なデータベースとAIの戦略および 科学データのサステナビリティへの提言





1. EMBL Science Al Strategy

欧州EMBLは生命科学の総本山ともいえる研究組織で、欧州バイオインフォマティクス研究所EBIもここに含まれる。2025年に公開された本レポートではEMBLが生命科学とAIにどのように取り組んでいくかについて、分析とビジョンが表明されている。ウェットな研究機関とドライのデータベース機関を擁するMBLでは、AI時代において、AIの理論的な基盤構築、AIレディなデータ整備、自動実験ワークフローによる研究の加速をつの主要な柱と位置づけている。

https://www.embl.org/news/connections/ai-in-the-life-sciences-philanthropy-fuels-embls-strategy/



2. How the National Library of Medicine should evolve in an era of artificial intelligence

米国国立医学図書館NLMはアメリカ国立生物工学情報センタ・NCBIを擁する世界最大の生命科学データリソースの拠点で、本レポートでは科学知識の真正性保証、科学進展の観測、そして知識のわかりやすい社会還元の重要性を論点としてNLMがAI時代にどのように変化していくべきかの考察が与えられている。とく「AI時代において科学的データの整合性の保護者としてデータの真正性安保証する機関となること、多様な利用者向けに科学情報を翻訳、変換、要約するハブとなることなどが謳われている。https://pmc.ncbi.nlm.nih.gov/articles/PMC12012362/

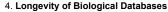


3. Exploring Determinants of Longevity of Biomedical Databases

Joseph Finkelstein et al., Joseph Stud Health Technol Inform, 6:290:135-139 (2022)

2009~2016年に論文発表された518データベースについて調査したところ2020年にはその35%にあたるデータベースがアクセスできなくなっていた。

https://pubmed.ncbi.nlm.nih.gov/35672986/



Teresa K Attwood et al., EMBnet.journal, e803 (2015)

この論文でも2015年までに発表された326データベースが18年間でどれくらい存続しているかを調査し、その60%が消失、14%はアーカイブされ更新がない状態であることを報告している。一方で研究機関からのサポートがあるデータベースは存続しているものが多いことを指摘し、現代の生命科学研究に必須なデータベースには長期的な戦略による予算とインフラの維持が求められることを訴えている。

https://journal.embnet.org/index.php/embnetjournal/article/view/803/1209



5. On the lifetime of bioinformatics web services

Fabian Kern et al., Nucleic Acids Research, 48(22):12523-12533 (2020) 2010年から2020年までに論文発表された生命情報科学のウェブツールは2010年のツールの50%から2019年と2020年のものの90%しか利用可能な状態を保っていなかった2週間以上にわたって利用不可能になっていた019年と2020年のツール47個のうち、著者に連絡を取ることで約半数は再び利用可能になったが、長期的な維持のためには十分な支援が必要と考えられる。

https://academic.oup.com/nar/article/48/22/12523/6018434



6. Biological databases in the age of generative artificial intelligence

Mihai Pop et al., Bioinformatics Advances, 5(1):vbaf044 (2025) 生命科学研究は公共データベースに依存しており、誤りの導入や伝播が研究の妨げとなる。 生成AIIにより大量の誤情報が拡散するリスクも高まっている。これを防ぐため、研究者教育の強化、誤りやデータ来歴に関する研究推進、データベース維持への資金支援が重要とされる。

https://academic.oup.com/bioinformaticsadvances/article/5/1/vbaf044/8088229



7. Global Biodata Coalition Highlights 2024

Global Biodata Coalition(GBC) は、生命科学研究に不可欠な国際的データベースやバイオデータ資源の持続可能な運営を支援・調整するために設立された国際組織で、各国の研究資金機関が連携し、データ基盤の長期的な維持と国際協力を推進しているDBJも含まれる52のGBCメンバーからなるGCBRフォーラムのレポートで下記のような記載がある12ページ):

Forum members highlighted key challenges, including ... concerns over Al-generated data potentially compromising trust in curated knowledgebase フォーラム参加者は、主要な課題としてAI生成データがキュレーション済み知識ベースの信頼性を損なう可能性への懸念などを挙げた。

https://globalbiodata.org/wp-content/uploads/gbc-highlights-brochure-2024.pdf

