

Understanding NGS raw data: Fastq format, quality checking

Overview

- Why this bioinformatics session?
- Basic NGS terminology
- NGS Pipelines
- Fastq format
- Fastq quality control

Why this session?

A typical RNA-Seq analysis at the Genomics Core delivers:

KU LEUVEN GENOMICS CORE
Herestraat 49 PO box 610,
3000 Leuven
tel. +32 16 33 08 21
mail: genomicscore@kuleuven.be

Project identification

Project Type RNA-seq
Number of Samples 10
Number of Conditions 2
Condition 1 MMP-9_KO_wt
Condition 2 MMP-9_KO_DS

Used Read Information

Sample Code	Condition	Used Read Count	Size Factor
CC104962	MMP-9_KO_wt	496208	0.48115132775072
CC104963	MMP-9_KO_wt	716403	1.27715467581361
CC104964	MMP-9_KO_wt	639245	1.064991387054
CC104965	MMP-9_KO_wt	409410	0.8332272277977
CC104970	MMP-9_KO_DS	686207	0.95136671272552
CC104971	MMP-9_KO_DS	297470	0.5916486132944
CC104972	MMP-9_KO_DS	729486	1.23614468091
CC104973	MMP-9_KO_DS	629628	1.1107879446147
CC104975	MMP-9_KO_wt	502251	1.2000610225668
CC104979	MMP-9_KO_wt	331976	0.5433710430181
CC104980	MMP-9_KO_wt	570708	1.2166626828147
CC104981	MMP-9_KO_wt	681343	1.45000741122796
CC104986	MMP-9_KO_DS	374337	0.81302637799265
CC104987	MMP-9_KO_DS	182259	1.4367640661185
CC104988	MMP-9_KO_DS	756891	1.3523966296006
CC104989	MMP-9_KO_DS	462772	1.284138308557

Sample Relations

Data quality assessment and quality control are essential steps of any data analysis. Here we define the term quality a Broom for purpose. Our purpose is the detection of differentially expressed genes, and we are looking in particular to samples whose experimental treatment suffered from an abnormality that renders the data points obtained from this particular samples detrimental to our purpose.

Variance stabilised data is used to create sample-to-sample distances. With these distances sample clustering becomes possible. The clustering should reflect the experimental design correctly, samples are more similar when they have the same treatment. The heatmap with true should show this effect (Figure 8). The Principal Component Analysis plot (PCA-plot) (Figure 9) is a 2 dimensional version of these analysis. Expected is that samples with a same treatment cluster together. Outliers and possible bias are easy to detect.

Top 30 highly expressed genes

By taking a look at the top 30 highly expressed genes, a first impression of the data can be made. The heatmap below shows this expression data (Figure 8). The data is normalised by using the variance stabilisation transformation. Samples should cluster together according to the experimental factor. However if this is not the case, this doesn't imply that there is no difference. This plot shows only the highest expressed genes, not the differentially expressed genes.



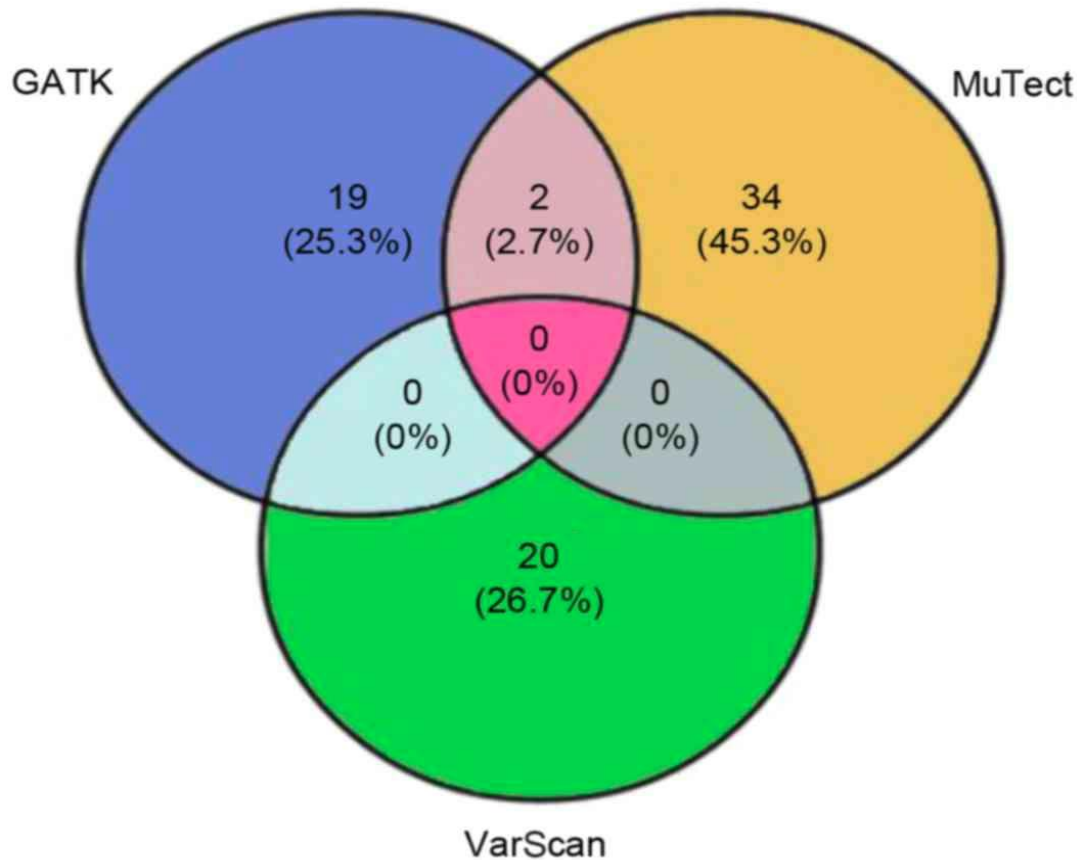
- 10 report files:
differentially expressed
genes, reads quality
report, counting report,...

3GB of raw and analyzed data:
Fastq, BAM, counts,
normalized counts, gene
expression files, heatmaps...

1. Results Interpretation & Analysis



2. Data Re-analysis

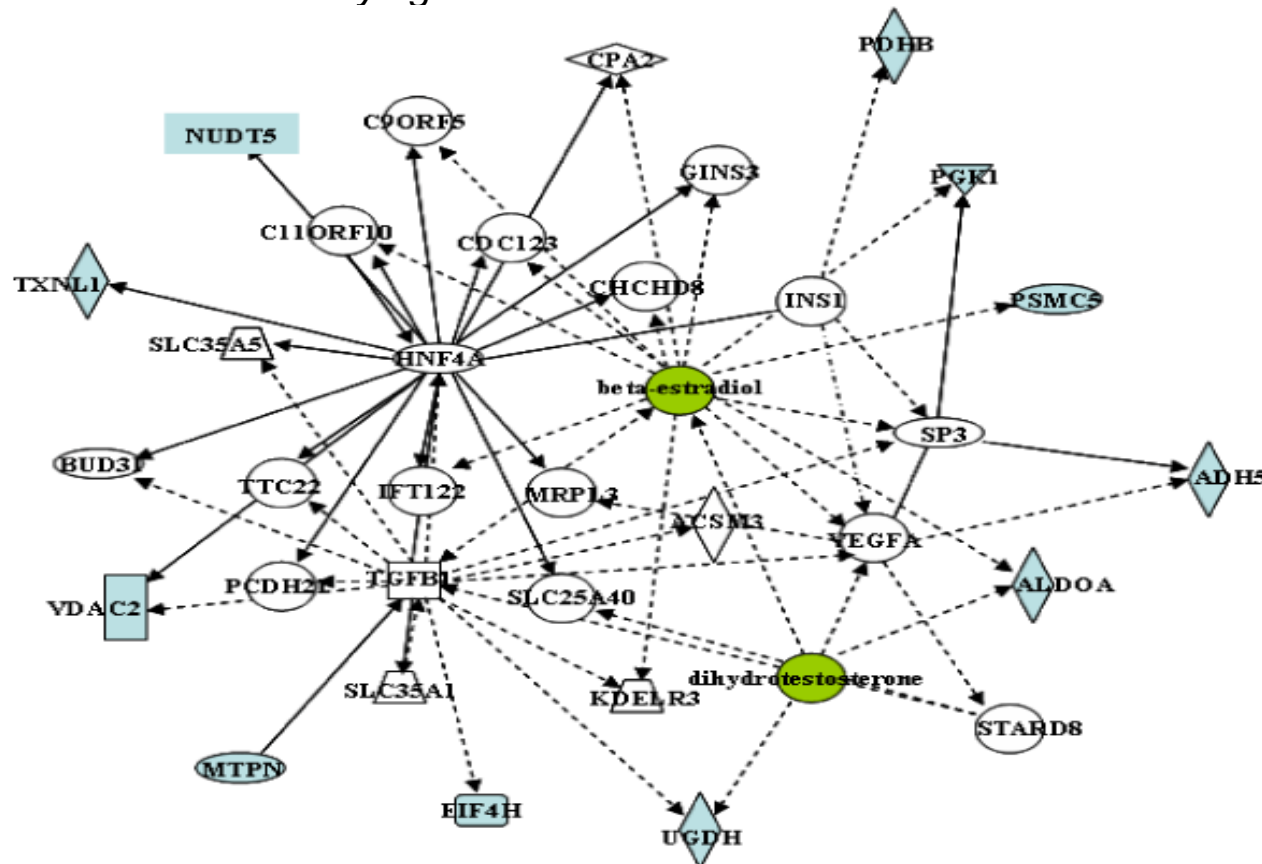


A three-caller pipeline for variant analysis of cancer whole-exome sequencing data

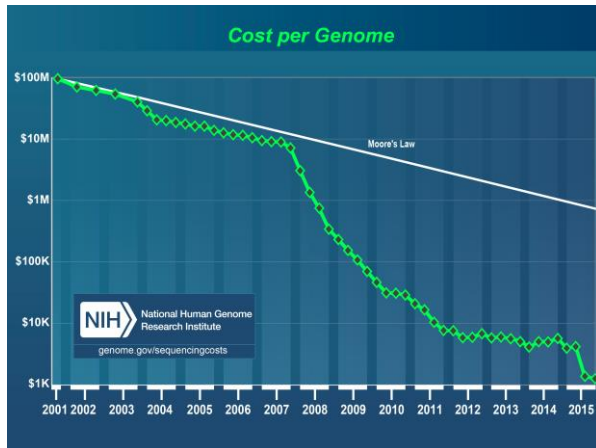
Authors: Ze-Kun Liu, Yu-Kui Shang, ✉ Zhi-Nan Chen, ✉ Huijie Bian

3. Downstream Analysis

Reviewer: “The transcriptional and proteomic profiling experiments propose some interesting followup pathways for further analysis that may shed light on the mechanisms underlying these conditions.”



4. Costs Saving



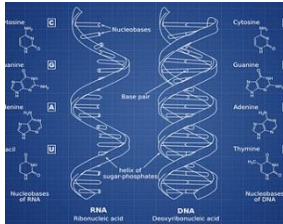
- **\$1.000 sequencing and \$5.000 bioinformatics costs?!?!**

Overview

- Why this bioinformatics session?
- Basic NGS terminology
- NGS Pipelines
- Fastq format
- Fastq quality control

NGS Bioinformatics

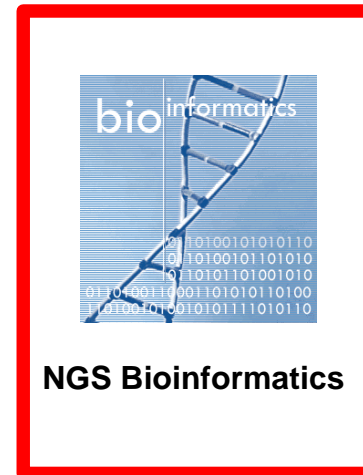
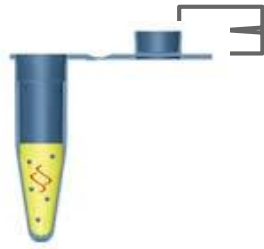
- **Experimental Design**



- **Sequencing**



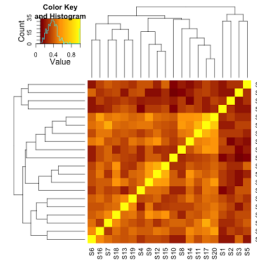
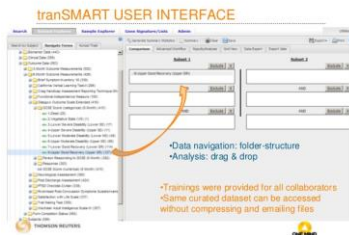
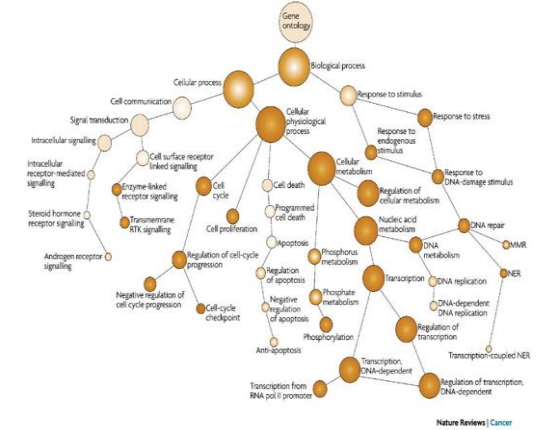
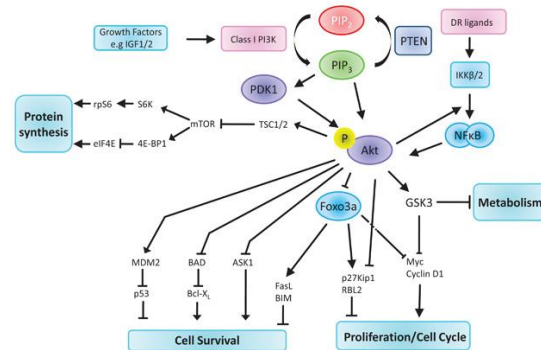
- **Library Preparation**



- **Follow up and support**

Bioinformatics...

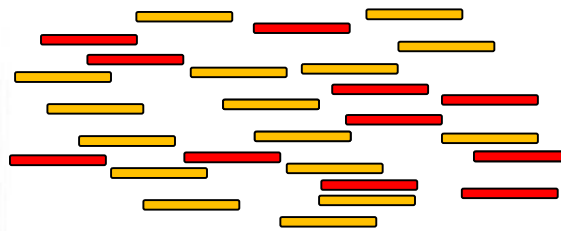
...is a broad field



```
a.length; c++) {  
    &b.push(a[c]);  
function h() {  
#user_logged").a(), a = q  
place(/+(?=)/g, ""); a.le  
) , b = [], c = 0; c < a.le  
    0 == r[a[c], b] && b.  
    c = c++; c.f - a.le  
    c = b.length - 1;  
    }  
}
```

<http://www.computerhope.com>

NGS Bioinformatics



NGS bioinformatics: interpretation and analysis
of **NGS data** using informatics tools

A day in NGS Bioinformatics

- Poly-A tail

Fastq

- Ontology

- Fragments cluster

- Xenograft Murine infiltration

- Hadoop

- Flow Cell

- Demultiplexing

- VCF File

- Ortholog gene

BAM

- Negative binomial distribution

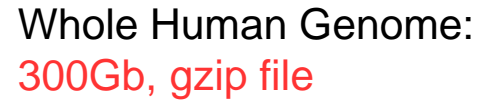
- Cluster

- snRNA

- Bonferroni Correction

- Barcode

NGS data means big data...means big computing power...by now.



RNA-Seq
1GB



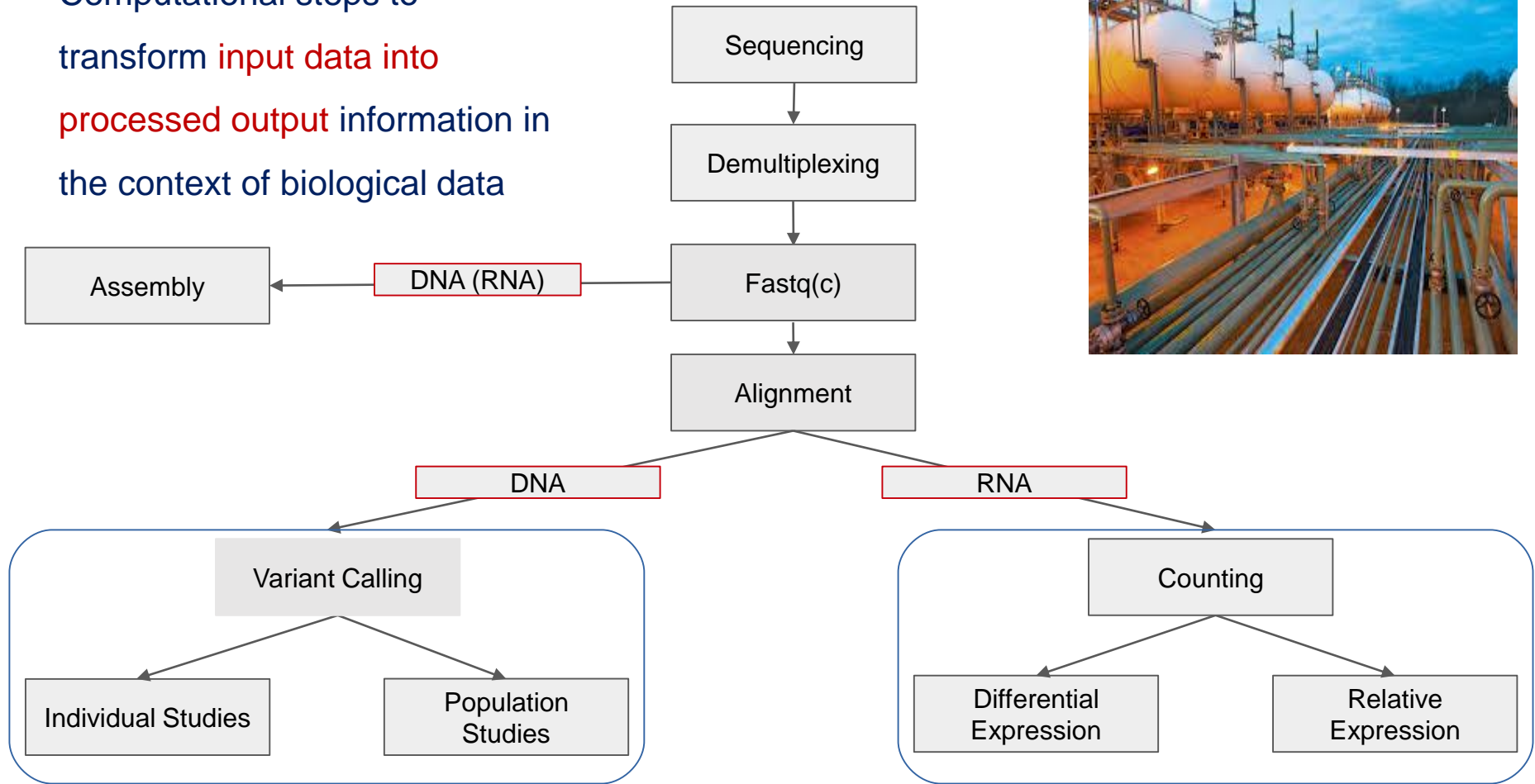
- UZ Leuven: Hydra, Google Cloud
- KU Leuven: VSC Flemish Super Computer

Overview

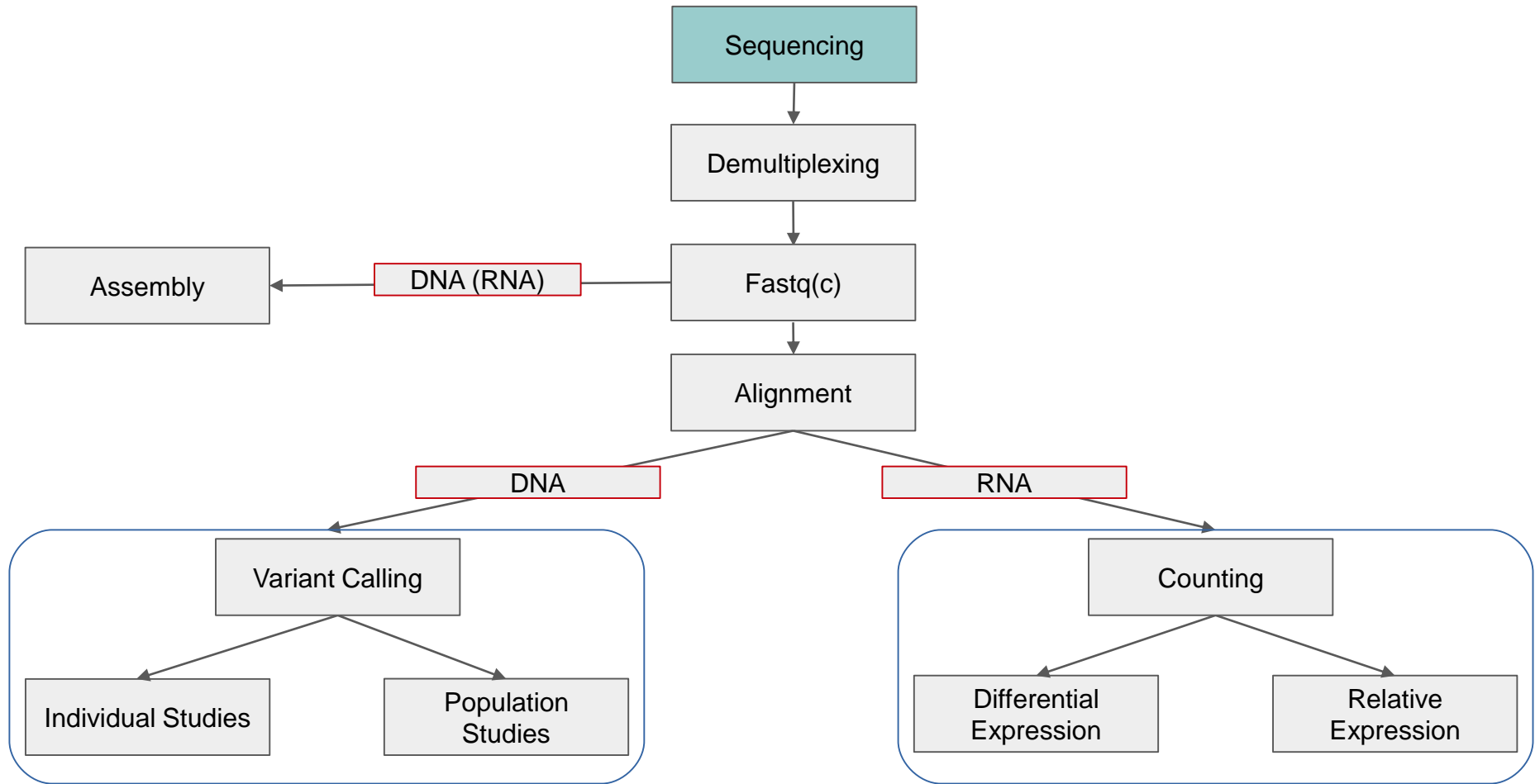
- Why this bioinformatics session?
- Basic NGS terminology
- **NGS Pipelines**
- Fastq format
- Fastq quality control

NGS Common Pipelines

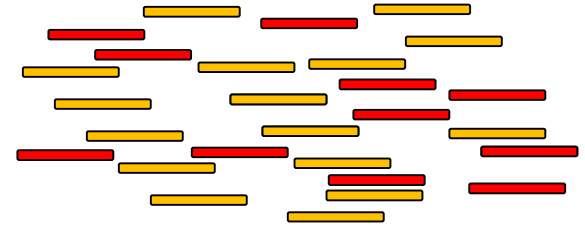
- Computational steps to transform **input data** into **processed output** information in the context of biological data



Sequencing



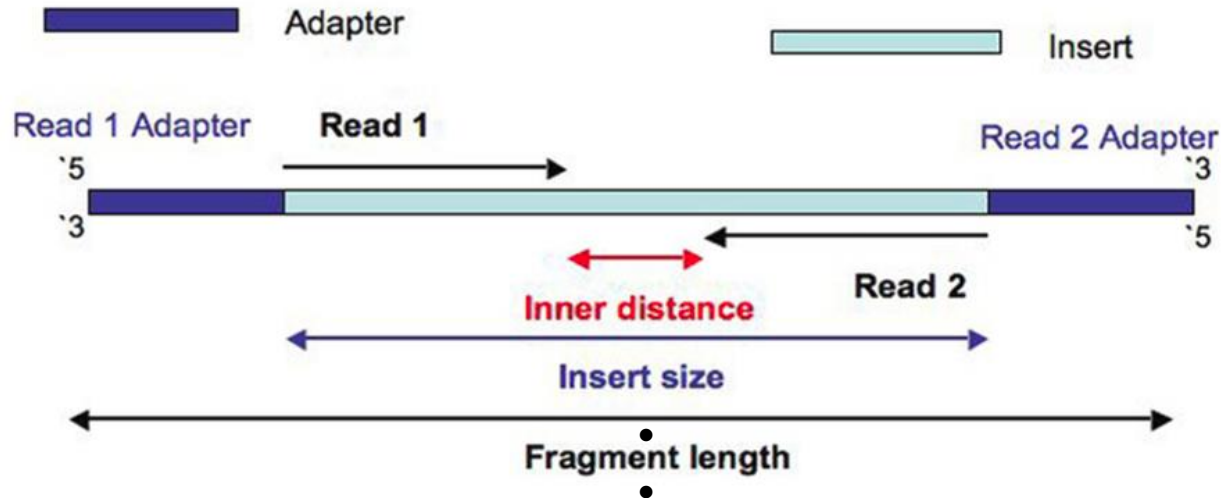
From DNA fragments to Digital reads



DNA Library

Reads

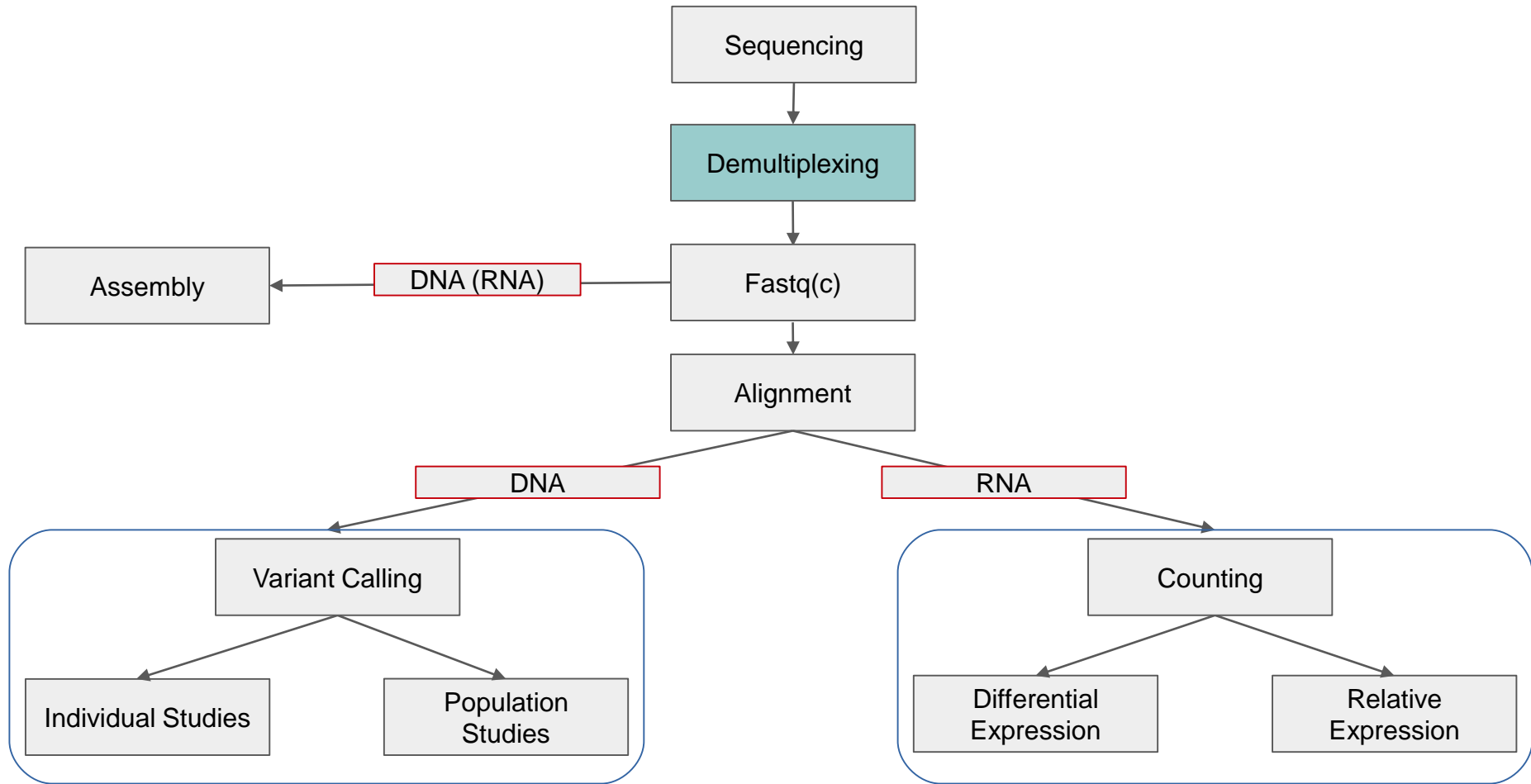
(c)DNA Fragments and Reads



Fragment: biological entity
Read: bioinformatics concept

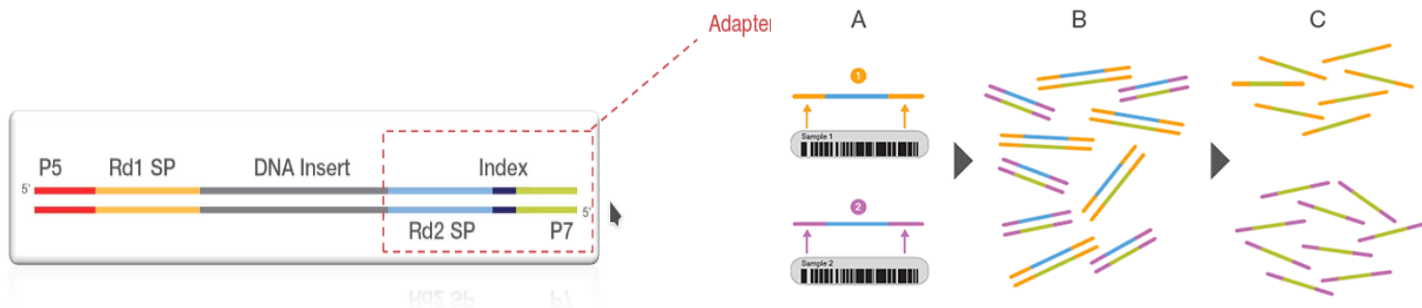
- **Fragment:** the DNA template + adapters that were loaded on the sequencing machine (is not completely sequenced)
- **Read:** a raw sequence originating from a sequencing machine
- **Single Read:** Sequencing only from one end
- **Paired-end:** Sequencing starting from both ends of the insert

Demultiplexing

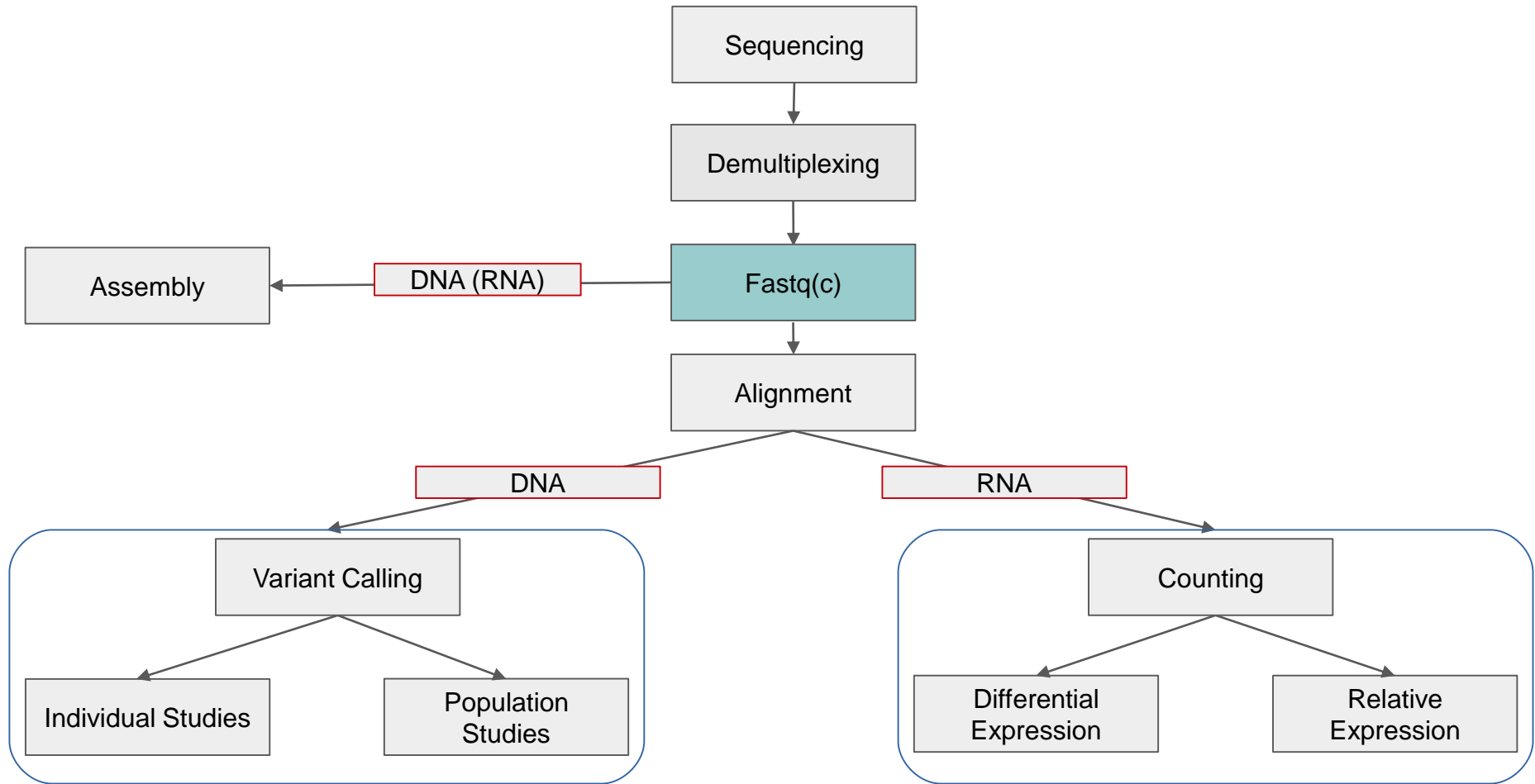


(De)Multiplexing

Multiple samples can be *pooled together or multiplexed* into one or more flowcells



Fastq Files



Fastq Files

The result of demultiplexing is one or more **fastq** files containing **raw reads**

Fastq files are:

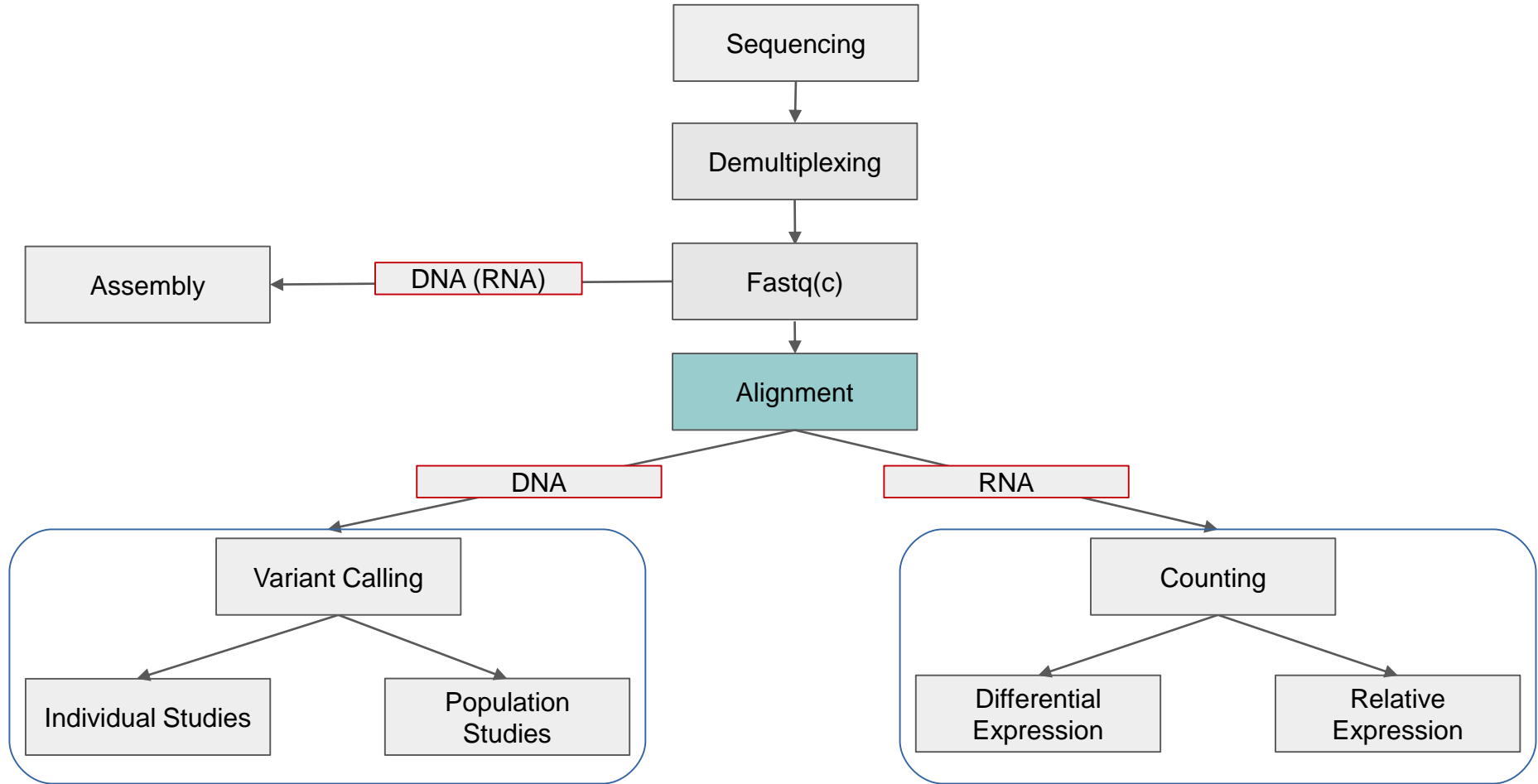
- Human **readable** (not binary) text files.
- Referred as **raw** data.
- Real **diamonds** NGS bioinformatics project
- **Not ordered** wrt originating (c)DNA
- Often compressed using **gzip**

```
-rwxr-xr-x 1 vsc31439 lp_biogenomics 32G 17 janv. 15:05 GC036462.R1.fastq.gz
-rwxr-xr-x 1 vsc31439 lp_biogenomics 39G 17 janv. 15:49 GC036462.R2.fastq.gz
-rwxr-xr-x 1 vsc31439 lp_biogenomics 31G 17 janv. 16:25 GC036463.R1.fastq.gz
-rwxr-xr-x 1 vsc31439 lp_biogenomics 37G 17 janv. 17:07 GC036463.R2.fastq.gz
-rwxr-xr-x 1 vsc31439 lp_biogenomics 31G 17 janv. 17:41 GC036464.R1.fastq.gz
-rwxr-xr-x 1 vsc31439 vsc31439 37G 17 janv. 18:21 GC036464.R2.fastq.gz
```

```
vsc31420@hpc-p-login-1 /staging/leuven/stg_00019/full_genomes/test_ws 11:48 $ ls -lha GC036463.R1.fastq
-rwxr-xr-x 1 vsc31420 vsc31420 134G 29 mars 11:24 GC036463.R1.fastq
```



Mapping and Alignment



Mapping and Alignment to Reference

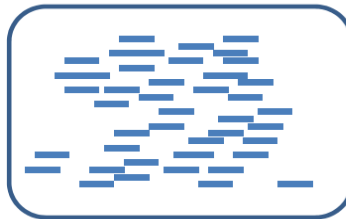
Mapping refers to the process of **aligning short reads to a reference** sequence, whether the reference is a complete genome, transcriptome, or *de novo* assembly.



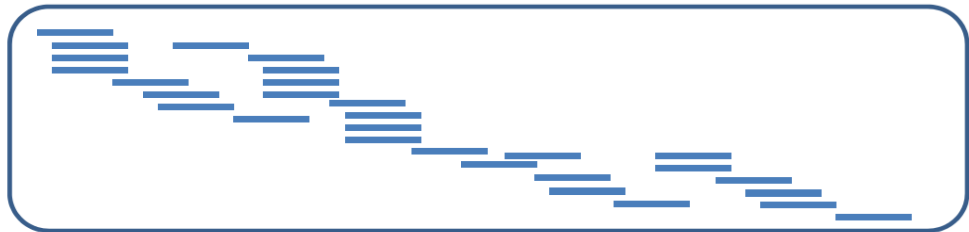
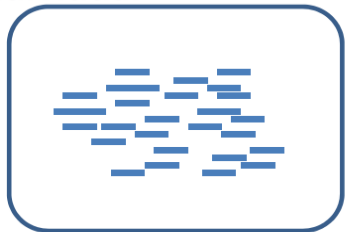
Sequencing Reads

Reference Genome

Individual A



Individual B



Reference Sequence in Fasta Format

- genome.fa **human-readable** nucleotide sequence
- **Species** dependent: Mouse genome: 2.6GB
- Evolves as it is **updated**

•
AATAAGTCAATGGCCTTTCTCTACACAAAGAATAAACAGGCTGAGAAAGAAATTAGGGAA
ACAACACCCTTCTCAATAGTCACAAATAATATAACATATCTCGGCGTGACTCTAACTAAG
GAAGTGAAAGATCTGTATGATAAAAACTTCAAGTCTCTGAAGAAAGAAATTAAAGAAGAT
CTCAGAAGATGGAAAGATCTCCCATGCTCATGGATTGGCAGGATCAATATTGTAAAAATG
GCTATCTTGCCAAAAGCAATCTACAGATTCAATGCAATCCCCATCAAATTCCAACCTCAA
TTCTTCAACGAATTAGAAGGAGCAATTTGCAAATTCATCTGTAATAACAAAAAACCTAGG
ATAGCAAAAAGTCTTCTCAAGGATAAAAAGAACCTCTGGTGGAATCACCATGCCTGACCTA
AAGCTTTACTACAGAGCAATTGTGGTAAAACTGCATGGTACTGGTATAGAGACAGACAA
GTAGACCAATGGAATAGAATTGAAGACCCAGAAATGAACCCACACACCTATGGTCACTTG
ATCTTCGACAAGGGAGCTAAAACCATCCAGTGGAAGAAAGACAGCATTTTCAACAAATGG
TGCTGGCACAACCTGGTTGTTATCATGTAGAAGAATGCGAATCGATCCATACTTATCTCCT
TGTAATAAGGTCAAATCTAAATGGATCAAAGAACTTCACATAAAACCAGAGACACTGAAA
CTTATAGAGGAGAAAGTGGGGAAAAGCCTTGAAGATATGGGCACAGGGGAAAAATTCCTG
AACAGAACAGCAATGGCTTGTGCTGTAAAGATTGAGAATTGACAAATGGGACCTAATGAAA
CTCCAAAGTTTCTGCAAGGCAAAAAGACACCGTCAATAAGAGAAAGAGACCACCAACAGAT
TGGGAAAGGATCTTTACCTATCCTAAATCAGATAGGGGACTAATATCCAACATATATAAA
GAACTCAAGAAGGTGGACTTCAGAAAATCAAACAACCCCATTAATAAATGGGGCTCAGAA
CTGAACAAAGAATTCTCACCTGAGTTATACCGAATGGCAGAGAAGCACCTGAAAAAATGC
TCAACATCCTTAATCATCAGGGAAATGCAAATCAAACAACCCCTGAGATTCCACCTCACA
CCAGTCAGAATGTCTAAGATCAAAAATTCAGGTGACAGCAGATGCTGGCGAGGATGTGGA
GAAAGAAGAACTCCTCCATTGTTGGTGGGATTGCAGGCTTGTACAACCACTCTGGAAA
TCCGTCTGGCGGTTCTCAGAAAATTGGACATAGTACTACCGGAGGATCCAGCAATACCT
CTCCTGGGCATATATCCAGAAGATGCCCAACTGGTAAGAAGGACACATGCTCCACTATG
TTCATAGCAGCCTTATTTATAATAGCCAGAAGCTGGAAAGAACCCAGATGCCCTCAACA
GAGGAATGGATACAGAAAATGTGGTACATCTACACAATGGAGTACTACTCAGCTATTAAA
AAGAATGAATTTATGAAATTCCTAGCCAAATGGATGGACCTGGAGGGCATCATCCTGAGT

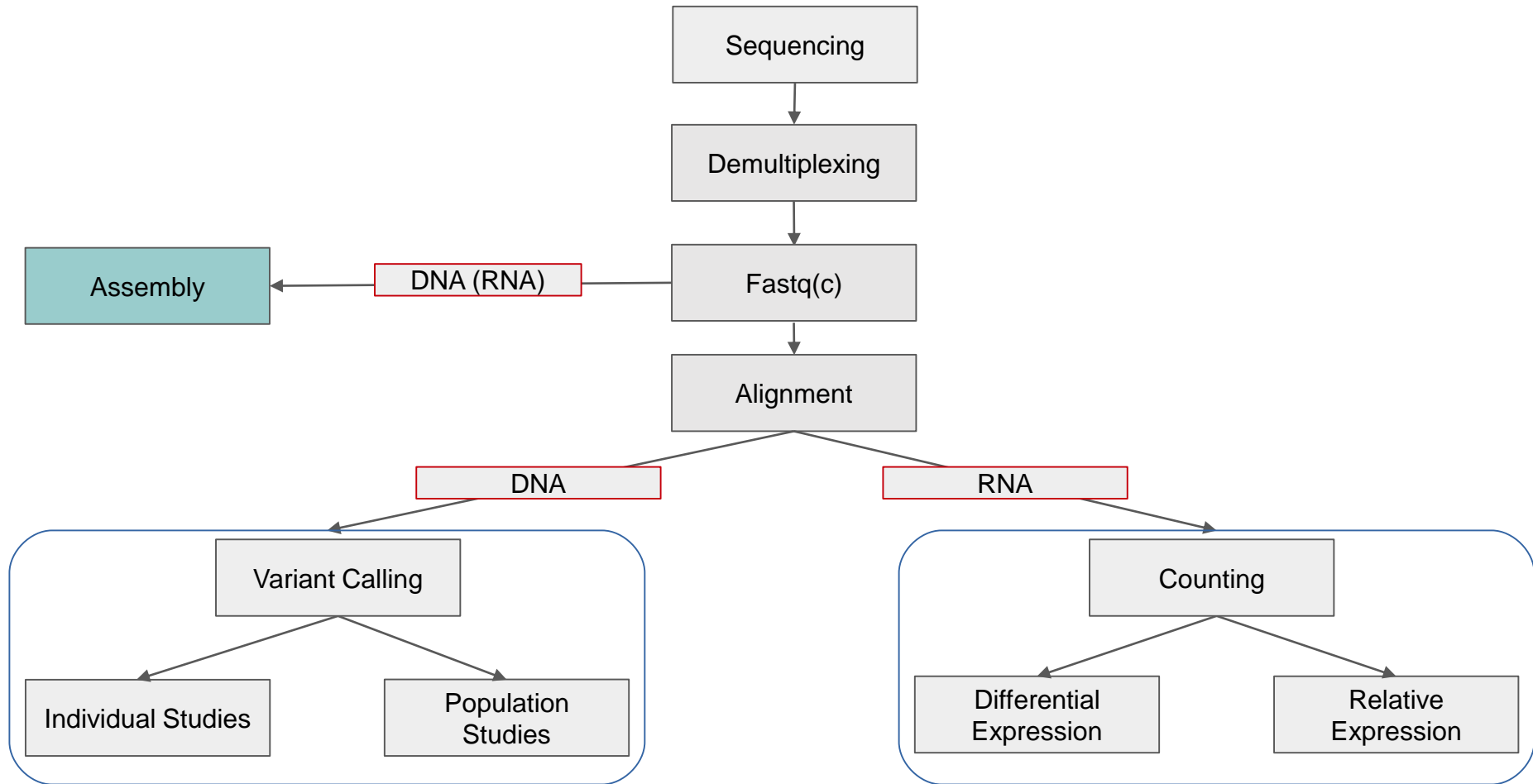
Mapping vs Alignment

- Reference **Ch1**:

1234567890123456
ATGGTTACACCATT
 - Read:

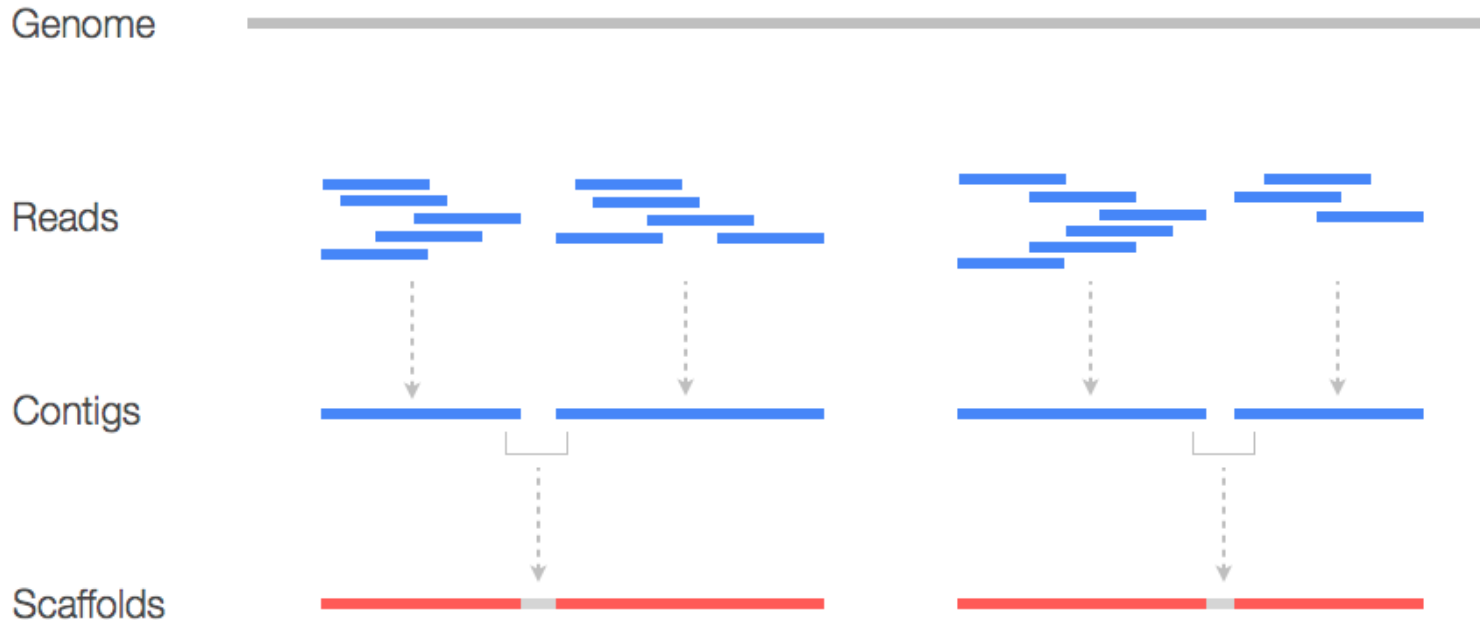
GGTTCA
 - Possible alignment:
ATGGTTACACCATT
GGTT-CA
 - Mapping: **Ch1-pos3**
- *“Also of note is that by this time the terms “read alignment” and “read mapping” had become interchangeable. The BWA and Bowtie papers both used both terms, as did many other papers.”*
 - <https://liorpachter.wordpress.com/2015/11/01/what-is-a-read-mapping>

Assembly



Assembly

- The **generation of a reference**, from scratch (*de novo*) or reference assisted.
- Overlapping reads are merged to **contigs** (smallest unitable unit without unknown bases)
- **Contigs** that belong together, but where the connecting sequence is unknown, can be connected to scaffolds, inserting N's for the unknown bases



Overview

- Why this bioinformatics session?
- Basic NGS terminology
- NGS Pipelines
- Fastq format
- **Fastq quality control**

Fastq Format and Reads

- What is a “read”?
 - A raw sequence (ordered collection) of nucleotides names A,C,G,T, or N.
 - Read length is **experiment-dependent**
 - mRNA differential expression: 51bp SE
 - Whole-Genome sequencing: 150bp PE
- **Fastq** file?
 - Plain-text file, where each read and complementary information occupies 4 consecutive lines

[illegible]

Raw Reads

[illegible]

Line 1: Read identifier and is followed by a sequence that is

Unique for each read, platform dependent
Begins with a '@' character

@HISEQ:574:C6VG2ANXX:2:1110:1400:2194 1:N:0:ATCACG

Instrument Run, flowcellID, lane, tile, pos, read, filter, control Index

Line 2: Raw sequence of nucleotides

Line 3: begins with a '+' character and is optionally followed by the same sequence identifier.

Line 4: Quality values for the sequence in Line 2

Raw Reads: Base Calling Quality

```
@HISEQ:574:C6VG2ANXX:2:1110:1400:2194 1:N:0:ATCACG  
GGGGGATTCTCACTAGGTCTCAAGTCTCTCACTCTCGGTAGTGTTCCAG  
+  
CCCCCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG
```

- A quality score (**Q-score**) is a prediction of the probability of an error in base calling.
- It serves as a compact way to communicate very small **error** probabilities
 - $P = 10^{(-Q/10)}$
 - $Q = -10 \log_{10}(P)$

ASCII_BASE=33 Illumina, Ion Torrent, PacBio and Sanger											
Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII
0	1.00000	33 !	11	0.07943	44 ,	22	0.00631	55 7	33	0.00050	66 B
1	0.79433	34 "	12	0.06310	45 -	23	0.00501	56 8	34	0.00040	67 C
2	0.63096	35 #	13	0.05012	46 .	24	0.00398	57 9	35	0.00032	68 D
3	0.50119	36 \$	14	0.03981	47 /	25	0.00316	58 :	36	0.00025	69 E
4	0.39811	37 %	15	0.03162	48 0	26	0.00251	59 ;	37	0.00020	70 F
5	0.31623	38 &	16	0.02512	49 1	27	0.00200	60 <	38	0.00016	71 G
6	0.25119	39 '	17	0.01995	50 2	28	0.00158	61 =	39	0.00013	72 H
7	0.19953	40 (18	0.01585	51 3	29	0.00126	62 >	40	0.00010	73 I
8	0.15849	41)	19	0.01259	52 4	30	0.00100	63 ?	41	0.00008	74 J
9	0.12589	42 *	20	0.01000	53 5	31	0.00079	64 @	42	0.00006	75 K
10	0.10000	43 +	21	0.00794	54 6	32	0.00063	65 A			

Raw Reads: Base Calling Quality

@HISEQ:574:C6VG2ANXX:2:1110:1400:2194 1:N:0:ATCACG
GGGGATTCTCACTAGTCTCAAGGTCTCTCACTCTCGGTAGTGTTCCCAG
+
CCCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG

- Base: G
- Quality: C
- ASCII: 67
- Q: 34
- P: 0.00040

- Base: G
- Quality: G
- ASCII: 71
- Q: 38
- P: 0.00016

ASCII BASE=33 Illumina, Ion Torrent, PacBio and Sanger

Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII
0	1.00000	33 !	11	0.07943	44 ,	22	0.00631	55 7	33	0.00050	66 B
1	0.79433	34 "	12	0.06310	45 -	23	0.00501	56 8	34	0.00040	67 C
2	0.63096	35 #	13	0.05012	46 .	24	0.00398	57 9	35	0.00032	68 D
3	0.50119	36 \$	14	0.03981	47 /	25	0.00316	58 :	36	0.00025	69 E
4	0.39811	37 %	15	0.03162	48 0	26	0.00251	59 ;	37	0.00020	70 F
5	0.31623	38 &	16	0.02512	49 1	27	0.00200	60 <	38	0.00016	71 G
6	0.25119	39 '	17	0.01995	50 2	28	0.00158	61 =	39	0.00013	72 H
7	0.19953	40 (18	0.01585	51 3	29	0.00126	62 >	40	0.00010	73 I
8	0.15849	41)	19	0.01259	52 4	30	0.00100	63 ?	41	0.00008	74 J
9	0.12589	42 *	20	0.01000	53 5	31	0.00079	64 @	42	0.00006	75 K
10	0.10000	43 +	21	0.00794	54 6	32	0.00063	65 A			

Raw Reads: Base Calling Quality

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.....
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
.....IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII.....
.....JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.....
..LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.....
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMN...
|               |       |               |                               |
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)

Raw Reads: Base Calling Quality

Checking individual base
calling manually is
rarely done

Fastq file quality control

- How many reads do I have ?
- Is that enough ?
- How good is the quality of the data ?
- Did I sequence what I wanted to sequence ?
- Is pre-processing needed ?
-

Example 1 – targeted sequencing

- Finding the genetic cause of a disease

- ~ 6,000 genes
- Illumina PE 126bp

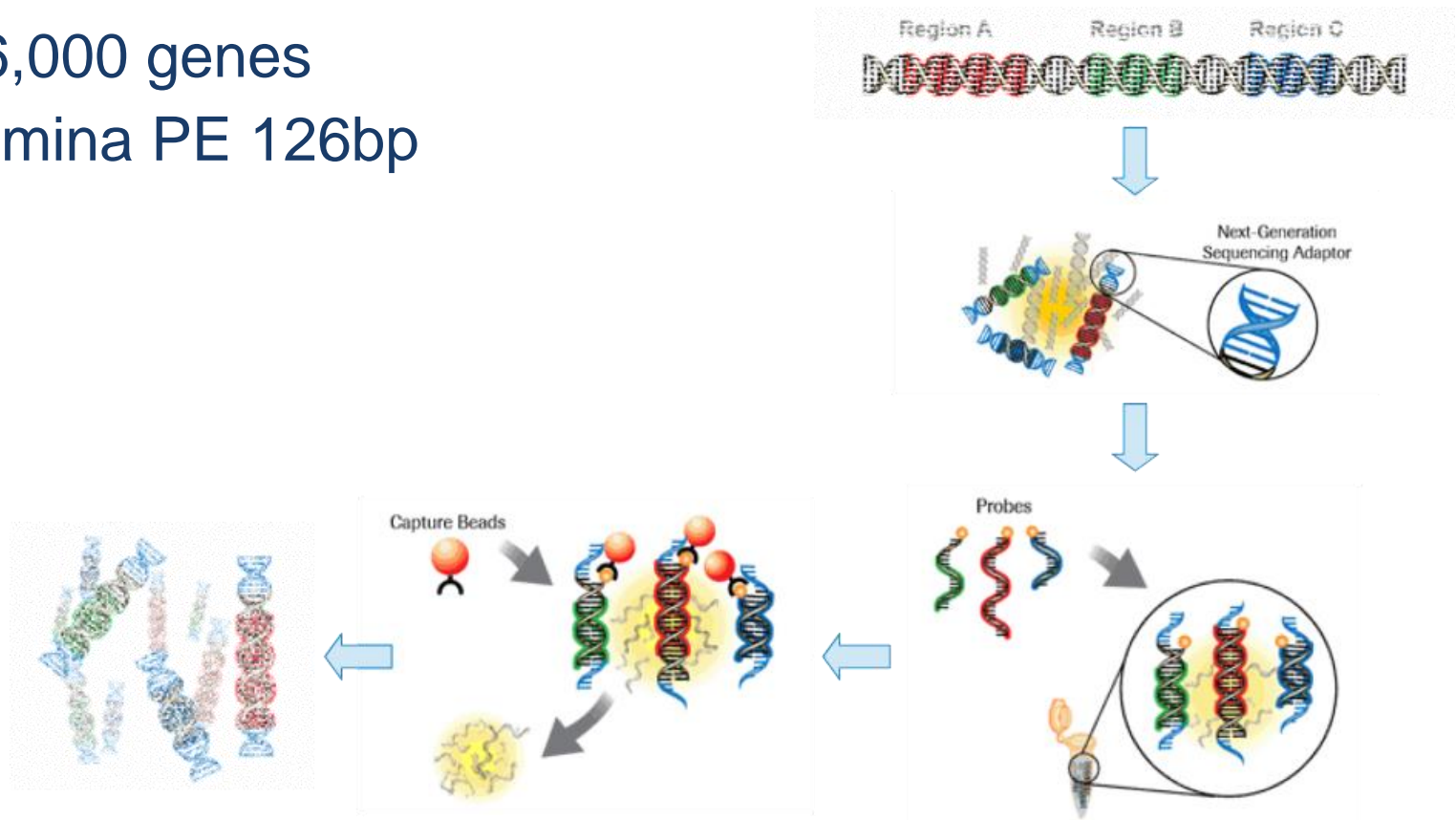
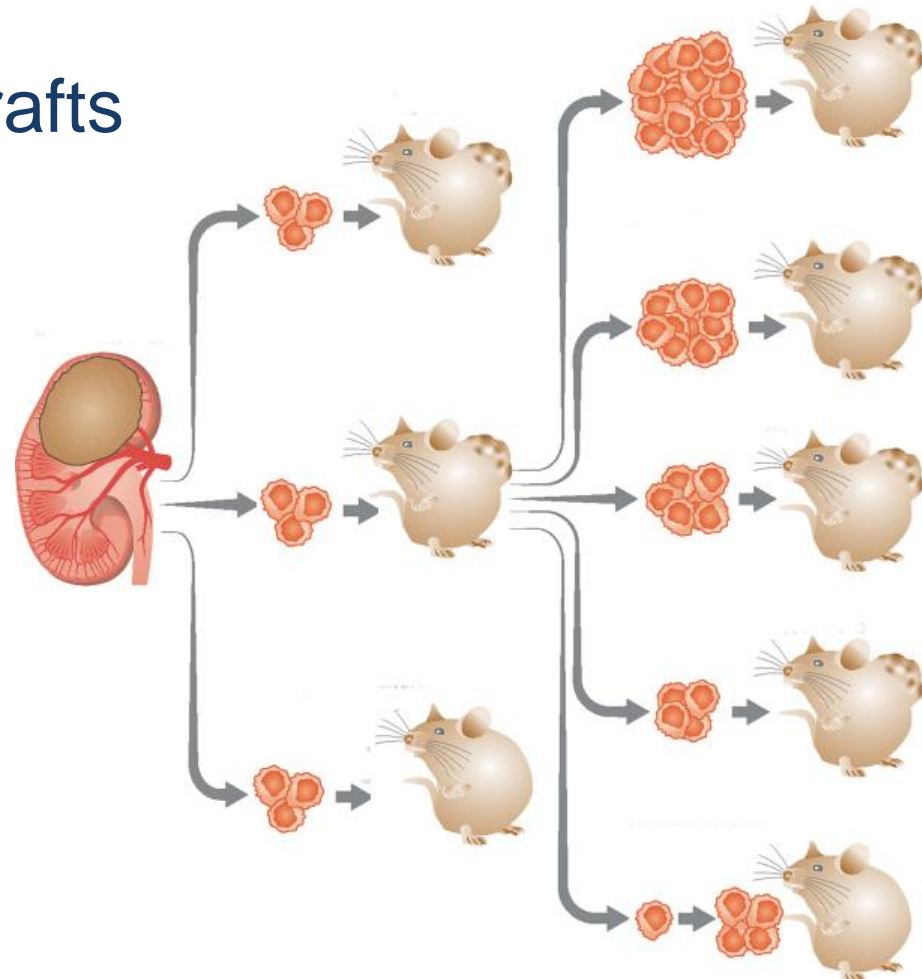


Figure adapted from Nimblegen

Example 2 – amplicon sequencing

- Fingerprinting of xenocrafts

- 31 SNPs
- Illumina PE 151bp

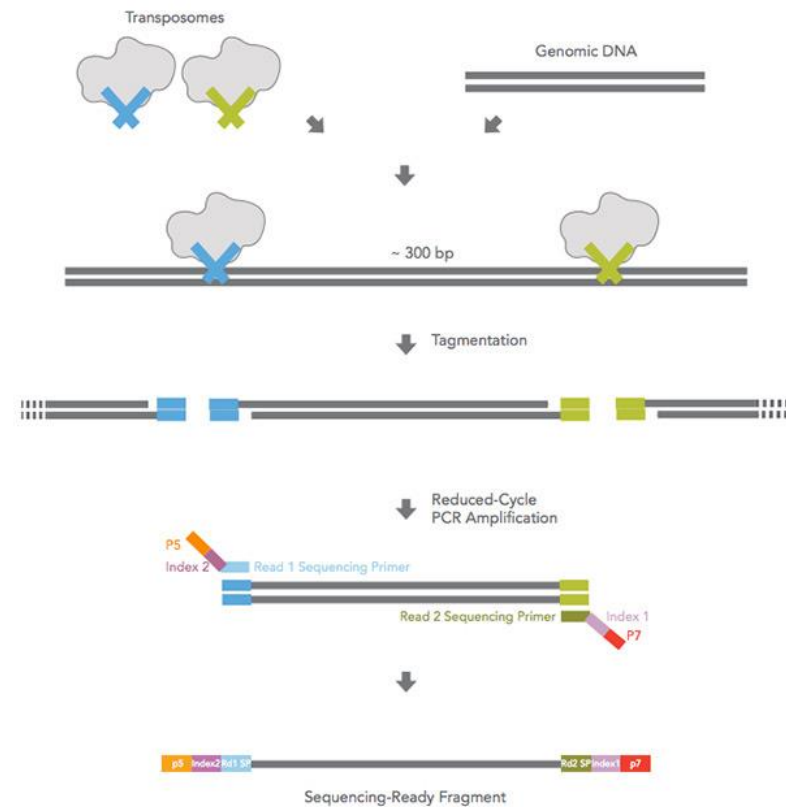


Trace platform, <http://www.uzleuven-kuleuven.be/lki/trace/>

Figure adapted from Peter Hohenstein, EMBO Molecular Medicine: 5 (1), 2013

Example 3 – whole genome sequencing

- Predicting bacterial resistance
 - Whole Genome Sequencing (WGS)
 - *Streptococcus pneumoniae*
 - *Mycobacterium tuberculosis*
 - Illumina PE 301 bp

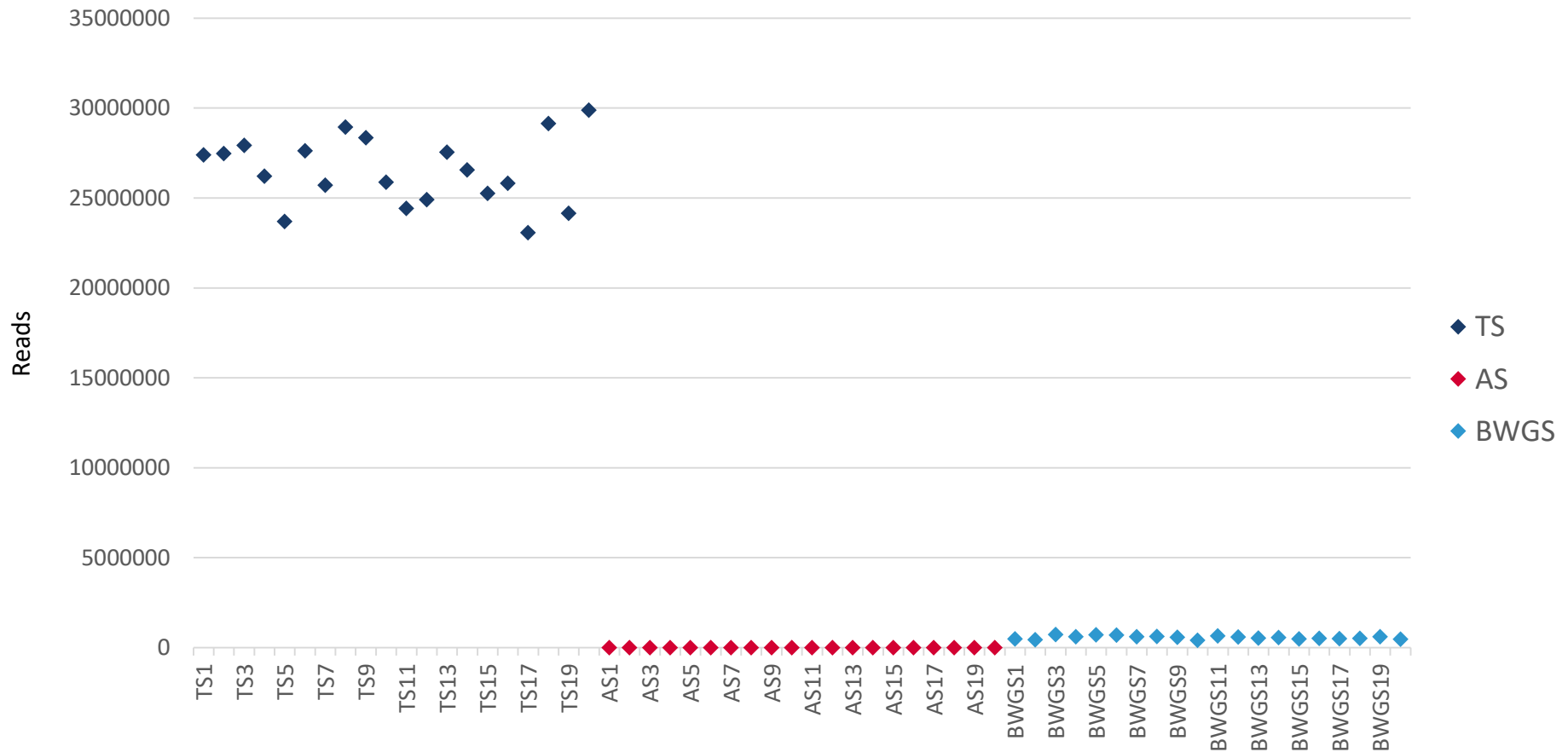


Number of reads

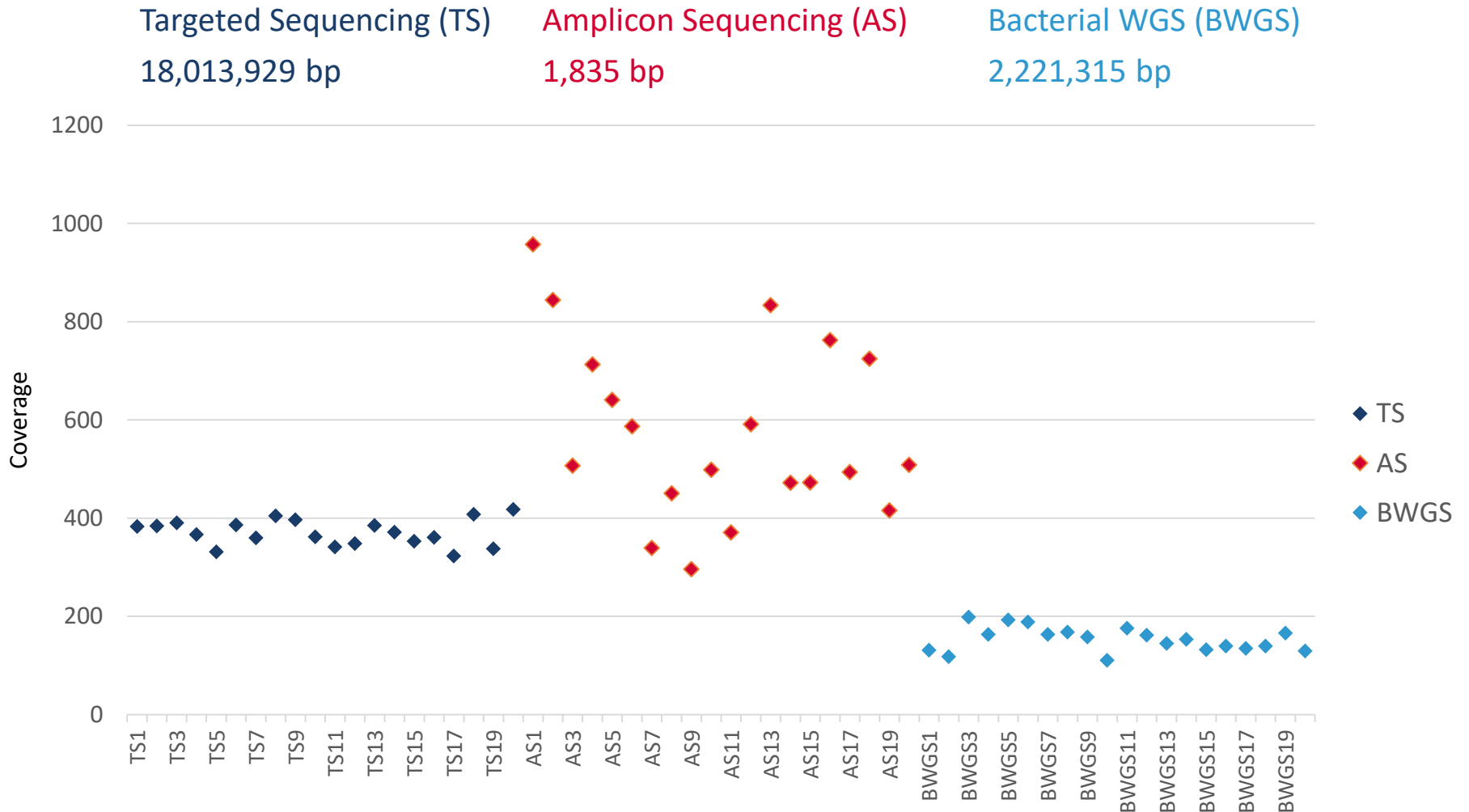
Targeted Sequencing (TS)
18,013,929 bp

Amplicon Sequencing (AS)
1,835 bp

Bacterial WGS (BWGS)
2,221,315 bp



Estimated coverage



Fastq files – QC

- FastQC
 - Check Phred quality scores
 - Check GC content
 - Check read content
 - ...

<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
<https://multiqc.info>

FastQC

- Summary report

```
fastqc -o result sample.R1.fastq.gz
```

```
fastqc -o result sample.R2.fastq.gz
```

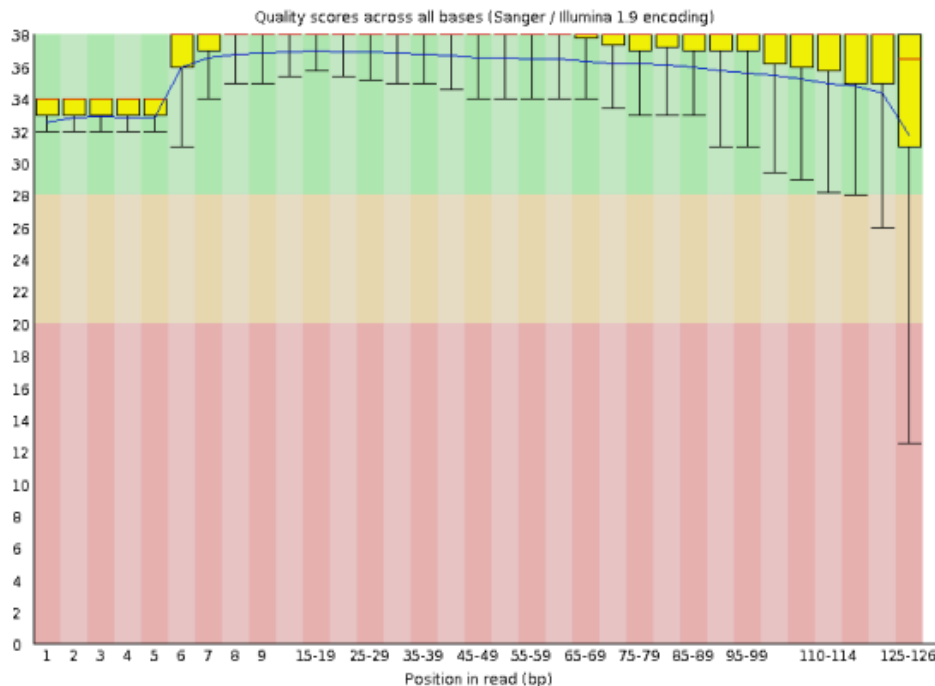
FastQC Report

Summary

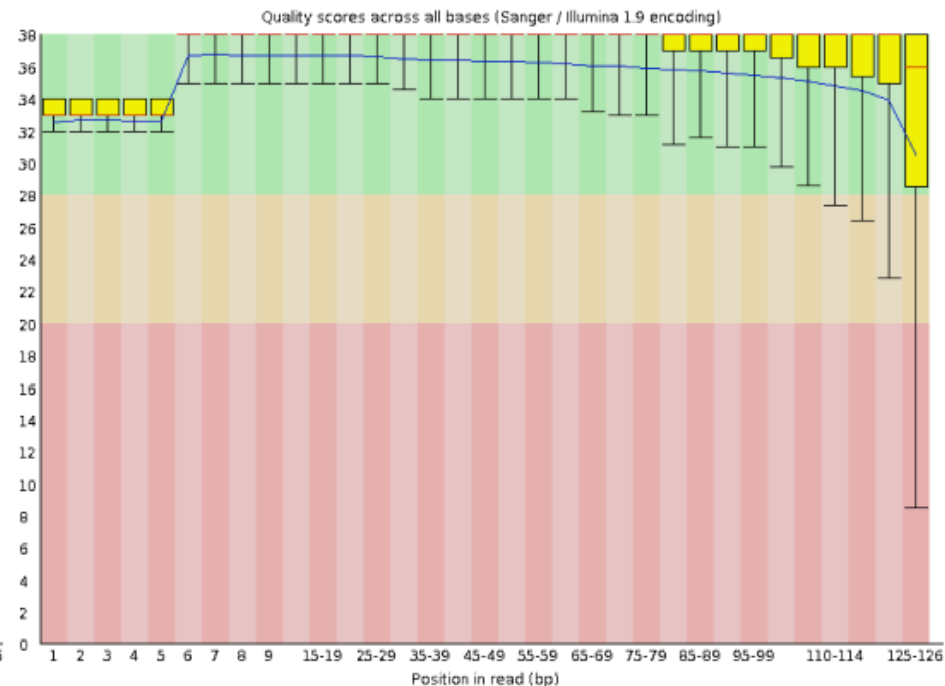
- ✓ [Basic Statistics](#)
- ✓ [Per base sequence quality](#)
- ✓ [Per sequence quality scores](#)
- ✓ [Per base sequence content](#)
- ✓ [Per base GC content](#)
- ✗ [Per sequence GC content](#)
- ✓ [Per base N content](#)
- ✓ [Sequence Length Distribution](#)
- ✓ [Sequence Duplication Levels](#)
- ✓ [Overrepresented sequences](#)
- ! [Kmer Content](#)

FastQC – Phred quality score by position

- Example 1 – targeted sequencing



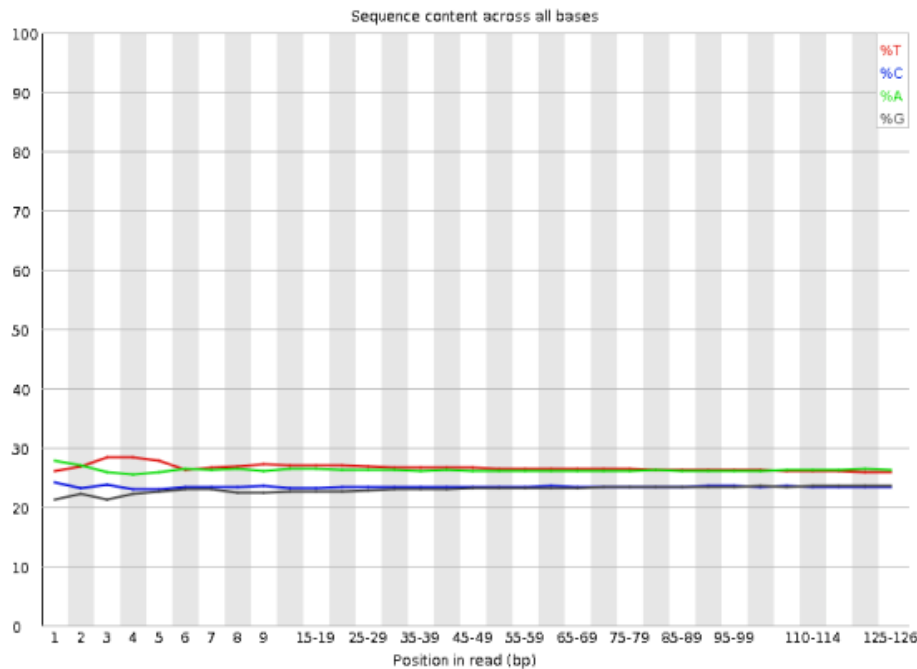
R1



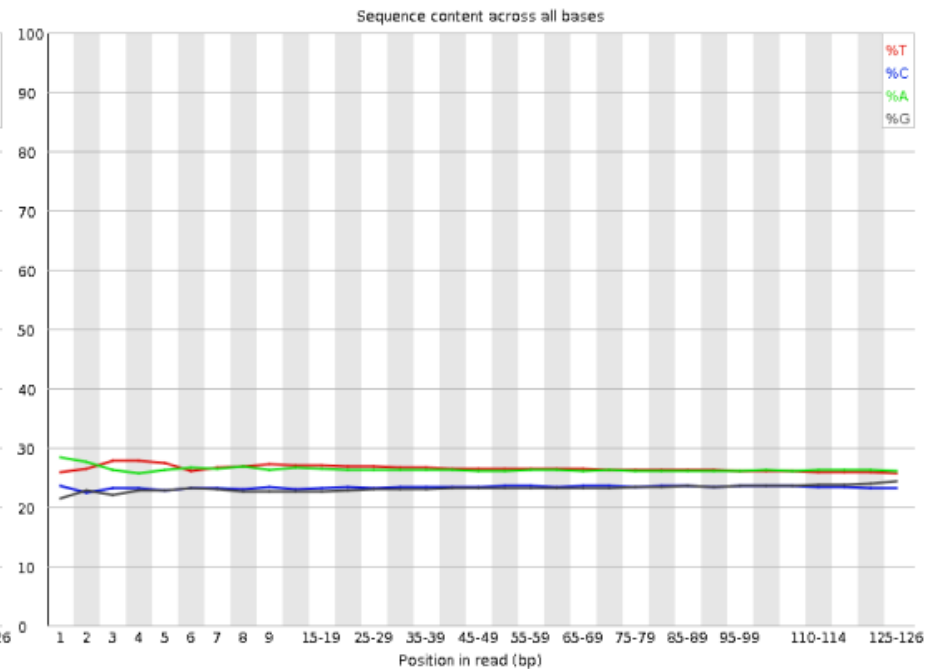
R2

FastQC – Base content by position

- Example 1 – targeted sequencing
 - G-C 25-26%
 - A-T 24-25%



R1



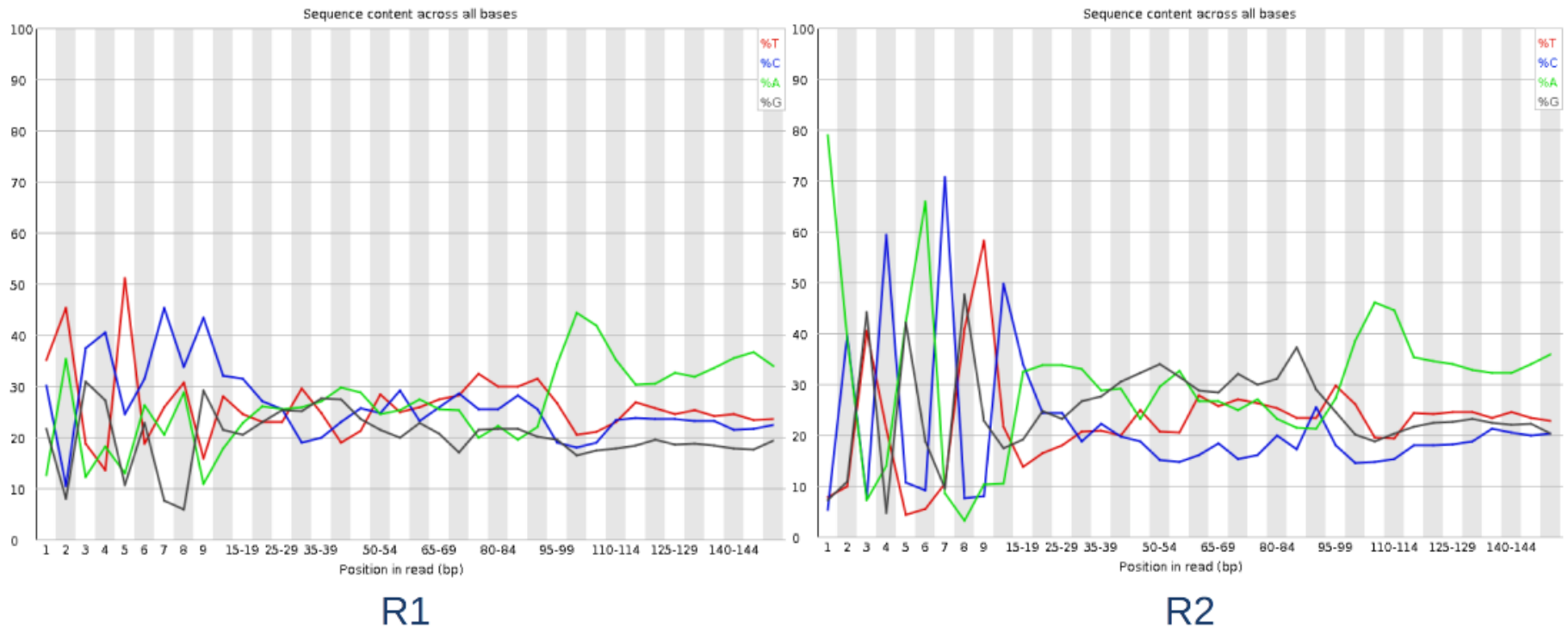
R2

FastQC – Base content by position

- Example 2 – amplicon sequencing

- G-C 25-27%

- A-T 24-24%



FastQC – Over-represented sequences

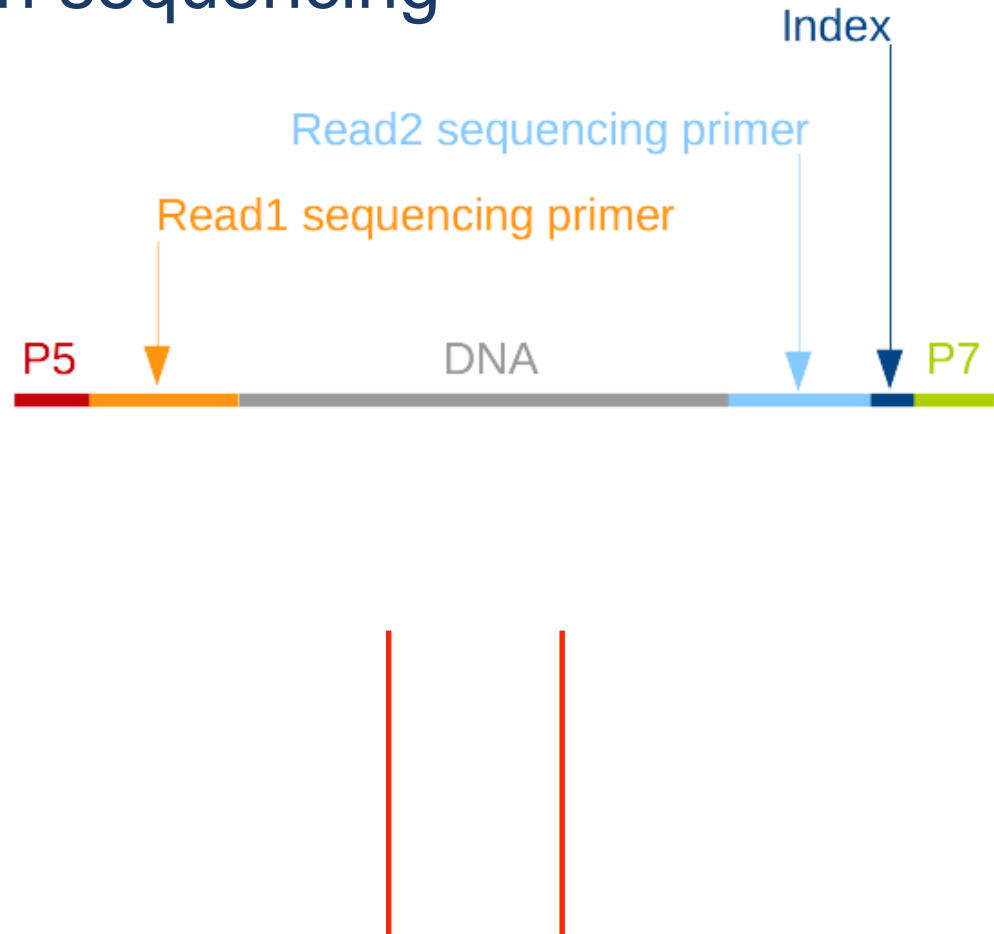
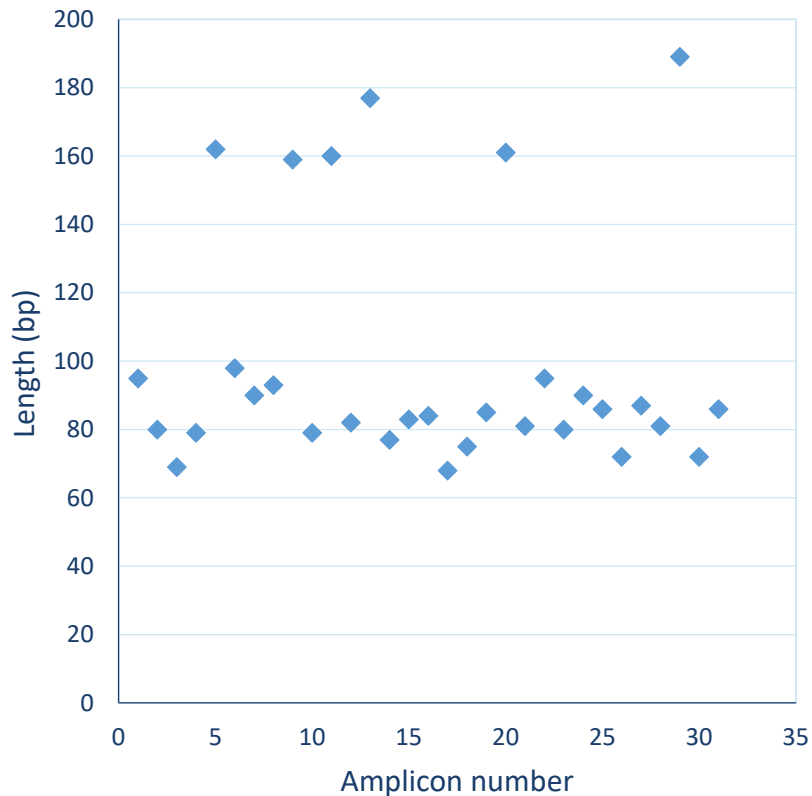
- Example 2 – amplicon sequencing

✖ Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GTGATCTCCAACCTTTGACCTGACCGTCGCTTAGATCGGAAGGACACAGT	92	5.2421652421652425	No Hit
CAGTGACACTAGTCTGCAACAACGCCACTTAGATCGGAAGGACACAGT	88	5.014245014245014	No Hit
TTCTCCTCAGTGCCTGACCGTCGCTTAGATCGGAAGGACACAGTCTGAA	52	2.9629629629629632	Illumina Multiplexing PCR Primer 2.01 (100% over 24bp)
TTCTCCTCAGTGCCTGCAACAACGCCACTTAGATCGGAAGGACACAGT	49	2.792022792022792	No Hit
CTCAGTACAGCTGACCTGACCGTCGCTTAGATCGGAAGGACACAGTCT	31	1.7663817663817662	Illumina Multiplexing PCR Primer 2.01 (100% over 21bp)
CAGCCTCTGCTCTCAGCTGACCTGACCGTCGCTTAGATCGGAAGGACACA	26	1.4814814814814816	No Hit
TTCTCCTCAGTGCCTGCACTCAATCATGCTCTCTAGATCGGAAGGACAC	25	1.4245014245014245	No Hit
TTCTCCTCAGTGCCTGCACTCTCTCAGCTCACCCTGCACTCTCTCTCAC	22	1.2535612535612535	No Hit
CTCAGTACAGCTGACACAAGTACGACCTAGGACCACTTGAATAGAGAGCTCAGT	22	1.2535612535612535	No Hit
TTCTCCTCAGTGCCTGCACTCTCTCAGCTCTCTCTGCACTCAATCATC	21	1.1965811965811968	No Hit
GACCAAGAAAGCTGACCTGACCGTCGCTTAGATCGGAAGGACACAGTCT	21	1.1965811965811968	No Hit
CTCAGTACAGCTGACCGTCGCTTAGATCGGAAGGACACAGTCTGAACT	19	1.0826210826210827	Illumina Multiplexing PCR Primer 2.01 (100% over 26bp)
TTCTCCTCTTGTATGTGACCGTCGCTTAGATCGGAAGGACACAGTCTGAA	14	0.7977207977207977	Illumina Multiplexing PCR Primer 2.01 (100% over 23bp)
GTACAGCTGACAGTGAAGATCGGAAGGACACAGTCTGAACTCCAGTCA	14	0.7977207977207977	Illumina Multiplexing PCR Primer 2.01 (100% over 33bp)
TTCTCCTCAGTGCCTGCACTCAATGAATGTTTTTATAAAAGGCTGTGGC	12	0.6837606837606838	No Hit
AGGTAAAGTGACAGTTTGTCTCATGGAAAGGAGATAGATCGGAAGGACACA	12	0.6837606837606838	No Hit
GTGATCTCCAACCTTTGACCGTCGCTTAGATCGGAAGGACACAGTCTGAA	11	0.6267806267806267	Illumina Multiplexing PCR Primer 2.01 (100% over 24bp)
CAAGAGCTCAGAGGAGGAGCTGTCAGAGATCGGAAGGACACAGTCTGAA	11	0.6267806267806267	Illumina Multiplexing PCR Primer 2.01 (100% over 23bp)
TTTGACTTGTACCTGACCGTCGCTTAGATCGGAAGGACACAGTCTGAA	11	0.6267806267806267	Illumina Multiplexing PCR Primer 2.01 (100% over 24bp)
CTCAGTACAGCTGACACTTAAAGTCGGGAGTCAGAAAGTACCAAGGAG	9	0.5128205128205128	No Hit
TTGGTGTACATGTGTTGTGTGTGTGTGTGGGGGAAAGTTGAGTAGATCG	9	0.5128205128205128	No Hit
TTCTCCTCAGTGCCTGCACTCGATAATTCAATACATAATATTCAATAATT	9	0.5128205128205128	No Hit
GTACAGCTGACAGTGAAGAACAGATCGGAAGGACACAGTCTGAACTCCAG	9	0.5128205128205128	Illumina Multiplexing PCR Primer 2.01 (100% over 30bp)
TTCTCCTCAGTGCCTGACCGTCGCTTAGATCGGAAGGACACAGTCTGAA	8	0.4558404558404558	Illumina Multiplexing PCR Primer 2.01 (100% over 24bp)
TTGGTGTACATGTGTTGTGTGTGTGTGGGGGAAAGTTGAGTAGATCGGAAG	8	0.4558404558404558	No Hit
CATCTGCATGGTATCTGGGCTCTGTAGTGGTGGCTGCAAGAGGTGCT	8	0.4558404558404558	No Hit
CATTTCATTGCAACCGAGTCCATTGTGACAGATGGAAGACAGCAGT	8	0.4558404558404558	No Hit
GATGTTCAAGGATCCCAAGTATGAGTAAACCCCTTATGATCAGTCACTAT	7	0.39886039886039887	No Hit
AGGTAAAGTGACAGTTTGTCTCAGGGAAGGTGAGATTGGATTCTTTAAAC	7	0.39886039886039887	No Hit
TGGCCTTGACAAACAGATCGGAAGGACACAGTCTGAACTCCAGTCAAG	7	0.39886039886039887	TruSeq Adapter, Index 2 (97% over 35bp)
TTCTCCTCAGTGCCTGCACTCTCTCAGCTCACCCTCATCAGCTCACC	7	0.39886039886039887	No Hit
ACACTGGGCTAGACACTCGTATGGTGTATGGGTTTCTCTTCTTAGAGA	7	0.39886039886039887	No Hit
TTTGACTTGTACCTGGGCGCATGTTCAATTTTCAGTTGTGGATAGCAC	7	0.39886039886039887	No Hit
AAGAGCTCGCTGACCGTCGCTTAGATCGGAAGGACACAGTCTGAACTC	6	0.3418803418803419	Illumina Multiplexing PCR Primer 2.01 (100% over 27bp)

FastQC – Over-represented sequences

- Example 2 – amplicon sequencing



Contamination check

- FastQ Screen

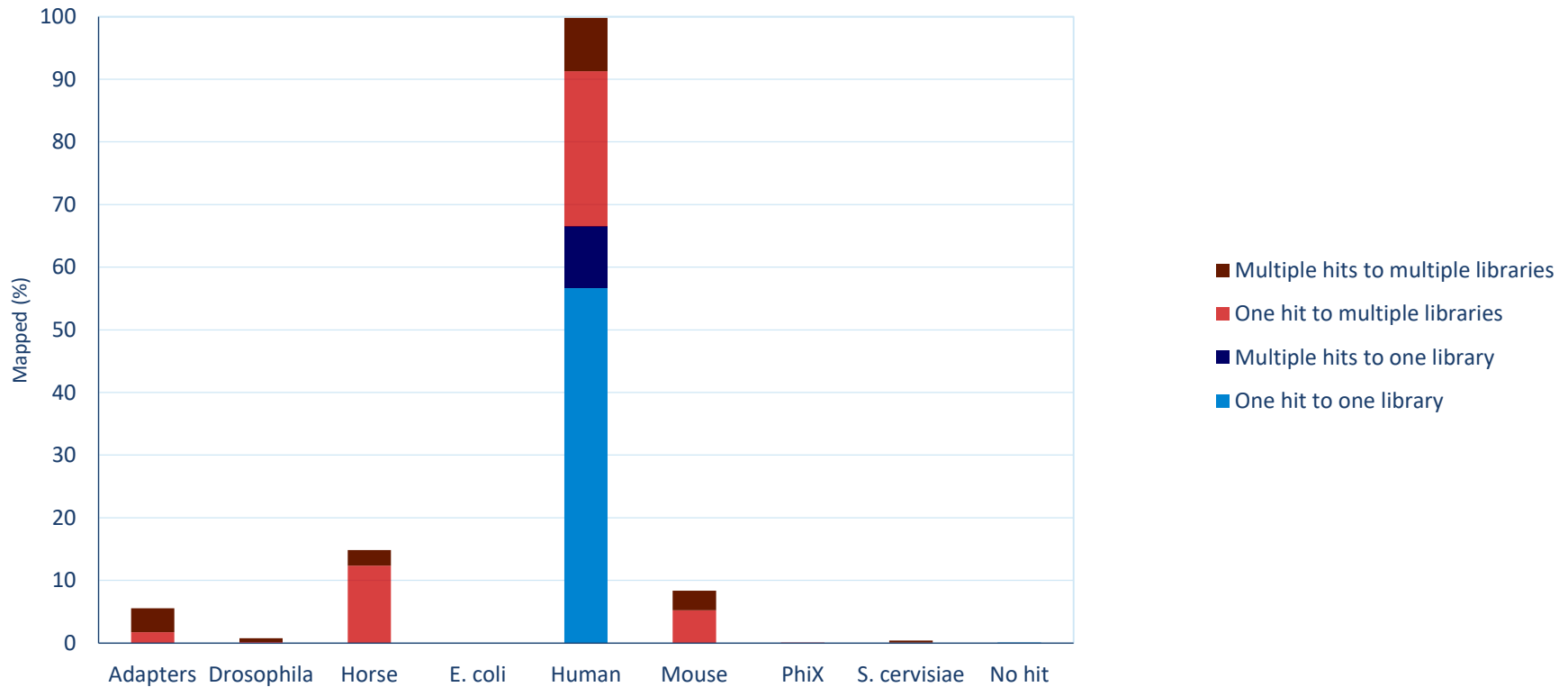
- Compare reads to various libraries
- Any library can be searched against
- Output proportion of reads with
 - One hit to one library
 - Multiple hits to one library
 - One hit to multiple libraries
 - Multiple hits to multiple libraries

```
fastq_screen --subset 100000 --conf fastq_screen.conf --aligner  
bowtie2 --outdir result --nohits sample.R1.fastq.gz
```

https://www.bioinformatics.babraham.ac.uk/projects/fastq_screen/

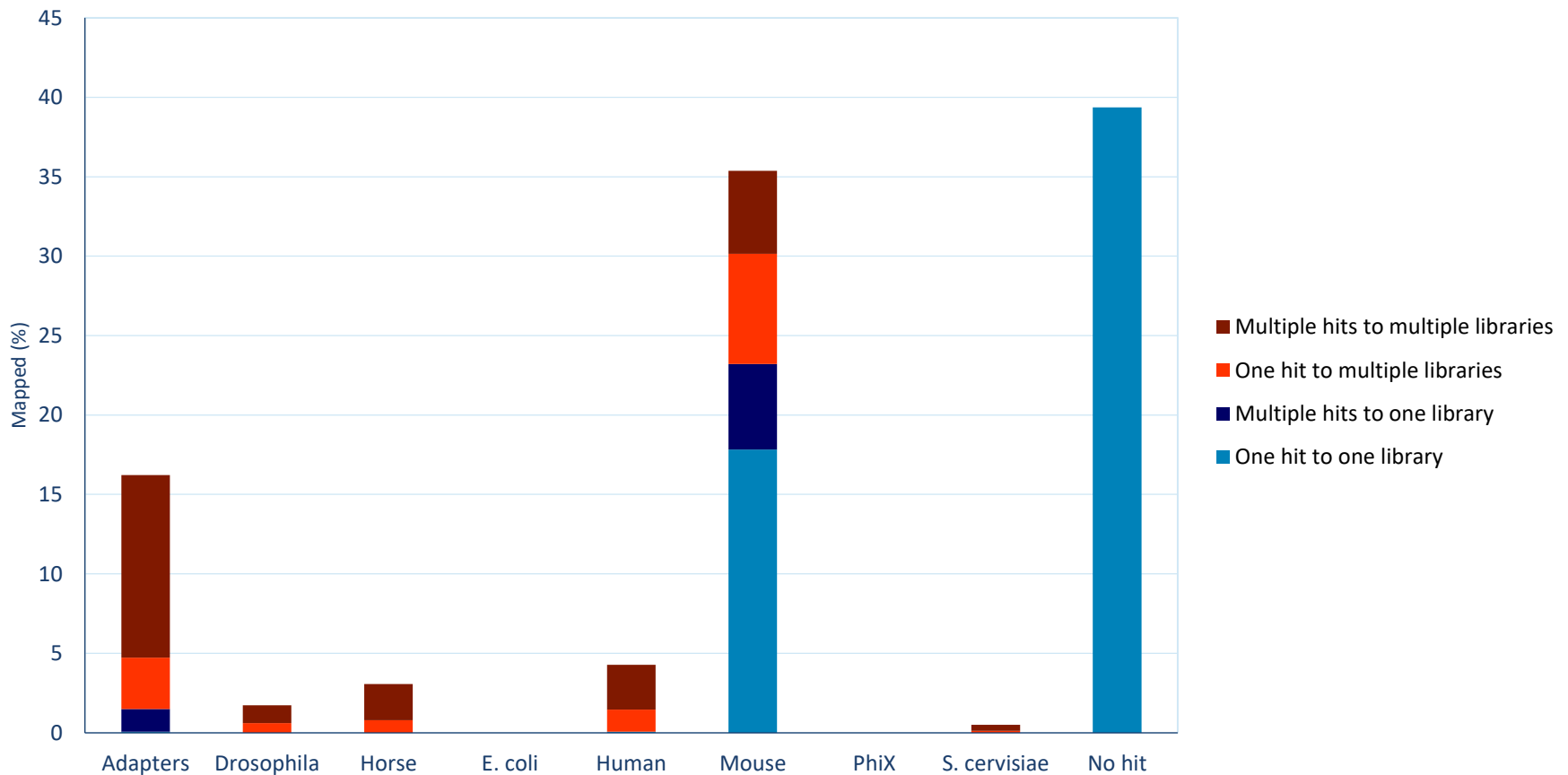
Contamination check – FastQ Screen

- Example 1 – targeted sequencing



Contamination check – FastQ Screen

- Example 2 – amplicon sequencing



Contamination check – kraken2

- Taxonomic sequence classification system
 - Build (custom) database
 - Compare k-mer from reads to database

```
kraken2-build --download-taxonomy --db database
kraken2-build --download-library library --db database
kraken2-build --build --db database --minimizer-spaces 0
kraken2 --db database --paired sample.R1.fastq.gz
sample.R2.fastq.gz -report kraken2Report.txt --use-names >
kraken2.output.txt
```

<https://ccb.jhu.edu/software/kraken2/>

Contamination check – kraken2

• Example 3 – bacterial WGS

Fragments covered by clade (%)	Fragments covered by clade	Fragments assigned to taxon	Rank code	NCBI taxonomic ID	scientific name
0.04	243	243 U		0	Unclassified
99.96	632527	0 R		1	root
99.96	632527	0 R1		131567	cellular organisms
99.96	632491	121 D		2	Bacteria
99.55	629954	1 D1		1783272	Terrabacteria group
99.53	629777	0 P		201174	Actinobacteria
99.53	629777	17 C		1760	Actinobacteria
99.52	629754	20 O		85007	Corynebacteriales
99.52	629733	74 F		1762	Mycobacteriaceae
99.51	629657	5911 G		1763	Mycobacterium
98.56	623663	5295 G1		77643	Mycobacterium tuberculosis complex
97.69	618141	613920 S		1773	Mycobacterium tuberculosis
0.04	227	78 S		78331	Mycobacterium canettii
0.01	37	35 S		1768	Mycobacterium kansasii
0.03	176	0 P		1239	Firmicutes
0.03	176	0 C		91061	Bacilli
0.03	172	0 O		186826	Lactobacillales
0.03	172	0 F		1300	Streptococcaceae
0.03	172	7 G		1301	Streptococcus
0.03	161	155 S		1313	Streptococcus pneumoniae
0	2	1 S		28037	Streptococcus mitis
0	1	0 S		257758	Streptococcus pseudopneumoniae

Pre-processing ?

- Process fastq files prior to further analysis
 - Remove reads from other species
 - Trim adapters
 - Clip low quality bases
 - Merge overlapping reads from same fragment
 - ...

Adapter clipping & trimming

- FastqMcf
 - Detect & remove sequencing adapters and primers
 - Detect & clip poor quality at the ends of reads
 - Remove low complexity reads
 - Detect & remove Ns from ends of reads
 - Keep PE reads in right order

```
fastq-mcf -H -X -o sample_filtered.R1.fastq.gz -o  
/sample_filtered.R2.fastq.gz adapters.fa sample.R1.fastq  
sample.R2.fastq
```

<https://github.com/ExpressionAnalysis/ea-utils/blob/wiki/FastqMcf.md>

Adapter clipping & trimming – FastqMcf

- Example 2 – amplicon sequencing
 - Input
 - 2 fastq files of 1,834 reads each
 - Outputs
 - 2 fastq files of 1,801 reads each
 - List of adapter found

Adapter TruSeq_Adapter,_Index_1 : counted **1038** at the 'end' of 'sample.R1.fastq'

...

Adapter Illumina_Single_End_Sequencing_Primer_3p : counted **1046** at the 'end' of 'sample.R2.fastq'

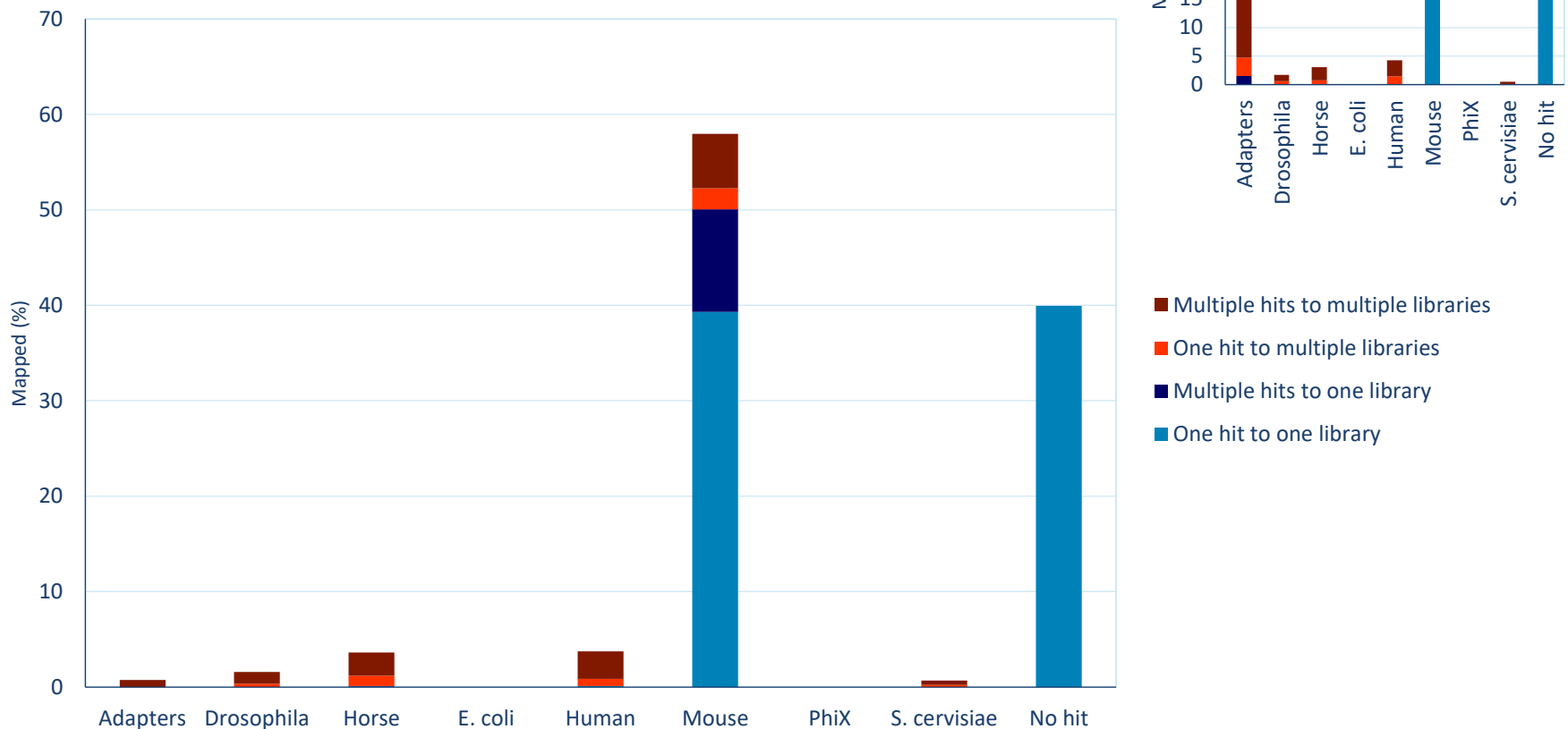
...

Total reads: 1801

Too short after clip: 33

Adapter clipping & trimming – FastqMcf

- Example 2 – amplicon sequencing



Read selection – Kraken 2 & seqtk

- Select reads from sequenced organism
 - Kraken2 output
 - Seqtk: toolkit for processing sequences in FASTA/Q formats

```
grep "organism" kraken2.output.txt | cut -f2 > reads.list  
seqtk subseq sample.R1.fastq.gz reads.list | gzip - >  
sample.selected.R1.fastq.gz  
seqtk subseq sample.R2.fastq.gz reads.list | gzip - >  
sample.selected.R2.fastq.gz
```

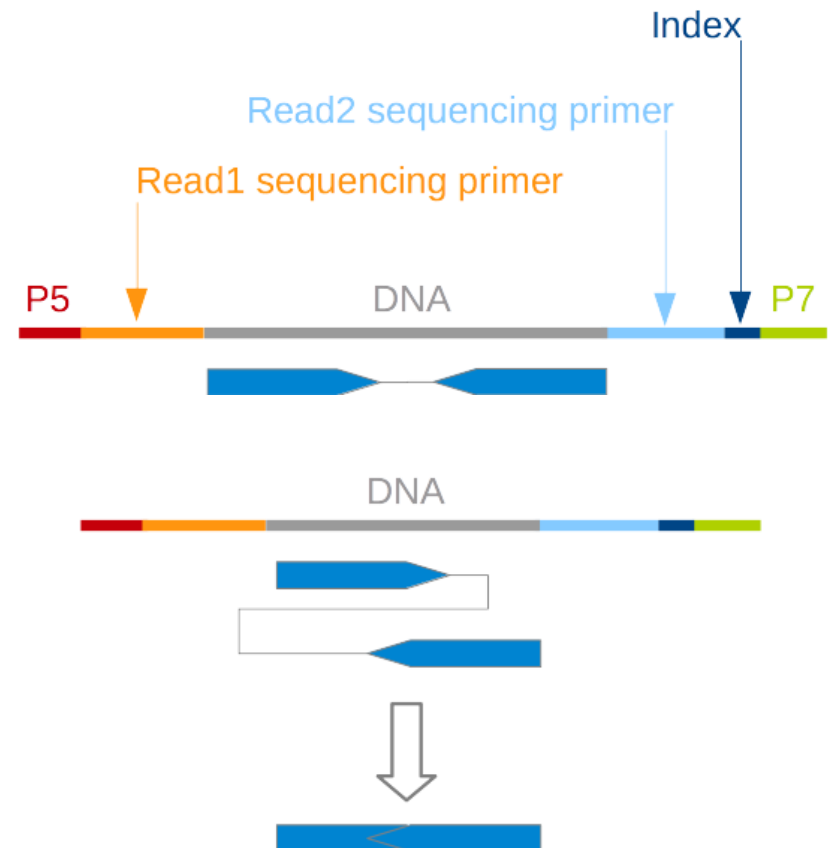
<https://ccb.jhu.edu/software/kraken2/>

<https://github.com/lh3/seqtk>

Merge overlapping reads – FLASH2

- FLASH (Fast Length Adjustment of SHort reads)

- Merge paired-end reads
- Keep DNA fragment only



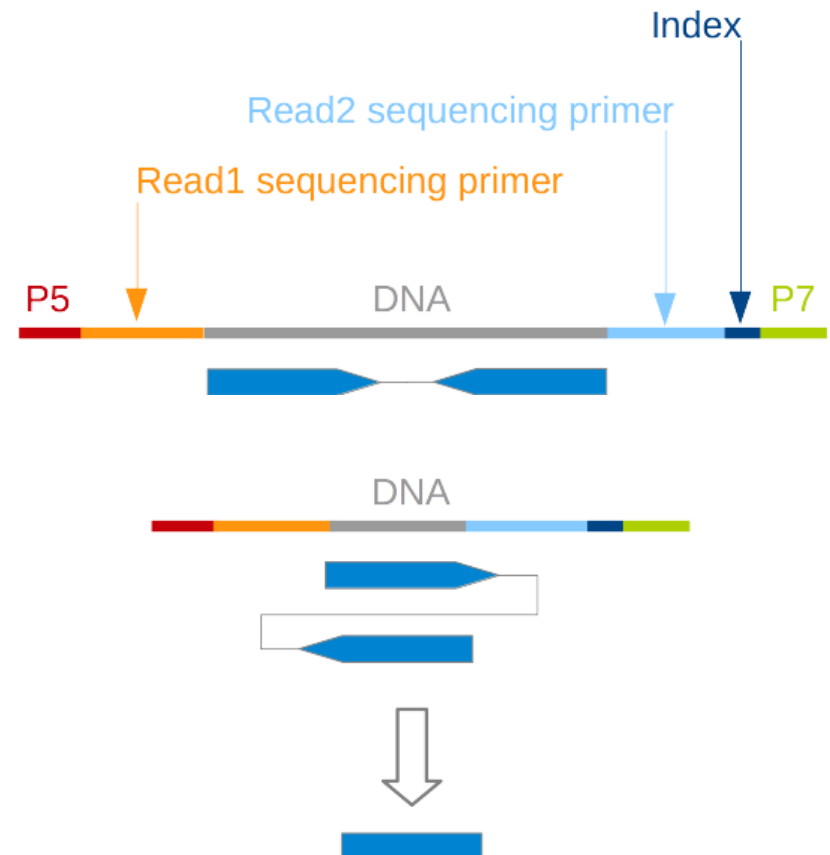
```
flash2 --max-overlap=250 --min-overlap=20 --allow-outies -d result -o  
sample.flashed sample.R1.fastq.gz sample.R2.fastq.gz > flash.log
```

<https://github.com/dstreett/FLASH2>

Merge overlapping reads – FLASH2

- FLASH (Fast Length Adjustment of SHort reads)

- Merge paired-end reads
- Keep DNA fragment only



```
flash2 --max-overlap=250 --min-overlap=20 --allow-outies -d result -o  
sample.flashed sample.R1.fastq.gz sample.R2.fastq.gz > flash.log
```

<https://github.com/dstreett/FLASH2>

A hand wearing a blue nitrile glove holds a rectangular microarray chip. The chip has a grid of small, light-colored spots. The background is a blurred laboratory with computer monitors and equipment.

Questions?