



Genômica Computacional

Introdução e processamento de dados de sequenciamento

Professor: Ricardo A. Vialle

CS31 - Genômica Computacional [14/07-17/07 - 10h00-12h00]

14 de Julho de 2025

Cronograma

Data	Tema
14-Jul	Introdução e processamento de dados de sequenciamento (teórico-prática)
15-Jul	Montagem de genomas (teórico-prática)
16-Jul	Anotação de genomas (teórico-prática)
17-Jul	Análise de variabilidade genética (teórico-prática)

Avaliação

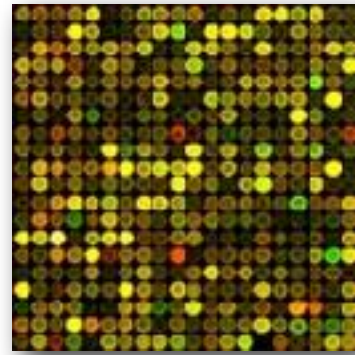
Nota final será determinada de acordo com presença

Genomics: shaped by technology



Sanger DNA
sequencing

1977-1990s



DNA Microarrays

Since mid-1990s



2nd-generation DNA
sequencing

Since ~2007



3rd-generation &
single-molecule
DNA sequencing

Since ~2010

These provide very high-resolution snapshots of the world of nucleic acids (not just DNA)

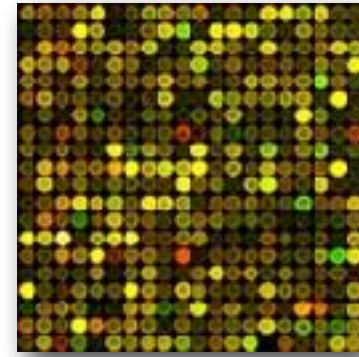
		Illumina					Pacific Biosciences	Oxford Nanopore
	Sanger	MiSeq	NextSeq	HiSeq	NovaSeq	Ion Torrent		
Throughput range per run (Gb)	<i>c.</i> 0.0005	10–15	10–120	1000–1800	2000–6000	1–15	0.5–10	0.1–1
Read length	Up to 1 kb	300	150	150	250	200–400	up to 60 kb	up to 100 kb
Read type	SR	SR, PE	SR, PE	SR, PE	SR, PE	SR	SR	SR
Error rate (%)	0.001–1	0.8	0.8	0.2	0.2–0.8	1–2	13	5–40
Error type	Substitutions	Substitutions	Substitutions	Substitutions	Substitutions	Indels, homopolymers	Indels	Indels, deletions
Advantages	Read accuracy and length	Read length	Throughput	Throughput, low error rate	High throughput	Speed, read length	Speed, read length	Read length, portability

Genomics technology



Sanger DNA
sequencing

1977-1990s



DNA Microarrays

Since mid-1990s



2nd-generation DNA
sequencing

Since ~2007

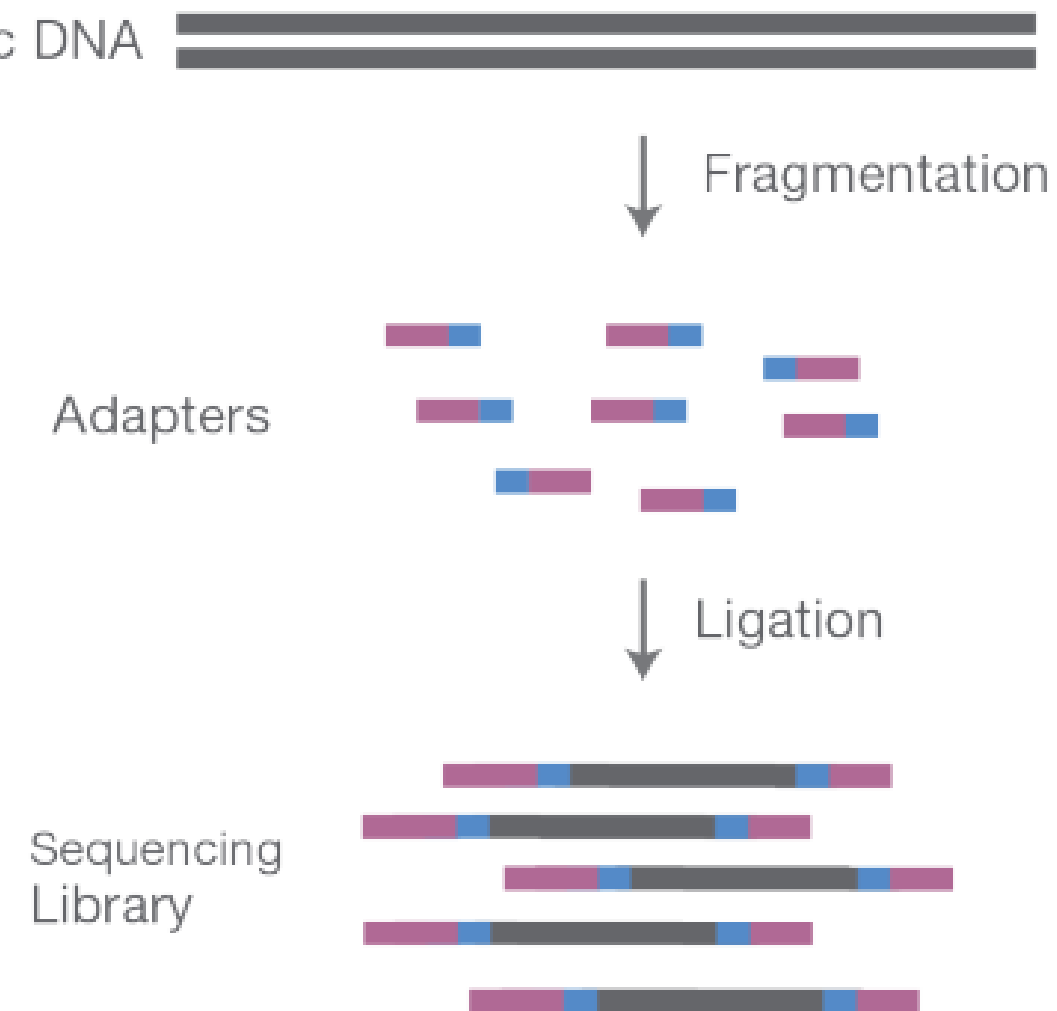


3rd-generation &
single-molecule
DNA sequencing

Since ~2010

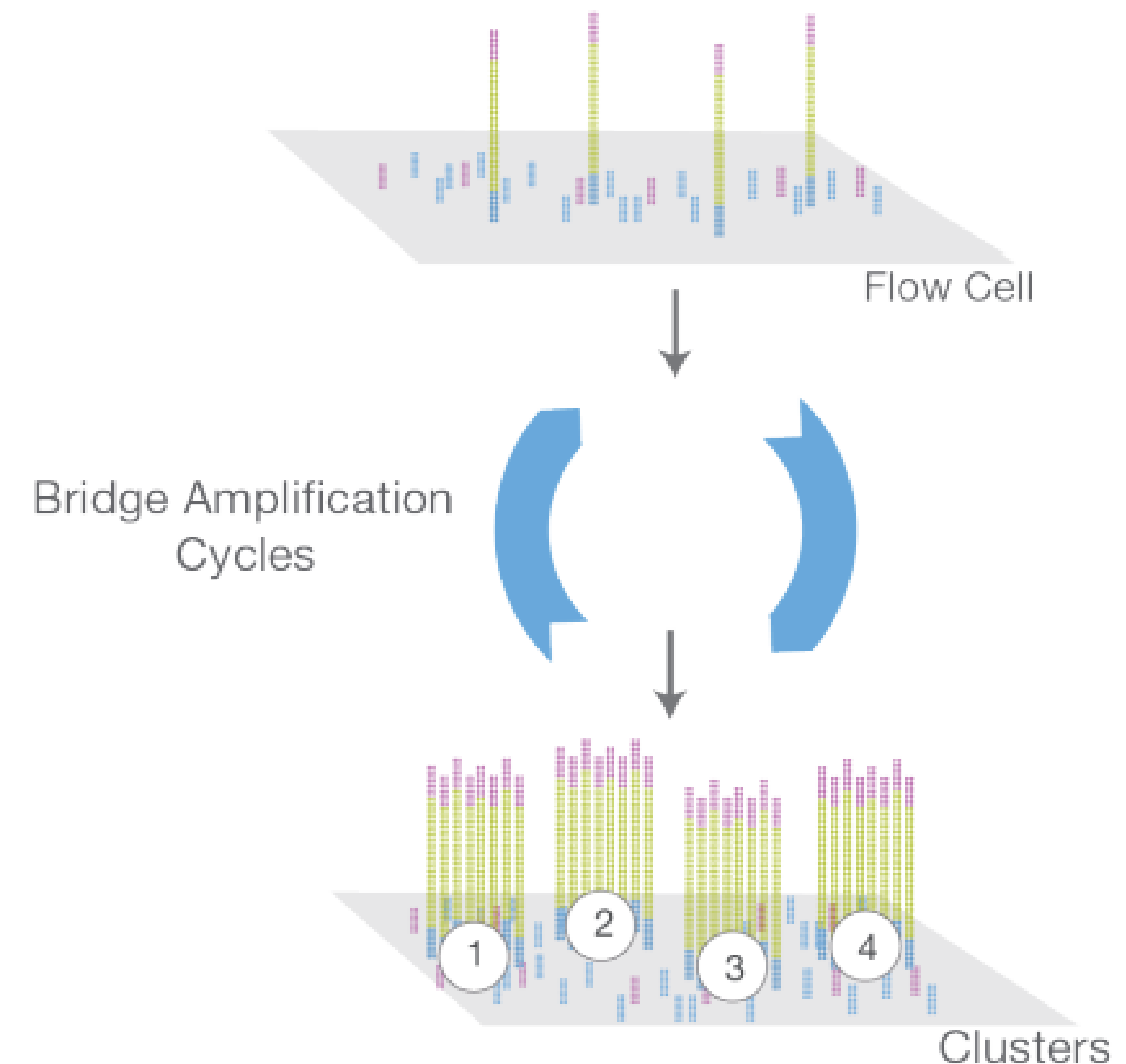


A. Library Preparation



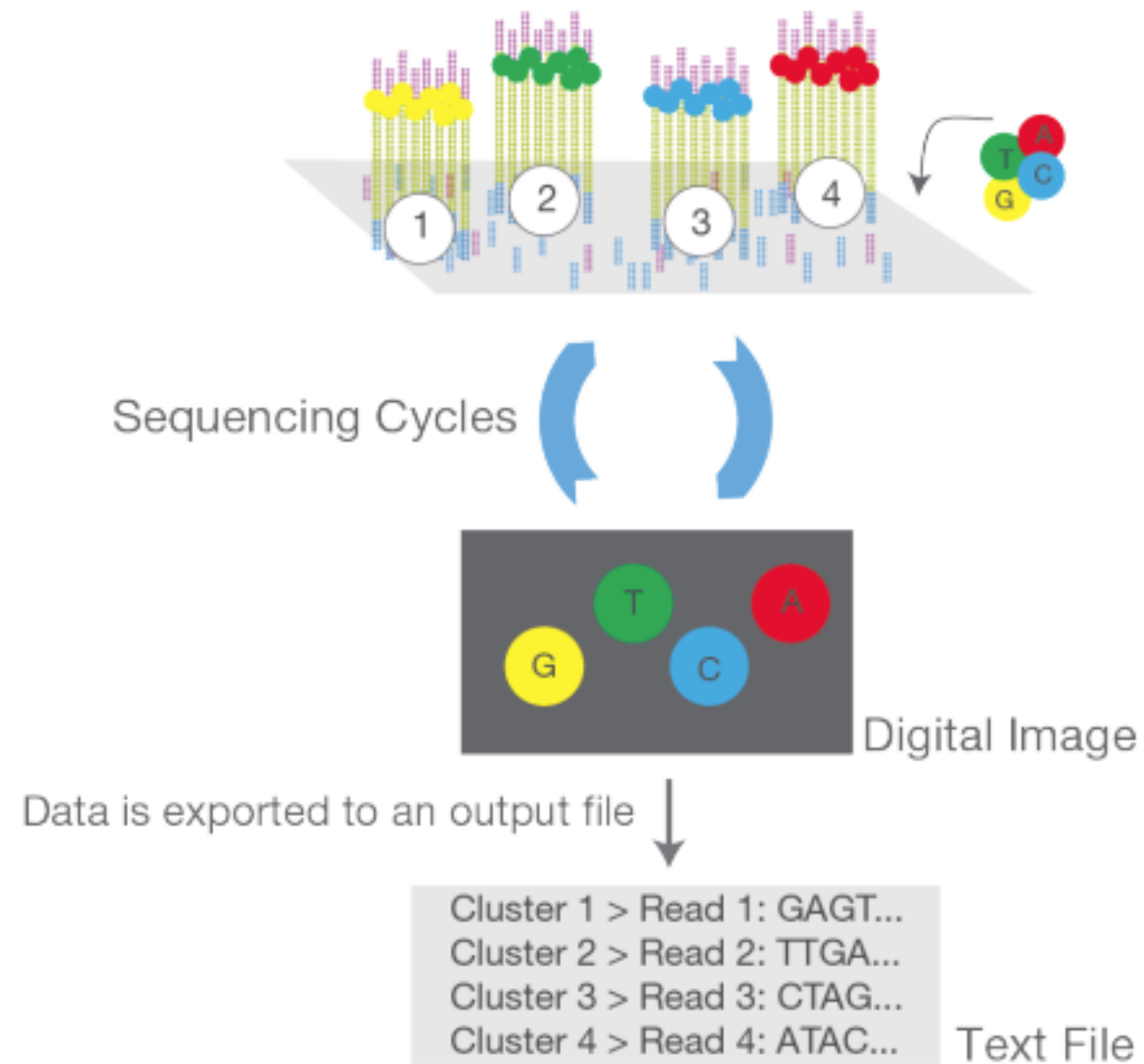
NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

B. Cluster Amplification



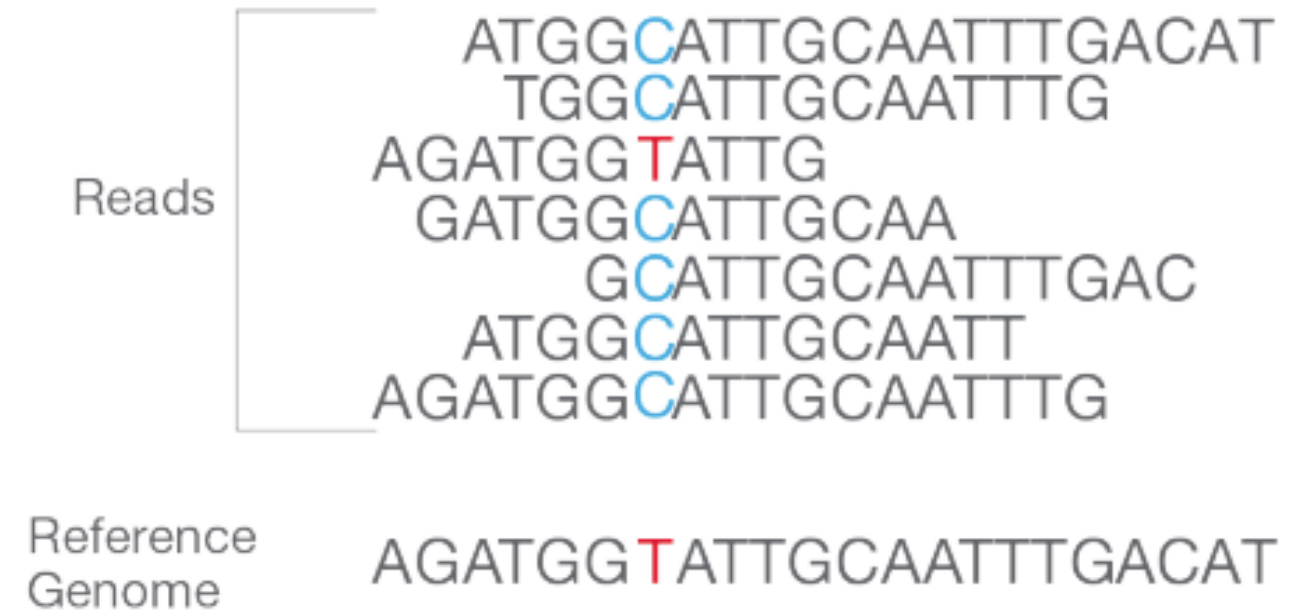
Library is loaded into a flow cell and the fragments are hybridized to the flow cell surface. Each bound fragment is amplified into a clonal cluster through bridge amplification.

C. Sequencing

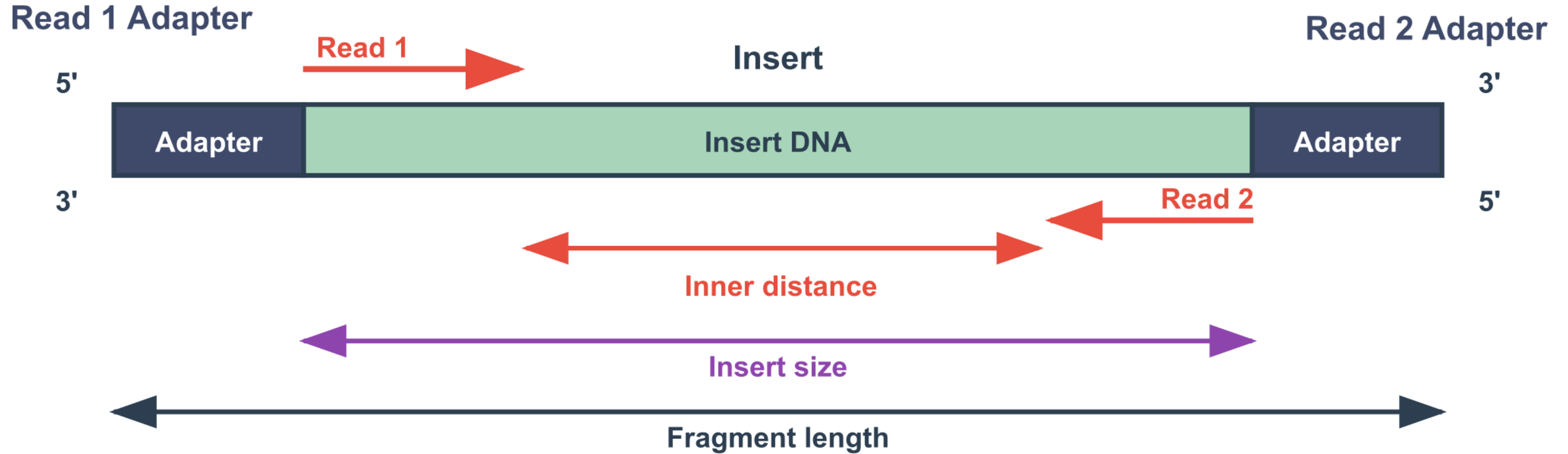


Sequencing reagents, including fluorescently labeled nucleotides, are added and the first base is incorporated. The flow cell is imaged and the emission from each cluster is recorded. The emission wavelength and intensity are used to identify the base. This cycle is repeated "n" times to create a read length of "n" bases.

D. Alignment and Data Analysis



Reads are aligned to a reference sequence with bioinformatics software. After alignment, differences between the reference genome and the newly sequenced reads can be identified.



- Fragment length: Total length including adapters
- Insert size: Length of DNA insert between adapters
- Inner distance: Distance between paired reads

A read in FASTQ format

Name @ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:3607 2:Y:18:ATCACG
Sequence ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAAT
(ignore) +
Base qualities ?@@FFBFFDDHHBCEAFGEGII DHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G

Always starts with “@”

ERR194146.1HSQ1008:141:D0CC8ACXX:3 – Machine, Run, Flowcell, Lane

1308:20201:3607 – Tile, X-pos, Y-pos

2:Y:18:ATCACG – Direction, Filtered?, Control bits, Index/Sample

Base qualities

Bases and qualities line up:

```
AGCTCTGGTGACCCATGGGCAGCTGCTAGGGA
| | | | | | | | | | | | | | | | | |
I H H H H H H H H H H H G C G C 5 F E F F F G H H H H H
```

Base quality is ASCII-encoded version of $Q = -10 \log_{10} p$

Base quality

Probability that base call is incorrect

$Q = 10 \rightarrow 1$ in 10 chance call is incorrect

$Q = 20 \rightarrow 1$ in 100

$Q = 30 \rightarrow 1$ in 1,000

ASCII TABLE

Decimal	Hexadecimal	Binary	Octal	Char	Decimal	Hexadecimal	Binary	Octal	Char	Decimal	Hexadecimal	Binary	Octal	Char
0	0	0	0	[NULL]	48	30	110000	60	0	96	60	1100000	140	`
1	1	1	1	[START OF HEADING]	49	31	110001	61	1	97	61	1100001	141	a
2	2	10	2	[START OF TEXT]	50	32	110010	62	2	98	62	1100010	142	b
3	3	11	3	[END OF TEXT]	51	33	110011	63	3	99	63	1100011	143	c
4	4	100	4	[END OF TRANSMISSION]	52	34	110100	64	4	100	64	1100100	144	d
5	5	101	5	[ENQUIRY]	53	35	110101	65	5	101	65	1100101	145	e
6	6	110	6	[ACKNOWLEDGE]	54	36	110110	66	6	102	66	1100110	146	f
7	7	111	7	[BELL]	55	37	110111	67	7	103	67	1100111	147	g
8	8	1000	10	[BACKSPACE]	56	38	111000	70	8	104	68	1101000	150	h
9	9	1001	11	[HORIZONTAL TAB]	57	39	111001	71	9	105	69	1101001	151	i
10	A	1010	12	[LINE FEED]	58	3A	111010	72	:	106	6A	1101010	152	j
11	B	1011	13	[VERTICAL TAB]	59	3B	111011	73	;	107	6B	1101011	153	k
12	C	1100	14	[FORM FEED]	60	3C	111100	74	<	108	6C	1101100	154	l
13	D	1101	15	[CARRIAGE RETURN]	61	3D	111101	75	=	109	6D	1101101	155	m
14	E	1110	16	[SHIFT OUT]	62	3E	111110	76	>	110	6E	1101110	156	n
15	F	1111	17	[SHIFT IN]	63	3F	111111	77	?	111	6F	1101111	157	o
16	10	10000	20	[DATA LINK ESCAPE]	64	40	1000000	100	@	112	70	1110000	160	p
17	11	10001	21	[DEVICE CONTROL 1]	65	41	1000001	101	A	113	71	1110001	161	q
18	12	10010	22	[DEVICE CONTROL 2]	66	42	1000010	102	B	114	72	1110010	162	r
19	13	10011	23	[DEVICE CONTROL 3]	67	43	1000011	103	C	115	73	1110011	163	s
20	14	10100	24	[DEVICE CONTROL 4]	68	44	1000100	104	D	116	74	1110100	164	t
21	15	10101	25	[NEGATIVE ACKNOWLEDGE]	69	45	1000101	105	E	117	75	1110101	165	u
22	16	10110	26	[SYNCHRONOUS IDLE]	70	46	1000110	106	F	118	76	1110110	166	v
23	17	10111	27	[ENG OF TRANS. BLOCK]	71	47	1000111	107	G	119	77	1110111	167	w
24	18	11000	30	[CANCEL]	72	48	1001000	110	H	120	78	1111000	170	x
25	19	11001	31	[END OF MEDIUM]	73	49	1001001	111	I	121	79	1111001	171	y
26	1A	11010	32	[SUBSTITUTE]	74	4A	1001010	112	J	122	7A	1111010	172	z
27	1B	11011	33	[ESCAPE]	75	4B	1001011	113	K	123	7B	1111011	173	{
28	1C	11100	34	[FILE SEPARATOR]	76	4C	1001100	114	L	124	7C	1111100	174	
29	1D	11101	35	[GROUP SEPARATOR]	77	4D	1001101	115	M	125	7D	1111101	175	}
30	1E	11110	36	[RECORD SEPARATOR]	78	4E	1001110	116	N	126	7E	1111110	176	~
31	1F	11111	37	[UNIT SEPARATOR]	79	4F	1001111	117	O	127	7F	1111111	177	[DEL]
32	20	100000	40	[SPACE]	80	50	1010000	120	P					
33	21	100001	41	!	81	51	1010001	121	Q					
34	22	100010	42	"	82	52	1010010	122	R					
35	23	100011	43	#	83	53	1010011	123	S					
36	24	100100	44	\$	84	54	1010100	124	T					
37	25	100101	45	%	85	55	1010101	125	U					
38	26	100110	46	&	86	56	1010110	126	V					
39	27	100111	47	'	87	57	1010111	127	W					
40	28	101000	50	(88	58	1011000	130	X					
41	29	101001	51)	89	59	1011001	131	Y					
42	2A	101010	52	*	90	5A	1011010	132	Z					
43	2B	101011	53	+	91	5B	1011011	133	[
44	2C	101100	54	,	92	5C	1011100	134	\					
45	2D	101101	55	-	93	5D	1011101	135]					
46	2E	101110	56	.	94	5E	1011110	136	^					
47	2F	101111	57	/	95	5F	1011111	137	_					

ASCII TABLE

These character
don't print

Decimal	Hexadecimal	Binary	Octal	Char	Decimal	Hexadecimal	Binary	Octal	Char	Decimal	Hexadecimal	Binary	Octal	Char
0	0	0	0	[NULL]	48	30	110000	60	0	96	60	1100000	140	`
1	1	1	1	[START OF HEADING]	49	31	110001	61	1	97	61	1100001	141	a
2	2	10	2	[START OF TEXT]	50	32	110010	62	2	98	62	1100010	142	b
3	3	11	3	[END OF TEXT]	51	33	110011	63	3	99	63	1100011	143	c
4	4	100	4	[END OF TRANSMISSION]	52	34	110100	64	4	100	64	1100100	144	d
5	5	101	5	[ENQUIRY]	53	35	110101	65	5	101	65	1100101	145	e
6	6	110	6	[ACKNOWLEDGE]	54	36	110110	66	6	102	66	1100110	146	f
7	7	111	7	[BELL]	55	37	110111	67	7	103	67	1100111	147	g
8	8	1000	10	[BACKSPACE]	56	38	111000	70	8	104	68	1101000	150	h
9	9	1001	11	[HORIZONTAL TAB]	57	39	111001	71	9	105	69	1101001	151	i
10	A	1010	12	[LINE FEED]	58	3A	111010	72	:	106	6A	1101010	152	j
11	B	1011	13	[VERTICAL TAB]	59	3B	111011	73	;	107	6B	1101011	153	k
12	C	1100	14	[FORM FEED]	60	3C	111100	74	<	108	6C	1101100	154	l
13	D	1101	15	[CARRIAGE RETURN]	61	3D	111101	75	=	109	6D	1101101	155	m
14	E	1110	16	[SHIFT OUT]	62	3E	111110	76	>	110	6E	1101110	156	n
15	F	1111	17	[SHIFT IN]	63	3F	111111	77	?	111	6F	1101111	157	o
16	10	10000	20	[DATA LINK ESCAPE]	64	40	1000000	100	@	112	70	1110000	160	p
17	11	10001	21	[DEVICE CONTROL 1]	65	41	1000001	101	A	113	71	1110001	161	q
18	12	10010	22	[DEVICE CONTROL 2]	66	42	1000010	102	B	114	72	1110010	162	r
19	13	10011	23	[DEVICE CONTROL 3]	67	43	1000011	103	C	115	73	1110011	163	s
20	14	10100	24	[DEVICE CONTROL 4]	68	44	1000100	104	D	116	74	1110100	164	t
21	15	10101	25	[NEGATIVE ACKNOWLEDGE]	69	45	1000101	105	E	117	75	1110101	165	u
22	16	10110	26	[SYNCHRONOUS IDLE]	70	46	1000110	106	F	118	76	1110110	166	v
23	17	10111	27	[ENG OF TRANS. BLOCK]	71	47	1000111	107	G	119	77	1110111	167	w
24	18	11000	30	[CANCEL]	72	48	1001000	110	H	120	78	1111000	170	x
25	19	11001	31	[END OF MEDIUM]	73	49	1001001	111	I	121	79	1111001	171	y
26	1A	11010	32	[SUBSTITUTE]	74	4A	1001010	112	J	122	7A	1111010	172	z
27	1B	11011	33	[ESCAPE]	75	4B	1001011	113	K	123	7B	1111011	173	{
28	1C	11100	34	[FILE SEPARATOR]	76	4C	1001100	114	L	124	7C	1111100	174	
29	1D	11101	35	[GROUP SEPARATOR]	77	4D	1001101	115	M	125	7D	1111101	175	}
30	1E	11110	36	[RECORD SEPARATOR]	78	4E	1001110	116	N	126	7E	1111110	176	~
31	1F	11111	37	[UNIT SEPARATOR]	79	4F	1001111	117	O	127	7F	1111111	177	[DEL]
32	20	100000	40	[SPACE]	80	50	1010000	120	P					
33	21	100001	41	!	81	51	1010001	121	Q					
34	22	100010	42	"	82	52	1010010	122	R					
35	23	100011	43	#	83	53	1010011	123	S					
36	24	100100	44	\$	84	54	1010100	124	T					
37	25	100101	45	%	85	55	1010101	125	U					
38	26	100110	46	&	86	56	1010110	126	V					
39	27	100111	47	'	87	57	1010111	127	W					
40	28	101000	50	(88	58	1011000	130	X					
41	29	101001	51)	89	59	1011001	131	Y					
42	2A	101010	52	*	90	5A	1011010	132	Z					
43	2B	101011	53	+	91	5B	1011011	133	[
44	2C	101100	54	,	92	5C	1011100	134	\					
45	2D	101101	55	-	93	5D	1011101	135]					
46	2E	101110	56	.	94	5E	1011110	136	^					
47	2F	101111	57	/	95	5F	1011111	137	_					

FASTQ

Read 1	Name
	Sequence
	(placeholder)
	Base qualities
Read 2	Name
	Sequence
	(placeholder)
	Base qualities
Read 3	Name
	Sequence
	(placeholder)
	Base qualities
Read 4	Name
	Sequence
	(placeholder)
	Base qualities
Read 5	Name
	Sequence
	(placeholder)
	Base qualities

[illegible]

Sample S1 L001 R1 001.fastq.gz

Sample: Sample name

S1: Sample number

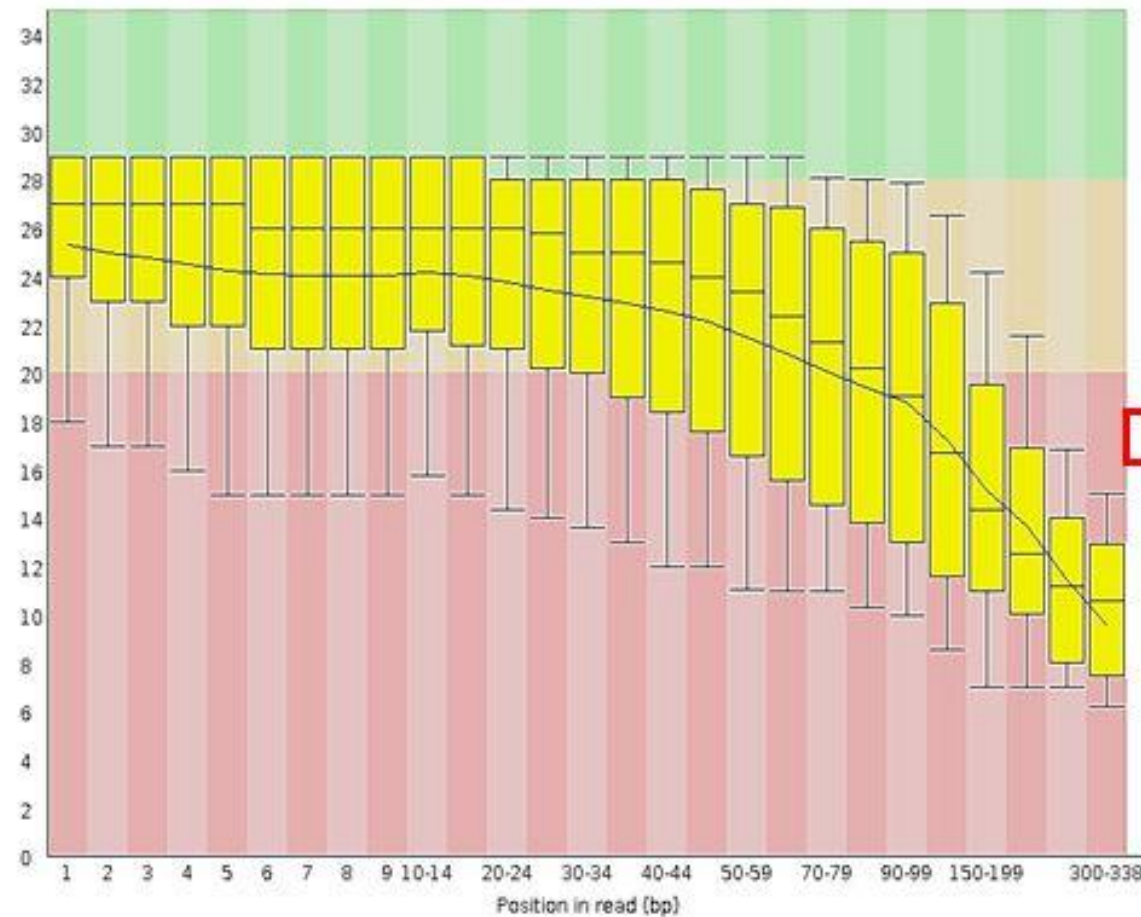
L001: Lane number

R1: Read direction (If starts with "I", Index read direction)

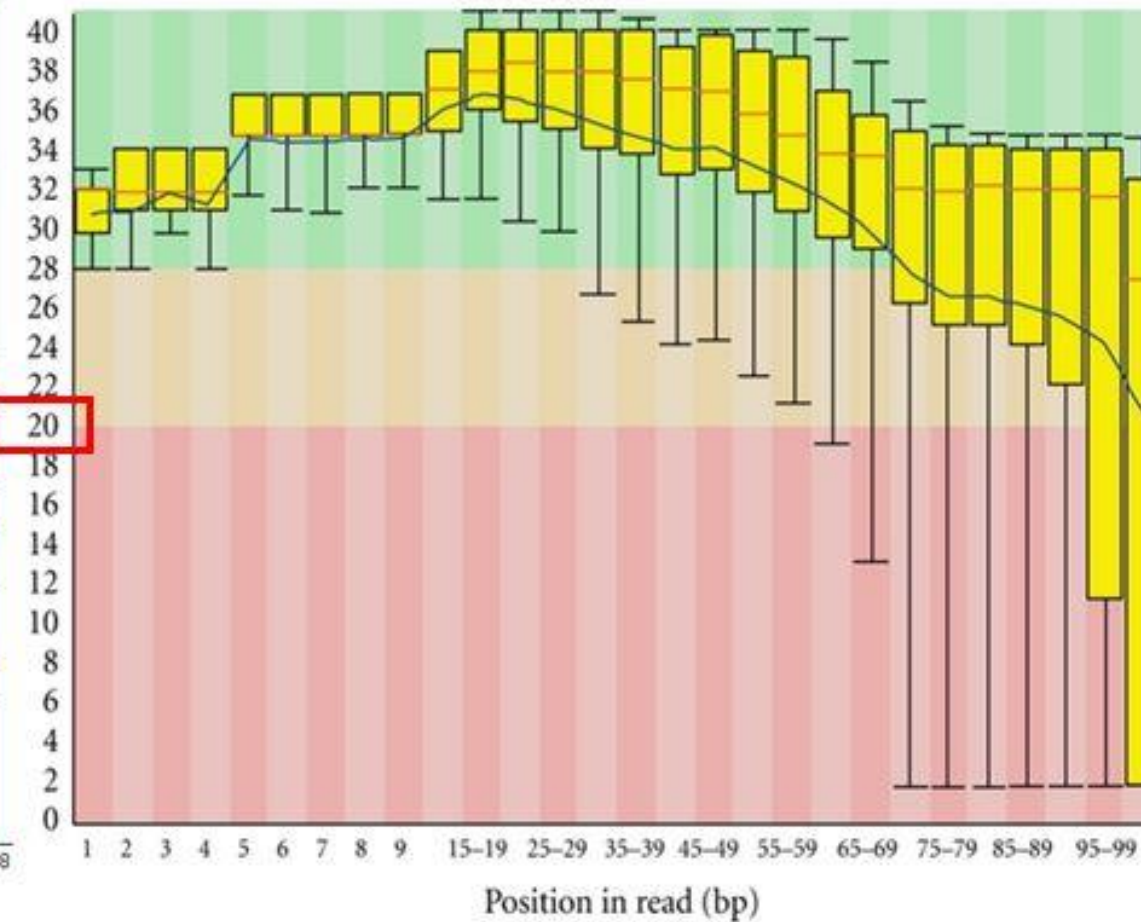
001: File number (always 001 on modern systems)

Comparison in sequencing quality

Ion Torrent PGM

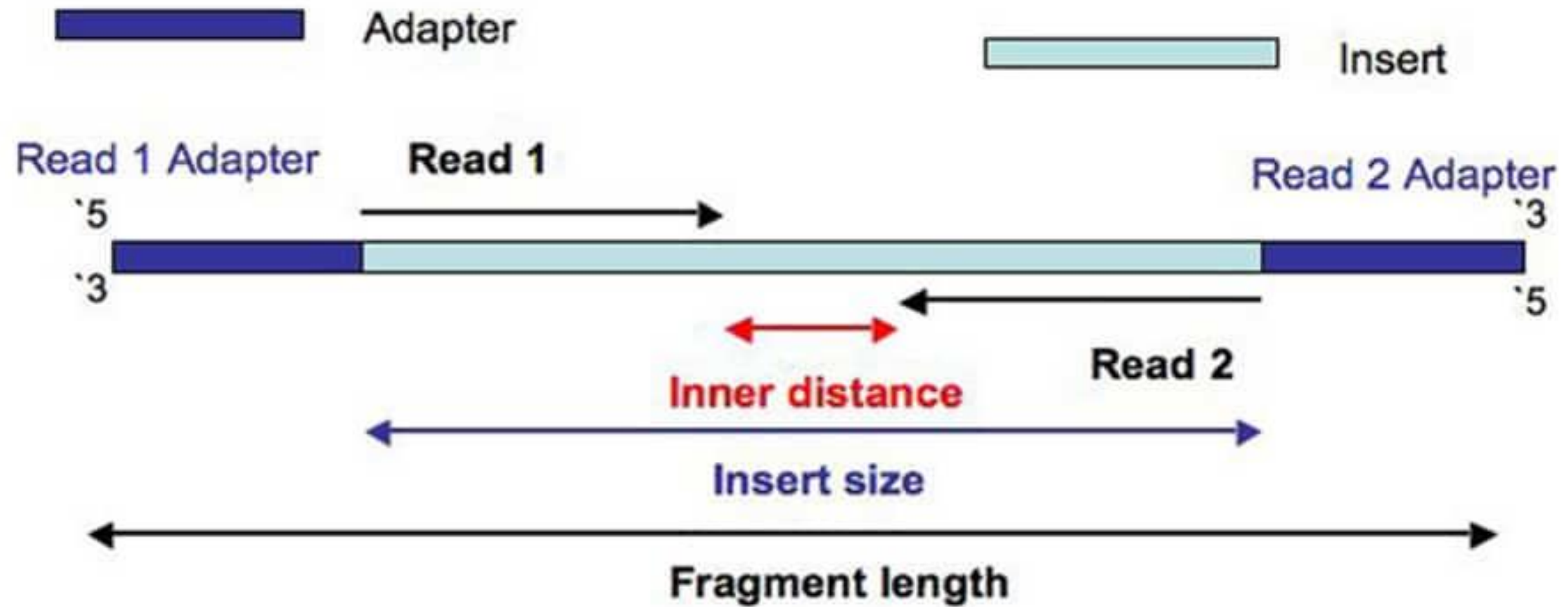


Illumina Hiseq 2000

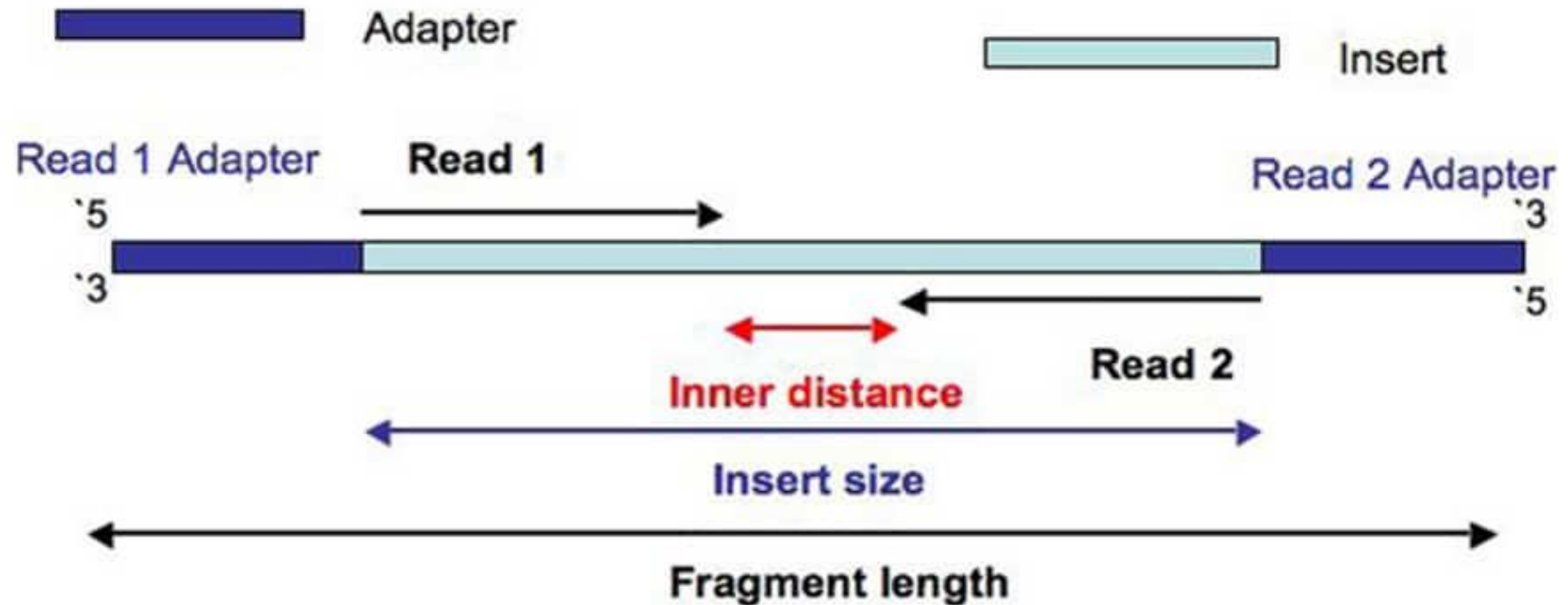


Per base sequence quality of samples generated by FASTQC. The yellow box show the base-calling quality scores across all sequencing reads. The blue line indicates the mean quality score. Q20=99% accuracy. Q30=99.9% accuracy...

- Trimming
 - Barcode & adapter sequences



- Trimming
 - Barcode & adapter sequences



- Poor quality sequence at the starts/ends of reads

Which trimming threshold? Examples

RNAseq

- Gentle trimming
- $Q > 5$ should be enough
- Too aggressive trimming => losing part of the dataset

SNP calling

- You need to be sure of the bases
- $Q > 20$

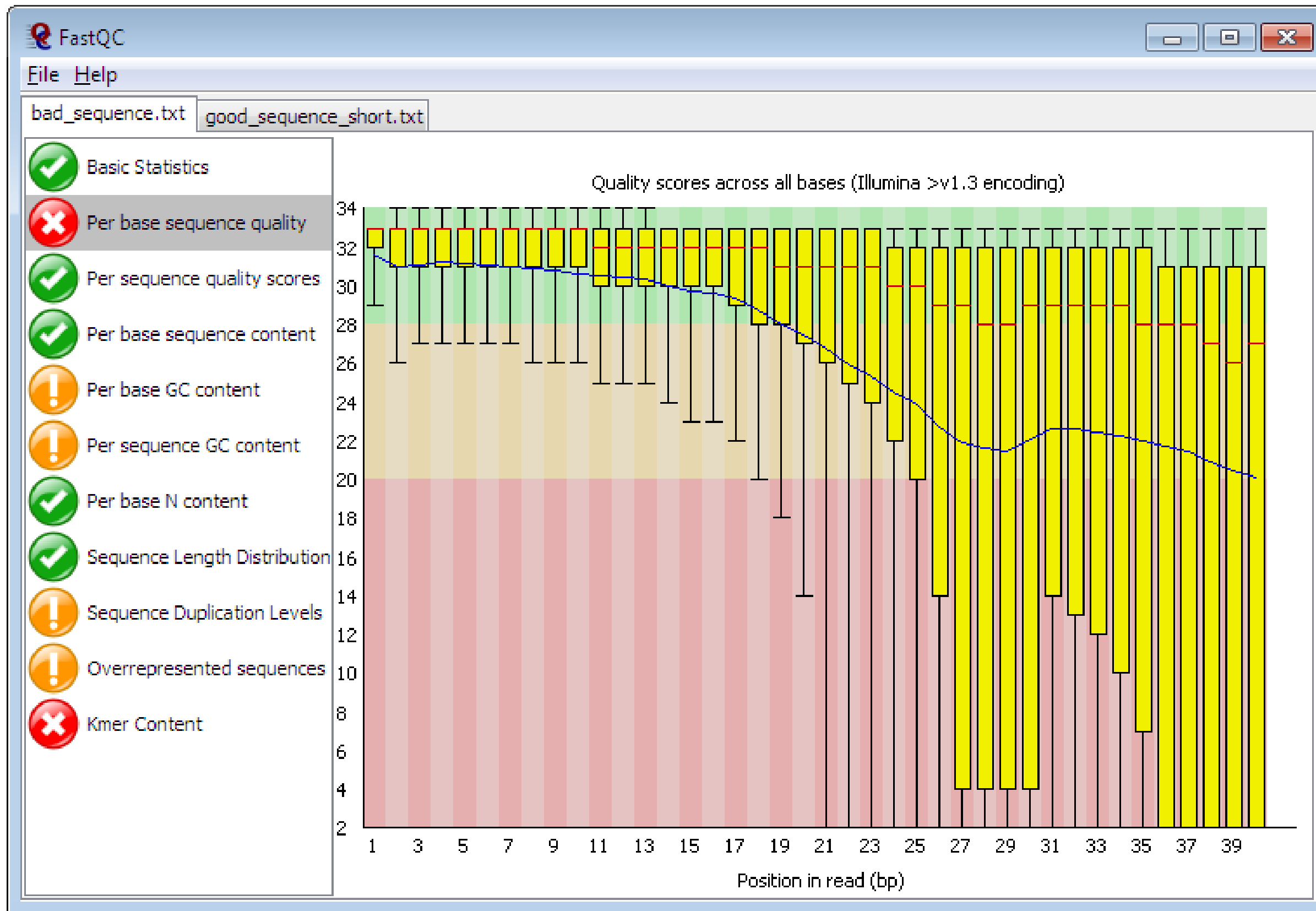
Adapter trimmer [[Scythe](#)]

Uses the quality information in a FASTQ entry and a prior to decide whether a 3' substring is adapter.

Very basically, it takes a naïve Bayesian approach to classifying 3'-end contaminants only. Because these are the most poor quality bases and most likely to be contaminated (especially as reads get longer and longer), Scythe is designed to specifically remove these contaminants.

Low quality trimmer [[Sickle](#)]

Sickle is a sliding window quality trimmer, designed to be used after Scythe. Unlike *cutadapt* and other tools, this pipeline remove adapter contaminants before quality trimming, as removing poor quality bases throws away any useful information that could be used in identifying a 3'-end adapter contaminant.



FastQC

[See each report detail here](https://www.bioinformatics.babraham.ac.uk/projects/fastqc/)

A modular tool to aggregate results from bioinformatics analyses across many samples into a single report.

Report generated on 2017-11-19, 21:42 based on data in: `/Users/hadrien/Documents/workspace/ar`

General Statistics

Copy table

Configure Columns

Plot

Showing 4/4 rows and 4/5 columns.

Sample Name	% Dups	% GC	Length	M Seqs
SRR957824_500K_R1	16.2%	49%	150 bp	0.5
SRR957824_500K_R2	7.2%	50%	150 bp	0.5
SRR957824_trimmed_R1	2.9%	51%	142 bp	0.4
SRR957824_trimmed_R2	2.7%	51%	136 bp	0.4

FastQC

[FastQC](#) is a quality control tool for high throughput sequence data, written by Simon Andrews at the Babraham Institute in Cambridge.

Sequence Quality Histograms

3

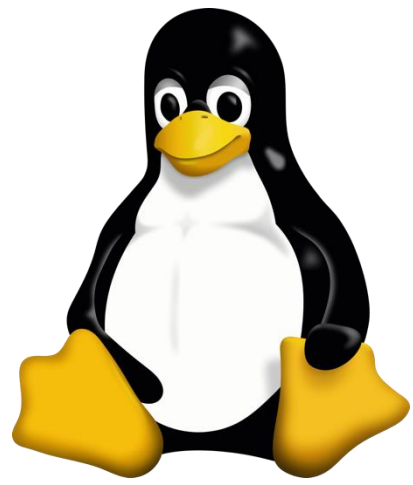
1

Improving quality: toolbox

- Trimmomatic
- Cutadapt
- Scythe
- Sickle
- Atropos

Fastp:

<https://github.com/OpenGene/fastp>



1. ls command

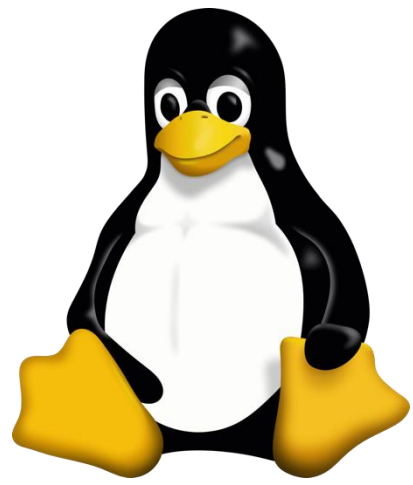
The **ls** command lists files and directories in your system. Here's the syntax:

```
ls [/directory/folder/path]
```

```
root@srv:/# ls /directory/folder/path  
file1.txt
```

If you remove the path, the **ls** command will show the current working directory's content. You can modify the command using these options:

- **-R** – lists all the files in the subdirectories.
- **-a** – shows all files, including hidden ones.
- **-lh** – converts sizes to readable formats, such as **MB**, **GB**, and **TB**.



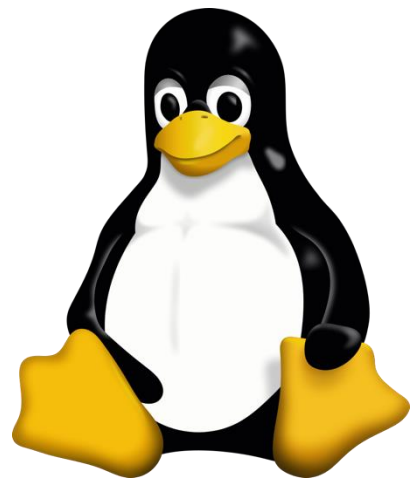
2. pwd command

The **pwd** command prints your current working directory's path, like **/home/directory/path**. Here's the command syntax:

```
pwd [option]
```

It supports two options. The **-L** or **--logical** option prints environment variable content, including symbolic links. Meanwhile, **-P** or **--physical** outputs the current directory's actual path.

```
root@srv:/directory/folder/path# pwd  
/directory/folder/path
```



3. cd command

Use the **cd** command to navigate the Linux files and directories. To use it, run this syntax with sudo privileges:

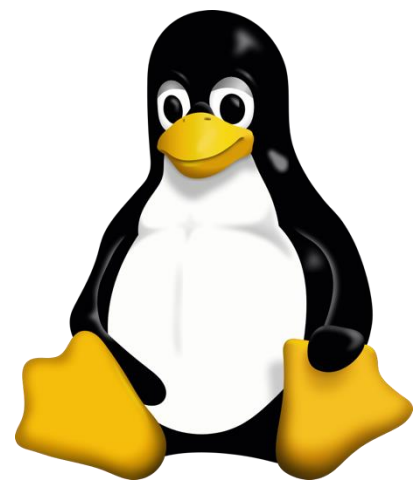
```
cd /directory/folder/path
```

```
root@srv:/# cd /directory/folder/path  
root@srv:/directory/folder/path#
```

Depending on your current location, it requires either the full path or the directory name. For example, omit **/username** from **/username/directory/folder** if you are already within it.

Omitting the arguments will take you to the home folder. Here are some navigation shortcuts:

- **cd ~[username]** – goes to another user's home directory.
- **cd ..** – moves one directory up.
- **cd-** – switches to the previous directory.



4. mkdir command

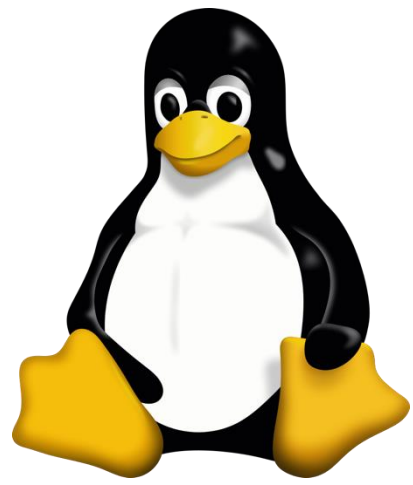
Use the **mkdir** command to create one or multiple directories and set their permissions. Ensure you are authorized to make a new folder in the parent directory. Here's the basic syntax:

```
mkdir [option] [directory_name]
```

To create a folder within a directory, use the path as the command parameter. For example, **mkdir music/songs** will create a **songs** folder inside **music**. Here are several common **mkdir** command options:

- **-p** – creates a directory between two existing folders. For example, **mkdir -p Music/2023/Songs** creates a new **2023** directory.
- **-m** – sets the folder permissions. For instance, enter **mkdir -m777 directory** to create a directory with read, write, and execute permissions for all users.
- **-v** – prints a message for each created directory.

```
root@srv:/# mkdir -v new-folder
mkdir: created directory 'new-folder'
```



6. rm command

Use the `rm` command to permanently delete files within a directory. Here's the general syntax:

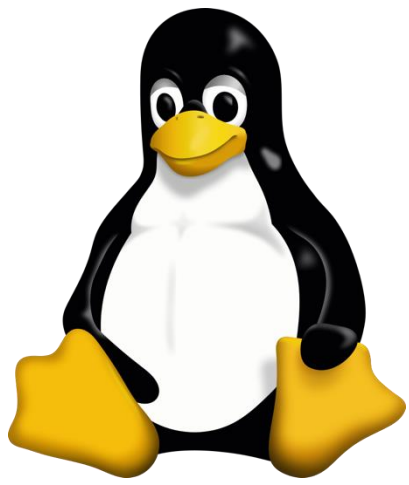
```
rm [filename1] [filename2] [filename3]
```

Adjust the number of files in the command according to your needs. If you encounter an error, ensure you have the **write** permission in the directory.

To modify the command, add the following options:

- **-i** – prompts a confirmation before deletion.
- **-f** – allows file removal without a confirmation.
- **-r** – deletes files and directories recursively.

! **Warning!** Use the `rm` command with caution since deletion is irreversible. Avoid using the `-r` and `-f` options since they may wipe all your files. Always add the `-i` option to avoid accidental deletion.



7. cp command

Use the **cp** command to copy files or directories, including their content, from your current location to another. It has various use cases, such as:

- Copying one file from the current directory to another folder. Specify the file name and target path:

```
cp filename.txt /home/username/Documents
```

- Duplicating multiple files to a directory. Enter the file names and the destination path:

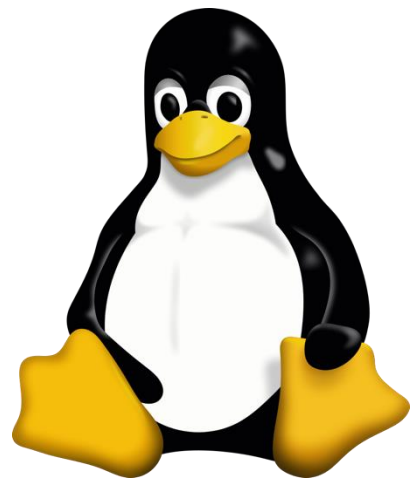
```
cp filename1.txt filename2.txt filename3.txt /home/username/Documents
```

- Copying a file's content to another within the same directory. Enter the source and the destination file:

```
cp filename1.txt filename2.txt
```

- Duplicating an entire directory. Pass the **-R** flag followed by the source and destination directory:

```
cp -R /home/username/Documents /home/username/Documents_backup
```



8. mv command

Use the **mv** command to move or rename files and directories. To move items, enter the file name followed by the destination directory:

```
mv filename.txt /home/username/Documents
```

Meanwhile, use the following syntax to **rename a file in Linux** with the **mv** command:

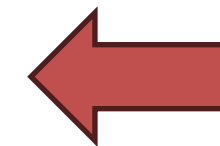
```
mv old_filename.txt new_filename.txt
```

CS31 - Genômica Computacional

Esse repositório contém materiais de aula para a disciplina de Genômica Computacional da IAMSPE.

Para facilitar execução dos tutoriais, utilizaremos o [Google Cloud Shell](#).

Aula	Data	Tema	Slides	Tutoriais
1	2025-07-14	Introdução e processamento de dados de sequenciamento	Slides	Fastq Quality-Control (QC) tutorial
2	2025-07-15	Montagem de genomas		
3	2025-07-16	Anotacao de genomas		
4	2025-07-17	Mapeando variantes		



https://github.com/RushAlz/IAMSPE-CS31-Genomica_Computacional

Open in Cloud Shell

You are about to clone the repo:

 https://github.com/RushAlz/IAMSPE-CS31-Genomica_Computacional.git

This repo is not officially maintained by Google and is considered untrusted by default.

☒ Trust repo



CANCEL

CONFIRM

Proxima aula...

Montagem de genomas (teórico-prática)

