# Data repositories and File formats

**Gustavo Adolfo Silva-Arias Ph.D**
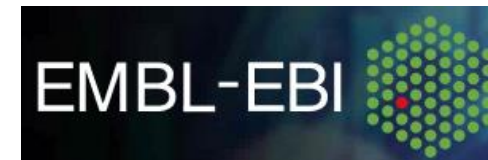
Technische Universität München
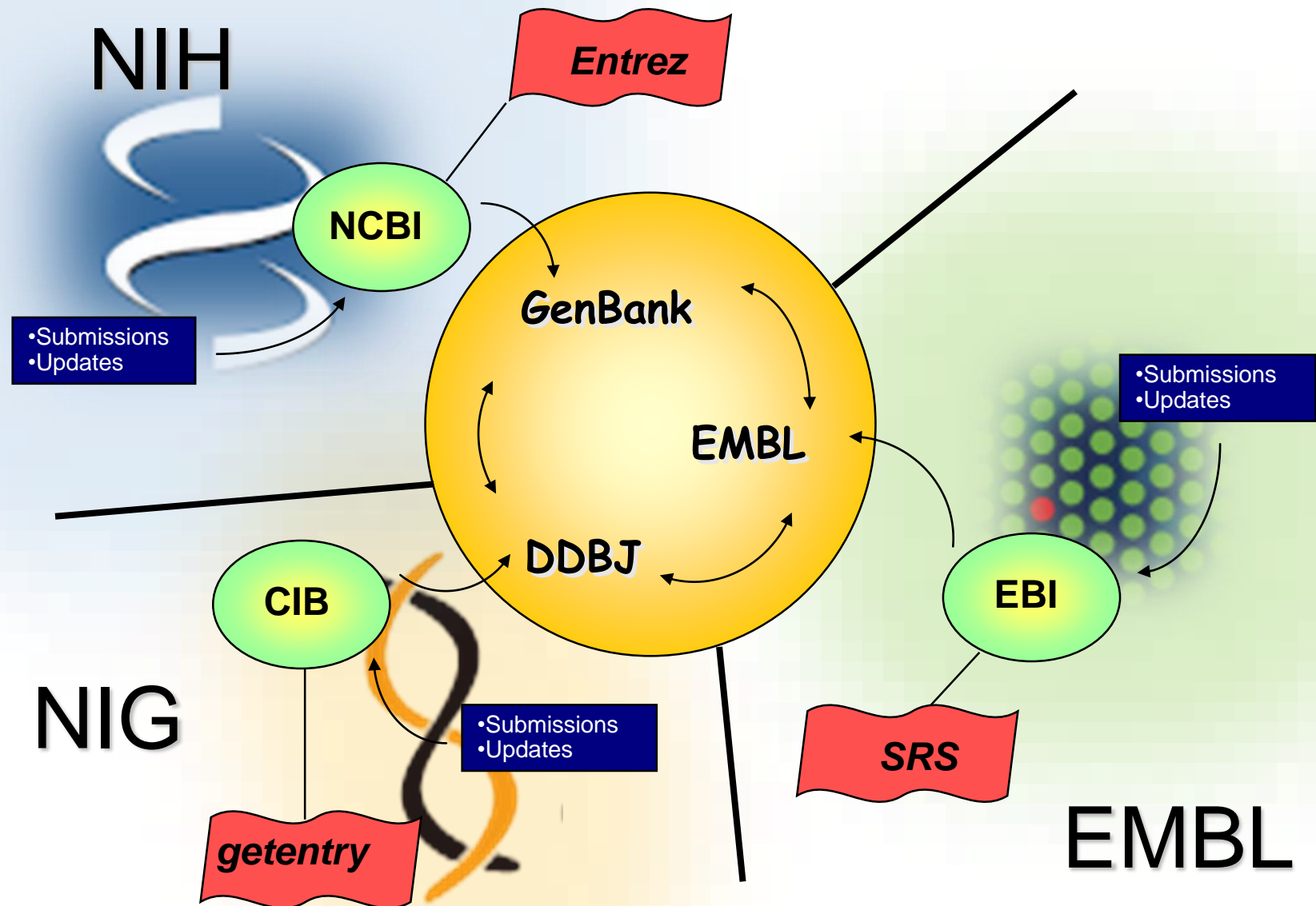
Bogotá, 24 Agosto 2021

# Overview

- Part 1
  - Associate any NGS related study publication with the databases
  - Link the raw data files with the study
  - Download all the relevant data

- Part 2
  - Familiarize with basic data file formats and understand the information

# Data repositories

- The National Center for Biotechnology Information (NCBI)
  - Genbank
  - Sequence Read Archive (SRA)

- EMBL-EBI
  - European Nucleotide Archive (ENA)
  - Ensembl
  - UniProt

- DNA DataBank of Japan (DDBJ)

- … among many others
- https://www.nature.com/sdata/policies/repositories
- https://en.wikipedia.org/wiki/List_of_biological_databases

# The International Sequence Database Collaboration

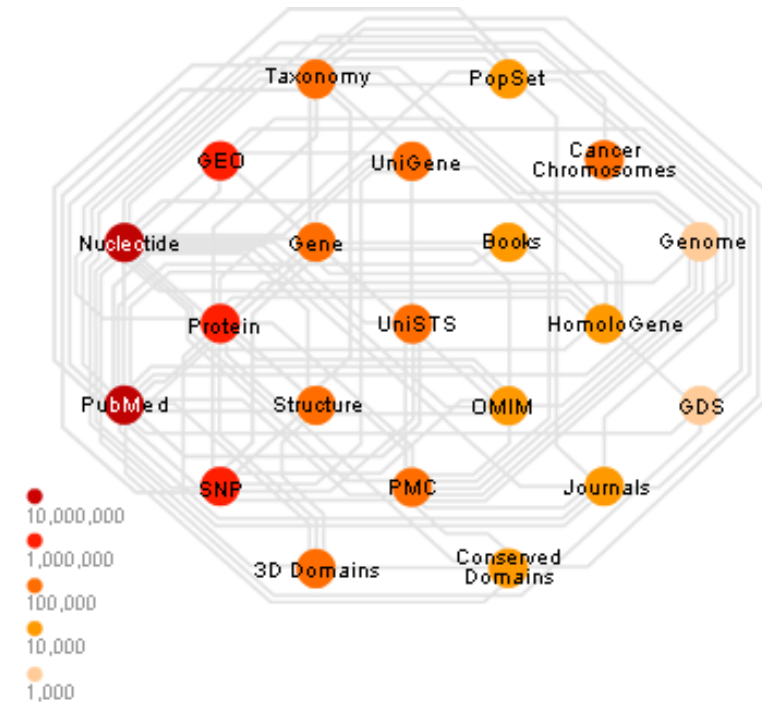# Data repositories - NCBI

Accepts submissions of:
- Bibliographic records (publication)
- Primary research data (nucleotide sequences for an organism/gene)

Organizes the information into databases, maintains them, makes them available to the world

Develops software to retrieve and analyze the data conducts basic research to make new biological discoveries
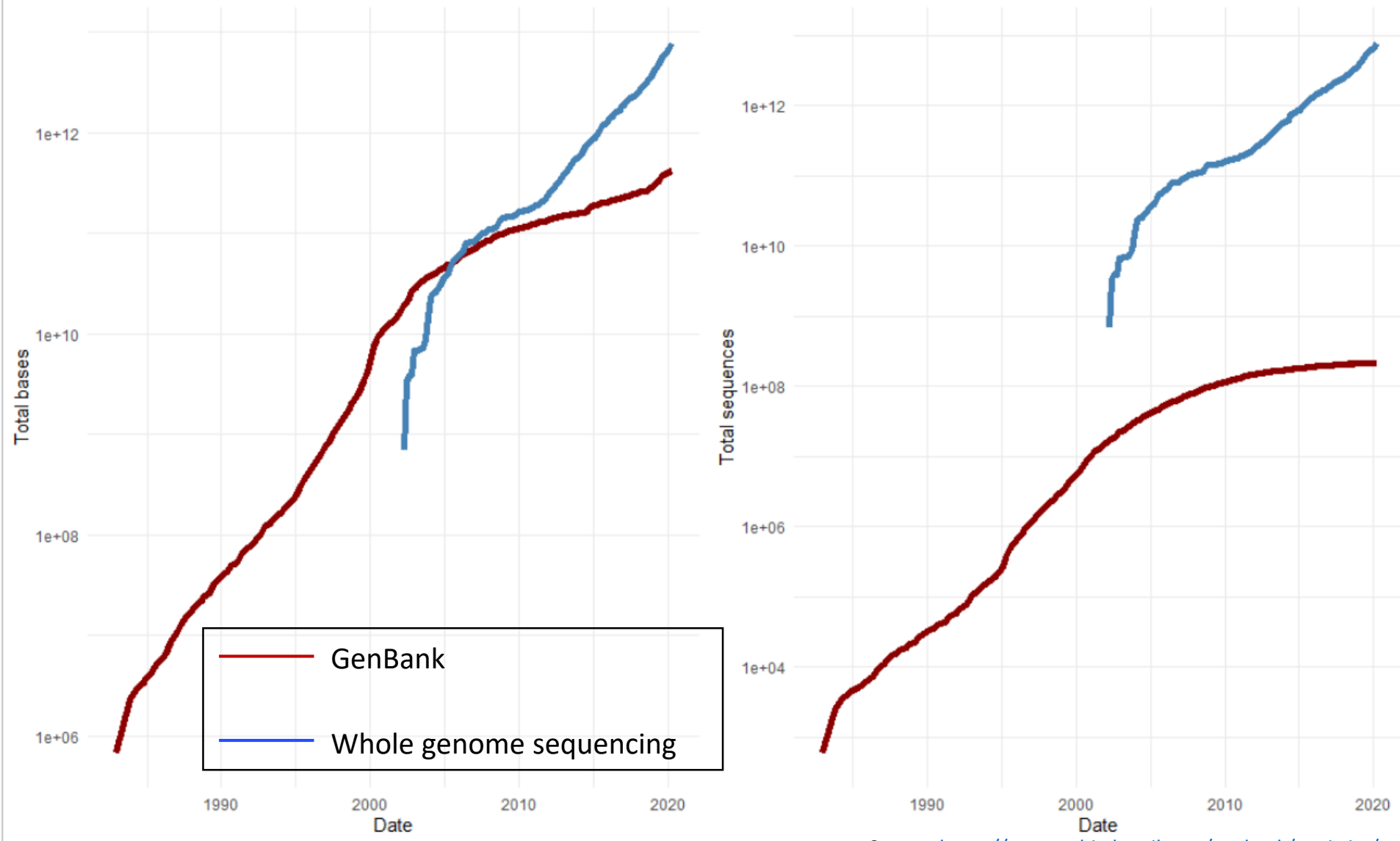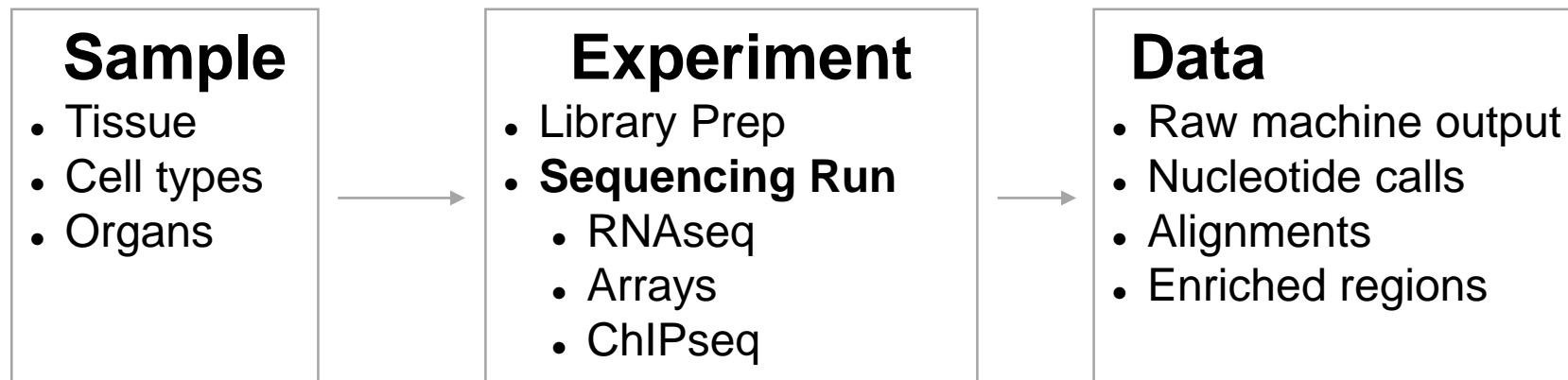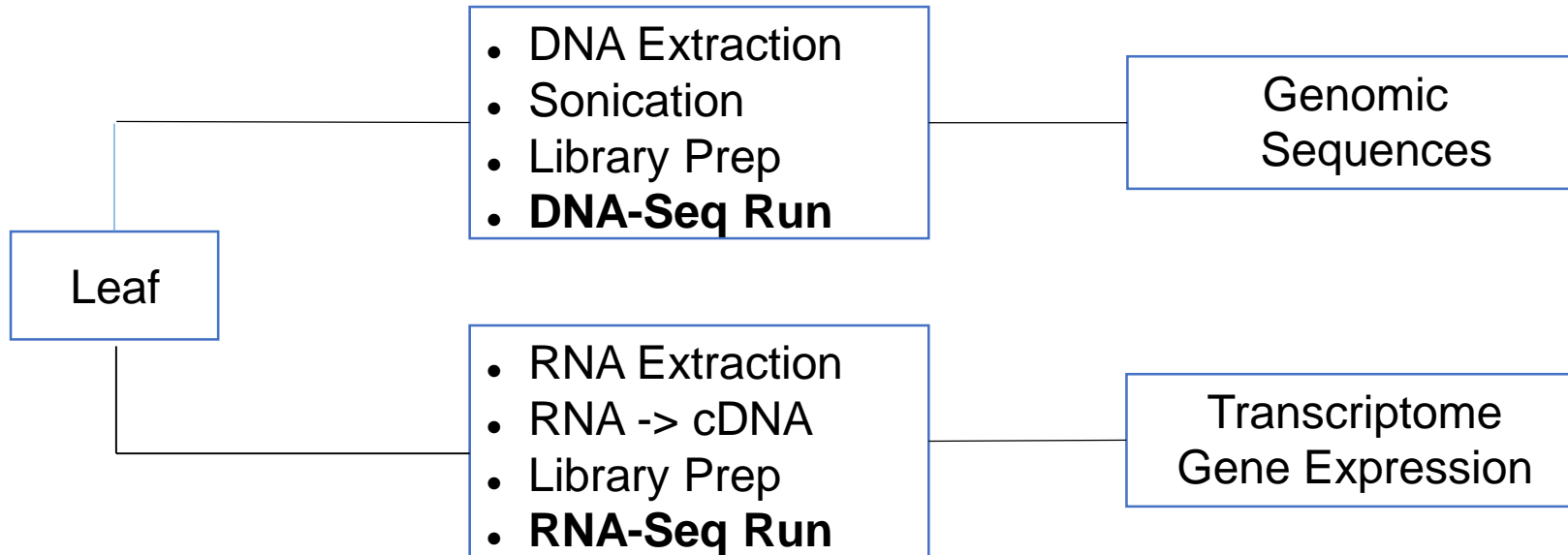
**NCBI databases**
https://www.ncbi.nlm.nih.gov/guide/all/



https://www.ncbi.nlm.nih.gov/Web/Search/entrezfs.html

# Data repositories - NCBI



Source: https://www.ncbi.nlm.nih.gov/genbank/statistics/

# Data repositories - Basic structure

**Sample**
- Tissue
- Cell types
- Organs

**Experiment**
- Library Prep
- **Sequencing Run**
  - RNAseq
  - Arrays
  - ChIPseq

**Data**
- Raw machine output
- Nucleotide calls
- Alignments
- Enriched regions

**S**ample ⟶ **E**xperiment ⟶ **D**ata

# Data structure - Basic projects



Leaf

- DNA Extraction
- Sonication
- Library Prep
- **DNA-Seq Run**

Genomic
Sequences

- RNA Extraction
- RNA -> cDNA
- Library Prep
- **RNA-Seq Run**

Transcriptome
Gene Expression

**S**ample ⟶ **E**xperiment ⟶ **D**ata

# Data structure - Basic projects



Healthy Leaf

Control Transcriptome

- RNA Extraction
- RNA → cDNA
- Library Prep
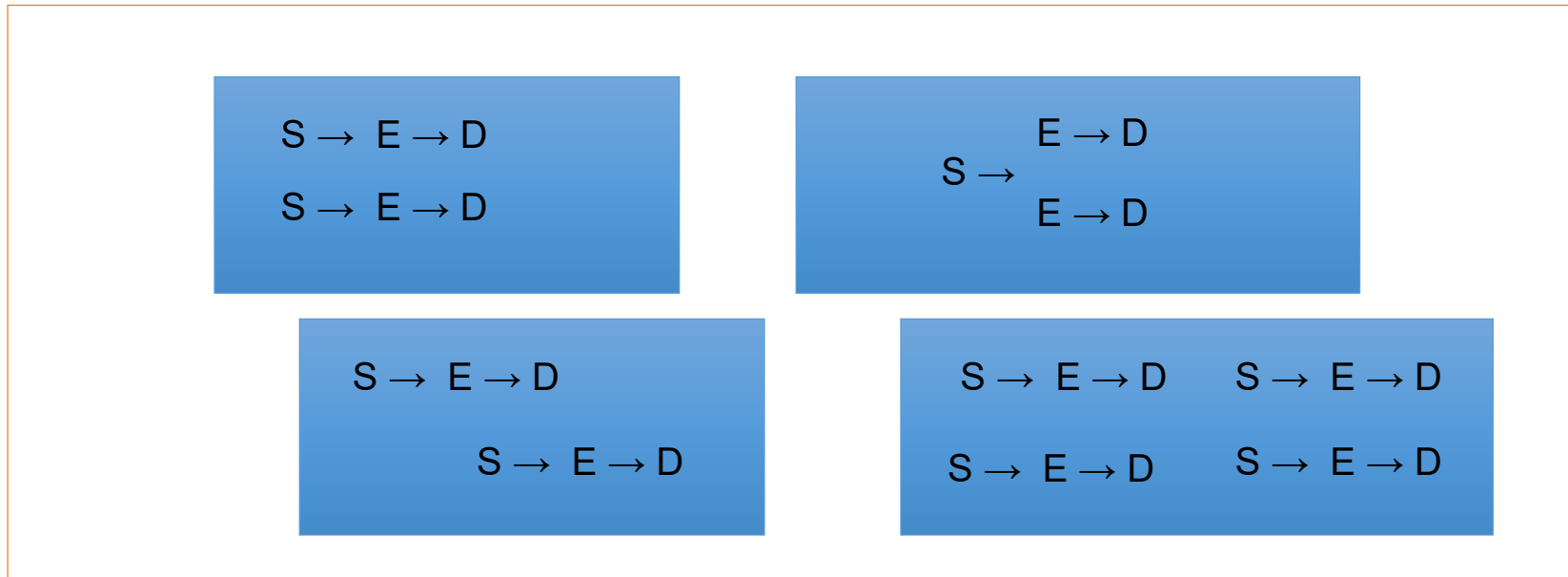- RNA-Seq Run

Infected Leaf

Case Transcriptome

**S**ample ⟶ **E**xperiment ⟶ **D**ata

# Data structure - Complex projects

BioProjects

- Initiative

- Organizations/Consortium

- Many studies in one big project

# Data structure - Complex projects
## 1000 genomes bioproject



Accession : PRJNA28889

http://www.1000genomes.org/

# Different SRA ID types

- **Study** (SRP)– A study is a set of experiments and has an overall goal.

- **Experiment** (SRX) – An experiment is a consistent set of laboratory operations on input material with an expected result.

- **Sample** (SRS)– 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** (SRR)– Results are called runs. Runs comprise the data gathered for a sample or sample bundle and refer to a defining experiment.

- **Submission** (SRA) – A submission is a package of metadata and/or data objects and a directive for what to do with those objects.

Source : http://www.ncbi.nlm.nih.gov/books/NBK47533/
Also:  http://www.ncbi.nlm.nih.gov/books/NBK56913/

# Data structure - Complex projects
## Tomato genome



BioSample: SAMN02981290

BioProjects:
    PRJNA66163 Solanum lycopersicum
                strain: Heinz 1706

    PRJNA119 Solanum lycopersicum
                cultivar:Heinz 1706



Source: https://www.sgn.cornell.edu/organism/Solanum_lycopersicum/genome

# Data structure - Complex projects
## Tomato genome

BioProject:
PRJNA119

**Project Data:**

| Resource Name | Number of Links |
|---|---:|
| SEQUENCE DATA | |
| Nucleotide (total) | 13 |
| WGS master | 1 |
| SRA Experiments | 11 |
| PUBLICATIONS | |
| PubMed | 2 |
| PMC | 1 |
| OTHER DATASETS | |
| BioSample | 12 |
| Assembly | 1 |

# Data structure - Complex projects
## Tomato genome

Experiment:
[SRX129876](#)

**SRX129876**: Tomato genome annotation using RNASeq data
1 ABI_SOLID (AB SOLiD System 3.0) run: 269.5M spots, 13.5G bases, 10.9Gb downloads

**Submitted by:** SISTEMAS GENOMICOS

**Study:** International Tomato Genome Sequencing Consortium - RNASeq in tomato var Heinz - SOLiD sequencing
PRJNA119 • SRP011485 • All experiments • All runs
show Abstract

**Sample:** International Tomato Genome Sequencing Consortium - RNASeq from tomato var Heinz
SAMN00828737 • SRS300638 • All experiments • All runs
*Organism:* Rubinisphaera brasiliensis

**Library:**
*Name:* Tomato Heinz
*Instrument:* AB SOLiD System 3.0
*Strategy:* RNA-Seq
*Source:* TRANSCRIPTOMIC
*Selection:* unspecified
*Layout:* SINGLE

**Spot descriptor:**

```
 1      forward
```

**Runs:** 1 run, 269.5M spots, 13.5G bases, 10.9Gb

| Run | # of Spots | # of Bases | Size | Published |
|-----|-----------|-----------|------|-----------|
| SRR445714 | 269,512,040 | 13.5G | 10.9Gb | 2012-05-31 |

# Data structure - Complex projects
## Tomato genome

Data:

SRR445714

**Tomato genome annotation using RNASeq data** (SRR445714)

| Metadata | Analysis | Reads | Data access |

| Run | Spots | Bases | Size | Published | Access Type |
|---|---|---|---|---|---|
| SRR445714 | 269.5M | 13.5Gbp | 11.7G | 2012-05-31 | public |

Quality graph (bigger)

This run has 1 read per spot:

L=50, 100%

Legend

| Experiment | Library Name | Platform | Strategy | Source | Selection | Layout |
|---|---|---|---|---|---|---|
| SRX129876 | Tomato Heinz | ABI Solid | RNA-Seq | TRANSCRIPTOMIC | unspecified | SINGLE |

| Biosample | Sample Description |
|---|---|
| SAMN00828737 (SRS300638) | RNASeq data from tomato (var Heinz). Equimolar amounts of total RNA from flowers at different developmental stages and fruit at different developmental s... sequencing. 50nt reads |

| Bioproject | SRA Study | Title |
|---|---|---|
| PRJNA119 | SRP011485 | International Tomato Genome Sequencing Consortium - RNASeq in tomato var Heinz - SOLiD sequencing |

Show abstract

# Data structure - Complex projects
## Tomato genome

SRA Run Selector

# Data structure - Complex projects
# Tomato genome

SRA Run Selector

# Data structure - Complex projects
## Tomato genome

SRA Run Selector

# Download NGS data

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

https://www.ncbi.nlm.nih.gov/home/download/

# Download NGS data

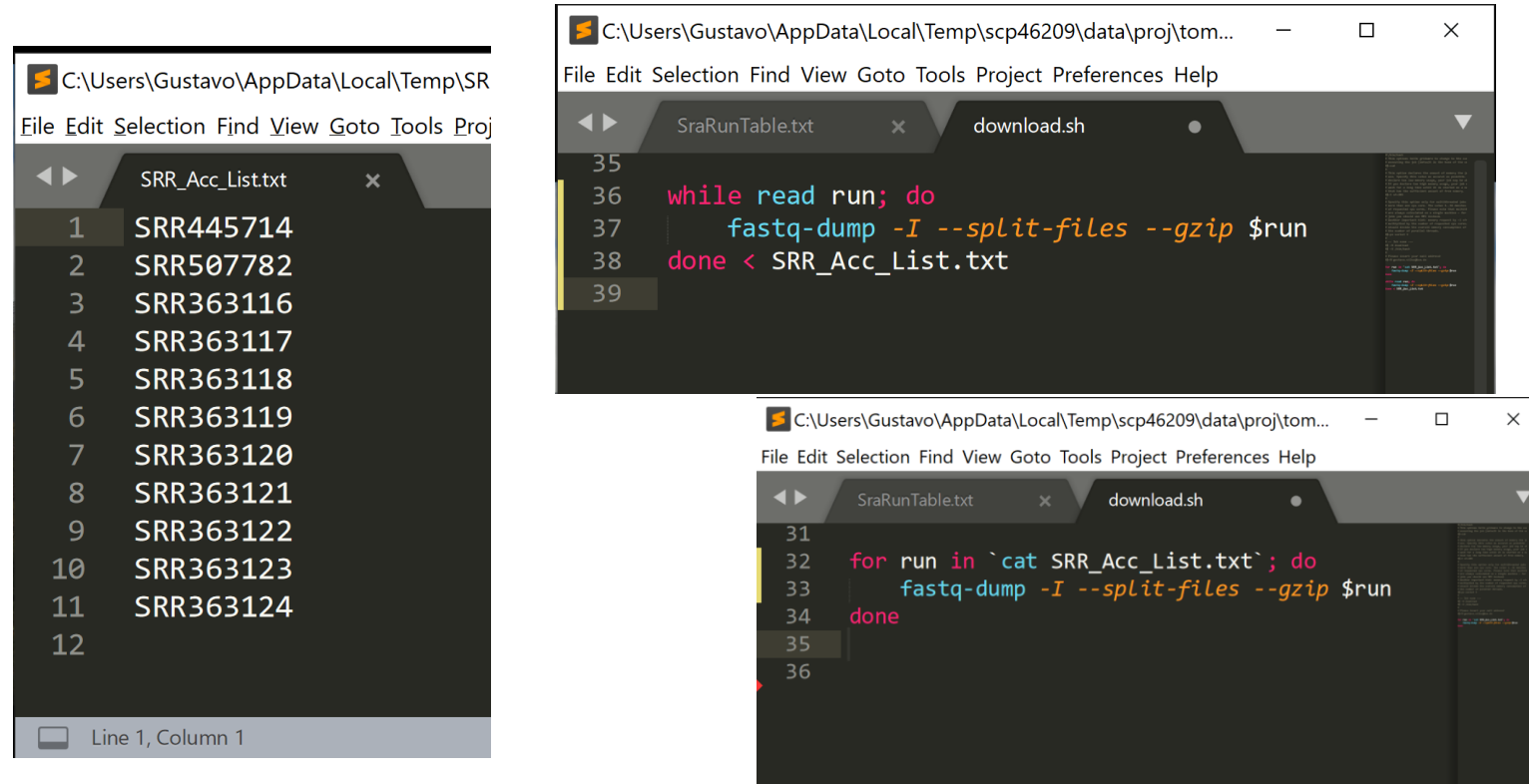- Direct Download (ftp, http, Aspera) (Browser or command line)

  wget ftp://ftp.sra.ebi.ac.uk/…/SRR4454118_1.fastq.gz

- Direct Download from EBI/DDBJ

- **sra-toolkit** software has a command fastq-dump

| Command | Argument | Input | Output |
|---------|----------|-------|--------|
| fastq-dump | -h | | Print help |
| fastq-dump | | SRR_ID | Download entire file |
| fastq-dump | -X <number> | SRR_ID | Download N spots |
| fastq-dump | --skip-technical | SRR_ID | Do not include technical reads |
| fastq-dump | -Z | SRR_ID | Print to terminal |
| fastq-dump | -F | SRR_ID | Get original id |
| fastq-dump | --split-files | SRR_ID | Print read pairs in separate files |

# Download NGS data

- **sra-toolkit**

Fastq-dump (download fastq files, pair end in separated files, compressed)

# Overview

- Part 1
  - Associate any NGS related study publication with the databases
  - Link the raw data files with the study
  - Download all the relevant data

- Part 2
  - Familiarize with basic data file formats and understand the information

# NGS data processing workflow

| Process | Output | File type |
|---|---|---|
| | Raw sequencing reads | fastq |
| Raw reads preprocessing | Clean sequencing reads | fastq |
| Raw reads mapping to the reference | | fasta / fastq |
| Post alignment processing | Alignment results | sam/bam |
| | Improved alignment results | bam |
| Variant calling | Raw variants | vcf |
| Variant filtering | Filtering variants | vcf |
| Variant annotation | Genome annotation | gff |
| | Annotated variants | vcf |

# Data files overview

**Fasta** – reference sequence
**Fastq** – unprocessed reads
**Sam** – aligned reads to the reference
**Bam** – binary (compressed) SAM file
**BED** – browser format (store genomic regions)
**GFF/GTF** – annotations
**VCF** – variant calls

# FASTA format

- Nucleotide or peptide sequence

- Simple structure
  - 2 lines per sequence

    > Header

    Sequence

- Multiple sequences per file

```
> H.sapiens chr17:126678768387-126787537
ACTGTCTCTGATTATTCTCTAGCTTCTAGCTATTCGATCGATTAGCTCTCGGATCGATCGATCTATGGGCG
ATTATATATCGGCTAGCTAGCTAGCTCTCATTCGCTAGCTAGCTAGCTAGCTATATCGATCGATCGATT
GCTCTAG

>the random protein sequence I found this morning
MDSTGEFCWICHQPEGPLKRFCGCKGSCAVSHQDCLRGWLETSRRQTCALCGTPYSMKWKTKPLREWTWGE
EEVLAAMEACLPLVLIPLAVLMIVMGTWLLVNHNGFLSPRMQVVLVVIVLLAMIVFSASASYVMVEGPGCL
DTCTAKNSTVTVNSIDEAIATQQPTKTDLGLARETLSTRFRRGKCRSCCRLGCVRLCCV
```

# FASTQ format

Standard format for high-throughput sequencing instruments

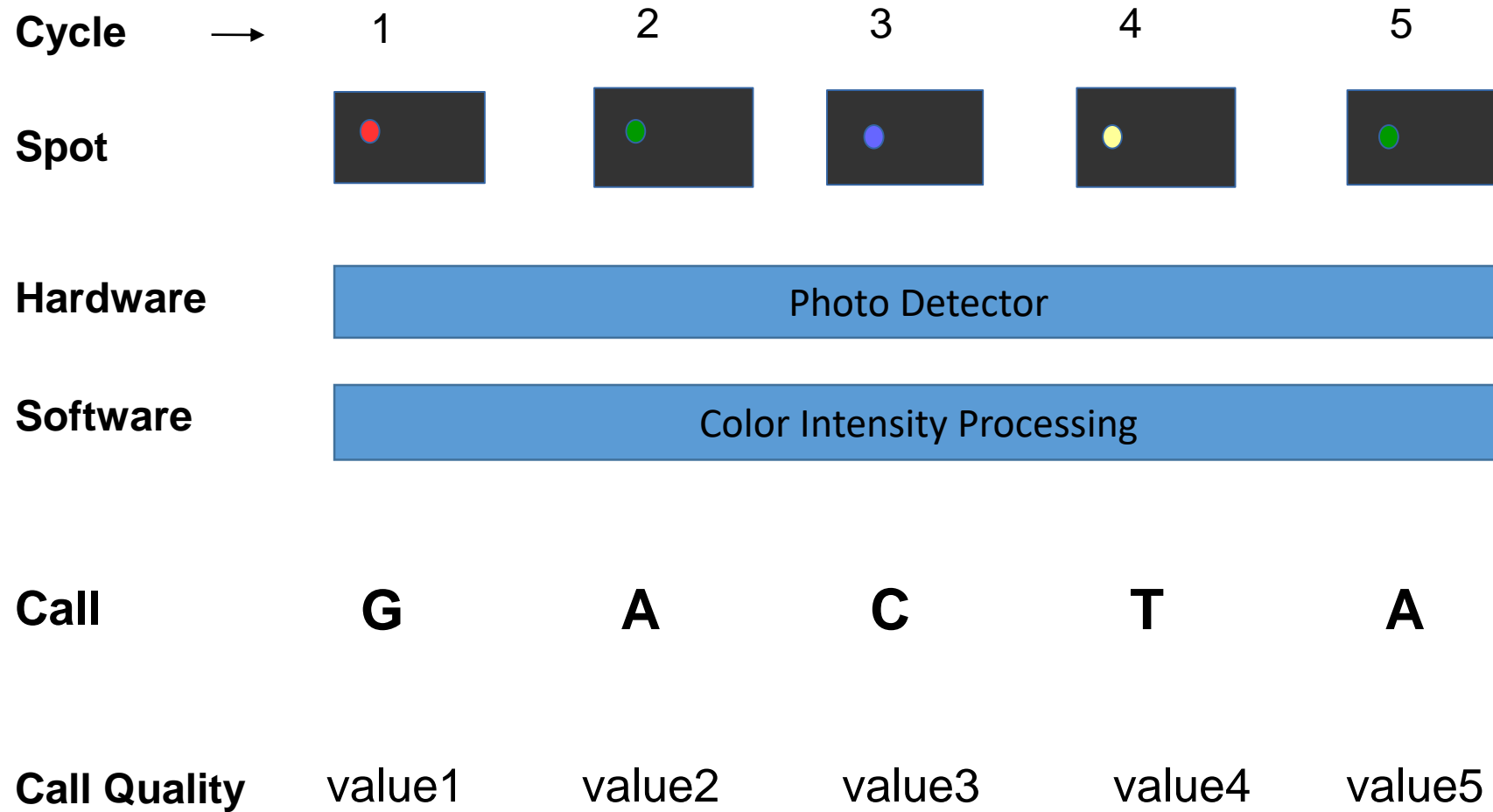4 lines per sequence (read)

@Header

Sequence

+

quality

Multiple sequences per file

```
@GWNJ-0850:627:GW190820000:5:1101:12033:1450 1:N:0:NCTCCTGA+NGGCTATA
CTTTTTCCTCGAGTATCTTTTGGAGGCGATTCTTTTTTTGAACTTGCTTTTTTTTTTGAGATCTACACGGTAGATTCAA
+
?@;DDDDDHDC<D<AEEHIGIII+<B@F?@FFGEHGIIII(77.7=7=AEHBBBB?=B8823(>A>(985++5(:@AC4
```
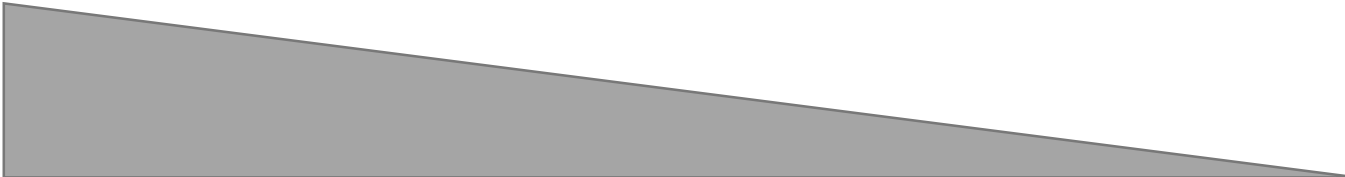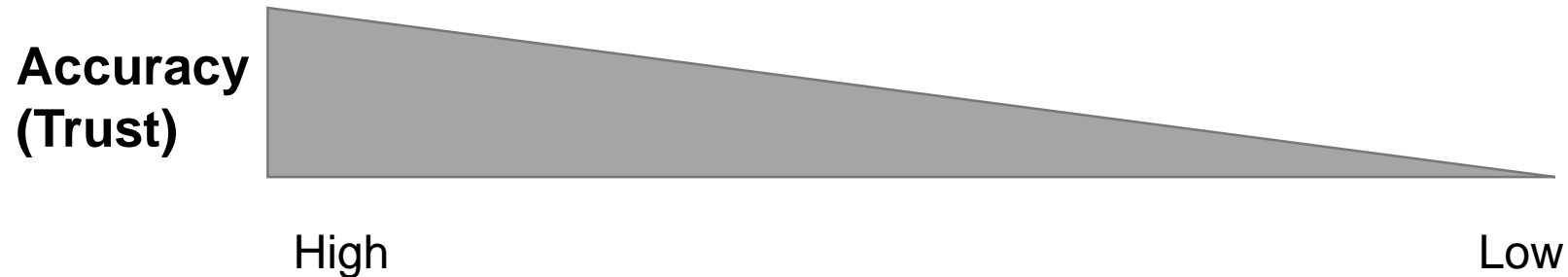
# FASTQ format
## Quality data

| Cycle → | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|



| | | | | | |
|---|---|---|---|---|---|
| **Spot** | | | | | |

| **Hardware** | Photo Detector |
|---|---|

| **Software** | Color Intensity Processing |
|---|---|

| **Call** | G | A | C | T | A |
|---|---|---|---|---|---|
| **Call Quality** | value1 | value2 | value3 | value4 | value5 |

# FASTQ format
## Quality data

| Call | G | A | C | T | A |
|---|---|---|---|---|---|
| **Call Quality** | value1 | value2 | value3 | value4 | value5 |
| **What User Wants** | Awesome | Very good | Good | OK | BAD |
| **Accuracy** | | | | | |

# FASTQ format
## Quality data

**Accuracy (Trust)**

High                                                                    Low

- Higher the value higher the trust
- **Higher the value higher the probability that call is correct**
- Amenable to statistical and probabilistic methods
- Common across all studies/platforms/machines
- Universally accepted
- Easily encoded/printed in a file

# Phred Score
## Quality data

- Denoted by letter **Q**

- $Q = -10 \log_{10} P$

- **P**: probability of error or the call being wrong

**Phred quality scores are logarithmically linked to error probabilities**

| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
| --- | --- | --- |
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |
| 50 | 1 in 100,000 | 99.999% |
| 60 | 1 in 1,000,000 | 99.9999% |

Phred scores: Phil Green's group, originally for Sanger reads. Ewing et al. (1998) Genome Res. 8:175-186

https://en.wikipedia.org/wiki/Phred_quality_score
https://www.illumina.com/documents/products/technotes/technote_Q-Scores.pdf

Source: Wikipedia

# Sequence data and Phred scores together

- Encoding ~ printing the phred scores along with base calls in a file.

- Nucleotides are typically available as a fasta file

- Quality scores could be added to the fasta file?

- Cumbersome and space consuming

```
>read1
ATGC
>read1
10 20 30 40
```

# Sequence data and Phred scores together

- ... better solution

- Put calls and quality scores and one below another

```
>read1
ATGC
10 20 30 40
```

**Encode ~ Encrypt**

```
10 = +
20 = 5
30 = ?
40 = I
```

⟶

```
>read1
ATGC
+5?I
```

# ASCII code

- **Decimal**

- 10 12 34 39 40 23 4 7 17 22 19 20 35 12 3 18 29 30 11 5 18 22


- Add 33 : 43 45 67 72 73 56 37 40 50 55 52 53 45 36 51 62 63 44 38 51 55

- **ASCII**

- +-CHI8%(2745-$3>?,&37

| Dec | Hex | Char | Dec | Hex | Char | Dec | Hex | Char | Dec | Hex | Char |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 00 | Null | 32 | 20 | Space | 64 | 40 | @ | 96 | 60 | ` |
| 1 | 01 | Start of heading | 33 | 21 | ! | 65 | 41 | A | 97 | 61 | a |
| 2 | 02 | Start of text | 34 | 22 | " | 66 | 42 | B | 98 | 62 | b |
| 3 | 03 | End of text | 35 | 23 | # | 67 | 43 | C | 99 | 63 | c |
| 4 | 04 | End of transmt | 36 | 24 | $ | 68 | 44 | D | 100 | 64 | d |
| 5 | 05 | Enquiry | 37 | 25 | % | 69 | 45 | E | 101 | 65 | e |
| 6 | 06 | Acknowledge | 38 | 26 | & | 70 | 46 | F | 102 | 66 | f |
| 7 | 07 | Audible bell | 39 | 27 | ' | 71 | 47 | G | 103 | 67 | g |
| 8 | 08 | Backspace | 40 | 28 | ( | 72 | 48 | H | 104 | 68 | h |
| 9 | 09 | Horizontal tab | 41 | 29 | ) | 73 | 49 | I | 105 | 69 | i |
| 10 | 0A | Line feed | 42 | 2A | * | 74 | 4A | J | 106 | 6A | j |
| 11 | 0B | Vertical tab | 43 | 2B | + | 75 | 4B | K | 107 | 6B | k |
| 12 | 0C | Form feed | 44 | 2C | , | 76 | 4C | L | 108 | 6C | l |
| 13 | 0D | Carriage return | 45 | 2D | - | 77 | 4D | M | 109 | 6D | m |
| 14 | 0E | Shift out | 46 | 2E | . | 78 | 4E | N | 110 | 6E | n |
| 15 | 0F | Shift in | 47 | 2F | / | 79 | 4F | O | 111 | 6F | o |
| 16 | 10 | Data link escape | 48 | 30 | 0 | 80 | 50 | P | 112 | 70 | p |
| 17 | 11 | Device control 1 | 49 | 31 | 1 | 81 | 51 | Q | 113 | 71 | q |
| 18 | 12 | Device control 2 | 50 | 32 | 2 | 82 | 52 | R | 114 | 72 | r |
| 19 | 13 | Device control 3 | 51 | 33 | 3 | 83 | 53 | S | 115 | 73 | s |
| 20 | 14 | Device control 4 | 52 | 34 | 4 | 84 | 54 | T | 116 | 74 | t |
| 21 | 15 | Neg. acknowledge | 53 | 35 | 5 | 85 | 55 | U | 117 | 75 | u |
| 22 | 16 | Synchronous idle | 54 | 36 | 6 | 86 | 56 | V | 118 | 76 | v |
| 23 | 17 | End trans. block | 55 | 37 | 7 | 87 | 57 | W | 119 | 77 | w |
| 24 | 18 | Cancel | 56 | 38 | 8 | 88 | 58 | X | 120 | 78 | x |
| 25 | 19 | End of medium | 57 | 39 | 9 | 89 | 59 | Y | 121 | 79 | y |
| 26 | 1A | Substitution | 58 | 3A | : | 90 | 5A | Z | 122 | 7A | z |
| 27 | 1B | Escape | 59 | 3B | ; | 91 | 5B | [ | 123 | 7B | { |
| 28 | 1C | File separator | 60 | 3C | < | 92 | 5C | \ | 124 | 7C | | |
| 29 | 1D | Group separator | 61 | 3D | = | 93 | 5D | ] | 125 | 7D | } |
| 30 | 1E | Record separator | 62 | 3E | > | 94 | 5E | ^ | 126 | 7E | ~ |
| 31 | 1F | Unit separator | 63 | 3F | ? | 95 | 5F | _ | 127 | 7F | ⌂ |

# Phred to ASCII

- Depends on encoding

- Sanger Encoding
  - Add 33 to the phred score and convert the number to character
  - Subtract 33 from the ascii code of the character

- Illumina encoding < 1.8 add 64

- Illumina encoding 1.8+ add 33

- **Software like FASTQC will tell you the encoding**

# Phred to ASCII

- Encoding

# Phred to ASCII

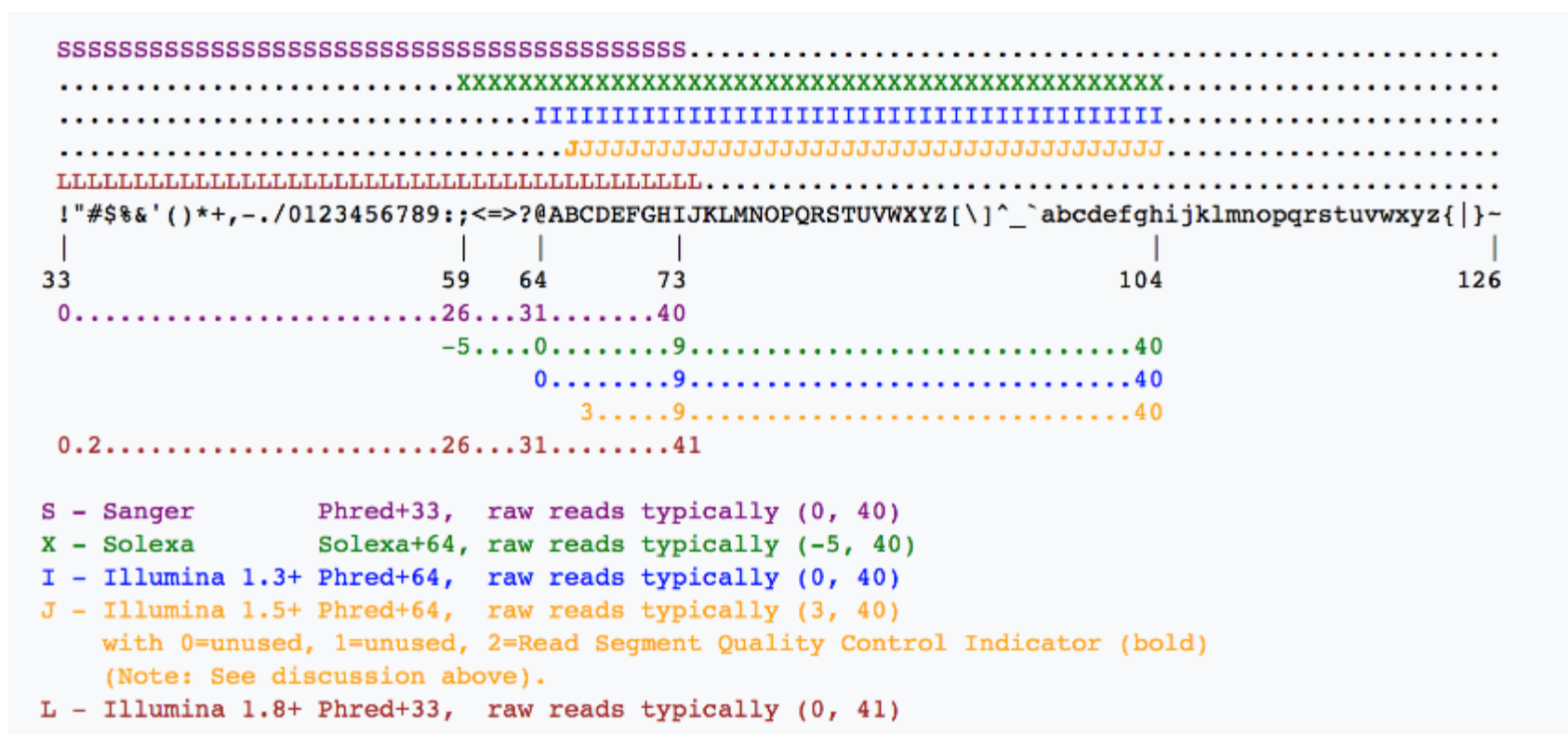- Encoding in different platforms

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.........................................
.............................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.............
...........................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII..............
.....................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.................
LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.........................................
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
|                              |    |          |                        |          |
33                            59   64         73                       104        126

0.........................26...31.......40
                          -5....0......9............................40
                              0........9............................40
                                  3.....9............................40
0.2.......................26...31........41
```

```
S - Sanger        Phred+33,   raw reads typically (0, 40)
X - Solexa        Solexa+64,  raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64,   raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64,   raw reads typically (3, 40)
    with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
    (Note: See discussion above).
L - Illumina 1.8+ Phred+33,   raw reads typically (0, 41)
```

# Phred to ASCII

- Illumina 1.9 uses **ASCII-33**, i.e. Illumina **quality score of 40** becomes

- 40 + 33 = **73** : "**I**"

# FASTQ format

- Each read is **4** lines

- Read starts with a character @ followed by the read descriptor

- Sequence follows in the second line

- Third line is reserved for additional info

- Fourth line is the Phred score encoding


- Read pairs are typically in different files

```
@GWNJ-0850:627:GW190820000:5:1101:12033:1450 1:N:0:NCTCCTGA+NGGCTATA
CTTTTTCCTCGAGTATCTTTTGGAGGCGATTCTTTTTTTGAACTTGCTTTTTTTTTTGAGATCTACACGGTAGATTCAA
+
?@;DDDDDHDC<D<AEEHIGIII+<B@F?@FFGEHGIIII(77.7=7=AEHBBBB?=B8823(>A>(985++5(:@AC4
```

# FASTQ format

- Store calls (ATGC … )

- Store Phred scores (Encoded)


- Store Machine make/ID

- Store Flowcell id for each spot

- Store coordinates of the spot

- Store additional info (Seq names …. )

- Easily parsed and stored.

```
@GWNJ-0850:627:GW190820000:5:1101:12033:1450 1:N:0:NCTCCTGA+NGGCTATA
CTTTTTCCTCGAGTATCTTTTGGAGGCGATTCTTTTTTTGAACTTGCTTTTTTTTTTGAGATCTACACGGTAGATTCAA
+
?@;DDDDDHDC<D<AEEHIGIII+<B@F?@FFGEHGIIII(77.7=7=AEHBBBB?=B8823(>A>(985++5(:@AC4
```

# FASTQ format

- Header



@HWI-ST863:211:C1MHVACXX:3:1101:1245:1869 1:N:0:TGGTTGTT — Identifier

NCATACCAGCGACGACGAGGACGGGGATGAAGACCCTGAGAGTGCACAGACGTCGTGTCGGCCCTCTATGGAGAATCCTATTTCG ATGACGACTGCCCA — Sequence

#1=DDFFEHADHGIGHIEHIIIGHIIHHIIIG=CHE@BDFCE3>ACCCCCBB5;B=B<?<@57;?ACCCDCCBA(:@CCCCCEDC @CD@:?@BBBCCB4 — Quality string

| HWI-ST863 | Instrument name |
|-----------|-----------------|
| 211 | Run id |
| C1MHVACXX | Flowcell id |
| 3 | Flowcell lane |
| 1101 | Tile number of flowcell lane |
| 1245 | 'x'-coordinate of the cluster within the tile |
| 1869 | 'y'-coordinate of the cluster within the tile |
| 1 | the member of a pair, 1 or 2 |
| N | Y if the read is filtered, N otherwise |
| 0 | 0 when none of the control bits are on, otherwise it is an even number |
| TGGTTGTT | Index sequence |

# FASTQ format

- Old header format

@HWUSI-EAS100R:6:73:941:1973#0/1

| HWUSI-EAS100R | the unique instrument name |
|---|---|
| 6 | flowcell lane |
| 73 | tile number within the flowcell lane |
| 941 | 'x'-coordinate of the cluster within the tile |
| 1973 | 'y'-coordinate of the cluster within the tile |
| #0 | index number for a multiplexed sample (0 for no indexing) |
| /1 | the member of a pair, /1 or /2 (paired-end or mate-pair reads only) |

Source: Wikipedia

# SAM format

- **SAM** : Sequence Alignment/Map

- **BAM** : Binary Alignment/Map (binary SAM)

- Used for: aligned reads

- Multiple *tab* delimited columns

- It is flexible enough to store **all the alignment information** generated by various alignment programs

- It allows most of the operations on the alignment to work on a stream without loading the whole alignment into memory

- It allows the file to be **indexed by genomic position** to efficiently retrieve all reads aligning to a locus

# SAM format

**Two sections**

- Header section, each line begins with "@"–Several record types

  5 fields

- Alignment section

  11 mandatory fields (columns)

# SAM format

- **Header section**

- @HD The header line. The first line if present.
  - VN* Format version

- @SQ Reference sequence dictionary. The order of @SQ lines defines the alignment sorting order
  - SN* Reference sequence name
  - LN* Reference sequence length

- @RG Read group. Unordered multiple @RG lines are allowed.
  - ID* Read group identifier

- @PG Program
  - ID* Program record identifier

- @CO One-line text comment

# SAM format

- Header section + Alignment section (first line)

```
@HD     VN:1.5  GO:none SO:coordinate
@SQ     SN:Spenn-ch01   LN:109333515
@SQ     SN:Spenn-ch02   LN:59803892
@SQ     SN:Spenn-ch03   LN:75414019
@SQ     SN:Spenn-ch04   LN:77197300
@SQ     SN:Spenn-ch05   LN:77991103
@SQ     SN:Spenn-ch06   LN:60730942
@RG     ID:LA2932_28.CA1PNANXX.3.CAGATC PU:CA1PNANXX.3.CAGATC   LB:LA2932_28
SM:LA2932_28    PL:ILLUMINA
@PG     ID:MarkDuplicates       VN:1.119 CL:picard.sam.MarkDuplicates …
@PG     ID:bwa  VN:0.7.12-r1039 CLbwa mem -M -t 24 Spenn2.fa …
@PG     ID:GATK IndelRealigner  VN:3.7-0-gcfedb67       …
HISEQ:202:CA1PNANXX:4:1311:19476:42830  163     Spenn-ch01      13      0
4S99M23S        =       124     228
TTATGGCCAACCGGATGCATAGACAAGGTCTTGACGGACGTCCACAAAAAAATTTGCCATTTTTGATGTCGGAATCCGG
ATCACCCAGAAAATGGTTTGCTATGTCACACGGAAATCGTTAAAATG
BBBBB<FFFFFFFFFBFFFFFFFFFBFFFF<F<FFFF</FFBBFFFFFFFFFFFFFBF/BFF/<F/FF<</<<FBBB/BFFFF
FFFFFFFFFFB<FFFBFFFFFFFF<7/7//<7/7F//777/BFFFFB  MC:Z:117M4S
MD:Z:46G1G34A15 PG:Z:MarkDuplicates     RG:Z:LA2932_28.CA1PNANXX.3.CAGATC
NM:i:3  MQ:i:0  AS:i:84 XS:i:90

…
```

Source: Wikipedia

# SAM format

- Alignment section

```
HISEQ:202:CA1PNANXX:4:1311:19476:42830
```
Query name: shared by pair-end mates

```
163
```
Flag value: Decimal > Binary > Multiple True/False values

```
Spenn-ch01
```
Chromosome/Contig where the read aligned

```
13
```
Position on chromosome/contig where the read aligned

```
0
```
Alignment confidence (Phred)

```
4S99M23S
```
CIGAR string

```
=
```
Chomosome/Contig where mate aligned (= if same)

```
124
```
Position on chromosome/contig where the mate aligned

```
228
```
Lenght of the reference sequence read aligned to

```
TTATGGCCAACCGGATGCATAGACAAGGTCTTGACG...
```
Read sequence

```
BBBBB<FFFFFFFFBFFFFFFFFFBFFFF<F<FFFF<...
```
Read quality score (same as fastq file)

```
        MC:Z:  MD:Z:  PG:Z:  RG:Z:
   NM:i:3  MQ:i:0  AS:i:84 XS:i:90
```
Optional tags: http://samtools.github.io/hts-specs/SAMtags.pdf

# SAM format

# SAM format

- bitwise flag

```
HISEQ:202:CA1PNANXX:4:1311:19476:42830    163    Spenn-ch01    13    0
4S99M23S           =        124      228
TTATGGCCAACCGGATGCATAGACAAGGTCTTGACGGACGTCCACAAAAAATTTGCCATTTTTGATGTCGGAATCCGG
ATCACCCAGAAAATGGTTTGCTATGTCACACGGAAATCGTTAAAATG
BBBBB<FFFFFFFFBFFFFFFFFBFFFF<F<FFFF</FFBFFFFFFFFFFFFBF/BFF/<F/FF<</<<FBBB/BFFFF
FFFFFFFFFB<FFFBFFFFFFFF<7/7//<7/7F//777/BFFFFB   MC:Z:117M4S
MD:Z:46G1G34A15 PG:Z:MarkDuplicates      RG:Z:LA2932_28.CA1PNANXX.3.CAGATC
NM:i:3   MQ:i:0   AS:i:84 XS:i:90
```

1 + 2 + 32 + 64 = **99**

1 + 2 + 16 + 128 = **147**

**Decoding SAM flags** https://broadinstitute.github.io/picard/explain-flags.html

# SAM format

- CIGAR string – Describes how the read align to the reference

```
HISEQ:202:CA1PNANXX:4:1311:19476:42830  163      Spenn-ch01         13      0
4S99M23S              =         124      228
TTATGGCCAACCGGATGCATAGACAAGGTCTTGACGGACGTCCACAAAAAAATTTGCCATTTTTGATGTCGGAATCCGG
ATCACCCAGAAAATGGTTTGCTATGTCACACGGAAATCGTTAAAATG
BBBBB<FFFFFFFFFBFFFFFFFFFBFFFF<F<FFFF</FFBBFFFFFFFFFFFFFBF/BFF/<F/FF<</<<FBBB/BFFFF
FFFFFFFFFFB<FFFBFFFFFFFFF<7/7//<7/7F//777/BFFFFB  MC:Z:117M4S
MD:Z:46G1G34A15 PG:Z:MarkDuplicates      RG:Z:LA2932_28.CA1PNANXX.3.CAGATC
NM:i:3  MQ:i:0  AS:i:84 XS:i:90
```

| Op | BAM | Description | Consumes query | Consumes reference |
|----|-----|-------------|----------------|--------------------|
| M | 0 | alignment match (can be a sequence match or mismatch) | yes | yes |
| I | 1 | insertion to the reference | yes | no |
| D | 2 | deletion from the reference | no | yes |
| N | 3 | skipped region from the reference | no | yes |
| S | 4 | soft clipping (clipped sequences present in SEQ) | yes | no |
| H | 5 | hard clipping (clipped sequences NOT present in SEQ) | no | no |
| P | 6 | padding (silent deletion from padded reference) | no | no |
| = | 7 | sequence match | yes | yes |
| X | 8 | sequence mismatch | yes | yes |

Source: https://samtools.github.io/hts-specs/SAMv1.pdf

# SAM format

- CIGAR string – Describes how the read align to the reference

```
Reference: ATGAAGGATAGTGATACTCTAGAGGG
Read: ACGAATAGTGATACTCGGGTAGAGGG
```

| Op | BAM | Description | Consumes query | Consumes reference |
|----|-----|-------------|----------------|--------------------|
| M | 0 | alignment match (can be a sequence match or mismatch) | yes | yes |
| I | 1 | insertion to the reference | yes | no |
| D | 2 | deletion from the reference | no | yes |
| N | 3 | skipped region from the reference | no | yes |
| S | 4 | soft clipping (clipped sequences present in SEQ) | yes | no |
| H | 5 | hard clipping (clipped sequences NOT present in SEQ) | no | no |
| P | 6 | padding (silent deletion from padded reference) | no | no |
| = | 7 | sequence match | yes | yes |
| X | 8 | sequence mismatch | yes | yes |

Source: https://samtools.github.io/hts-specs/SAMv1.pdf

# SAM format

- CIGAR string – Describes how the read align to the reference

```
Reference:    ATGAAGGATAGTGATACTC---TAGAGGG
Read:         ACGAA---TAGTGATACTCGGGTAGAGGG
CIGAR:        5M3D11M3I7M
```

| Op | BAM | Description | Consumes query | Consumes reference |
|----|-----|-------------|----------------|--------------------|
| M | 0 | alignment match (can be a sequence match or mismatch) | yes | yes |
| I | 1 | insertion to the reference | yes | no |
| D | 2 | deletion from the reference | no | yes |
| N | 3 | skipped region from the reference | no | yes |
| S | 4 | soft clipping (clipped sequences present in SEQ) | yes | no |
| H | 5 | hard clipping (clipped sequences NOT present in SEQ) | no | no |
| P | 6 | padding (silent deletion from padded reference) | no | no |
| = | 7 | sequence match | yes | yes |
| X | 8 | sequence mismatch | yes | yes |

Source: https://samtools.github.io/hts-specs/SAMv1.pdf

# SAM format

- CIGAR string – Describes how the read align to the reference

```
Reference:    ATGAAGGATAGTGATACTC---TAGAGGG
Read:         ACGAA---TAGTGATACTCGGGTAGAGGG
CIGAR:        5M   3D 11M          3I 7M
```

| Op | BAM | Description | Consumes query | Consumes reference |
|----|-----|-------------|----------------|--------------------|
| M | 0 | alignment match (can be a sequence match or mismatch) | yes | yes |
| I | 1 | insertion to the reference | yes | no |
| D | 2 | deletion from the reference | no | yes |
| N | 3 | skipped region from the reference | no | yes |
| S | 4 | soft clipping (clipped sequences present in SEQ) | yes | no |
| H | 5 | hard clipping (clipped sequences NOT present in SEQ) | no | no |
| P | 6 | padding (silent deletion from padded reference) | no | no |
| = | 7 | sequence match | yes | yes |
| X | 8 | sequence mismatch | yes | yes |

Source: https://samtools.github.io/hts-specs/SAMv1.pdf

# BAM format

- Binary SAM

- Used for: aligned reads

- **25% of the size**

- **SAMtools** to convert

- **.bai** = BAM index

# VCF format

- Variant call format

# GFF/GTF format

- Genome annotation

# Hands-on ...

NGS wiki

https://github.com/gsilvaarias/NGS2021-AGROSAVIA/wiki/02.-Bases-de-datos-y-formatos-de-archivos-NGS

**1st part**: Download sequence data and familiarize with SRA data repository

**2nd part**: Dive into different file formats, read on your own, try to dissect information (focus on fastq and SAM)