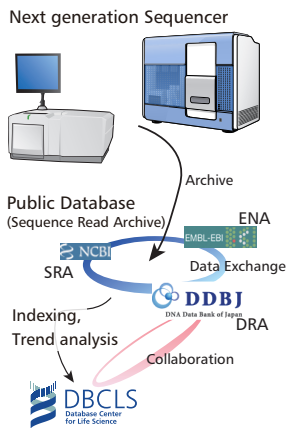


# DBCLS SRA: Functional mining and characterization of public NGS data



**Takeru Nakazato\***, **Tazro Ohta**, **Hidemasa Bono**  
Database Center for Life Science (DBCLS), Research Organization of Information and Systems (ROIS), JAPAN  
\*: nakazato@dbcls.rois.ac.jp

## Backgrounds

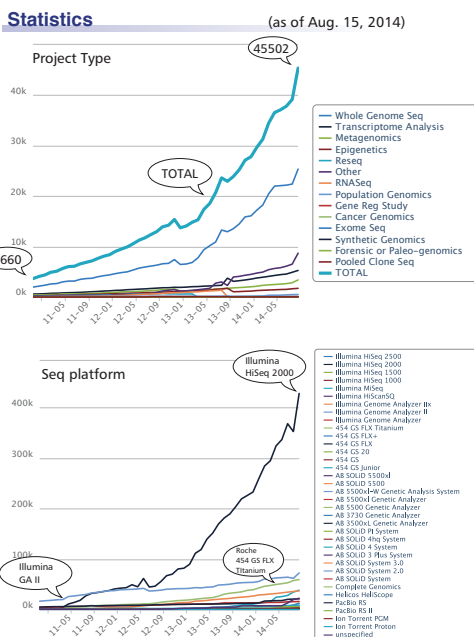


The next-generation sequencing (NGS) data is archived in public database, namely the sequence read archive (SRA), and the data is collaboratively maintained by DDBJ, EBI, and NCBI. In Japan, Database Center for Life Science (DBCLS) has developed infrastructure for researchers to access and re-use these data easily by providing index and stats pages and constructing a portal site for life science databases and tools in collaboration with DDBJ.

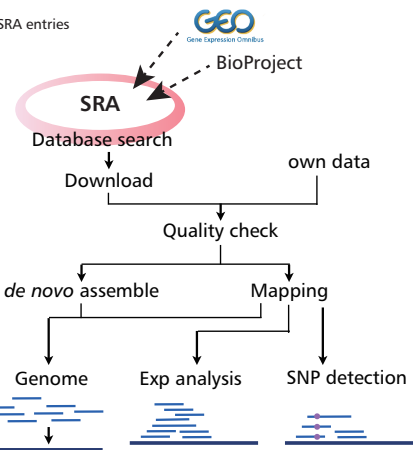
## Solutions

**DBCLS SRA**  
<http://sra.dbcls.jp/>

Exp-design based characterization and quality check of public NGS data  
Search engine for public NGS data  
Also available publication list referring to SRA entries



FREE!

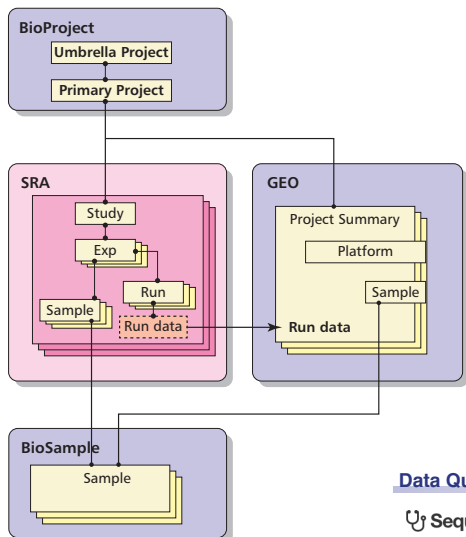


## Publication List

to get NGS data with sufficient good quality to publish articles

Study Type	Filter	Platform	Species	Search
Total: 6198				
SRA ID	SRA Title	Disease	疾病名	PMID
SRA026055	DNA sequencing of a cytogenetically normal acute myeloid leukaemia genome	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	19867736
SRA026055	DNA sequencing of a cytogenetically normal acute myeloid leukaemia genome	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	19657110
SRA026055	Recurring Mutations Found by Sequencing an Acute Myeloid Leukemia Genome	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	19867736
SRA026055	Recurring Mutations Found by Sequencing an Acute Myeloid Leukemia Genome	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	19657110
SRA009887	In-depth characterization of the microRNA transcriptome in a leukemia progression model	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	18849523
SRA029797	Exome Sequencing Identifies Somatic Mutations in Acute Monocytic Leukemia	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	21396634

## The data structure of SRA



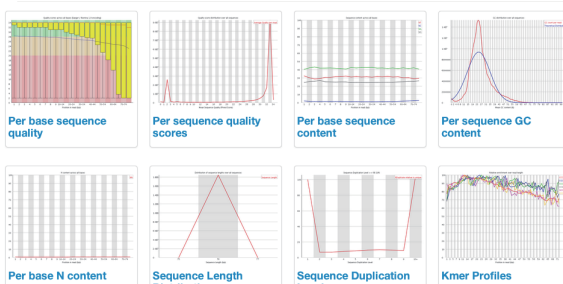
Species (top 15)	Count
Homo sapiens	262584
Mus musculus	50126
human gut metagenome	21731
human metagenome	18150
Plasmodium falciparum	17468
Staphylococcus aureus	15843
Streptococcus pneumoniae	15374
Saccharomyces cerevisiae	12545
soil metagenome	11750
Drosophila melanogaster	10728
Danio rerio	10293
Anopheles gambiae	8231
Mycobacterium tuberculosis	8086
rhizosphere metagenome	7727
Caenorhabditis elegans	7392
Total	759598 (experiments)

## NGS data relevant to diseases

SRA ID	SRA Title	Disease	疾病名	PMID
SRA026055	DNA sequencing of a cytogenetically normal acute myeloid leukaemia genome	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	19867736
SRA026055	DNA sequencing of a cytogenetically normal acute myeloid leukaemia genome	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	19657110
SRA026055	Recurring Mutations Found by Sequencing an Acute Myeloid Leukemia Genome	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	19867736
SRA026055	Recurring Mutations Found by Sequencing an Acute Myeloid Leukemia Genome	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	19657110
SRA009887	In-depth characterization of the microRNA transcriptome in a leukemia progression model	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	18849523
SRA029797	Exome Sequencing Identifies Somatic Mutations in Acute Monocytic Leukemia	Leukemia, Myeloid, Acute	白血球急性骨髄性白血病	21396634

## Data Quality

### Sequence Quality Statistics



- The results of FastQC is pre-calculated and presented in detail view.  
- Researchers are not only spared the problem of downloading low quality data but also able to skip checking process of data quality.

Recently, some NGS run data is not deposited to SRA but GEO and GenBank. Moreover, project and sample information is captured in external other databases: BioProject and BioSample. This complexity of data structure prevents researchers from retrieving the NGS data of interests.

Reference:  
Experimental design-based functional mining and characterization of high-throughput sequencing data in the Sequence Read Archive.  
Nakazato T., Ohta T., Bono H., PLOS One, 8 (10): e77910 (2013) PMID: 24167589

**ECCB'14**  
(European Conference on Computational Biology)  
Strasbourg Convention Centre, Strasbourg, France  
Sep. 7-10, 2014