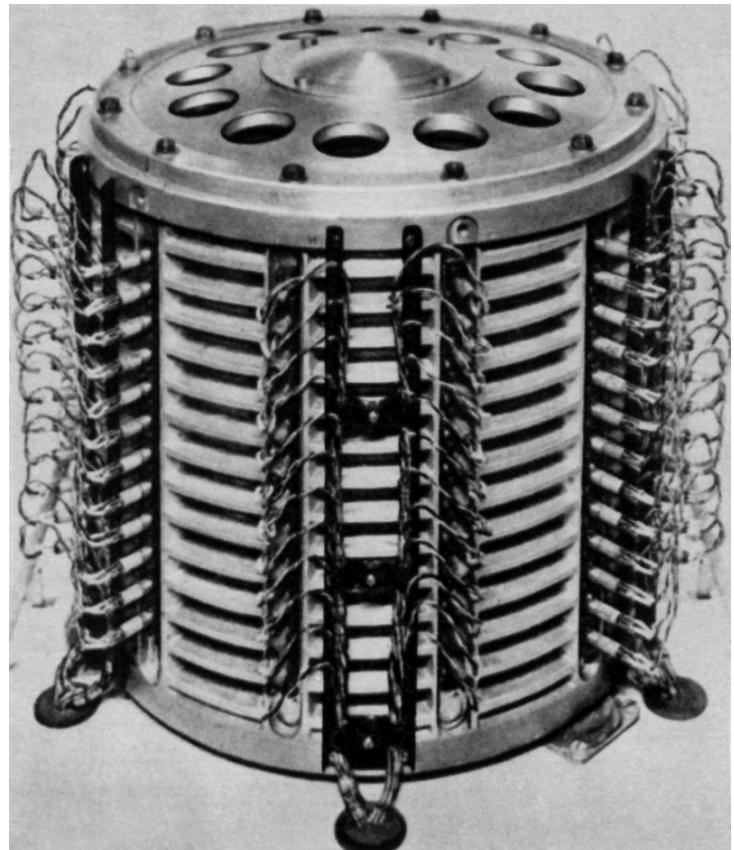


# Online Resources for Nucleotide Data

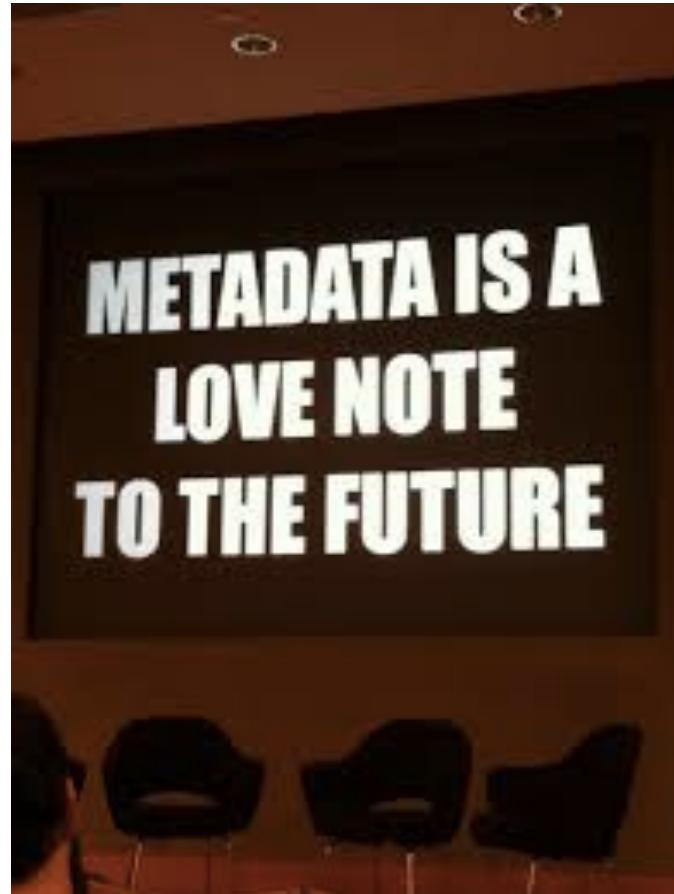
# Why do we need databases?

- To archive and preserve information
- To put all the things in one place (facilitate discovery)
- To enforce and maintain format standards
- To allow reuse of data (its expensive, use it more than once!)
- To prevent fraud in research
- To have reproducibility of research
- To store metadata



# Metadata

- Metadata is data about data
- Where did the data come from?
  - Organism or substrate, experimental conditions, location, time and date, tissue
- How was the data collected?
  - Field methods, lab methods, instruments, calibration
- How has the data been processed?
  - Normalization, removal of “bad” data, any processing at all
- User rights and management for the data
  - Is this open for additional publication or is it embargoed?
  - Does it carry a license?



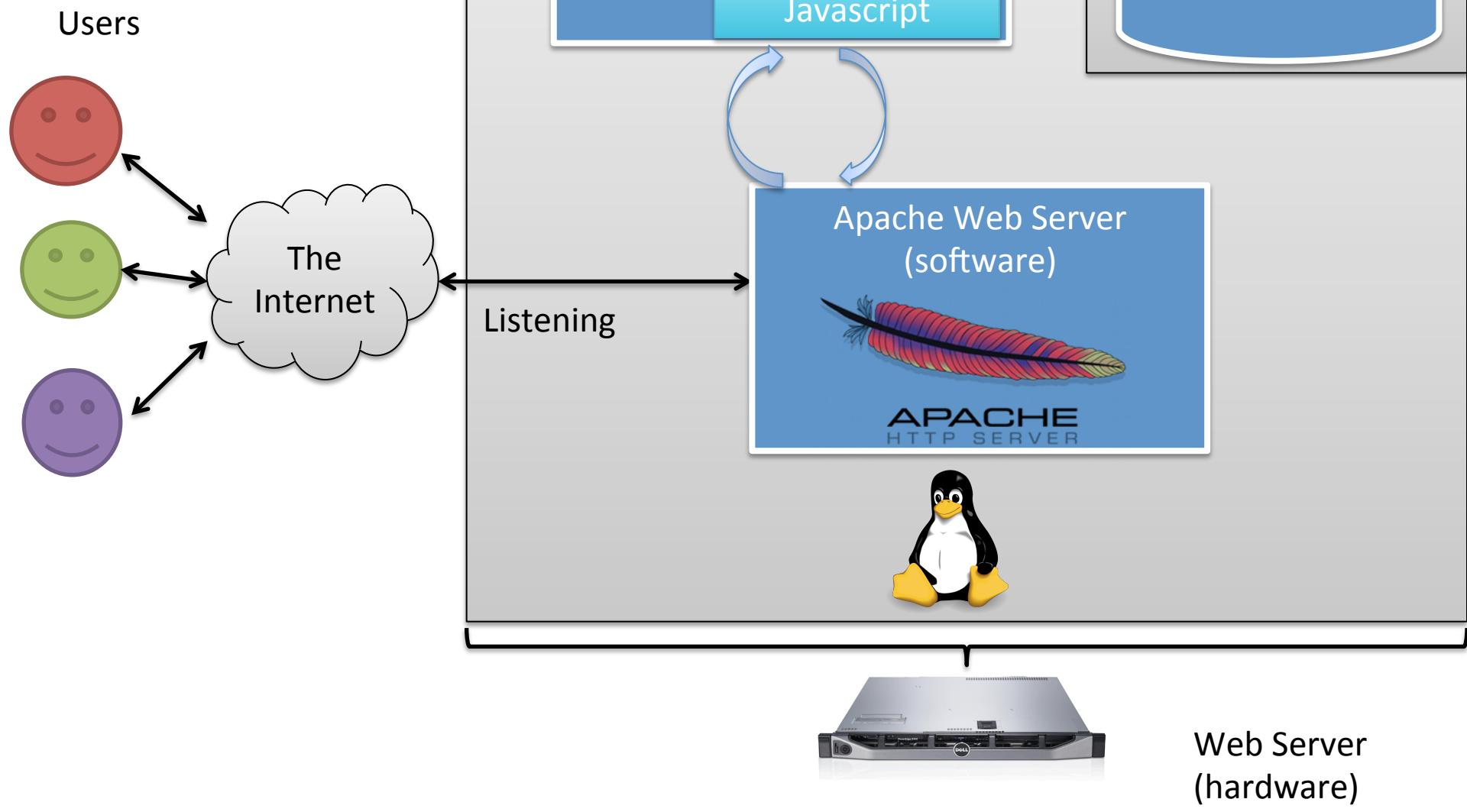
# A Few Types of Databases

- I. International, Primary Repositories
- II. Protein DB resources
- III. Community DBs

These are quite possibly totally unrelated to your data of interest. Publications and internet searches will help you identify the right database for you.

Sometimes there just isn't a home. Non-human metabolomics data?

# The Structure of a Website and Database



# International Nucleotide Sequence Database (INSD)

- Consists of the following 3 dbs:
  - DDBJ (DNA Data Bank of Japan)
  - EMBL (European Molecular Biology Laboratory)
  - NCBI (National Center for Biotechnology Information)
- repositories for nucleotide sequence data from all organisms
- all three databases accept nucleotide sequence submissions, and then exchange new and updated data on a daily basis
- Primary database = house original sequence data



# NCBI

- Discover
- Download
- Submit
- Analyze

The screenshot shows the NCBI homepage. At the top left is the NCBI logo and the text "National Center for Biotechnology Information". A dropdown menu "All Databases" is open. On the left is a vertical navigation bar with a blue header "NCBI Home" containing a list of links: Resource List (A-Z), All Resources, Chemicals & Bioassays, Data & Software, DNA & RNA, Domains & Structures, Genes & Expression, Genetics & Medicine, Genomes & Maps, Homology, Literature, Proteins, Sequence Analysis, Taxonomy, Training & Tutorials, and Variation. To the right of the navigation bar is the main content area. At the top right of the content area is the heading "Welcome to NCBI". Below it is a brief description: "The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information." Below this are links to "About the NCBI", "Mission", "Organization", and "NCBI News". The main content area is divided into several sections: "Submit" (Deposit data or manuscripts into NCBI databases, with an upward arrow icon), "Download" (Transfer NCBI data to your computer, with a downward arrow icon), "Learn" (Find help documents, attend a class or watch a tutorial, with a book icon), "Develop" (Use NCBI APIs and code libraries to build applications, with a square icon), "Analyze" (Identify an NCBI tool for your data analysis task, with a scatter plot icon), and "Research" (Explore NCBI research and collaborative projects, with a microscope icon).

<http://www.ncbi.nlm.nih.gov/>

# Discover “Entrez”

<http://www.ncbi.nlm.nih.gov/gquery/>

**Search NCBI databases** [Help](#)

Fraxinus [X](#) [Search](#)

Results found in 22 databases for "Fraxinus"

Literature			Genes		
Books	7	books and reports	EST	12,100	expressed sequence tag sequences
MeSH	2	ontology used for PubMed indexing	Gene	5	collected information about gene loci
NLM Catalog	0	books, journals and more in the NLM Collections	GEO DataSets	4	functional genomics studies
PubMed	654	scientific & medical abstracts/citations	GEO Profiles	0	gene expression and molecular abundance profiles
PubMed Central	1,060	full-text journal articles	HomoloGene	0	homologous gene sets for selected organisms
<b>Health</b>			PopSet	240	sequence sets from phylogenetic and population studies
ClinVar	0	human variations of clinical significance	UniGene	0	clusters of expressed transcripts
dbGaP	0	genotype/phenotype interaction studies	<b>Proteins</b>		
GTR	0	genetic testing registry	Conserved Domains	0	conserved protein domains
MedGen	6	medical genetics literature and links	Protein	809	protein sequences
OMIM	0	online mendelian inheritance in man	Protein Clusters	0	sequence similarity-based protein clusters
PubMed Health	0	clinical effectiveness, disease and drug reports	Structure	0	experimentally-determined biomolecular structures
<b>Genomes</b>			<b>Chemicals</b>		

# Download

<http://www.ncbi.nlm.nih.gov/home/download.shtml>

## Download

The majority of NCBI data are available for downloading, either directly from the NCBI FTP site or by using software tools to download custom datasets.



## **ADDITIONAL LINKS**

## How to download custom data sets

Large Data Download Best Practices

SRA Download Reference

FTP

Download data from the NCBI  
FTP site

Aspera

High-speed downloads provided by Aspera software

## Download Tools

## Tools and APIs for downloading customized datasets





# aspera

- Private software owned by IBM
- Free for clients
- Can be hundreds of times faster than http and ftp
- For NCBI, you need to download and install a ~~work~~ browser plug in, Aspera Connect



## Fast Aspera Download [How to setup Aspera.](#)

Please ensure you are running a current version of AsperaConnect. It is available at [Aspera Connect](#) under the "RESOURCES" tab.

Set your bandwidth rate and continue increasing it until the data transfer rate plateaus. Many sites can transfer data at 200-500Mbps. and nearly all sites can transfer at faster than 10Mbps.

Please refer to [Aspera Transfer Guide](#) and [Aspera's documentation](#) for more information.

[Collapse tree](#)

Name	Total size	Content	Last update
└ SRR292241	956.32 Mb	1 file	2015-06-28 01:33
└ SRR292241.sra	956.32 Mb		2015-06-28 01:33



# Sequence Read Archive

- GenBank was the original name for the database to store all sequence reads
- GenBank now encompasses the Sequence Read Archive (SRA), which accepts next generation sequence data
  - Raw sequencing data
  - Alignment information
- <http://www.ncbi.nlm.nih.gov/sra>

# Emphasis on metadata

- This is a new paradigm from the old Trace Archive
  - **Study (BioProject)** – A study is a set of experiments and has an overall goal.
  - **Experiment** – An experiment is a consistent set of laboratory operations on input material with an expected result.
  - **Sample** – An experiment targets one or more samples. Results are expressed in terms of individual samples or bundles of samples as defined by the experiment.
  - **Run** – Results are called runs. Runs comprise the data gathered for a sample or sample bundle and refer to a defining experiment.

# Hierarchical Design

Study(BioProject)

Experiment

Experiment

Sample

Sample

Sample

Sample

Sample

Run

Run

Run

Run

Run

Run

# BioProjects

A BioProject is a collection of biological data related to a single initiative, originating from a single organization or from a consortium. A BioProject record provides users a single place to find links to the diverse data types generated for that project.

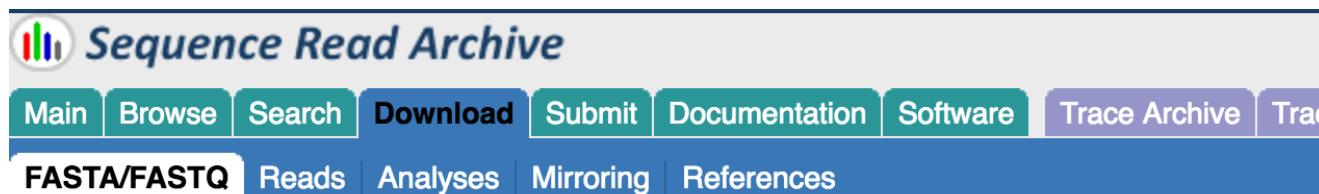
- Biosamples
- Raw reads
- Genome assembly
- Transcriptome assembly
- Genome annotation
- Markers

- You can create a BioProject page at the very beginning of the project (without data)
- Provide ongoing updates

# Sequence Read Archive Format

- Their own format : SRA format
- (this is why a lot of our lessons use FASTQ files sourced from EMBL)
- there is a web tool for downloading fastq files if you have a list of accessions and want to do this over the web:

[https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=search\\_seq\\_name](https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=search_seq_name)



## Downloading SRA data in either fasta or fastq format

Experiment(s):  Show Runs

?

What can be entered in this field?

# SRA Toolkit

- CLI tool for downloading and converting to/  
from SRA format
- Most important command:
  - **fastq-dump**: Download SRA data from the  
internet and convert into fastq format

# Submitting to the SRA

- Journals will require that you submit all NGS data to SRA (and most other ‘omic data forms somewhere!)
- Collect all data while the experiment is being done
- Know what data you need – they have spreadsheets!
- Start the submission process early (especially if you have a lot of data)

# NCBI Data Submission: Earlier is Better!

For submitting any processed data such as a genome sequence: NCBI does a lot of contamination screening, so plan ahead!

1. submit and make sure you pass their QA screens
2. Receive an accession number but keep the data private
3. Do downstream analysis
4. When you are ready to publish, make data public



If you don't do this, then they may ask you to completely change your data, which will (possibly) invalidate your downstream analysis.

# Submitting to the SRA

## Before you begin

### Gather information

#### Why did you perform your analysis?

- Project title and abstract
- Aims and objectives
- Organism(s) sequenced
- Optional: Funding sources, publications, etc.

#### What did you sequence?

- Descriptive sample information
- Tabular format is ideal
- Examples: Organism(s), age(s), gender(s), location data, cell line(s), etc.

#### How did you sequence your samples?

- Sequencing methods
- Kits used
- Instrument model(s)

#### What is your data file format?

- Files in acceptable format(s): BAM, FASTQ, etc.
- MD5 checksum for each file
- Minimum of 1 unique dataset per sample

## Register metadata

### BioProject

- 
- A description of the research effort
  - "Why" you sequenced your samples

### BioSample

- 
- A description of biologically or physically unique specimens
  - "What" you sequenced

## Provide technical details

### SRA Study

### SRA Experiment

- A description of a sample-specific sequencing library
- "How" you performed the sequencing
- Multiple Experiments can "point" to a single Sample, but not vice-versa

## Transmit data files

### SRA Sample

### SRA Run

- All files linked to a Run are "merged" into a single dataset
- Files are converted to SRA format
- Files submitted by FTP or Aspera once steps 1 and 2 are complete

**Castanea mollissima strain:Vanuxem Targeted Locus (Loci)**

The integrated genetic and physical map for Chinese chestnut was utilized to identify bacterial artificial clones (BACs) located in the three previously identified QTL regions conferring blight resistance. [More...](#)

[See Genome Information for Castanea mollissima](#)

**Related Resources:**

- [Link to assembly results](#)

**Project Data Type:** Targeted Locus (Loci)

**Attributes:** Scope: Monoisolate; Material: Genome; Capture: Targeted Locus Loci; Method type: Sequencing

**Relevance:** Environmental

**Project Data:**

Resource Name	Number of Links
SEQUENCE DATA	
SRA Experiments	8
OTHER DATASETS	
BioSample	1

**▼ SRA Data Details**

Parameter	Value
Data volume, Gbases	7
Data volume, Mbytes	7147

**Lineage:** Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliophyta; eudicotyledons; Gunneridae; Pentapetalae; rosids; fabids; Fagales; Fagaceae; Castanea; Castanea mollissima [Taxonomy ID: 60419]

**Submission:**

Registration date: 8-Nov-2014

[University of Tennessee Institute of Agriculture](#)

[NAVIGATE ACROSS](#)

2 additional projects are related by organism.

# European Nucleotide Archive

- [www.ebi.ac.uk/ena](http://www.ebi.ac.uk/ena)
- Mirrors all data in NCBI
- Can submit/download from here
- They use fastq format – no conversion needed!
- Can be slower for large datasets b/c they are far away



# Lets go look at NCBI and download some data!

When The Sequencing Data Comes In  
#WhatShouldWeCallGradSchool



<http://whatshouldwecallgradschool.tumblr.com/post/127656695585/when-the-sequencing-data-comes-in>