

支援事業

・データ共有支援事業

データベース名	研究分野、コミュニティ	活動内容
<生命科学データベースの統合利活用> ・TogoDX (生命科学データ俯瞰探索) ・TogoVar (日本人ゲノム多様性データベース) ・RDF ポータル (生命科学 RDF データ一覧) ・以上を支える RDF データ検索支援・作成支援などのツール群	生命科学、医学研究	データベース統合利用のための環境構築、データベース統合化のための基盤技術開発及び国際標準整備
<極域科学の公開用データベースシステム> ・学術データベース ・北極南極データアーカイブシステム (ADS) ・超高層大気観測研究ネットワークデータシステム (IUGONET)	極域・地球環境	南北両極域から得られた科学データの公開と共同利用、データサイエンスを推進し地球環境研究に貢献する
<社会状況に関するデータの提供・共同利用> ・社会調査の個票データ 日本人の国民性調査 意識の国際比較調査 ・公的統計のマイクロデータ (オンサイト解析室による提供)	社会科学全般 (経済学、社会学、政治学、行政学、国際関係等)、世論調査、政府統計、ソーシャルビッグデータ	社会状況に関するデータ (社会調査、公的統計のマイクロデータ、ソーシャルビッグデータ) 有効利用のための大学間連携ネットワーク基盤形成及び地域社会への貢献
<データサイエンスに基づく人文学 (人文情報学) 分野の創生> ・日本古典籍データセット ・日本古典籍くずし字データセット ・江戸料理レシピデータセット ・顔貌コレクション ・江戸ビッグデータ ・歴史資料情報共有データベース	人文情報学、機械学習、美術史、古気候学、日本文化研究	情報学・統計学の最新技術を用いて人文学資料 (史料) を分析する「データ駆動型人文学」、人文学研究の成果に基づき構築したデータセットを超学際的に活用する「人文学ビッグデータ」など、オープンサイエンス時代の新しい人文学研究を展開

・データ解析支援事業

解析対象	研究分野、コミュニティ	活動内容
<さまざまな生物種のゲノムデータ解析> ・新規ゲノム決定 ・ゲノム再シーケンス ・トランスクリプトーム解析 ・メタゲノム解析	ゲノム生物学 ゲノム医学 ゲノム創薬	次世代型 DNA シーケンサーから得られる大量の配列データに基づいた多様な生命科学研究を対象に、情報科学的な解析支援を実施
<数値シミュレーション全般の理論・応用支援> ・人流シミュレーション ・宇宙機シミュレーション ・沿岸海洋モデル ・核融合プラズマの統合輸送シミュレーション ・磁気圏電離圏モデル	交通工学 宇宙工学 地球物理学 核融合学 ほかデータ中心科学の考え方・手法を用いる分野全般	データ同化やエミュレータなどの手法を活用する研究相談・支援及び共同研究 上記に関連するハンズオンによる講習会や体験学習の実施

データサイエンス共同利用基盤施設

〒190-0014 東京都立川市緑町10-3 データサイエンス棟
<https://ds.rois.ac.jp/>



ロゴマーク紹介

データサイエンスの略である DS の文字で地球 (地軸) の傾きを、そして周囲は 4 つの研究所を表現しています。
 全体では、枠に収まらずに形を変えながら発展するさまを表すとともに、親しみのあるシルエットにより社会貢献を意味する組織をイメージしています。



お問い合わせ窓口：データサイエンス推進室 電話 042-512-9254 E-mail ds_suishin@rois.ac.jp

大学共同利用機関法人 情報・システム研究機構

データサイエンス 共同利用基盤施設

Joint Support-Center for Data Science Research (ROIS-DS)

ROIS-DS

一連携、協働、そして発展へー
データ駆動型研究で研究者を支援する
分野融合的な共同利用・共同研究拠点です。

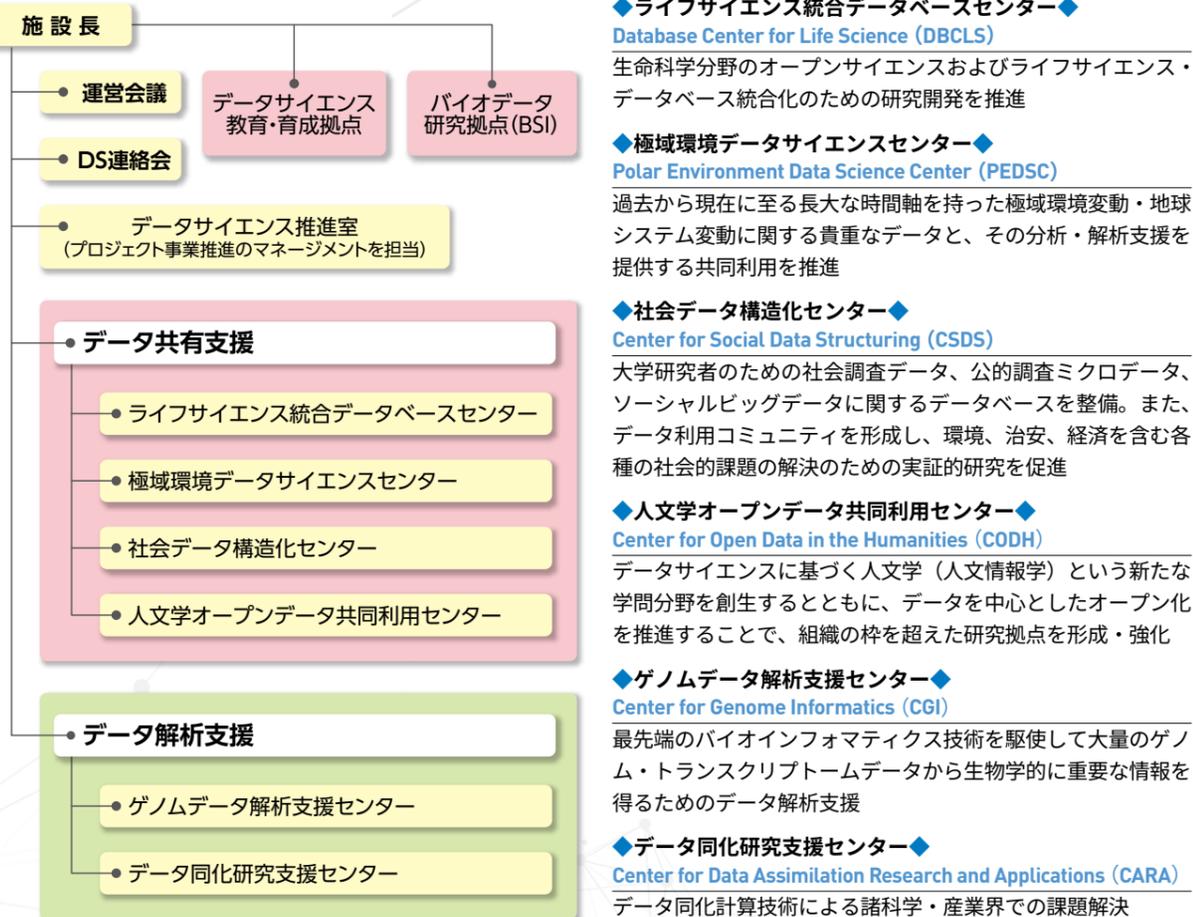


施設長 荒木 弘之

データサイエンスの推進拠点

情報・システム研究機構「データサイエンス共同利用基盤施設（DS施設）」は、大規模データの高度な解析により科学や社会の課題を解決する「データサイエンス」を全国規模で融合的に推進するための共同利用・共同研究拠点です。データサイエンス（データ駆動型研究）を合い言葉に大学等との連携・協働を強化する目的で、2016年に情報・システム研究機構により設置されました。2023年現在、ライフサイエンス統合データベースセンター（DBCLS）、極域環境データサイエンスセンター（PEDSC）、社会データ構造化センター（CSDS）、人文学オープンデータ共同利用センター（CODH）、ゲノムデータ解析支援センター（CGI）、データ同化研究支援センター（CARA）の計6センターで構成され、大学等の研究力強化に貢献しています。当施設が扱うデータは、ゲノムや遺伝子に関する大量の生命情報データや大気レーダーデータ等による観測データから、古典籍や社会調査、公的マイクロデータまで極めて広範囲に及ぶため、情報・システム研究機構を構成する4つの研究所（国立極地研究所、国立情報学研究所、統計数理研究所、国立遺伝学研究所）や、他の大学共同利用機関法人の研究所等と協力して活動しています。

データサイエンス共同利用基盤施設(ROIS-DS)



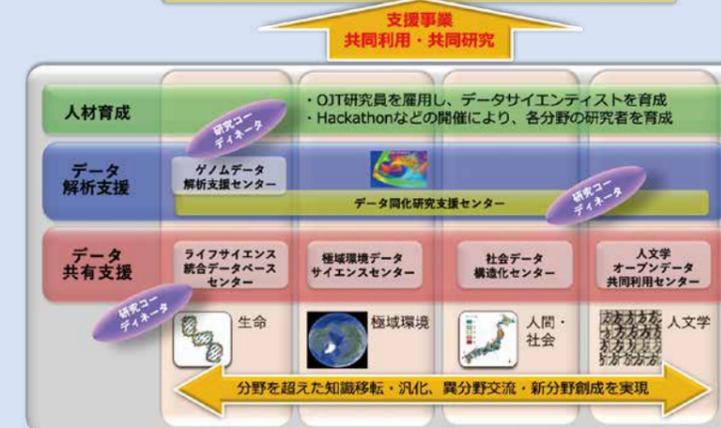
事業内容

- 支援事業（データ共有支援）
 - 生命科学分野におけるデータ共有支援事業
 - 極域環境科学分野におけるデータ共有支援事業
 - 人間・社会分野におけるデータ共有支援事業
 - 人文学オープンデータ共有支援事業
- 支援事業（データ解析支援）
 - ゲノムデータ解析支援事業
 - データ融合計算支援事業
- 人材育成事業（データサイエンティスト育成）
 - OJTによるデータサイエンス人材の育成
 - データサイエンス教育人材（DS教員）養成事業
- 公募型共同研究（ROIS-DS-JOINT）
 - 一般共同研究
 - 共同研究集会

支援事業の活動内容はP.12にあります。

公募型共同研究の説明はP.10にあります。

データ共有支援事業・解析支援事業および共同利用・共同研究を必要としている大学等のすべての研究者



拠点活動

- データサイエンス教育・育成拠点
 - DS施設と各研究所の連携により、様々なレベル・分野のデータサイエンス人材を育成
- バイオデータ研究拠点
 - DS施設と遺伝研等との連携・センター統合および「データの収集・整理・標準化等一体化」による世界レベルの研究効率化の支援
- 日本文化ビッグデータ研究ハブ
 - DS施設と人間文化研究機構との連携により、日本文化をデータ駆動型方法論で分析

設置：2016（平成28）年4月1日
 所在：東京都立川市緑町10-3 立川キャンパス内データサイエンス棟、および各研究所
 人員：65人（研究系42人、事務系23人）
 予算：5.86億円 2022年度（令和4年度） ※ DBCLS DB 統合推進事業及び SIP 受託事業等に係る予算を除く

ライフサイエンス統合データベースセンター



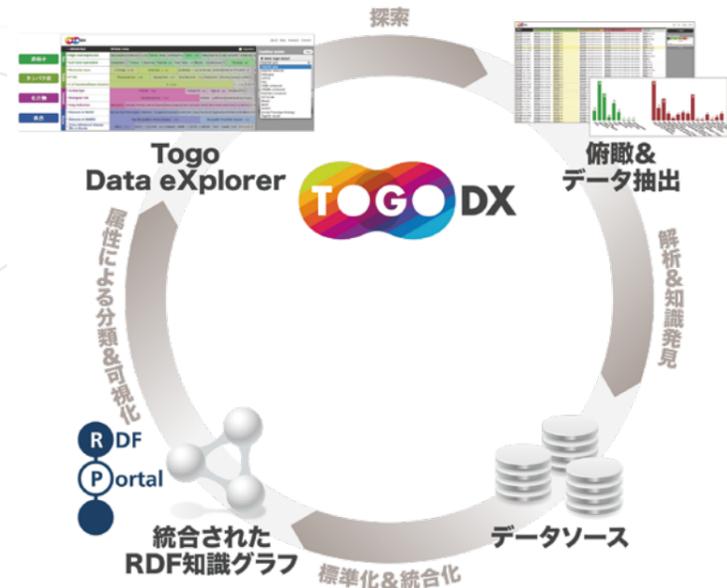
データベース統合化を通じた生命科学分野のオープンサイエンスの推進

本センターは生命科学分野のオープンサイエンスを目指し、全国の大学、研究機関などが所有・生産する多様かつ急速に増加するデータベースを一元的に活用するための「データベース統合化」に関する研究開発を行っています。データの記述に用いる用語とその分類体系(オントロジー)を標準化した知識グラフの構築、構築された知識グラフから必要な情報を効果的に抽出するための技術開発、文献情報の有効利用などに力を入れています。また、世界各国のデータベース開発機関からエキスパートを招き、バイオハッカソンなどの開発者会議を毎年開催し、統合化のための技術開発と標準化を国際的に主導しています。



小原 雄治センター長

▶ウェブ上に分散しているデータベース(DB)を統合的に利用できる環境を構築するために、生命科学DBの知識グラフ(Knowledge graph)化を目指し、各種データソースのResource Description Framework(RDF)化支援とRDFデータの集積を行っています。また、知識グラフ基盤を活用するための様々なアプリケーションを開発しています。例えば、様々な情報を統合的に探索するためのフレームワークTogoDXを開発し、ヒトに関する情報をワンストップで探索することができるアプリケーションTogoDX/Human (<https://togodx.dbcls.jp/human/>)を提供しています。研究対象の絞り込み、実験結果の解釈や考察の一助として利用できるほか、絞り込んだ結果を統合解析に応用可能です。今後はヒトに関するデータをさらに充実させるとともに、TogoDXを他の生物種へ応用するなどの展開をしていきます。



様々なデータベースを統合した知識グラフとデータを統合的に探索するためのフレームワークTogoDXによるデータ解析プラットフォーム

▶統合的利用の実現には、国内外のDB関係機関との連携が必要です。そのために約1週間の合宿形式で共同開発作業を集中的に行う国際版「BioHackathon」を始めとする各種のハッカソンを10年以上開催してきました。コロナ禍で開催が難しい中でも、国内版バイオハッカソン(年1回)、Togothon(旧 SPARQLthon、毎月1回連続2日間でデータ統合に関する議論や開発の相談を行う)の開催を継続し、2023年6月には4年ぶりの国際版バイオハッカソンを再開できました。これらにより、データ相互利用のためのルール、仕組みやツールができ、国内外の機関で採用されてきており、国際的な標準化を進めています。



2023年6月開催の国際版バイオハッカソン集合写真@小豆島

■ ROIS-DS-JOINT ライフサイエンス統合データベースセンター支援対象

大学、研究所、企業等に所属の方で、以下の技術およびその関連手法の開発・応用を考えている方
キーワード：【生命科学データ統合】【知識グラフ】【大規模データ解析】

詳しくはこちら [ROIS-DS-JOINT これまでの採択課題一覧](#)

極域環境データサイエンスセンター



南北両極域から得られた様々な科学データの公開と共同利用、データサイエンスを推進し地球環境研究に貢献することを目指しています。

当センターでは、国立極地研究所を中心とする極域科学研究コミュニティが、南極域、北極域での観測・研究活動により取得した、多分野かつ多種多様な科学データの公開と共同利用を促進し、大学等外部コミュニティとの連携をさらに強化し、より多くの研究成果の創出と、極域科学研究の価値のさらなる向上に、データ活動面からの支援を行うことを目的とした活動を行っています。また、国際的には、日本の極域観測・研究のナショナルデータセンターとしての役割も果たすことを目指しています。このように、極域科学におけるデータ活動の中核を担うとともに、データに基づく新しい極域科学を創生し、地球環境研究に貢献することを目的としています。



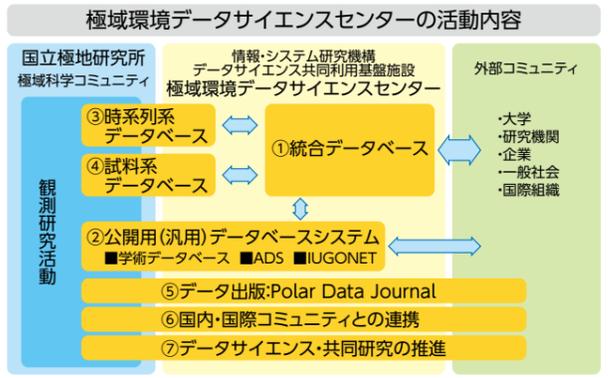
門倉 昭センター長

▶極域環境データサイエンスセンターが取り扱うデータは、南極域、北極域での科学観測・研究活動によって取得された全ての分野のデータになります。両極域では、国立極地研究所(極地研)を中心に、超高層、大気、海洋、雪氷、地学、生物など様々な分野の観測・研究が行われていて、様々な媒体に記録されたデジタルデータや、採取されて保管された試料系データなど、多種多様なデータが取得されています。それらのデータは、取得後に、様々な処理、解析、分析をされ、物理的に意味のあるデータとなった後に、それらを基にした科学的な成果が生み出されます。信頼される科学的成果を生み出すためには、データの信頼性が確保される必要があります。そのためには、データが確実に保管されていて、失われたり、劣化したり、改ざんされたりしないこと、そのデータが誰にでも利用可能で、同じ科学的成果の再現性が保たれること、などが求められます。



極域環境データサイエンスセンターでは、南北両極域での科学観測・研究活動によって取得された全ての分野のデータを扱います。

▶また一方で、地球環境変動のような研究では、多分野の多種多様なデータを同時に用いることによって全く新しい成果が生み出される、ということもあります。その場合は、様々なデータの所在情報、属性情報などのメタ情報(メタデータ)を統一的に扱う必要があります。また、ある分野のデータが、予想も出来ない分野に応用され、予想もされない新しい成果や価値が生み出される、ということもあります。そのためには、そのデータの公開性や所在の分かり易さが重要になります。当センターでは、こうした、極域科学データの、処理、解析、保管、共有、公開、共同利用、についての活動支援を行っています。



極域環境データサイエンスセンターは、極域科学コミュニティと外部コミュニティとの間の様々なデータ活動の橋渡しをします。
<http://pedsc.rois.ac.jp/ja/activity>

■ ROIS-DS-JOINT 極域環境データサイエンスセンター支援対象

大学、研究所、企業等に所属の方で、以下のデータやデータベースの利用を考えている方
キーワード：【極域科学データ】【学術データベース】【IUGONET】

詳しくはこちら [ROIS-DS-JOINT これまでの採択課題一覧](#)

社会データ構造化センター



社会を対象として得られる様々なデータの整備と 共同利用を通じて社会的課題の解決に貢献

社会データ構造化センターは、社会で生じる様々な事象を測定・計測することで得られる様々なデータ個人や組織を対象とする社会調査の実施を通じて得られる社会調査データ、官庁等の公的セクションが実施する統計調査の結果としての公的マイクロデータ、様々な機器を通じて社会活動をリアルタイムで計測するソーシャルビッグデータ、等一を整備すること、それらを広範な活用に供することで、各種の社会的課題の解決のための実証的学問を促進し、実証的データに基づく政策立案の実現のための研究基盤を発展させることを目標に活動しています。データの整備・活用に資する基盤的な技術開発も、本センターの活動の目標の一部となっています。



前田 忠彦センター長

▶本センターは、データが得られる分野に対応して、主に3つのグループに分かれて、次のプロジェクト・事業を展開しています。

社会調査関連事業

全国共同調査ネットワーク形成によるデータ収集、及び社会調査データの整備と公開を進めます。統計数理研究所から継承した大規模学術調査データの整備と公開、他機関の研究者と共同で実施する調査の企画やデータ共有、社会調査の実施に伴うコンプライアンスに関する課題の研究と普及などのプロジェクトを推進しています。

公的マイクロデータ事業

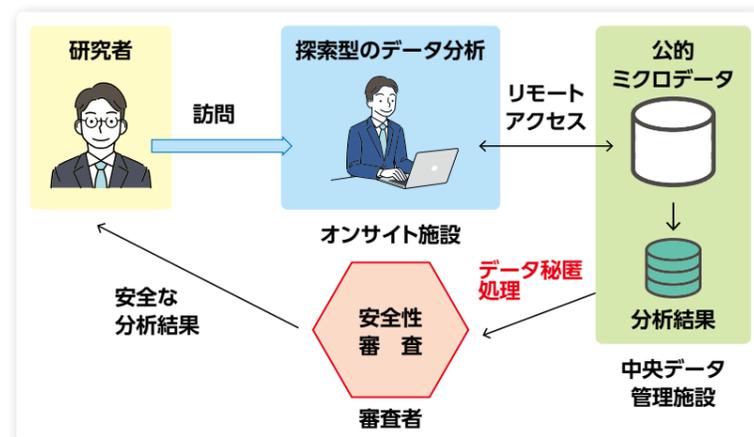
公的統計データの整備と共同利用システムの発展、及びオンラインのデータ解析システムの研究開発、オンサイト施設の運営などを担当します。公的マイクロデータの安全な公開に関わるマイクロデータ・セキュリティ、公的統計二次利用の推進、経済・金融分野でのリスク管理に関わる高度信用リスク、政府統計を用いたEBPM（証拠に基づく政策決定）、等の研究プロジェクトを推進しています。公的統計マイクロデータ研究コンソーシアムの事務局機能も担っています。

ソーシャルビッグデータ事業

異なる組織の研究者間によるソーシャルビッグデータを用いる研究活動で必要となる共同利用データの管理方法ならびにプラットフォームの整備を進めています。また、ソーシャルビッグデータを用いた共同研究も実施します。実社会データ共有基盤の開発を目指すプロジェクトでは、道路や交通などの社会インフラで、常に変動する実社会の状況の効率的な収集、状況把握、分析を可能にするデータ共有基盤システムを、自治体等と連携して実証的に開発します。



データライフサイクル



オンサイト解析のプロセス（施設利用についてはP.11を参照ください）

■ ROIS-DS-JOINT 社会データ構造化センター支援対象

大学、研究所、企業等に所属の方で、以下のデータに関わる整備と共同利用、共同利用に至る諸プロセスに関わる基盤的な技術の開発・研究を考えている方
キーワード：【社会調査データ】【公的マイクロデータ】
【ソーシャルビッグデータ】

詳しくはこちら [ROIS-DS-JOINT これまでの採択課題一覧](#)

人文学オープンデータ共同利用センター



人文学におけるオープンサイエンスと デジタル変革の推進： データ駆動型人文学と人文学ビッグデータの展開

人文学オープンデータ共同利用センター（CODH）は、人文学分野におけるオープンサイエンスとデジタル変革の推進を目指し、2つのテーマに注力して研究を進めています。第一に、情報学・統計学における最新のデータ駆動型技術の導入により人文学の研究手法を変革する「データ駆動型人文学」の研究です。日本古典籍に対するAIくずし字認識の研究やIIIF (International Image Interoperability Framework) を活用した美術史研究など、新しい技術を用いて人文学データから新たな知を創出します。第二に、人文学分野で生み出されたビッグデータを他分野の研究に活用する「人文学ビッグデータ」の研究です。江戸時代の日記から天気データを取り出して古気候復元に用いる歴史ビッグデータの研究などを進めています。



北本 朝展センター長

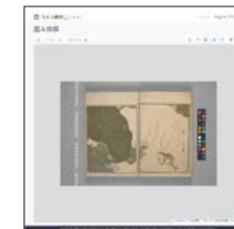
AIくずし字認識アプリ「みを」

スマホでくずし字画像を撮影すると、AIが数秒で現代日本語文字に変換して表示するアプリを、iOSおよびAndroidで無料公開しています。



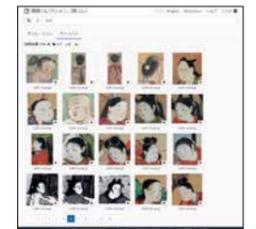
IIIF Curation Platform

IIIFの世界に、キュレーションという新しいコンセプトを導入し、利用者主導型のIIIFプラットフォームを実現します。



顔貌コレクション（顔コレ）

IIIF Curation Viewerを用いて、美術作品に出現する顔の部分を取り取って収集し、美術史研究に活用します。



江戸ビッグデータ

都市「江戸」に関する地理情報と紐づく、商業や観光・人物などの歴史ビッグデータを統合し、データを基に過去の世界を分析します。



武鑑全集

江戸時代の200年続いたベストセラー『武鑑』を網羅的に解析し、大名家や幕府役人に関する中核的情報プラットフォームを構築します。



れきすけ

歴史資料を利用した研究のために、歴史資料に関する知識や経験を、様々な分野の研究者で共有します。



■ ROIS-DS-JOINT 人文学オープンデータ共同利用センター支援対象

大学、研究所、企業等に所属の方で、以下の技術およびその関連手法の開発・応用を考えている方
キーワード：【日本古典籍・くずし字】【IIIF】【人文学ビッグデータ】【デジタルヒストリー】

詳しくはこちら [ROIS-DS-JOINT これまでの採択課題一覧](#)

ゲノムデータ解析支援センター

大量のゲノム・トランスクリプトームデータから生物学的に重要な情報を抽出するための情報科学的解析を支援します。



次世代シーケンシング (NGS) 技術の発展に伴い、さまざまな生命科学研究の分野で新規ゲノムシーケンスやリシーケンス、トランスクリプトーム解析などのNGSを用いた解析が広く行われるようになってきました。しかし、NGSデータはあくまで塩基配列の断片データでありデータ量も膨大なため、これらを効率的に解析し目的に合う結果を得るには生物学の知識に加えてバイオインフォマティクスの知識と技術が不可欠です。

ゲノムデータ解析支援センターでは、大量のゲノムデータを迅速かつ高精度に解析するための情報科学技術の研究開発や、最先端の手法を用いた実データの解析支援、またそのための人材の育成を行なっています。



野口 英樹センター長

▶ゲノムデータ解析支援センターでは、大学・研究機関等の研究者を対象にさまざまな種類のゲノムデータ解析支援を行なっています (図1)。2016年度から2021年度 (平成28~令和3年度) までの6年間には、32の大学・研究機関 (代表者の所属) からの依頼で合計51課題の解析支援を実施しました。

▶取り扱うゲノムデータは主にNGSの配列データですが、研究の目的や実験条件などは研究ごとに大きく異なります。また、対象の生物種も哺乳類やその他脊椎動物を中心に、昆虫、植物、真菌、原核生物などさまざま、ゲノムサイズや構造、進化的背景などに応じて適切な解析手法を選択する必要があります。当センターでは豊富な解析経験を活かして研究目的に応じた柔軟で高精度な解析支援を行なっています。

▶また、解析支援を円滑に実施するためにゲノムアノテーションパイプライン (図2) やゲノム再シーケンスパイプラインなどの各種解析パイプラインの開発を行なっているほか、遺伝子予測手法やRNA-seqアセンブラ、メタゲノムの種分類手法などの新規解析手法の開発も行ない、最先端の解析手法を提供できる環境を整えています。

解析内容	
<ul style="list-style-type: none"> De novo ゲノムシーケンス De novo ゲノムアセンブリ 参照配列のない新規生物種のゲノムを、NGSデータを用いて構築します。 ゲノムアノテーション ゲノム配列上の遺伝子の位置やエキソン・イントロン構造を同定し、注釈付けします。 ゲノムリシーケンス 全ゲノムリシーケンス 全ゲノムの配列リードを参照ゲノム配列と比較し、SNVや構造多型を検出します。 ターゲットゲノムリシーケンス・エピジェネティクス解析 エキソーム、RAD-seq、ChIP-seq、HiC-seq等々。 	<ul style="list-style-type: none"> トランスクリプトーム解析 遺伝子構造・発現解析 RNA-seqデータのde novoアセンブル・マッピングを通して、遺伝子構造同定や発現量解析を行います。 non-coding RNA解析 RNA2次構造予測、miRNAのターゲット検索など。 メタゲノム解析 メタゲノムアセンブリ メタゲノム配列リードをde novoでアセンブルします。 種分類、遺伝子予測 メタゲノム配列のクラスタリング、遺伝子予測、パスウェイ解析など。

図1 当センターで行なっているゲノムデータ解析

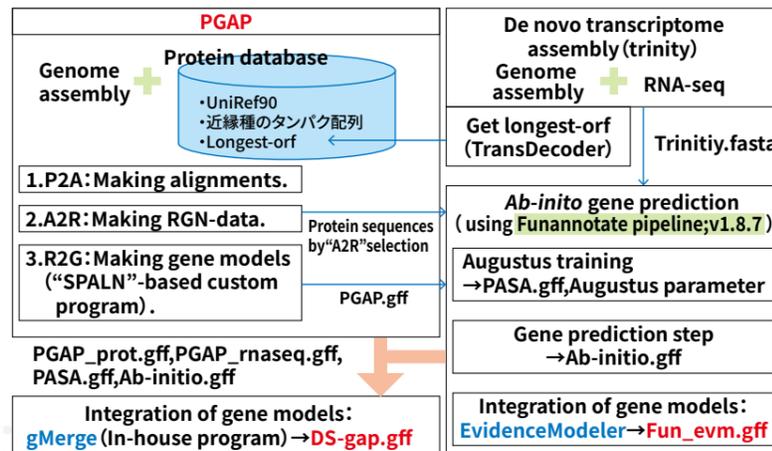


図2 ゲノムアノテーションパイプライン

ROIS-DS-JOINT ゲノムデータ解析支援センター支援対象

大学、研究所、企業等に所属の方で、以下の技術およびその関連手法の開発・応用を考えている方
 キーワード：【バイオインフォマティクス】【次世代シーケンシング (NGS)】【ゲノムDNA・RNAデータ解析】

詳しくはこちら ROIS-DS-JOINT これまでの採択課題一覧

データ同化研究支援センター

シミュレーションと観測データの統合による問題解決：合わないシミュレーションをまだ続けますか？観測データを予測に生かすには？



データ同化とは、観測データと数値シミュレーションを統合する方法です。データ同化により、高精度の予測が可能なシミュレーションである「データ同化システム」や、計算時間を大幅に短縮できるシミュレーションである「エミュレータ」を開発できます。本センターは、データ同化研究の相談窓口を開いており、面談で助言や技術指導を提供し、問題解決の支援を行っています。統計科学を基盤とするデータ同化の方法から、データ同化を応用する現場の観点からの相談対応、ならびに共同研究が可能です。データ同化を新しく導入したいがどうすればよいのか、データ同化の計算を完了したがこの先どうすればよいのか、など、ご相談をお待ちしております。



上野 玄太センター長

【面談現場の実況中継】

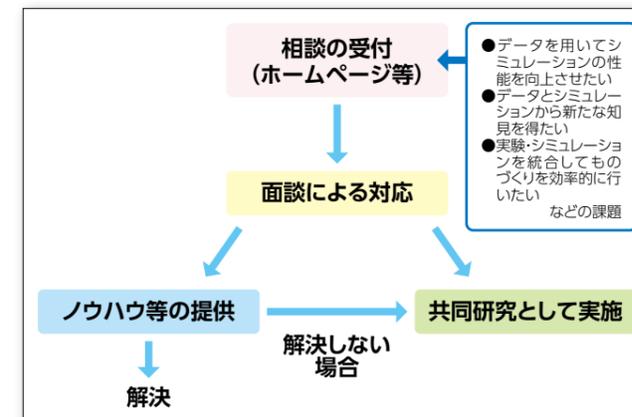
立川までおいでいただき、ありがとうございます。本日は、2時間程度、よろしくお願ひします。

データ同化を導入したいということですね。そちらの分野ではまだ使われていない。「データ同化」という言葉はどこで聞かれました？ 私のセミナーをお聞きに。ありがとうございます。それで、シミュレーションはお持ちで、観測データもあると。はい、ではご説明、よろしくお願ひします。(説明中)なるほど、ありがとうございました。

ではデータ同化の説明に入ります。説明いただいた資料の、この変数を全部まとめてx、この方程式を全部まとめてfと書きたいのですが、よいでしょうか。ホワイトボードに書きますね。はい、xは巨大なベクトルになります。時間ステップが違う変数は別にまとめます。データのほうも同じように、全部まとめてyとします。

これからがデータ同化の本題です。xとyが合わないことが問題なので、データ同化では、yを使ってxを修正するのは。修正といっても、単にxの値をyの値で置き換えてもうまくいかないところが難しいのです。うまく修正するための、パラメータがいくつも出てきます。感じをつかむために、簡単な例題を実装してみるのがよいです。

サンプルプログラムですか？ 今後を考えると、紹介した本を参考に、最初からプログラムを書いたほうがよいです。シミュレーションがFORTRANなら、データ同化もFORTRANで始めるのがよいです。頑張ってくださいね！



研究相談は、ホームページをご覧のうえメールにてお申し込みください。



センター長との初回相談が着落し、相談内容 (ホワイトボード) を写真撮影

ROIS-DS-JOINT データ同化研究支援センター支援対象

大学、研究所、企業等に所属の方で、以下の技術およびその関連手法の開発・応用を考えている方
 キーワード：【データ同化】【シミュレーション】【エミュレータ】

詳しくはこちら ROIS-DS-JOINT これまでの採択課題一覧

