

大量高速シーケンサー由来の配列データベース

454、Illumina、SOLiD などからのデータを登録



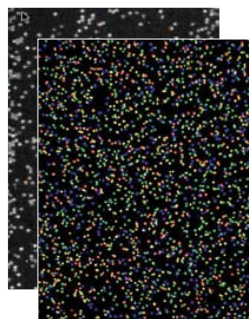
454 (Roche)



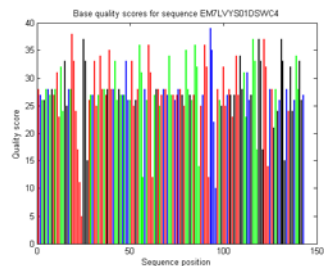
Solexa (Illumina)



SOLiD (ABI)



画像データ



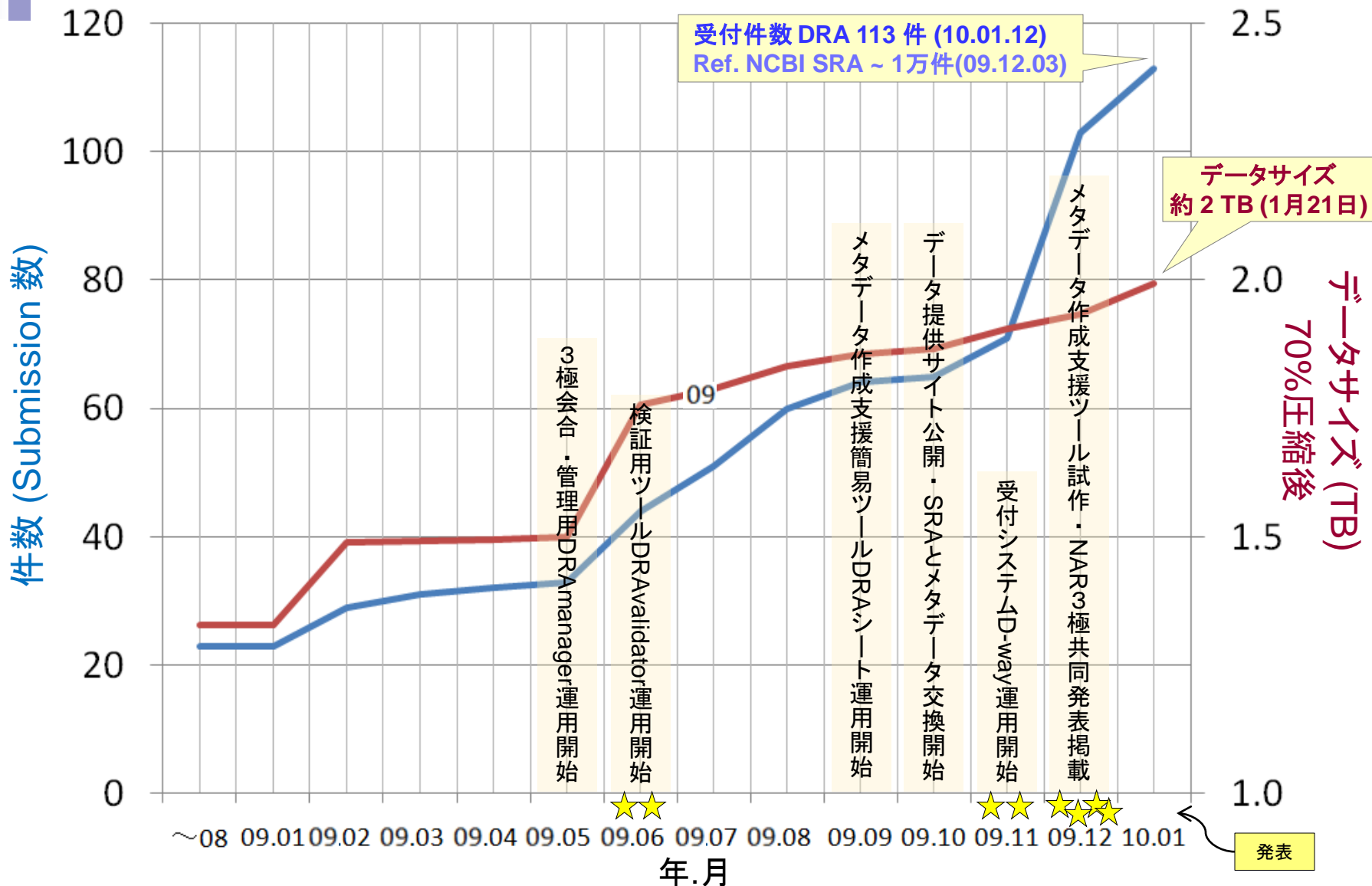
計測データ
 シグナル強度
 ノイズ
 ベースコール
 Quality value 等

登録対象

```
@SRR001654.1 9460:7:1:830:763 length=36
GTCAATATTAATCATACCAATATACTCAAAAAATAA
+SRR001654.1 9460:7:1:830:763 length=36
I+&*4)%+5#%#)&$%$#%#&%%"$%#%!"
@SRR001654.2 9460:7:1:402:781 length=36
GGTCTAAAAGCAAATTCAGTCTTCAAATAATTC
+SRR001654.2 9460:7:1:402:781 length=36
II+(%$+%&+*-0+/*("%&+"*%#"%%&$
@SRR001654.3 9460:7:1:433:775 length=36
GTGCTTTTTTTTTCCAGGAAGTTGTCTCCTCTATC
+SRR001654.3 9460:7:1:433:775 length=36
II3DI>IIIIIIIB7.,&%&)." +,,$&$&"%#
```

fastq データ
 塩基配列
 Quality Value

DRA開発と受付の年間動向



登録機関数 13機関
国内 11機関
国外 2機関(中国、タイ)

公開件数 18件
国内 17件
国外 1件(中国)

(注)ここまでの件数は、メタデータの中でいうところのSubmissionが単位。
メタデータの中でいうところのrunを単位とすると、登録件数が510件、
公開件数が91件となる。

2009年度はDDBJ Read Archive(DRA)の運用を開始した。すなわち、第2世代シーケンサ由来のデータを受付けし、受付けたデータにIDを発行し、さらに、利用者へ提供する枠組みを実現した。また、NCBIのSRAならびにEBIのERAとの連携も具体化した。

2010年度は、このDRAの受付システム、公開システムならびに管理システムの高度化を進める。また、SRAならびにERAと協調して、研究社会へのデータ登録への理解を広げるべく努め、さらに、第3世代シーケンサからのデータやそれを利用した新たな研究分野の成果への対応についても取り組む。

これによって、技術進歩の歩みが止まらないシーケンサから生成され続ける塩基配列リードに対応し、また、高速に入手できる大量のデータを活用して進化し続ける研究動向に対応可能なDRA基盤を開発しつつ安定に運用する。

① 受付システムの高度化と運用

DRAの受付システムは、受付け窓口に対応するD-way、データの正当性を検証するDRAvalidator、メタデータ作成を補助するGUIならびにショートリードを解析するパイプラインからの出力を受付けるインターフェースといったサブシステムで構成される。2009年度に構築したこれらのサブシステムにおいて、検証機能を向上させ、登録者に対する利用者インターフェースを直感的に分かりやすいものにし、さらに、受付査定業務の効率化を図る。ま

た、NCBIで開発が進められているSRA toolkitとの整合性も確保し、同時に、fastqデータより大規模になるが豊富な情報を含むSRAファイルにも対応する。

② 公開システムの高度化と運用

2009年度のデータ提供は限られた項目で検索しヒットしたデータセットをダウンロード可能にしたものであったが、2010年度には、メタデータに含まれる多様な項目による検索を可能として、その検索結果にヒットしたデータセットをダウンロード可能とする。特に、SRAファイルの提供に取り組む。

③ 管理システムの高度化と運用

2009年度はログファイルなどから手作業で集計していた各種統計情報の自動生成、データのバージョン管理、DRA側でのデータ検証、データ更新の自動反映、公開日管理の機能をシステム化して、得られる統計情報をシステム改善に活かしていく。

④ 国際連携

国際塩基配列データベースの国際実務者会議の機会などを生かして、SRAならびにERAと意識と情報を共有することによって、日米欧3極で整合性のあるリードアーカイブ構築を目指す。また、塩基配列データ登録の事例に倣って、シーケンサからの出力データをいずれかのアーカイブに登録するルール採用を、3極共同で学術雑誌などに働きかける。

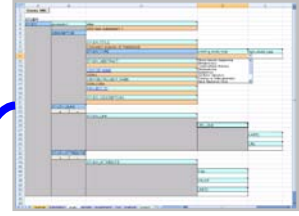
① 受付システムの高度化と運用

開発済

今年度からの
開発継続

2011年度開発

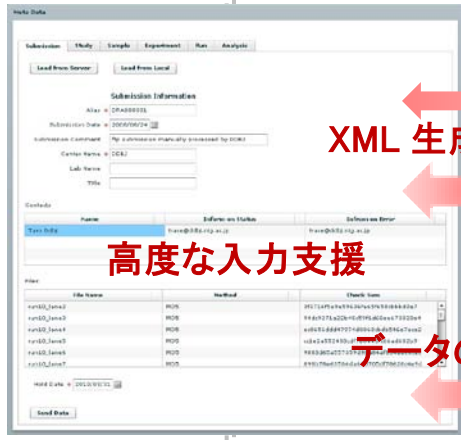
“DRA シート”



XML ファイル



“ウェブ GUI ツール”



結果
XML 生成/読み込み
チェック

“DRValidator”

“登録受付システム D-way”

mfujimot | New submission | Account | Password | Logout

D-way 登録履歴

DRA submission list for mfujimot

Submission ID	Accession	Study Title	Status	Creation Date
mfujimot-0005	---	---	new	2009-11-13
mfujimot-0004	---	High-resolution profiling of histone methylations in the human genome	sheet_uploaded	2009-11-12
mfujimot-0003	---	High-resolution profiling of histone methylations in the human genome	metadata_validated	2009-11-12
mfujimot-0002	---	High-resolution profiling of histone methylations in the human genome	metadata_validated	2009-11-12
mfujimot-0001	---	High-resolution profiling of histone methylations in the human genome	complete (public)	2005-10-14

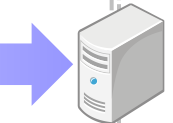
データの再利用

メタデータ

ランデータ

D-way サイトからアップロード

サーバにファイル転送



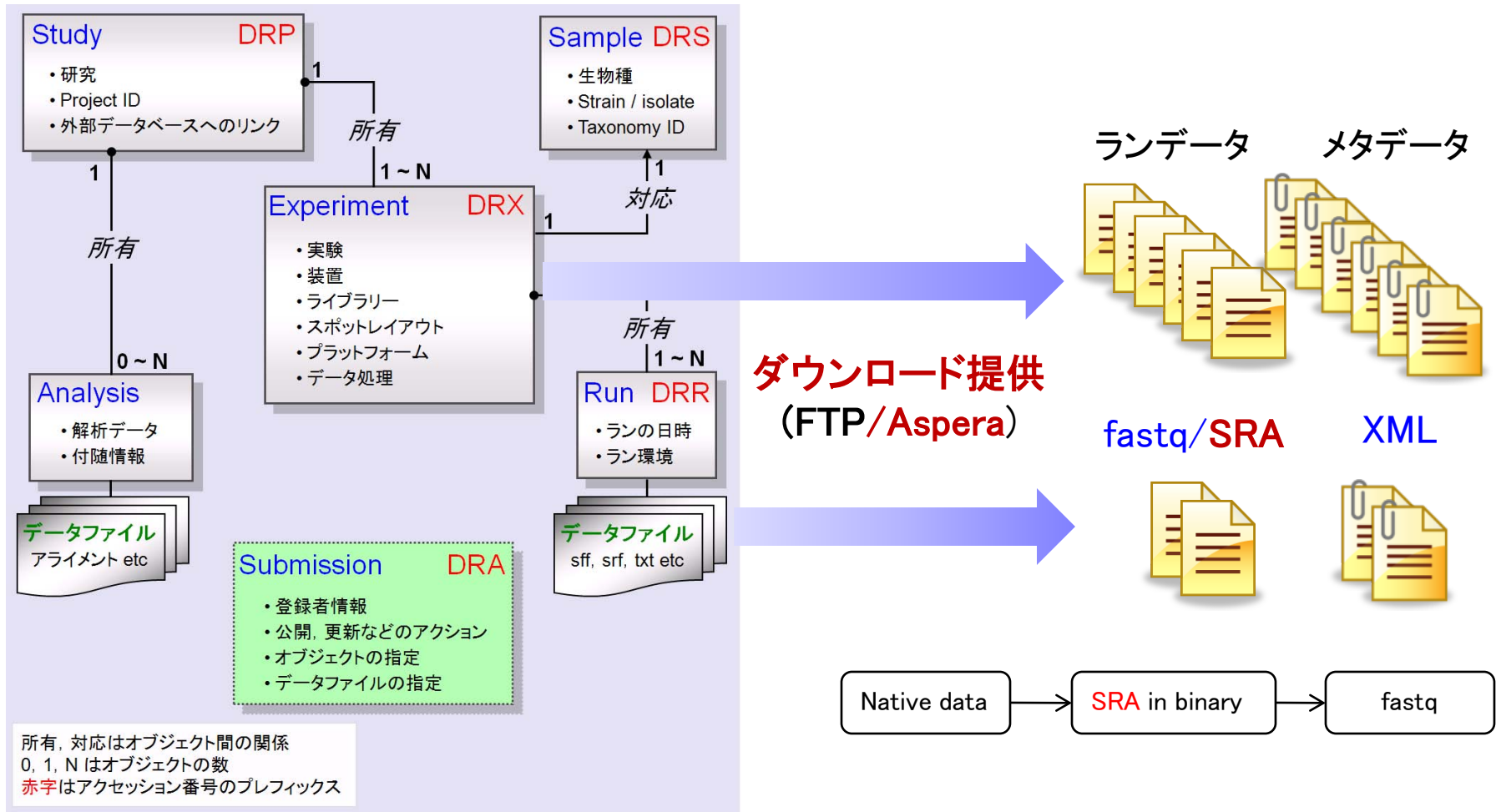
受付サーバ

“SRA Toolkit”
ファイル形式変換

解析パイプライン

② 公開システムの高度化と運用

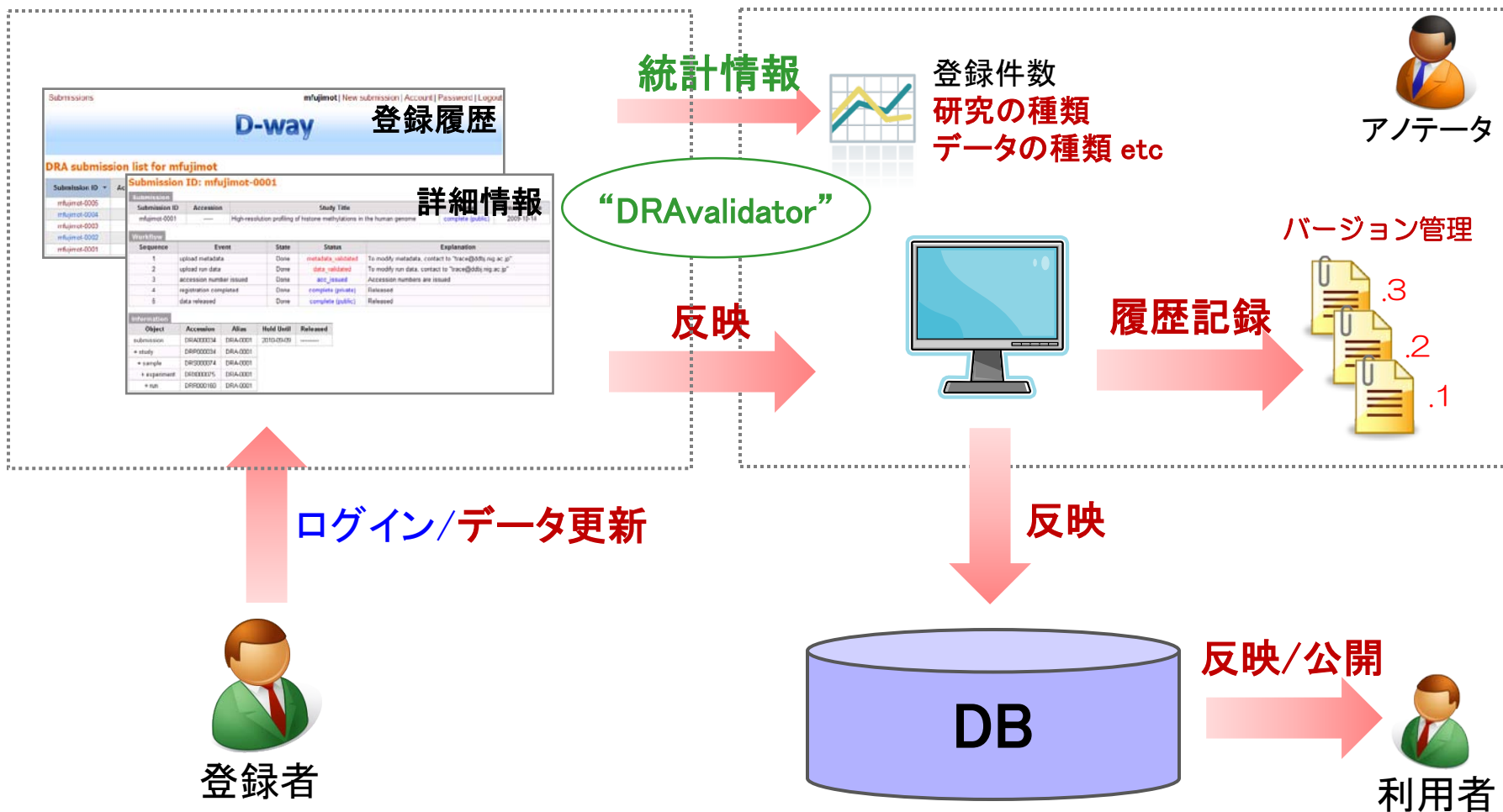
メタデータの検索



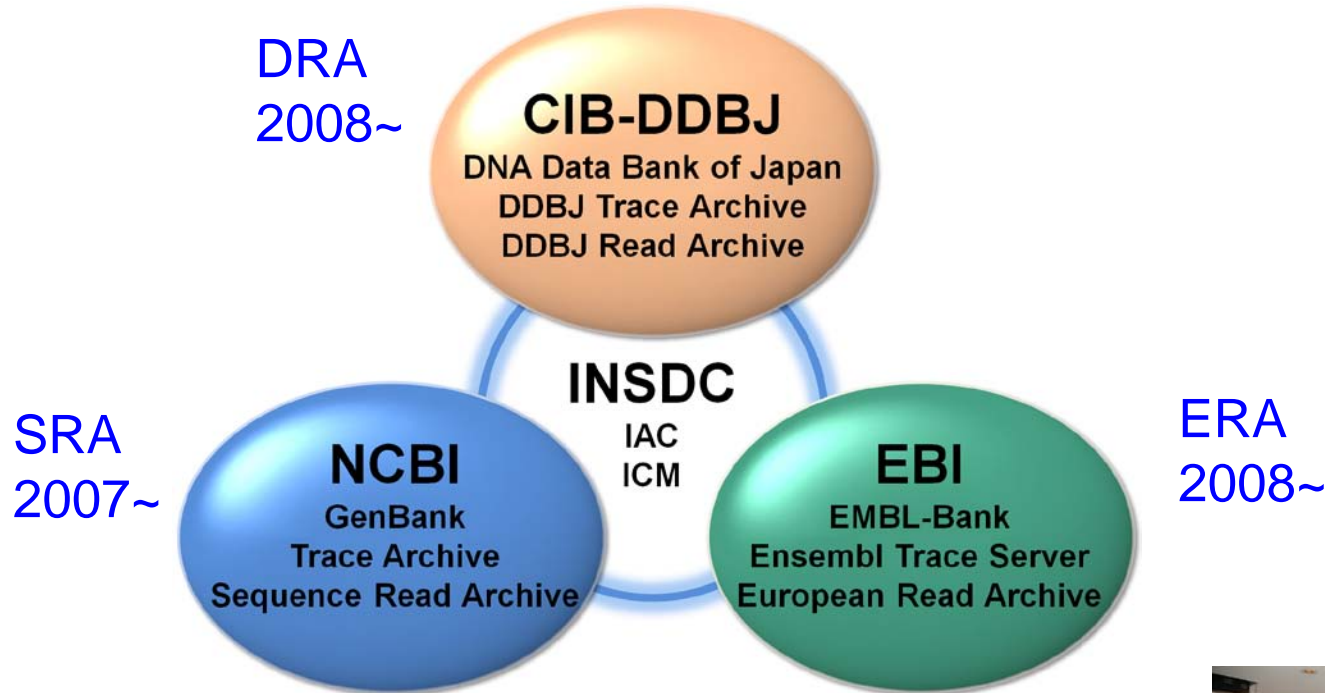
③ 管理システムの高度化と運用

“登録受付システム D-way”

“管理システム DRM”



④ 国際連携



第1回 DRA/ERA/SRA 国際会議 2009.5.14-15 at NCBI

DRA/ERA/SRA 3極共同で国際アーカイブを構築することで合意

Shumway M, Cochrane G, Sugawara H (2010)

“Archiving next generation sequencing data.” *Nucleic Acids Res* 38 (Database issue) :D870-871

第2回 2010.5.20-21 at EBI 予定



国内におけるシーケンサーの導入と稼働の状況を把握しつつ、ハードウェアとソフトウェアのスケールビリティを継続的に検証

おおよその国内設置台数(一部見込みを含む)

Illumina	～ 60台
Illumina新型機	～ 3台
SOLiD	～ 20台
454	～ 20台
Heliscope	～ 4台

DDBJにおけるDRAの位置付け

