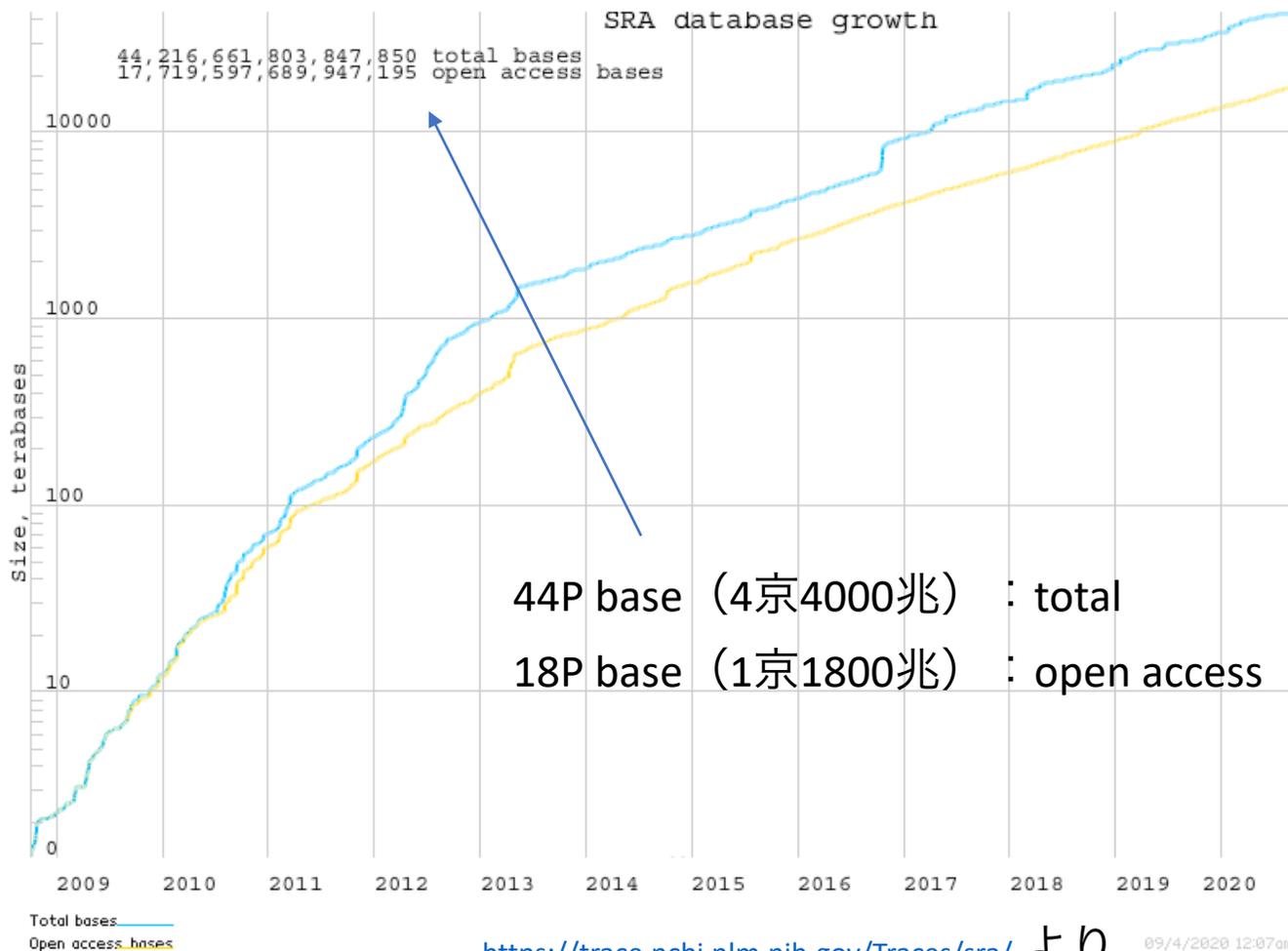○仲里　猛留、大田　達郎

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 ライフサイエンス統合データベースセンター
（DBCLS）

# NGSデータの現状：登録塩基数
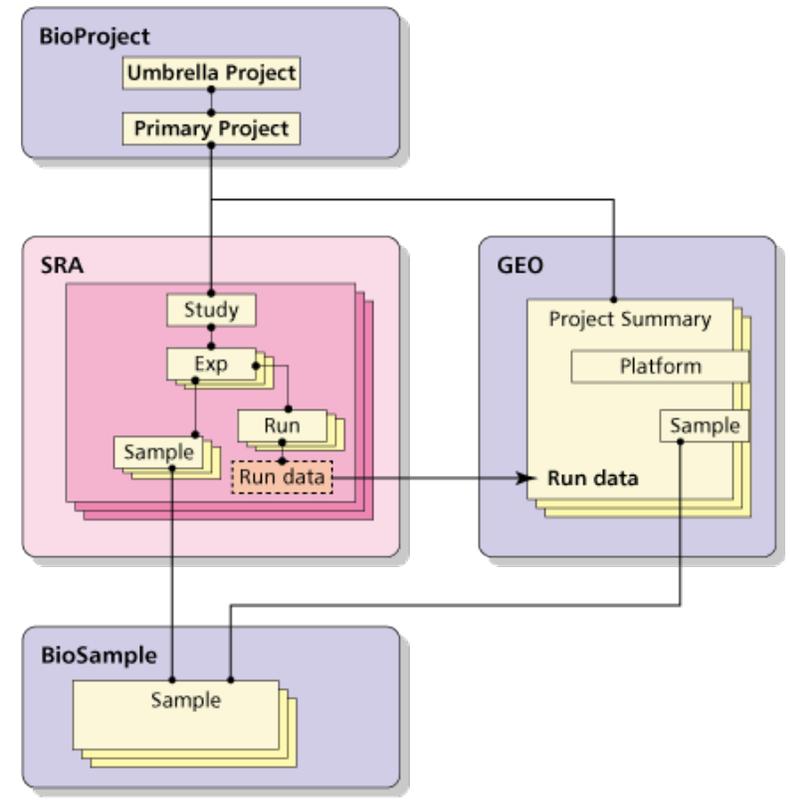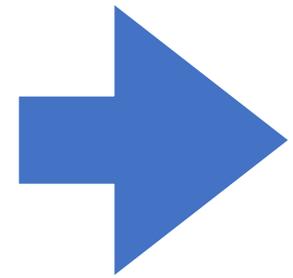
SRA database growth

44,216,661,803,847,850 total bases
17,719,597,689,947,195 open access bases

44P base（4京4000兆）：total
18P base（1京1800兆）：open access

Total bases
Open access bases

[https://trace.ncbi.nlm.nih.gov/Traces/sra/](https://trace.ncbi.nlm.nih.gov/Traces/sra/) より

09/4/2020 12:07am

2007年にNCBIがSRAを公開して公共NGSデータを
収集し始めて以来、NGSの普及とともに
データ量も爆発的に増加した。
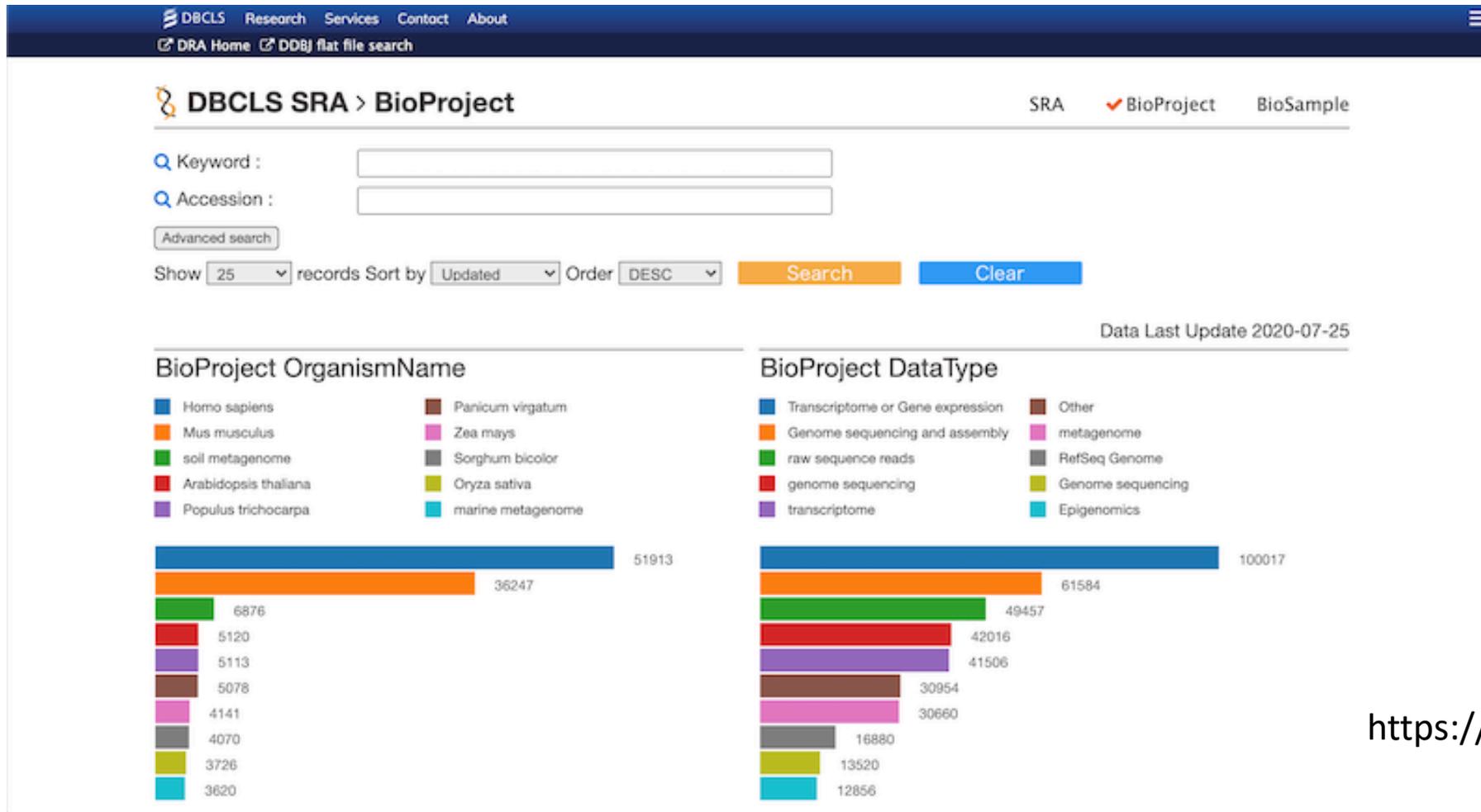
NGSデータはその目的の多様さからSRAだけでなく、GEO（transcriptome）やNucleotide（genomeなど）にも登録されるようになったため、BioProject/BioSampleのデータベースを作成し共通で参照できるようにした。

https://www.ddbj.nig.ac.jp/dra/submission.html を改

# DBCLS SRA

我々はこれまで
DBCLS SRAとしてSRA内の
公共NGSデータに対する
検索エンジンを提供して
きたが、今回 新たに
BioProjectやBioSampleの
データを検索できる
ようにした。

https://sra.dbcls.jp/

# DBCLS SRA・リスト画面

リスト画面



詳細画面

これまでのDBCLS SRAでの公共NGSデータ検索と同様にBioProjectデータを検索し、そのリスト表示や個々の登録の詳細の確認などを行うことができる。

# NGSデータの現状：目的別

各年の登録数（累積ではない！）

登録数の現状

BioProjectについての統計情報も提供している。

図は、目的別の各年と全体の目的別登録数。全体としてはtranscriptomeの登録が多い。

# NGSデータの現状：platform別

登録数の現状

SRAも公開から10年以上経ち、その間に多くの
シーケンサーが出現し、そして消えていった。
数年前まではIllumina HiSeq 2000がSRAデータで
1位を占めていたが、今はMiSeqが1位となっている。