# Using the NCBI Short Read Archive (SRA) as a Scientific Resource

Martin Shumway, Eugene Yaschenko, Kurt Rodarmer, Michael Kimelman, and James Ostell

National Center for Biotechnology Information (NCBI)
National Library of Medicine (NLM)
Bethesda MD USA

**NCBI**

*The Trace Archives*

## Current Holdings



**A** Coverage of 169 Complete Human Genomes (log scale)

**B** Bases Archived — Data through 31 Jan 2009. Includes deposits in European Read Archive mirrored to NCBI.

**C** Total Bases by Sequencing Study Type — Data through 31 Jan 2009. Includes deposits in European Read Archive mirrored to NCBI.

Human whole genome sequencing is the main activity.

- Whole Genome Sequencing
- Resequencing
- Epigenetics
- Transcriptome Analysis
- Other
- Metagenomics
- Gene Regulation Study
- Cancer Genomics

**D** Bases (species) by Taxonomic Division — Data through 31 Jan 2009. Includes deposits in European Read Archive mirrored to NCBI.

Diptera (Drosophila) (4) 52691194344 1%
human 4.69412E+12 96%

- human
- Diptera (Drosophila) (4)
- Bacteria + Archea (316)
- mouse
- Eukaryotes - Other Vertebrates (15)
- Fungi (3)
- Eukaryotes Other Invertebrates (15)
- Metagenomes (10)
- Protists (6)
- Caenorhabditis
- Viruses

**Figure 1**: Deposits in the SRA as of 31 Jan 2009. (A) lists 169 human genomes sequenced to 1X or greater. Most of these were generated by the 1000 Genomes Project. (B) shows the growth in tera bases of deposited sequence. (C, D) show the breakdown of sequencing by study type and taxonomy.

## Abstract

NCBI announces production release of the Short Read Archive (SRA), a repository of sequencing data from next generation sequencing platforms. Thus far the SRA has amassed over 4 TBp (Tera bases) of sequence from several thousand submissions executed via automated and interactive channels. The deposits are now indexed in the Entrez system, meaning that short read datasets can be retrieved via queries to PubMed, Genome Project, Taxonomy, and GEO. Short read data are provisioned in fastq format (bases and qualities), suitable for use with popular assemblers and aligners. The SRA regularly exchanges data with the European Read Archive (ERA) and the DNA Data Bank of Japan, providing users with worldwide scope for short read deposit and retrieval.

The SRA continues as a work-in-progress. Future developments include on-line sequence search capability, archival of short read alignment and assembly records, improved indexing and retrieval, a software library for interacting with archival data sets, and support for new platforms and sequencing technologies.
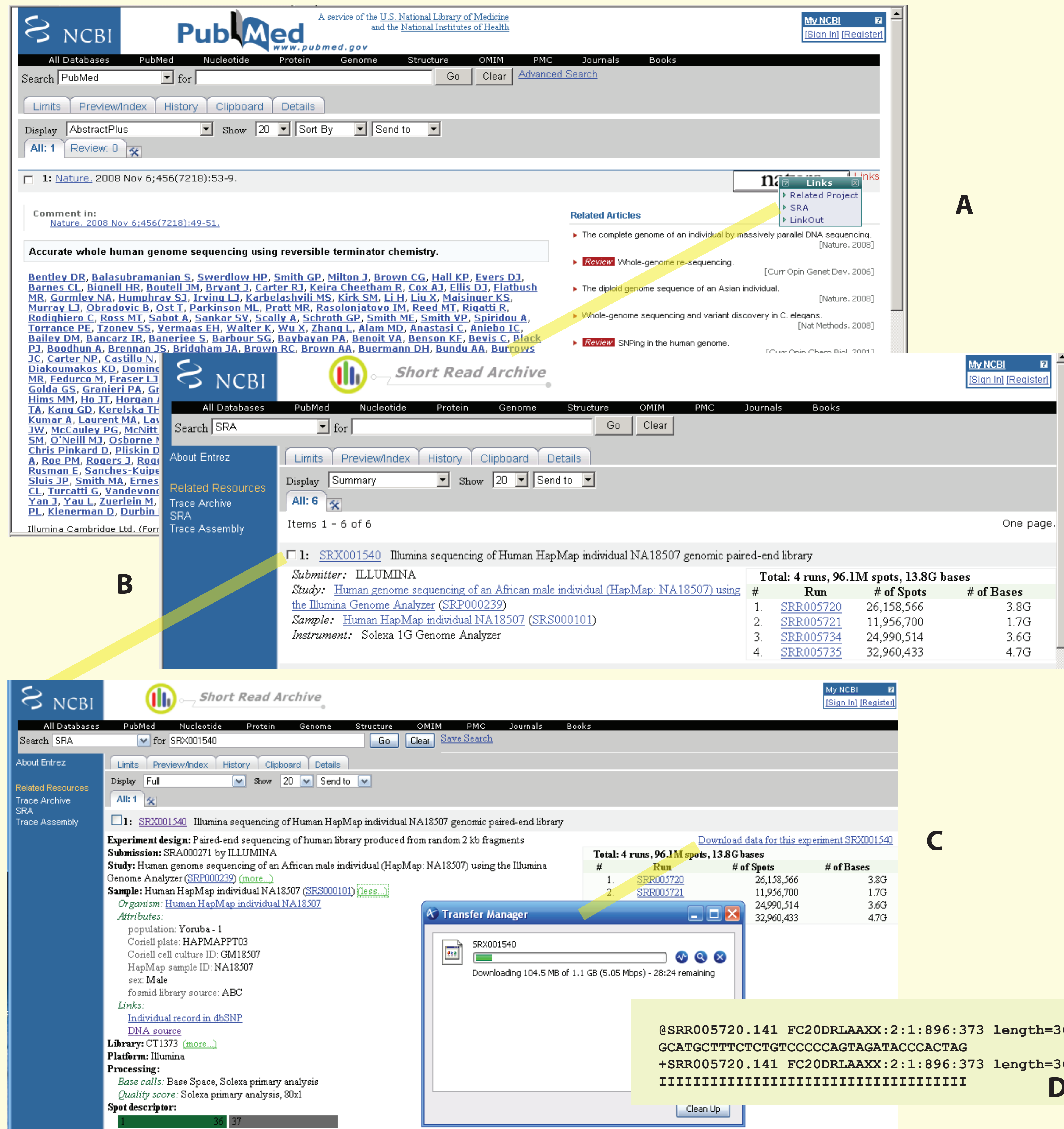
## Retrieval Workflow



**Figure 2**: A user finds an article of interest in PubMed. The user follows the SRA link (A), obtains a list of experiments run for the study (B), expands the detail of one such experiment (C), and downloads the data in fastq format (D).
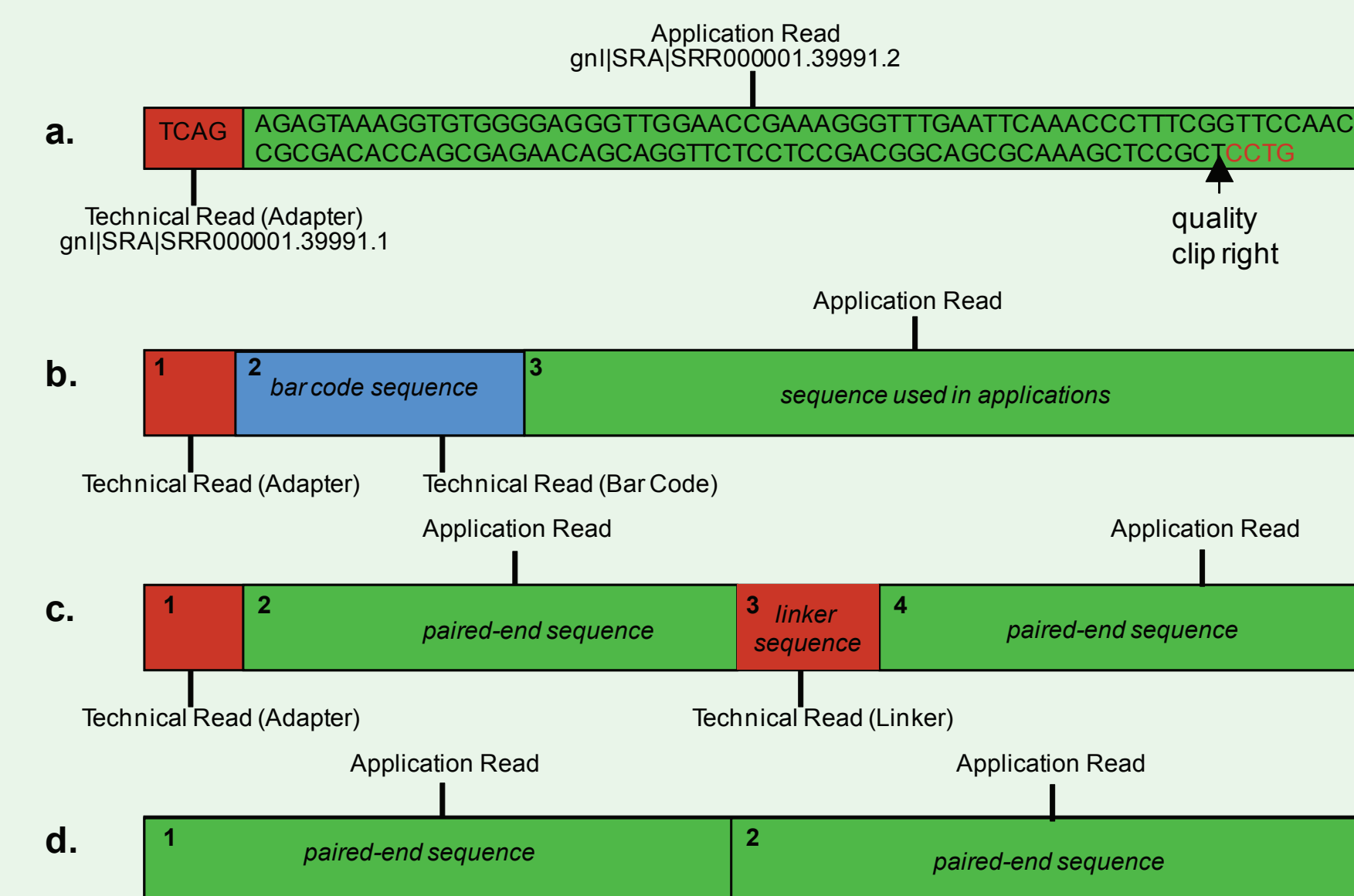
## Spot Abstraction



**Figure 3**: Short reads are extracted from an image cluster, or "spot", and partitioned during processing to present subsequences to downstream applications. The SRA presents the partition among "application" and "technical" reads of the spot's sequence. Annotations such as 3' quality clipping are stored as properties of the entire spot (**a**). Alternative spot layouts should be supported this way: barcoded reads (**b**), paired-end sequence with linker (**c**), and paired-end sequence without linker (**d**).

## Data Model



a. Main Objects    b. Helper Objects    c. Pooled Sample Scenario    d. Multiplexed Sample Scenario
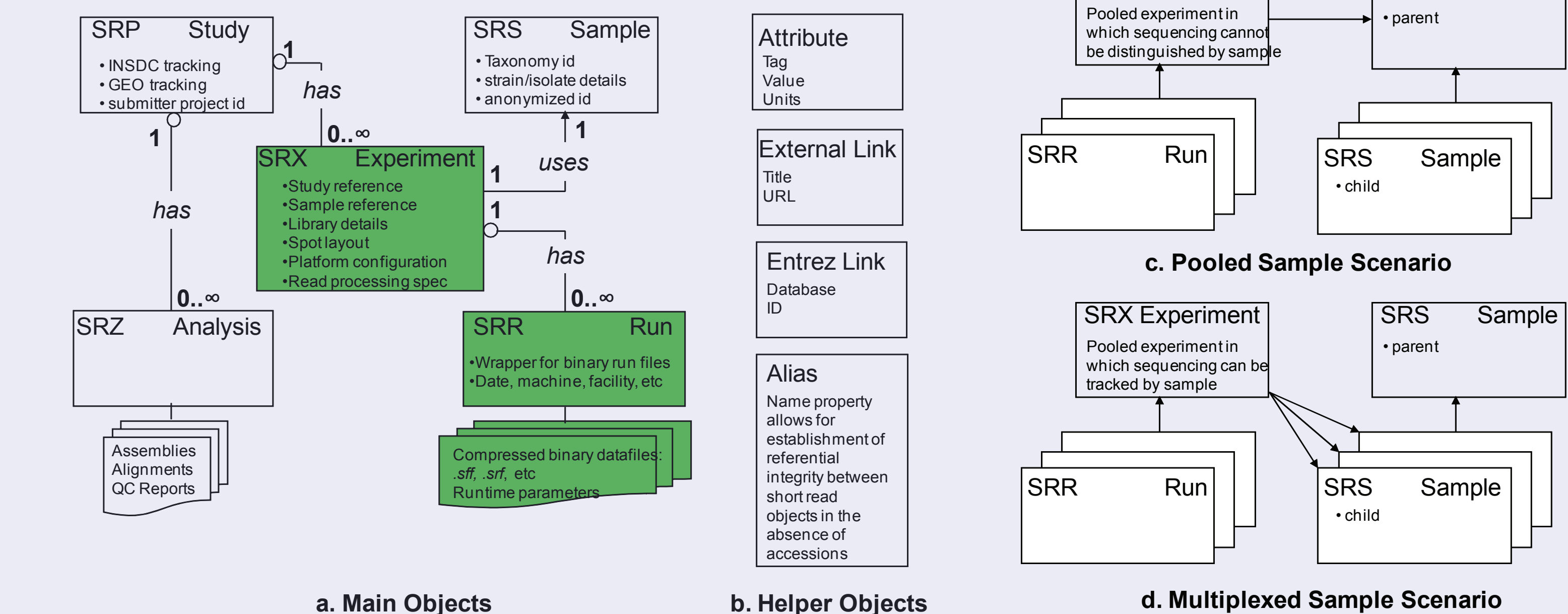
**Figure 4**: The SRA separates metadata from data (a), relates experiments to samples in flexible ways (c,d), and provides flexible data structures for decorating objects with properties, internal and external links, and names (b).

## Submitting to SRA

**Submissions Model**
- High Throughput Submission with XML
- Interactive Submissions Tool
- Consultative Submissions
- Submitter maintenance of metadata
- Sequencing Centers vs Individual Submitters

A distinct accession (SRA) is issued just for the submission session, separating submissions contact info, file manifests, exceptions from the data and metadata. Metadata can be submitted separately from data, and data and metadata for a given submission can arrive asynchronously.

**Hold Until Publish Feature**
- Hold for [days]
- Hold until [date]
- Hold [until released]
- Hold for broker [broker] = poll broker authority

**Restricted Access Feature**
- Metadata are public
- Run data are private
- Appropriate for patient data, personal genomics
- Brokered by dbGaP, which manages access

## Deposit Network



INSDC Shared Accession Space

DDBJ    SRA    ERA
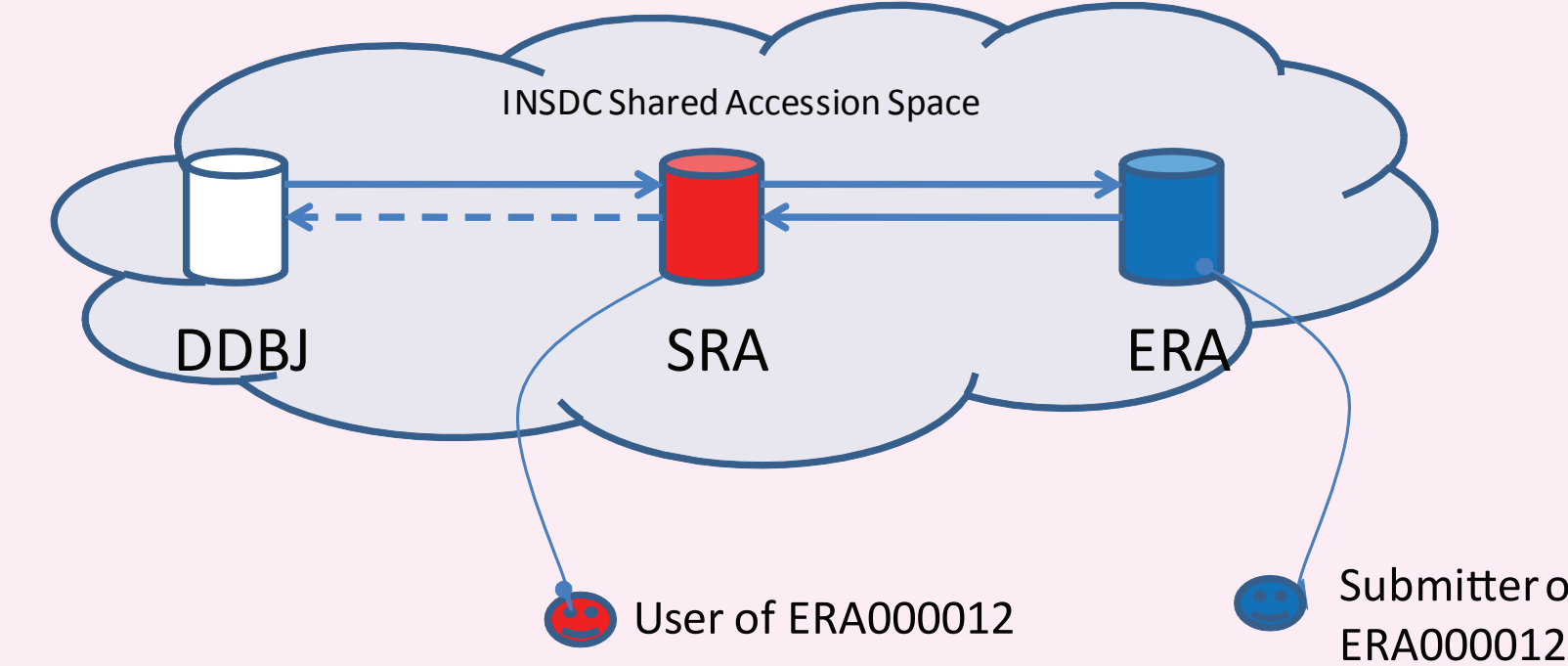
User of ERA000012    Submitter of ERA000012

**Figure 5**: Deposits are currently brokered by DNA Data Bank of Japan (DDBJ) to NCBI, and are mirrored between the SRA at NCBI and the European Read Archive (ERA) at European Bioinformatics Institute (EBI). This gives users across the world immediate access to deposits at the respective archives.

A global name space ensures that accessions generated by the individual archives are recognized by the collaborating archives.