

「使える」SRA データにすばやくたどりつくために

How to provide quick access to SRA entries of interest



*: nakazato@dbcls.rois.ac.jp

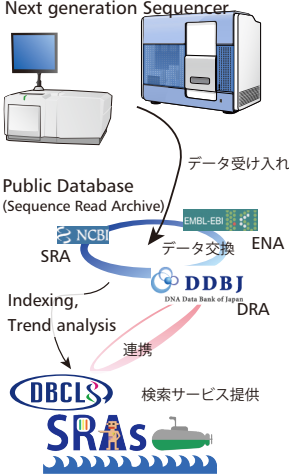
仲里 猛留*, 大田 達郎, 坊農 秀雅
Takeru Nakazato Tazro Ohta Hidemasa Bono

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
Database Center for Life Science (DBCLS), Res. Org. of Info. and Systems (ROIS)

※ デモしますので発表者 までお気軽に。

http://sra.dbcls.jp/

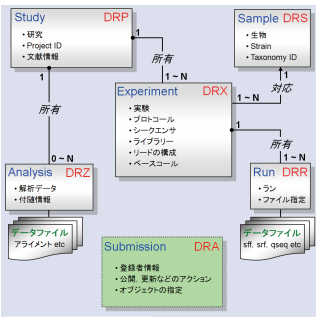
Backgrounds



SRA: the survey of read archives
http://sra.dbcls.jp/

NGS データも、マイクロレイのデータが GEO に登録されるのと同様に、公共データベースである Sequence Read Archive (SRA) に登録され、日本語の 3 局でデータ交換がなされている。その数は、プロジェクト数でおよそ 11300 (2012 年 5 月現在) に及んでいる。DBCLS では、DDBJ と連携して登録データに対して、目次作成、データの傾向分析を行い、NGS データの検索サイトを構築、提供している。

SRA のデータ構造



http://trace.ddbj.nig.ac.jp/dra/documentation.shtml より

Table with 5 columns: Submission, Study, Experiment, Run, Sample, Analysis. It lists various submission IDs and their corresponding counts across different levels, totaling 70003 submissions.

Results and Discussions

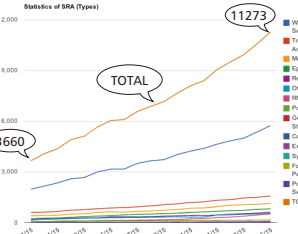
統計値より (2012-5-15 現在)

Big Project List

Table listing various Big Project List entries with columns for ID, SRA ID, Project Name, Type, Date, and other details. Includes projects like The Cancer Genome Atlas Project, Framingham Heart Study, etc.

Study Types

Table showing counts for various Study Types: Whole Genome Sequencing (5741), Transcriptome Analysis (1574), Metagenomics (1143), Epigenetics (819), Resequencing (613), etc.



Species of samples (top 15)

Table listing the top 15 species of samples, including Homo sapiens, Drosophila melanogaster, Escherichia coli, and various other organisms, with their respective sample counts.

文献より

Table of publications using NGS, showing PMID, Title, Journal, Date, and SRA ID. Includes titles like 'Efficient alignment of sequencing reads for re-sequencing applications' and 'A novel and well-defined benchmarking method for second generation read mapping'.

Publications using NGS

Corresponding NGS data

疾患より

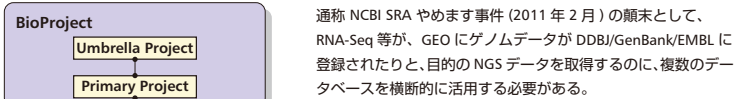
Table of diseases, showing SRA ID, SRA Title, Disease, and PMID. Includes titles like 'DNA sequencing of a cytogenetically normal acute myeloid leukaemia genome' and 'Recurring Mutations Found by Sequencing an Acute Myeloid Leukemia Genome'.

データのクオリティより

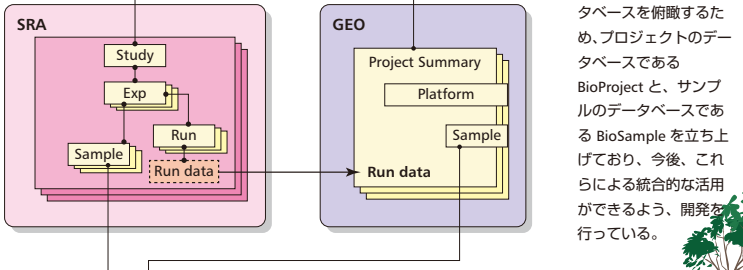
O12 : 公共 NGS データベース Sequence Read Archive におけるシーケンスクオリティによるデータ検索 (大田) 参照

Future works

NGS もデータベースの横断的な活用へ



通称 NCBI SRA やめず事件 (2011 年 2 月) の顛末として、RNA-Seq 等が、GEO にゲノムデータが DDBJ/GenBank/EMBL に登録されたりと、目的の NGS データを取得するのに、複数のデータベースを横断的に活用する必要がある。



3 局では、複数のデータベースを俯瞰するため、プロジェクトのデータベースである BioProject と、サンプルのデータベースである BioSample を立ち上げられており、今後、これらによる統合的な活用ができるよう、開発を行っている。

NGS 現場の会 第 2 回研究会
ホテル板倉エキスポパーク (大阪)
平成 24 年 5 月 23 日~25 日