

平成19年度  
委託業務成果報告書

「塩基配列アーカイブのデータベース構築と統合への貢献」

大学共同利用機関法人情報・システム研究機構  
国立遺伝学研究所

本報告書は、文部科学省の科学技術試験研究委託事業による委託業務として、大学共同利用機関法人情報・システム研究機構が実施した平成19年度「塩基配列アーカイブのデータベース構築と統合への貢献」の成果を取りまとめたものです。

従って、本報告書の著作権は、文部科学省に帰属しており、本報告書の全文又は一部の無断複製等の行為は、法律で認められたときをのぞき、著作権の侵害にあたるので、これらの利用行為を行うときは、文部科学省の承認手続きが必要です。

## 目 次

### 「塩基配列アーカイブのデータベース構築と統合への貢献」

	頁
業務題目-----	1
1. 委託業務の目的-----	1
2. 平成19年度（報告年度）の実施内容-----	1
2.1 実施計画-----	1
2.2 実施内容（成果）-----	2
1 公開FTPサイトとWWWサイトに関する開発-----	2
2 ゲノムネットワーク利用システムの開発-----	3
2.3 成果の外部への発表-----	4
2.4 活動（運営委員会等の活動等）-----	4
2.5 実施体制-----	4

業務題目：「塩基配列アーカイブのデータベース構築と統合への貢献」

担当者：（所属）大学共同利用機関法人情報・システム研究機構国立遺伝学研究所

（氏名）五條堀 孝

## 1. 委託業務の目的

各種生物の遺伝子やゲノムの塩基配列を決定するいわゆるシーケンシングセンターにおいては、その配列決定の原データになる波形データなどのいわゆるアーカイブ(Trace archive; 以下 Trace データという。)を有している。この Trace データは、品質管理やそれを基にして行う配列決定アルゴリズムの改良ならびに配列断片を連結するアセンブリにおいて大変重要で貴重な情報である。そして、これらの Trace データは、シーケンシングセンターの活動がその支持母体のプロジェクトの完了とともに終了すると、原則的には完全に消滅してしまう可能性が極めて高い状況にある。また、454 や Solexa 又は ABI-SOLiD といった次世代の超高速の塩基配列決定装置の登場により、その Trace データの量は飛躍的に巨大化してきており、シーケンシングセンター自身でもそのデータハンドリングを含めて保存はもちろんのこと対処が非常に困難な状況になっている。

この状況の理解の下、わが国における塩基配列決定における Trace データの保存と有効利用を目的として、当機関である情報・システム研究機構国立遺伝学研究所生命情報・DDBJ 研究センターの DDBJ が、Trace データのデータベース構築事業とデータ提供の事業を実施する。

具体的には、同機関の同研究所の生物遺伝資源情報総合センターアカデミア DNA シーケンシングセンターが保有する Trace データを当初の対象データとして、データベース構築を行なうことを目標とする。さらに、先行して Trace データを蓄積してきた米国 NCBI 及び欧州 EBI における Ensembl とはデータを共有して、DDBJ からの国際貢献とする。

そして最も重要なことに、このデータベースを統合データベースに提供し、統合データベースが必要とする基本的なアノテーションの実施やソフトウェアの開発を連携して行うものである。

## 2. 平成 19 年度（報告年度）の実施内容

### 2.1 実施計画

統合データベースにおける情報資源のひとつとして利用できるような基礎データアーカイブ作成を目指して以下のような計画で本年度の作業を進める。

#### 公開 FTP サイトと WWW サイトに関する開発

平成 20 年度以降の独立行政法人理化学研究所、財団法人かずさディー・エヌ・エー研究所、独立行政法人製品評価技術基盤機構、独立行政法人農業生物資源研究所などの塩基配列シーケンシングセンターである国内主要機関を中核とするコミュニティとの連携を深めながら、Trace データの実データの登録開始を目標にして、必要となる基盤技術（大規模 Trace データのデータベース化に関わるデータ登録・アーカイブ処理・データ提供に関する技術）の整備を行う。上記の基盤技術を用いて、大規模データ転送（数十 GB から数百 GB）に必要な FTP サーバの構築と、個別データの閲覧、波形データの表示、統計処理などを行うための WWW システムの開発を進める。

#### 登録処理及び波形表示システムに関する開発

平成 19 年度に登録処理及び波形表示システムのプロトタイプを構築し、このプロトタイプをもとにして情報・システム研究機構国立遺伝学研究所の生物遺伝資源情報総合センターアカデミア DNA シーケンシングセンターの Trace データを対象としたデータベース構築のための基礎とする。

## 2.2 . 実施内容 ( 成果 )

### 公開 FTP サイトと WWW サイトに関する開発

統合データベースへの Trace データ統合のために必要となる機能に関して、独立行政法人理化学研究所、東京大学、財団法人かずさディー・エヌ・エー研究所、独立行政法人製品評価技術基盤機構、独立行政法人農業生物資源研究所などの塩基配列シーケンシングセンターである国内主要機関を中核とするコミュニティとの連携を深めながら、シーケンサーによる塩基配列決定の際、波形情報などの一次情報として得られる Trace データの実データの登録開始を目標にして、必要となる基盤技術 ( 大規模トレースデータのデータベース化に関わるデータ登録・アーカイブ処理・データ提供に関する技術 ) の整備を進めることを目的として、上記の基盤技術を用いて、大規模データ転送 ( 数十 GB から数百 GB ) に必要な FTP サーバの構築と、個別データの閲覧、波形データの表示、統計処理などを行うための WWW システムの開発を進めた。現在、データの転送 ( 長時間が必要 ) を半自律化 ( エラー時再試行、完了後レポート送付、など ) する仕組みはプロトタイプが完成し、試験中である。

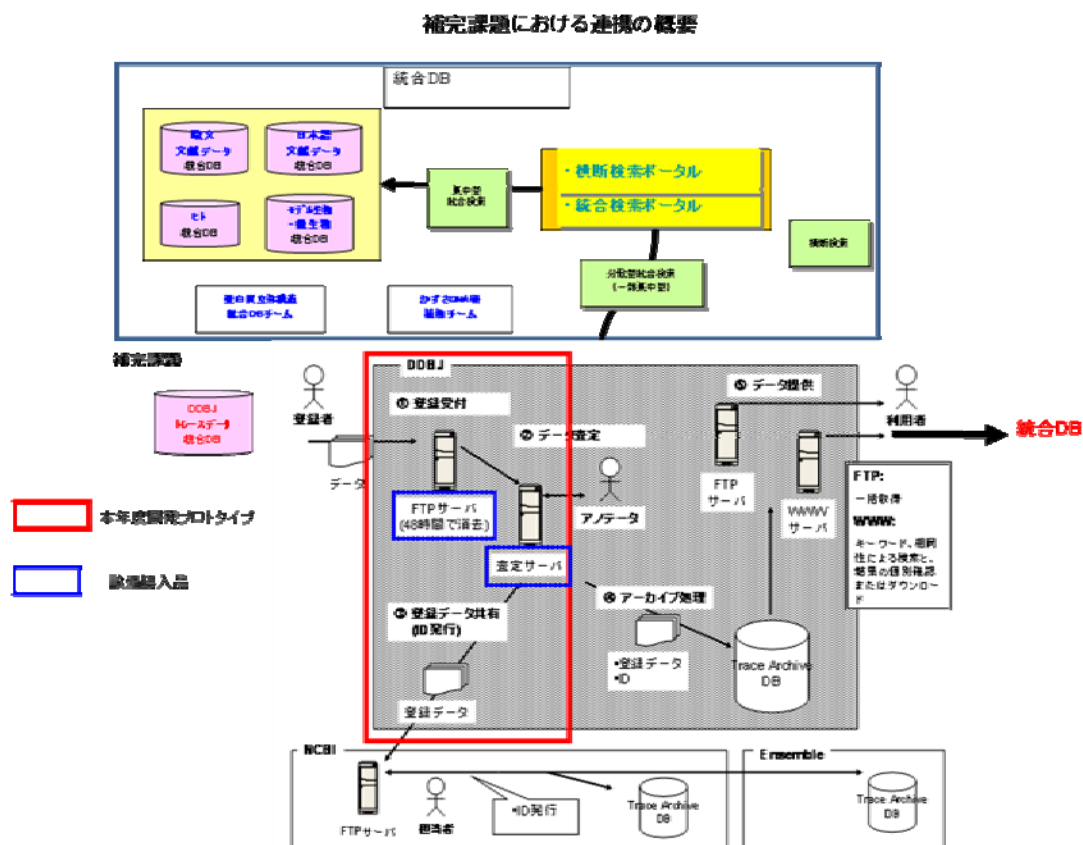


図1 統合との連携における DDBJ における Trace データの流れ

### ゲノムネットワーク利用システムの開発

統合データベースからの Trace データの利用実現のために、平成 19 年度にプロトタイプ ( NCBI / Ensembl の構成に基本的には準拠しながら、DDBJ 特有のデータベース管理システムやパイプ

ラインの効率化に合わせて、より国内の状況に応じたシステム)を構築し、このプロトタイプを情報・システム研究機構国立遺伝学研究所の生物遺伝資源情報総合センターアカデミア DNA シーケンシングセンターの Trace データを対象として適応し、データベースの構築を開始した。(1)登録データの受け付けから、受付時の査定、NCBI へのデータ提供、本年度は、登録データの機械的チェックのシステムや、トレースデータの形式の相互変換 (AB1 - SCF) など、最低必要な部分のシステムは一通りプロトタイプを完成した。(2)現在、国立遺伝学研究所の生物遺伝資源情報総合センターアカデミア DNA シーケンシングセンターより、最初のデータをテスト用に受領し、システムおよび運用の評価として、実際の処理を行っている。

また、データベース構築については、前述の通り、現段階では検索キーなどデータの論理設計の段階にある。

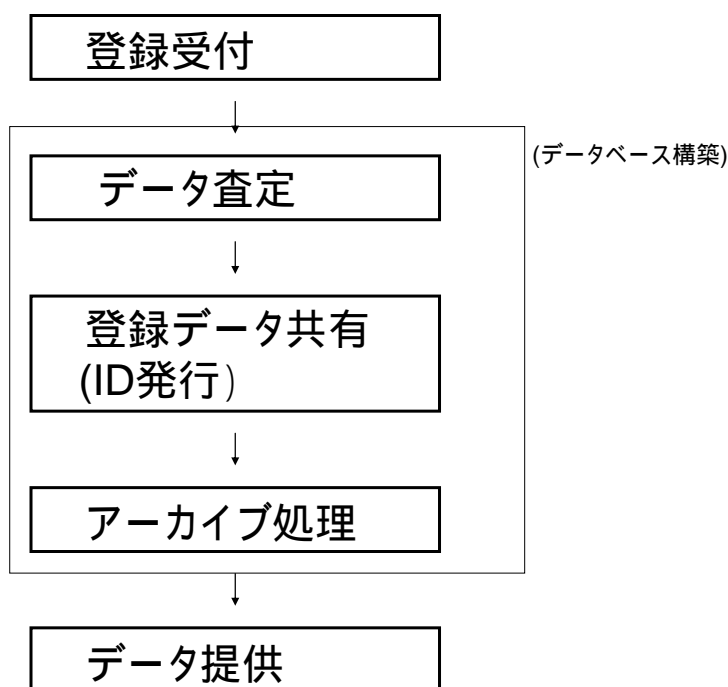


図2 データ受付のための業務フロー

### 2.3 成果の外部への発表

論文寄稿

07.10.25(Epub)

DDBJ with New System and Face (Tateno, Y., Sugawara, H., Ogasawara, O., Okubo, K. and Gojobori, T.) Nucleic Acids Res. 36(Database issue), 22-24

### 2.4 活動(運営委員会等の活動等)

2008年3月 国立遺伝学研究所生物遺伝資源情報総合センターアカデミア DNA シーケンシングセンターへのヒアリング

2008年3月 東京大学ヒトメタゲノムデータの調査

### 2.5 実施体制

別表1の通り。

別表1 平成18年度に於ける実施体制

研究項目	担当機関等	研究担当者
全体構成及び総括	国立遺伝学研究所 生命情報・DDBJ研究センター	五條 堀 孝
データ登録	国立遺伝学研究所 生命情報・DDBJ研究センター	菅原秀明
提供サービス内容	国立遺伝学研究所 生命情報・DDBJ研究センター 国立遺伝学研究所集団遺伝研究系	館野義男 斎藤成也
データ登録・ハードウェア	国立遺伝学研究所 生命情報・DDBJ研究センター	池尾一穂