
Cycle

« Analyse de données de séquençage à haut-débit »

Module 1/6: Analyses ADN

Pierre Pericard
Plateforme bilille - UMS PLBS

pierre.pericard@univ-lille.fr

Module 1/6: Analyses ADN

- NGS Introduction
- Reads Quality Control
- Reads Cleaning
- Aligning reads on reference → *Hélène Touzet*
- Alignment parameters → *Hélène Touzet*
- Reads duplicates
- Assembly → *Hélène Touzet*

Module 1/6: Analyses ADN

- NGS Introduction
- Reads Quality Control
- Reads Cleaning
- Aligning reads on reference → *Hélène Touzet*
- Alignment parameters → *Hélène Touzet*
- Reads duplicates
- Assembly → *Hélène Touzet*

Module 1/6: Analyses ADN

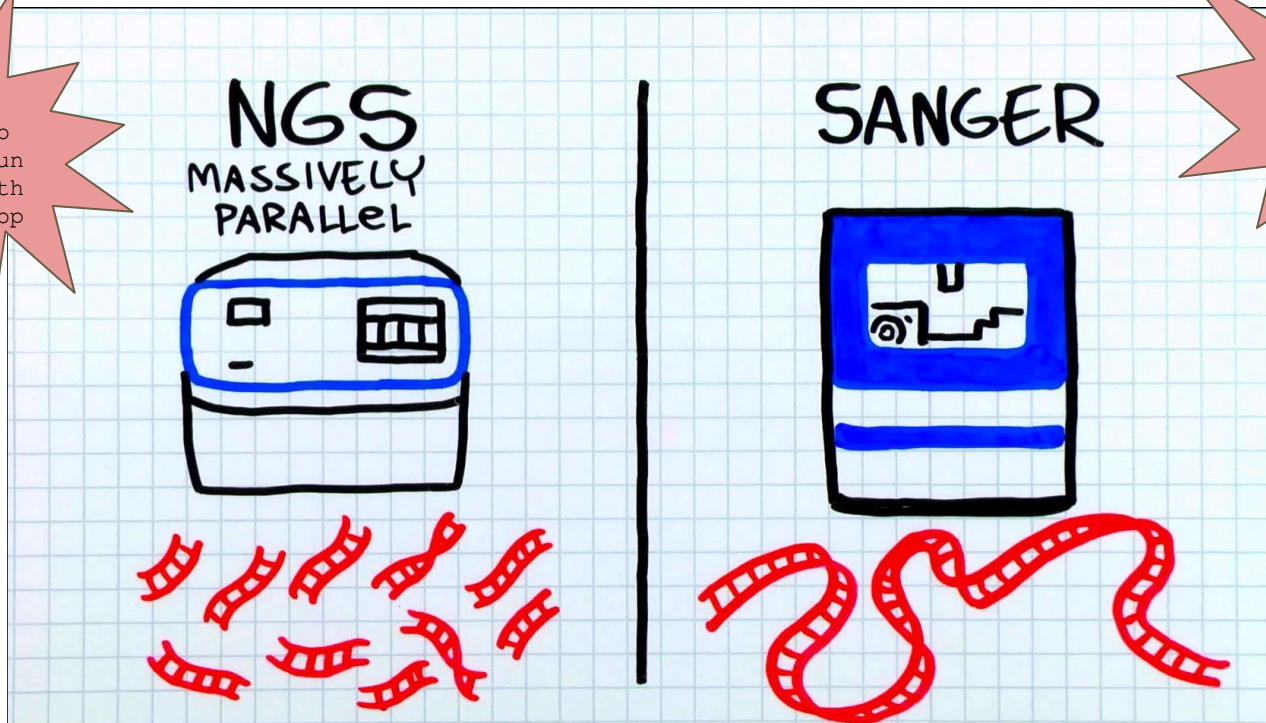
- NGS Introduction

- What is NGS?
- Sequencers
- Applications
- NGS workflow
- Output data

What is Next-Generation Sequencing (NGS)?

“Next-generation sequencing (NGS), also known as high-throughput sequencing, is the catch-all term used to describe a number of different modern sequencing technologies. These technologies allow for sequencing of DNA and RNA much more quickly and cheaply than the previously used Sanger sequencing, and as such revolutionised the study of genomics and molecular biology”

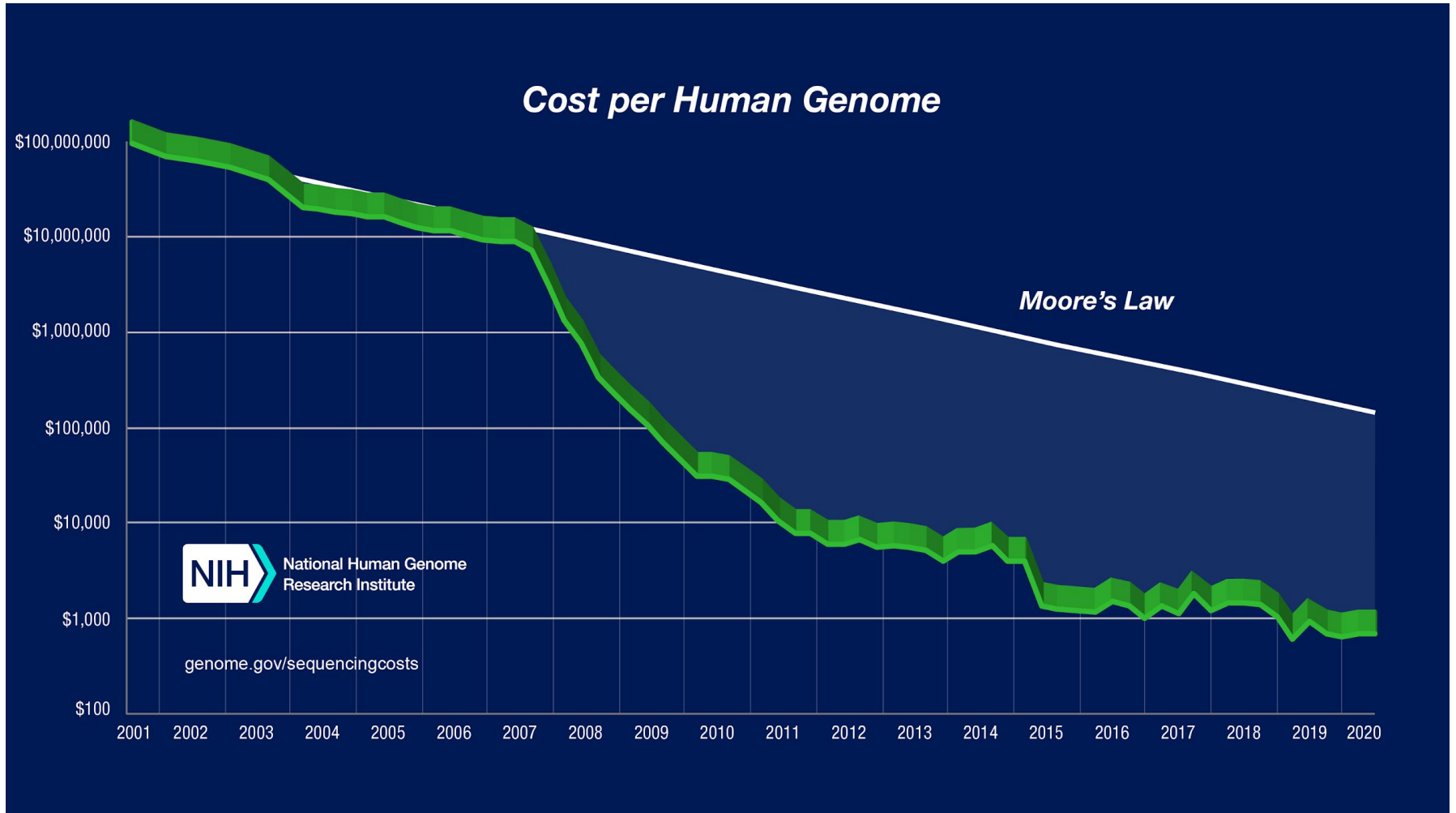
x **Gb**/run
~ 0.4 \$/Mb
~ 3 Day/run
read length
= 50~200 bp



x **Kb**/run
~ 5000 \$/Mb
~ 1 Day/run
read length
= 700 bp

The Human Genome Project was a **13-year-long** & cost **\$5 billion**

What is Next-Generation Sequencing (NGS)?



What is Next-Generation Sequencing (NGS)?

Illumina sequencing

Illumina sequencing works by simultaneously identifying DNA bases, as each base emits a unique fluorescent signal, and adding them to a nucleic acid chain

Ion Torrent: Proton / PGM sequencing (thermofisher)

Ion Torrent sequencing measures the direct release of H⁺ (protons) from the incorporation of individual bases by DNA polymerase and therefore differs from the previous two methods as it does not measure light.

illumina®

ThermoFisher
SCIENTIFIC

What is Next-Generation Sequencing (NGS)?

Illumina sequencing

Illumina sequencing works by simultaneously identifying DNA bases, as each base emits a unique fluorescent signal, and adding them to a nucleic acid chain

Ion Torrent: Proton / PGM sequencing (thermofisher)

Ion Torrent sequencing measures the direct release of H⁺ (protons) from the incorporation of individual bases by DNA polymerase and therefore differs from the previous two methods as it does not measure light.

Roche 454 pyrosequencing

illumina®

ThermoFisher
S C I E N T I F I C

454
SEQUENCING

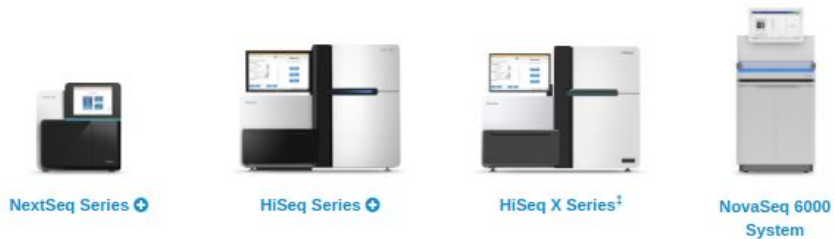
Sequencers – Illumina

	Benchtop Sequencers		Production-Scale Sequencers		
	iSeq 100	MiniSeq	MiSeq Series	NextSeq 550 Series	NextSeq 1000 & 2000
Popular Applications & Methods	Key Application	Key Application	Key Application	Key Application	Key Application
Large Whole-Genome Sequencing (human, plant, animal)					
Small Whole-Genome Sequencing (microbe, virus)	●	●	●	●	●
Exome & Large Panel Sequencing (enrichment-based)				●	●
Targeted Gene Sequencing (amplicon-based, gene panel)	●	●	●	●	●
Single-Cell Profiling (scRNA-Seq, scDNA-Seq, oligo tagging assays)				●	●
Transcriptome Sequencing (total RNA-Seq, mRNA-Seq, gene expression profiling)				●	●
Targeted Gene Expression Profiling	●	●	●	●	●
miRNA & Small RNA Analysis	●	●	●	●	●
DNA-Protein Interaction Analysis (ChIP-Seq)			●	●	●
Methylation Sequencing				●	●
16S Metagenomic Sequencing		●	●	●	●
Metagenomic Profiling (shotgun metagenomics, metatranscriptomics)				●	●
Cell-Free Sequencing & Liquid Biopsy Analysis				●	●
Run Time	9.5–19 hrs	4–24 hours	4–55 hours	12–30 hours	11–48 hours
Maximum Output	1.2 Gb	7.5 Gb	15 Gb	120 Gb	330 Gb*
Maximum Reads Per Run	4 million	25 million	25 million †	400 million	1.1 billion*
Maximum Read Length	2 × 150 bp	2 × 150 bp	2 × 300 bp	2 × 150 bp	2 × 150 bp

Sequencers – Illumina

	Production-Scale Sequencers		
	NextSeq 550 Series	NextSeq 1000 & 2000	NovaSeq 6000
Popular Applications & Methods	Key Application	Key Application	Key Application
Large Whole-Genome Sequencing (human, plant, animal)			●
Small Whole-Genome Sequencing (microbe, virus)	●	●	●
Exome & Large Panel Sequencing (enrichment-based)	●	●	●
Targeted Gene Sequencing (amplicon-based, gene panel)	●	●	●
Single-Cell Profiling (scRNA-Seq, scDNA-Seq, oligo tagging assays)	●	●	●
Transcriptome Sequencing (total RNA-Seq, mRNA-Seq, gene expression profiling)	●	●	●
Chromatin Analysis (ATAC-Seq, ChIP-Seq)	●	●	●
Methylation Sequencing	●	●	●
Metagenomic Profiling (shotgun metagenomics, metatranscriptomics)	●	●	●
Cell-Free Sequencing & Liquid Biopsy Analysis	●	●	●
Run Time	12–30 hours	11–48 hours	~13 - 38 hours (dual SP flow cells) ~13–25 hours (dual S1 flow cells) ~16–36 hours (dual S2 flow cells) ~44 hours (dual S4 flow cells)
Maximum Output	120 Gb	330 Gb*	6000 Gb
Maximum Reads Per Run	400 million	1.1 billion*	20 billion
Maximum Read Length	2 × 150 bp	2 × 150 bp	2 × 250**

Sequencers – Illumina (pre-2020)



Popular Applications & Methods	Key Application ■	Key Application ■	Key Application ■	Key Application ■
Large Whole-Genome Sequencing (human, plant, animal)	●	●	●	●
Small Whole-Genome Sequencing (microbe, virus)	●	●		●
Exome Sequencing	●	●		●
Targeted Gene Sequencing (amplicon, gene panel)	●	●		●
Whole-Transcriptome Sequencing	●	●		●
Gene Expression Profiling with mRNA-Seq	●	●		●
miRNA & Small RNA Analysis	●	●		●
DNA-Protein Interaction Analysis	●	●		●
Methylation Sequencing	●	●		●
Shotgun Metagenomics	●	●		●

Optimized NGS Sample Tracking and Workflows

See how BaseSpace Clarity LIMS (Laboratory Information Management System) enabled this large genomics lab to standardize lab procedures and cope with increasing sample volumes from diverse clients.

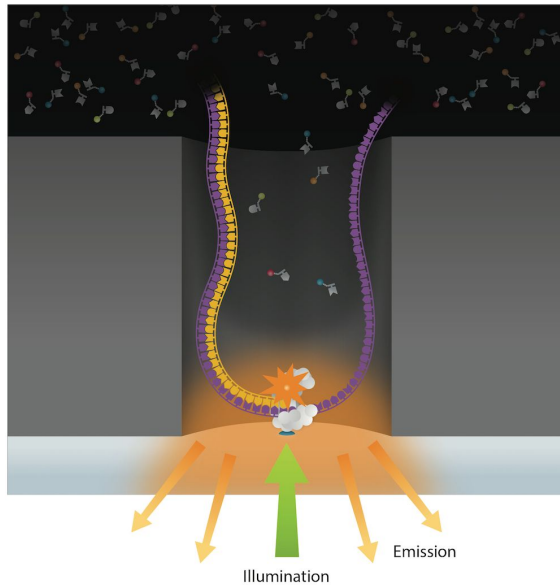
[Read Case Study >](#)

Run Time	12–30 hours	< 1–3.5 days (HiSeq 3000/HiSeq 4000) 7 hours–6 days (HiSeq 2500)	< 3 days	16–36 hours (Dual S2 flow cells) 44 hours (Dual S2 flow cells)
Maximum Output	120 Gb	1500 Gb	1800 Gb	6000 Gb ⁵
Maximum Reads Per Run	400 million	5 billion	6 billion	20 billion**
Maximum Read Length	2 × 150 bp	2 × 150 bp	2 × 150 bp	2 × 150 bp

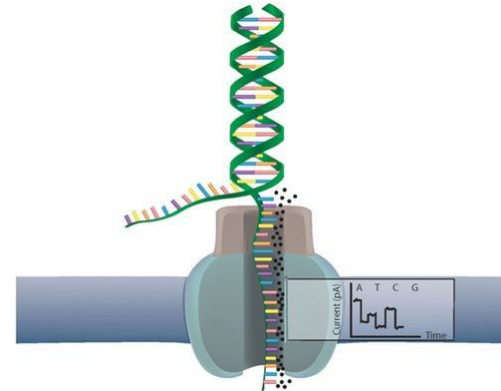
Sequencers – Thermo Fisher Scientific

Plateformes de séquençage	 <p data-bbox="647 378 879 478">Système Ion PGM™ pour le séquençage de nouvelle génération</p>	 <p data-bbox="879 378 1130 478">Système Ion S5™ pour le séquençage de nouvelle génération</p>	 <p data-bbox="1130 378 1431 478">Système Ion S5™ XL pour le séquençage de nouvelle génération</p>
<p data-bbox="511 506 647 578">Avantages</p> <p data-bbox="511 606 647 778">Rapidité : séquençage exécuté en 2 à 7 heures, selon la longueur de lecture et la sortie par la puce</p>	<p data-bbox="647 506 879 578">Évolutivité : de 30 Mo à 2 Go</p> <p data-bbox="647 606 879 778">Rapidité : séquençage exécuté en 2 à 7 heures, selon la longueur de lecture et la sortie par la puce</p>	<p data-bbox="879 506 1130 649">Simplicité : solutions de flux de travaux automatisé, de la préparation des échantillons à l'analyse</p> <p data-bbox="879 678 1130 749">Évolutivité : de 600 Mo à 15 Go</p> <p data-bbox="879 778 1130 921">Rapidité : séquençage effectué en 2,5 à 4 heures (quelle que soit la sortie par la puce)</p>	<p data-bbox="1130 506 1431 649">Simplicité : solutions de flux de travaux automatisé, de la préparation des échantillons à l'analyse</p> <p data-bbox="1130 678 1431 749">Évolutivité : de 600 Mo à 15 Go</p> <p data-bbox="1130 778 1431 921">Rapidité : de l'ADN aux données en 24 heures</p>
<p data-bbox="511 935 647 1006">Applications de séquençage</p> <p data-bbox="511 1006 647 1106">ARN ciblé ADN ciblé</p>	<p data-bbox="647 935 879 1106">ARN ciblé ADN ciblé Microbien</p>	<p data-bbox="879 935 1130 1106">ARN ciblé ADN ciblé Microbien Transcriptome Exome Séquençage de l'ARN</p>	<p data-bbox="1130 935 1431 1106">ARN ciblé ADN ciblé Microbien Transcriptome Exome Séquençage de l'ARN</p>

Third-generation sequencing



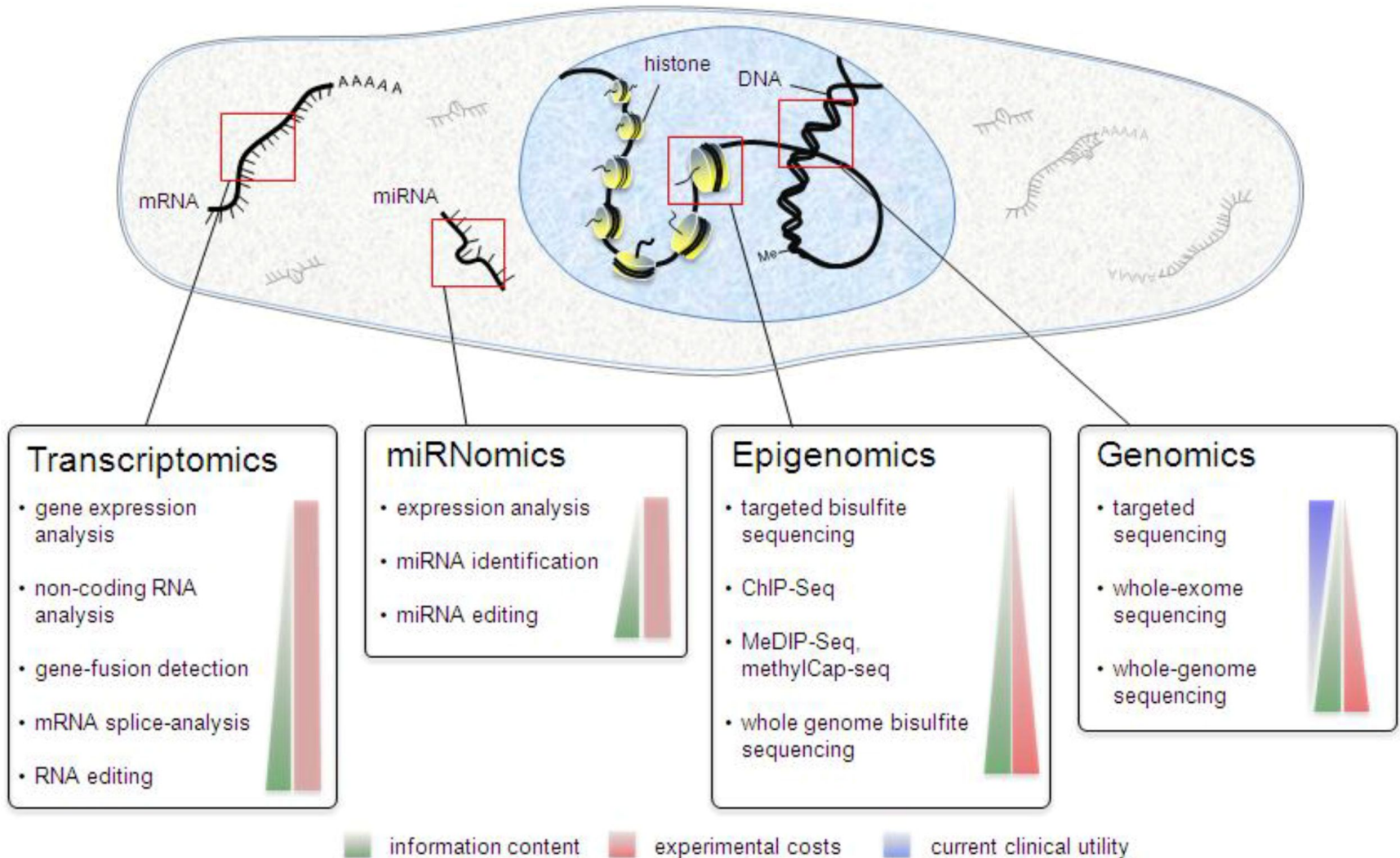
PacBio Sequencing



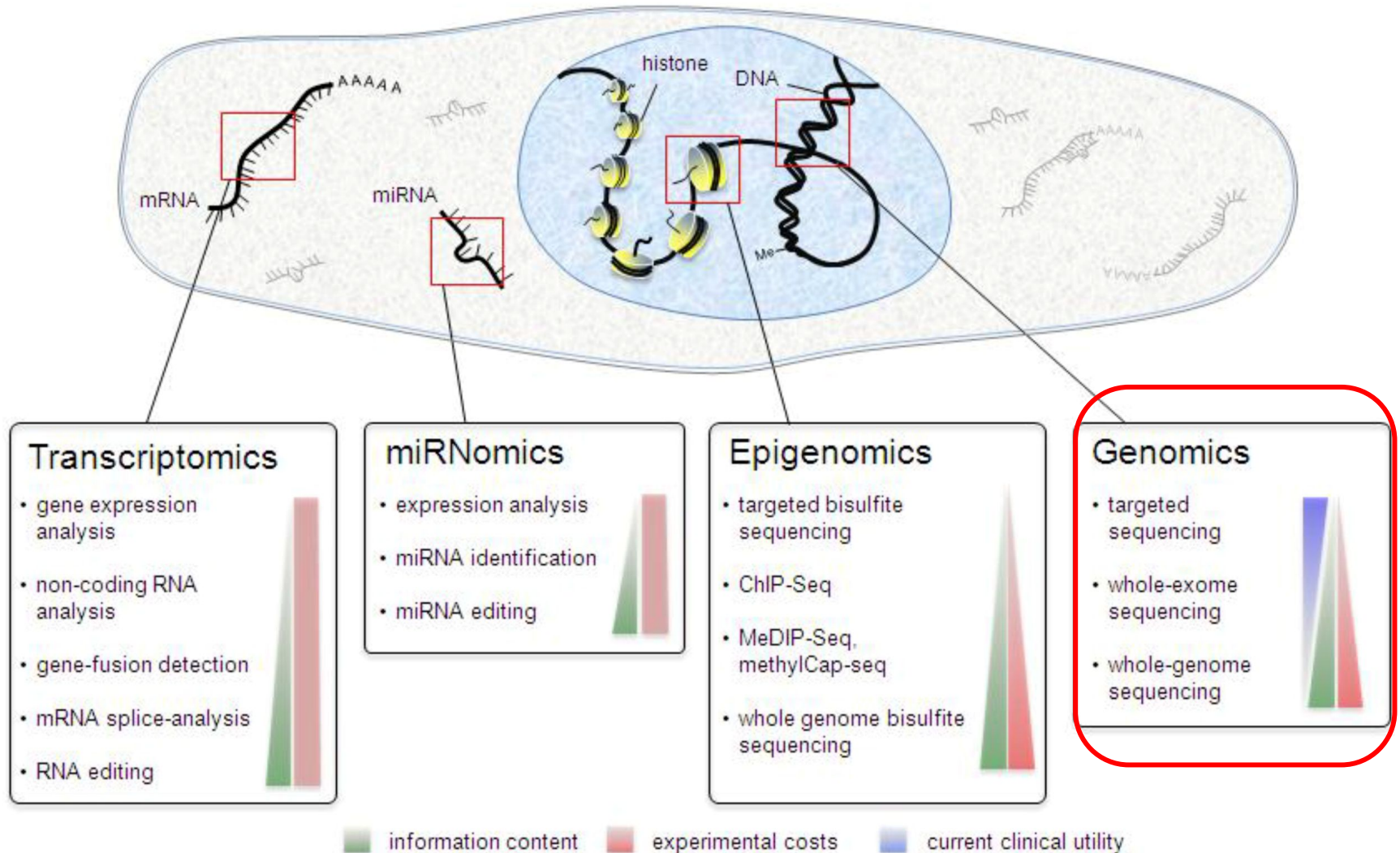
Nanopore technology

- much longer reads (> Kb)
- error rate (~ 10 → 40 %)

Applications

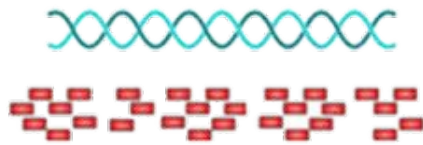


Applications



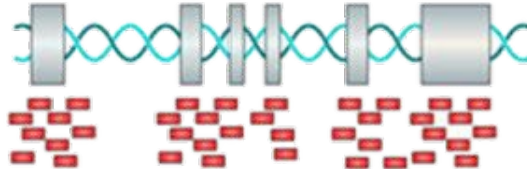
Applications : genomics (DNA-seq)

Whole genome sequencing



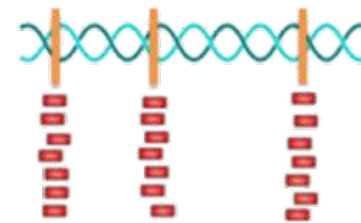
- Sequencing region : whole genome
- Sequencing Depth : >30X
- Covers everything – can identify all kinds of variants including SNPs, INDELs and SV.

Whole exome sequencing



- Sequencing region: whole exome
- Sequencing Depth : >50X ~ 100X
- Identify all kinds of variants including SNPs, INDELs and SV in coding region.
- Cost effective

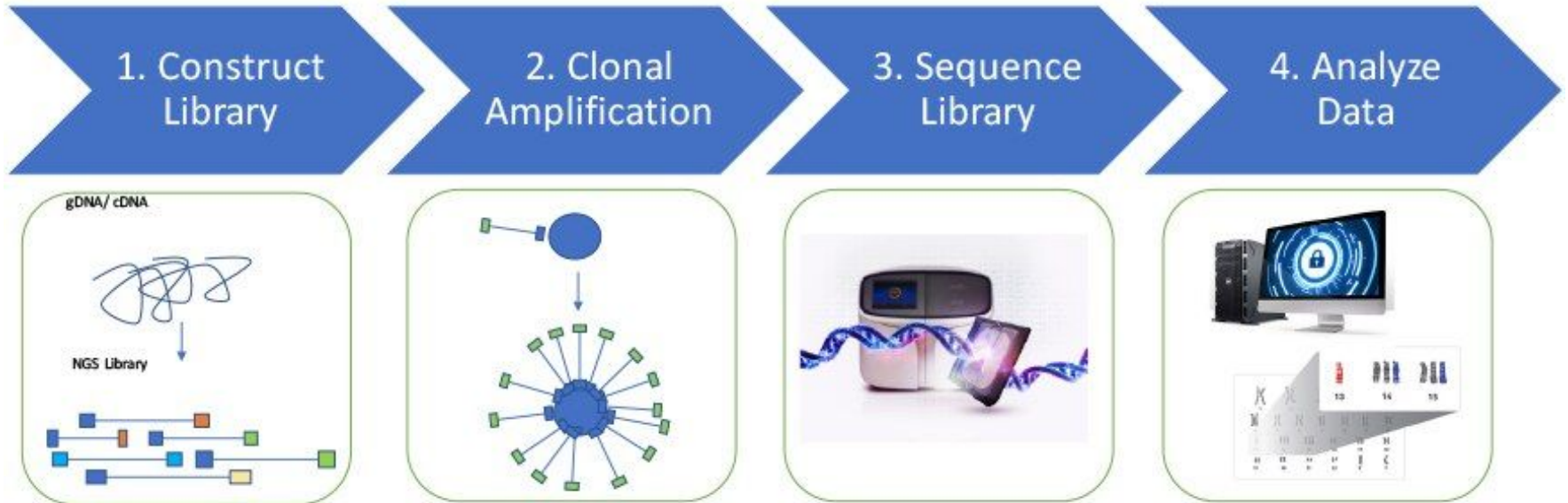
Targeted sequencing



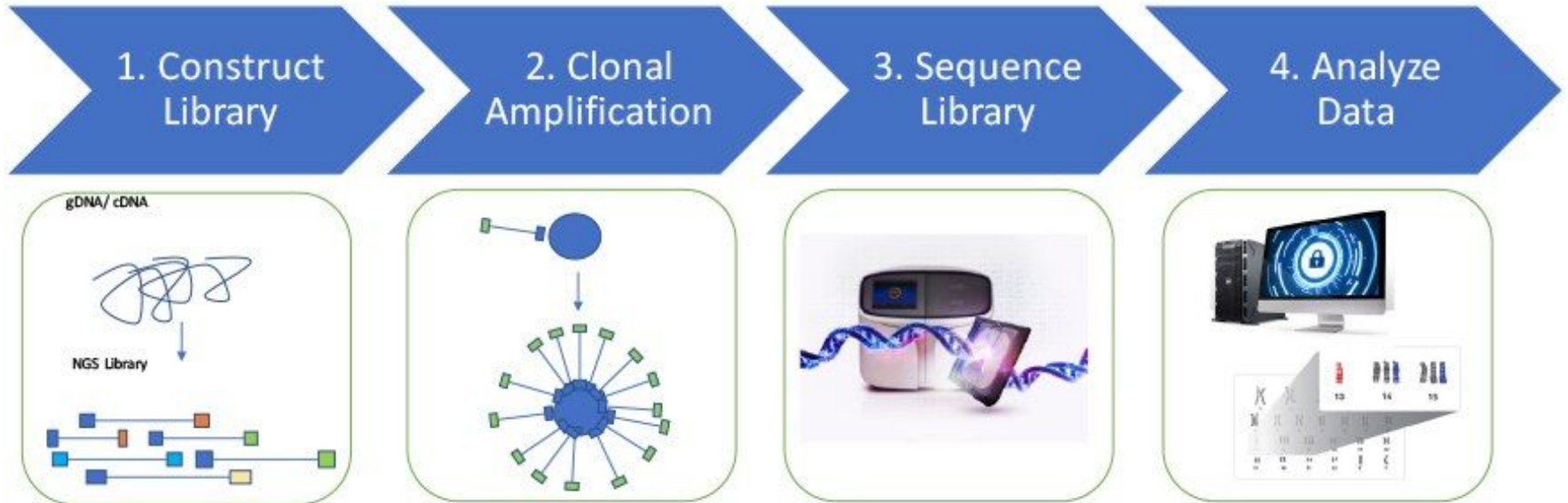
- Sequencing region: specific regions (could be customized)
- Sequencing Depth : >500X
- Identify all kinds of variants including SNPs, INDELs in specific regions
- Most Cost effective

- Targeted sequencing : rapid and cost-effective way to detect known and novel variants in selected sets of genes or genomic regions
- Whole exome sequencing : sequencing all of the protein-coding regions of genes in a genome (applications : discover rare-variants, adjacent splice-sites,...)
- Whole genome sequencing : alterations in regulatory sequences and non-coding regions, chromosomal rearrangements,

NGS workflow

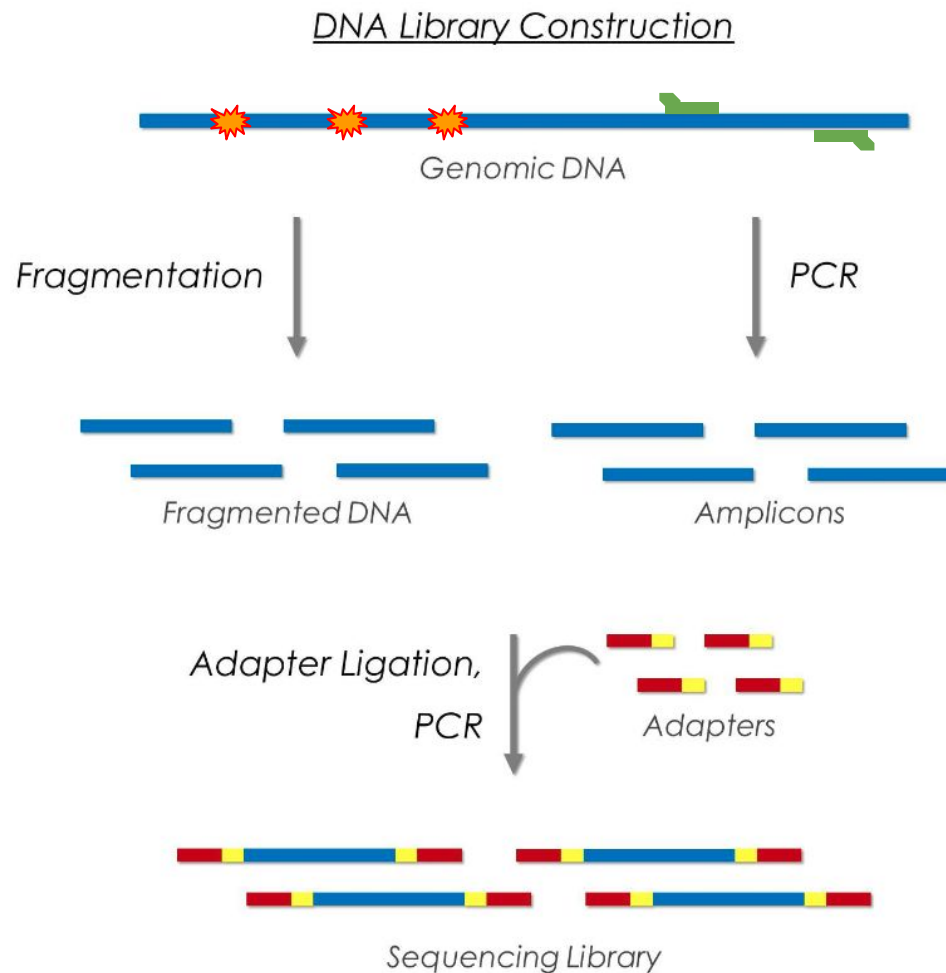


NGS workflow



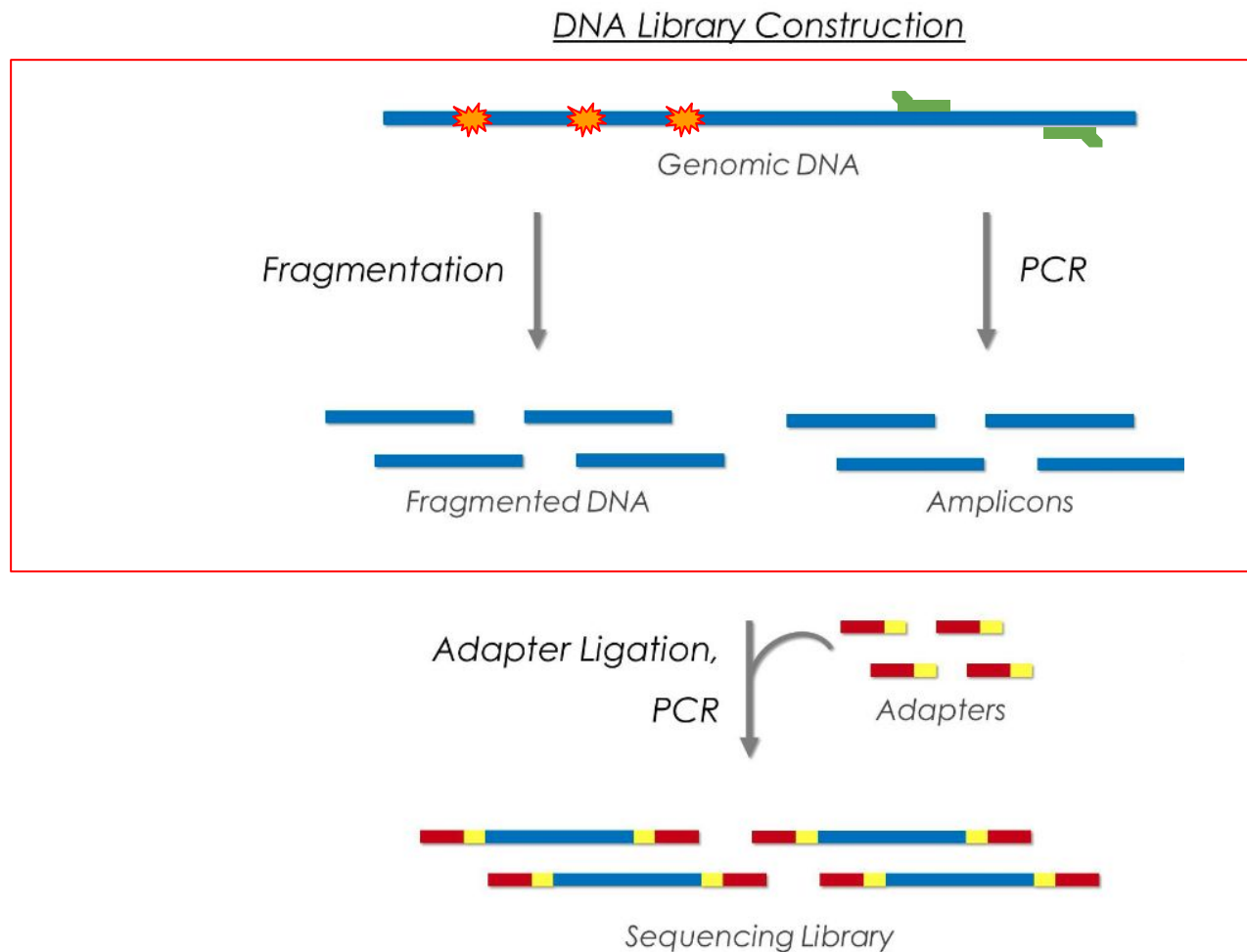
Library construction

A sequencing “library” must be created from the sample. The DNA (or cDNA) sample is processed into relatively short double-stranded fragments (100–800 bp)



Library construction

A sequencing "library" must be created from the sample. The DNA (or cDNA) sample is processed into relatively short double-stranded fragments (100–800 bp)

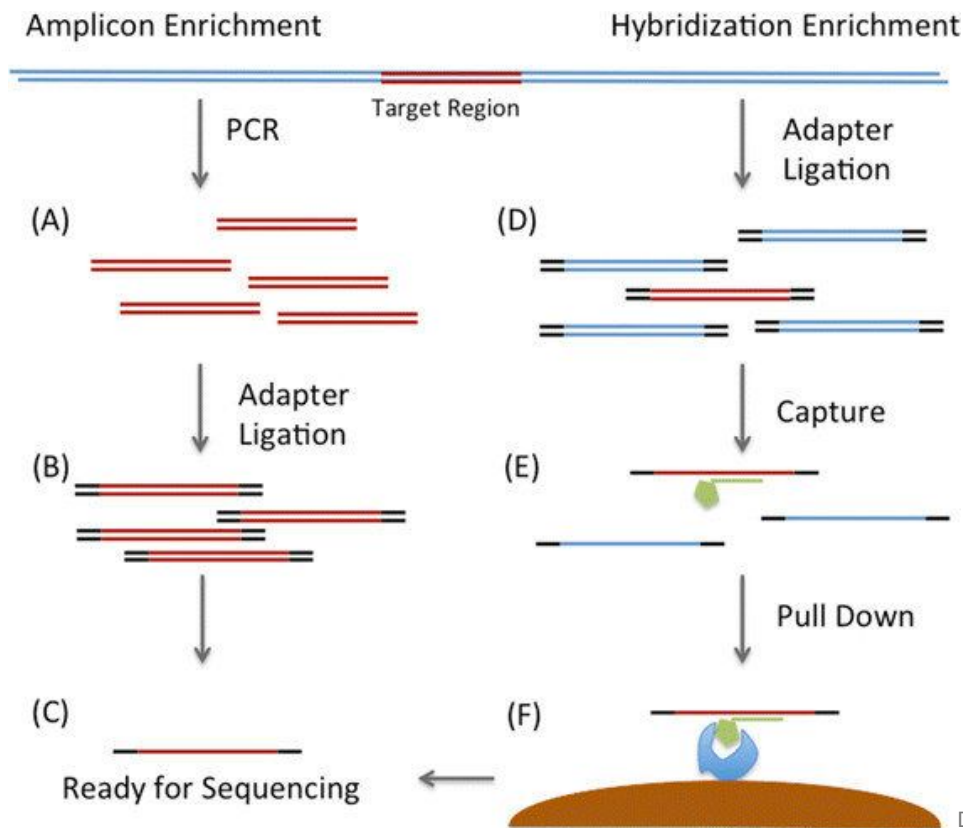


Library construction

Targeted sequencing : enrichment methods

Effective in enrichment and specificity
 Simple and fast protocol
 Target from Kb to Mb
 Low DNA input (100 ng)

HaloPlex
 AmpliSeq
 ...



Effective in enrichment and specificity
 Complex procedure
 Larger gene panels
 Higher DNA input (>1 µg)

Agilent's SureSelect
 Roche/Nimbelgen's SeqCap
 Illumina's TruSeq and Nextera
 ...

DOI: [10.1186/s13075-014-0490-4](https://doi.org/10.1186/s13075-014-0490-4)

The **BED** format is a text file format used to store genomic regions as coordinates and associated annotations

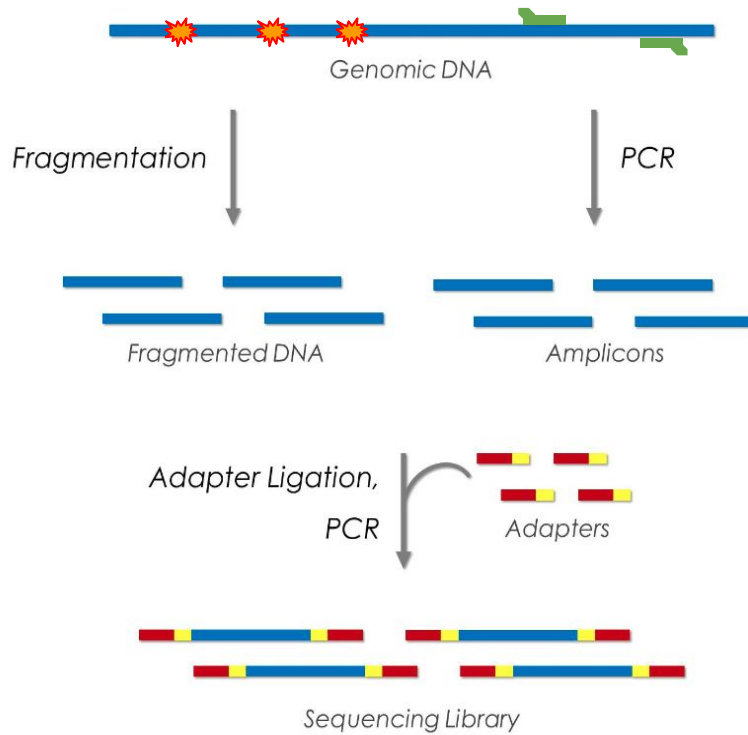
```
chr7 127471196 127472363
chr7 127472363 127473530
chr7 127473530 127474697
```

Library construction

Multiplex sequencing using DNA barcoding

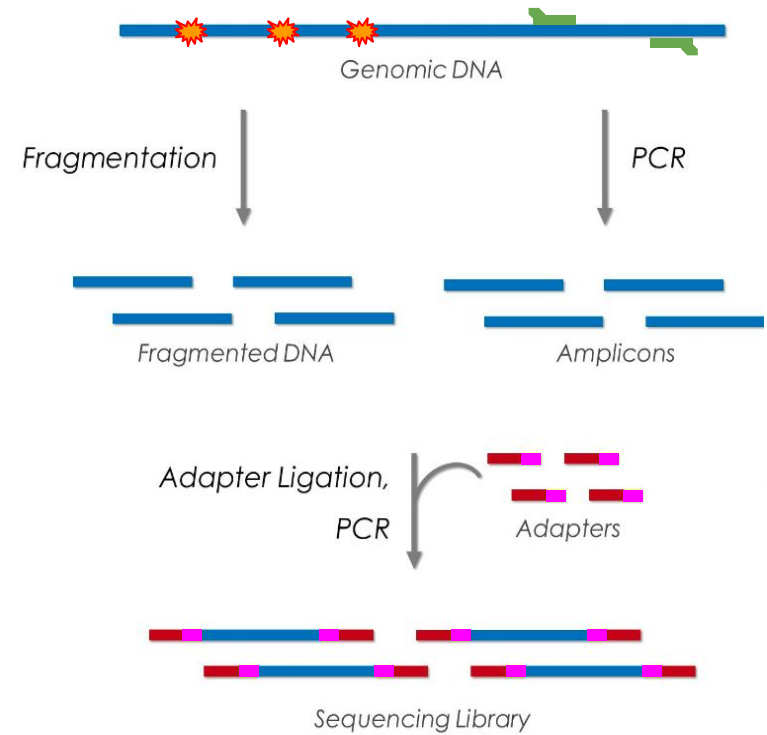
SAMPLE 1

DNA Library Construction



SAMPLE 2

DNA Library Construction



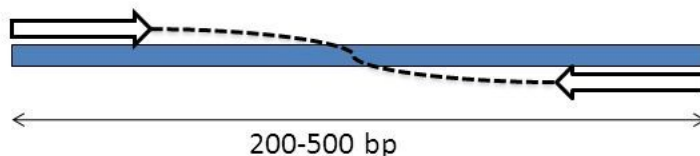
Single-end vs paired-end

- **Single-End Read:** When sequencing process only occurs in 1 direction
- **Paired-End Read:** When sequencing process occurs in both directions
- **Mate-pair Read:** Short fragments consisting of two segments that originally had a separation of several kilobases in the genome.

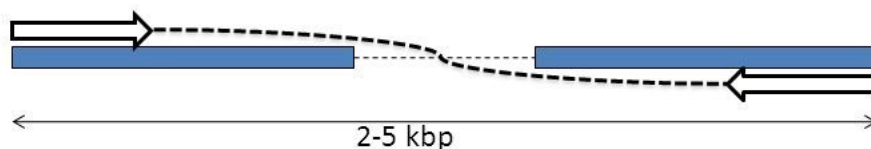
Single-End Reads - 5' or 3' (random)



Paired-End Reads - 5' and 3'



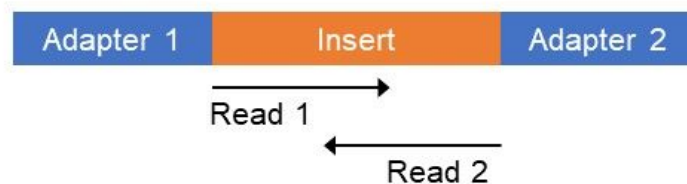
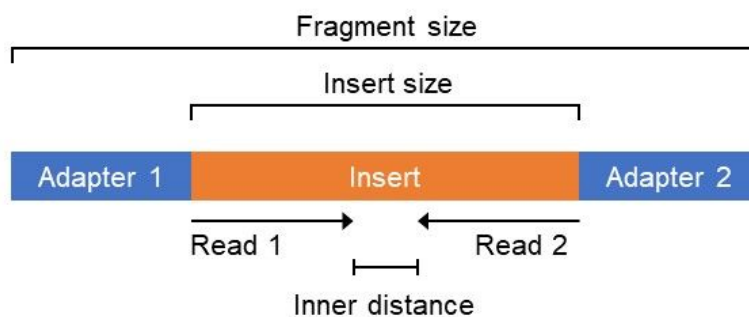
Mate-Pair Reads - 5' and 3'



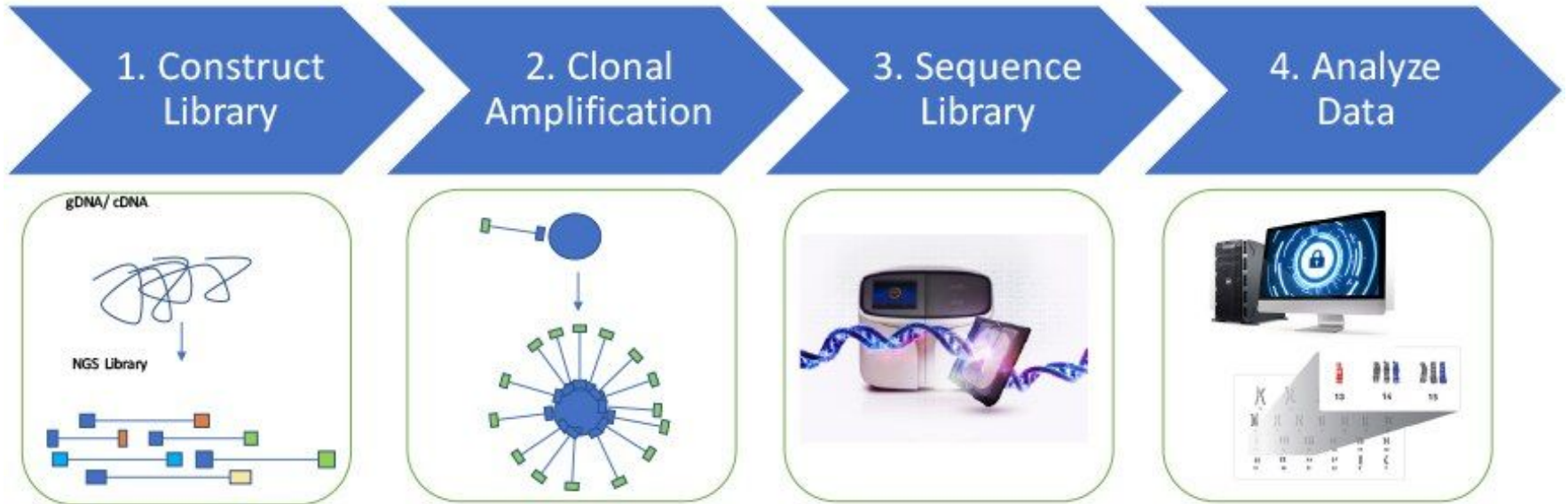
source: <http://slideplayer.com/slide/7847747/25/images/7/Types+of+Sequencing+Libraries.jpg>

Paired-end

- **The insert size** is the size of the piece of DNA of interest, without the adapters.

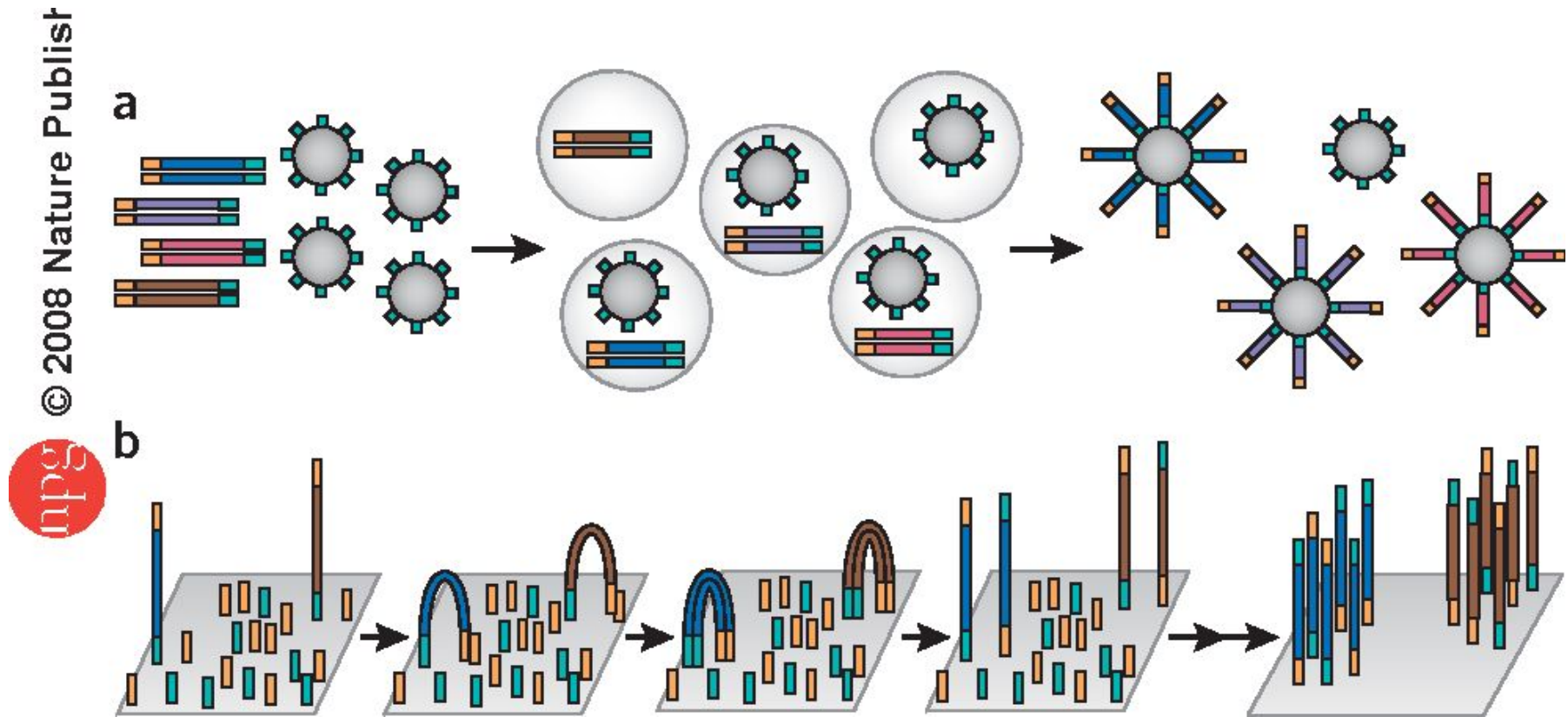


NGS workflow



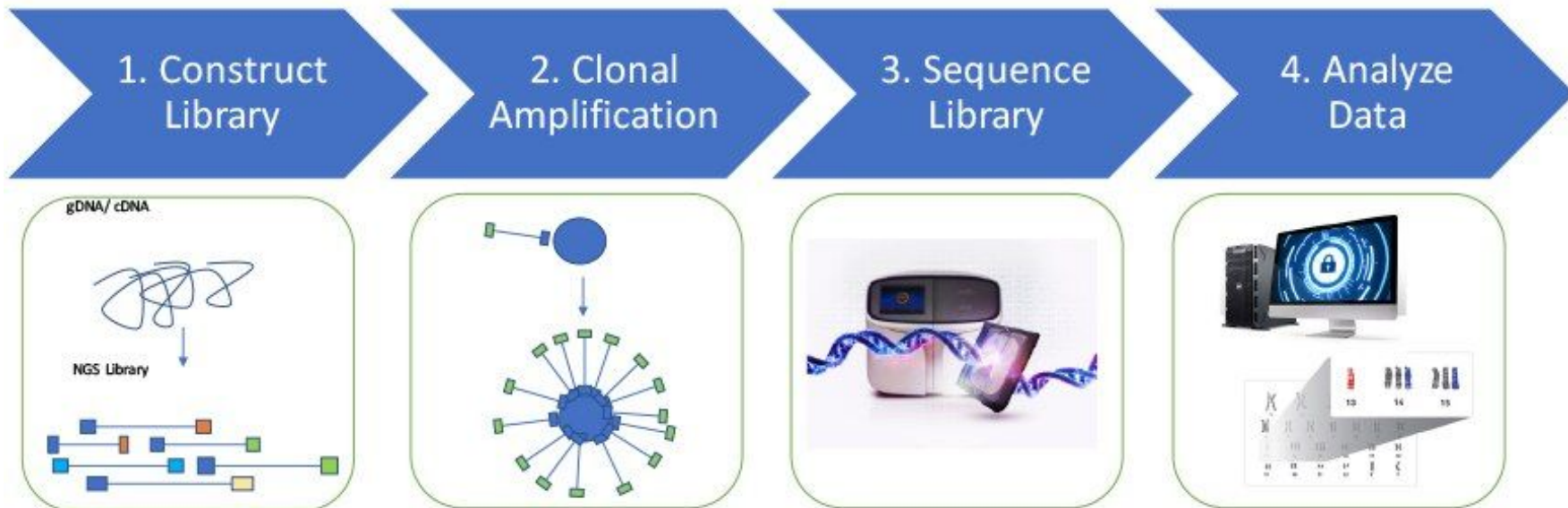
Clonal amplification

Prior to sequencing, the DNA library must be attached to a solid surface and clonally amplified to increase the signal that can be detected from each target during sequencing.



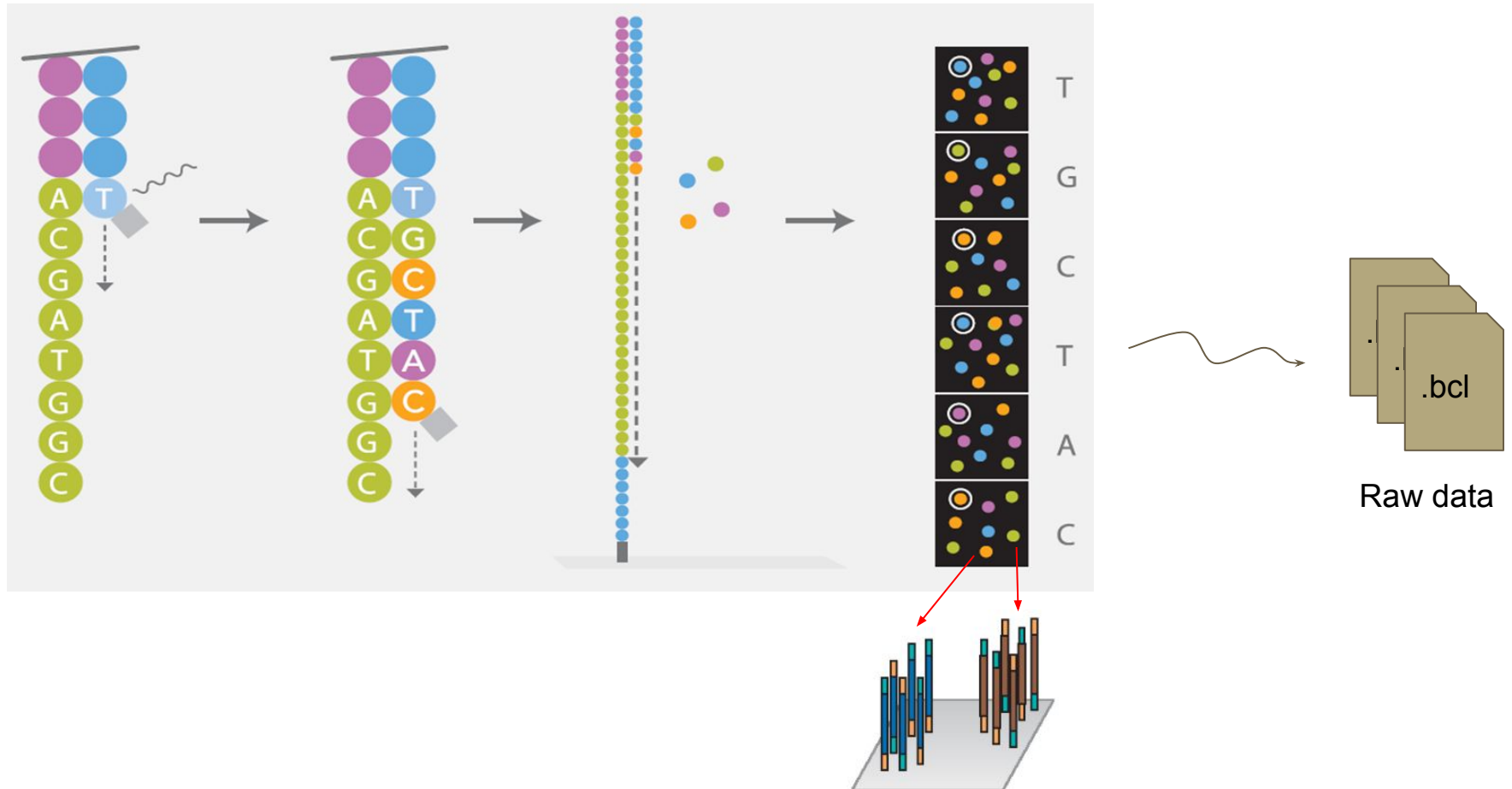
(a) thermofisher platforms rely on emulsion to amplify clonal sequencing features. (b) The Illumina technology relies on bridge PCR^{21,22} (aka 'cluster PCR') to amplify clonal sequencing features.

NGS workflow

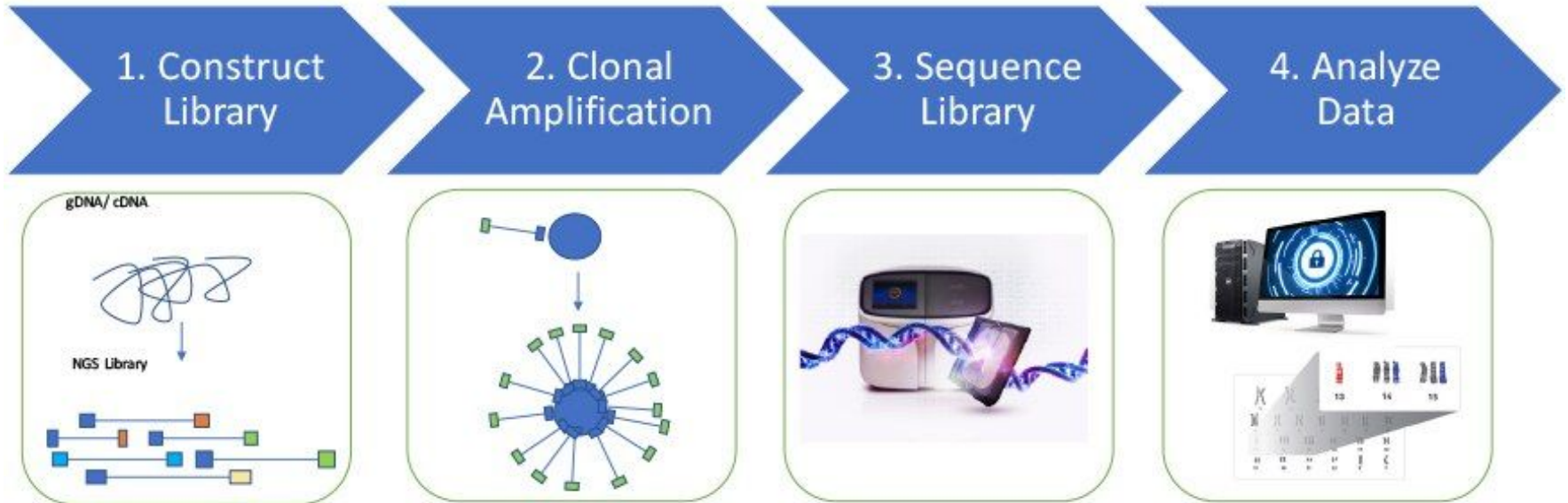


Sequencing

Illumina technology

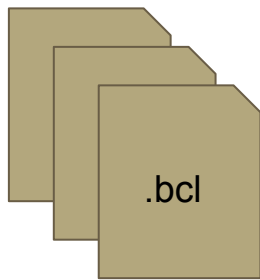


NGS workflow



Data analyses

Extracting reads, Demultiplexing



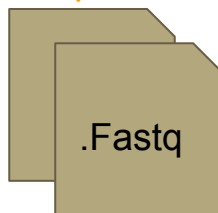
+

Sample Sheet

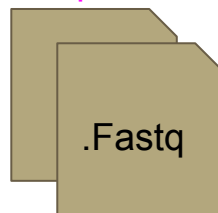
Sample Sheet							
[Header]							
IEMFileVersion	4						
Experiment Name	Project1						
Date	4/16/2016						
Workflow	GenerateFASTQ						
Application	NextSeq FASTQ Only						
Assay	TruSeq LT						
Description							
Chemistry	Default						
[Reads]							
	151						
	151						
[Settings]							
Adapter	AGATCGGAAGAGCACACGTCTGAACTCCAGTCA						
AdapterRead2	AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT						
[Data]							
Sample_ID	Sample_Name	Sample_Plate	Sample_Well	I7_Index_ID	index	Sample_Project	Description
Sample_1				A002	CGATGT		
Sample_2				A004	TGACCA		
Sample_3				A005	ACAGTG		
Sample_4				A006	GCCAAT		

bcl2fastq

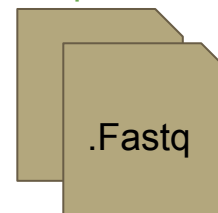
sample 1



sample 2

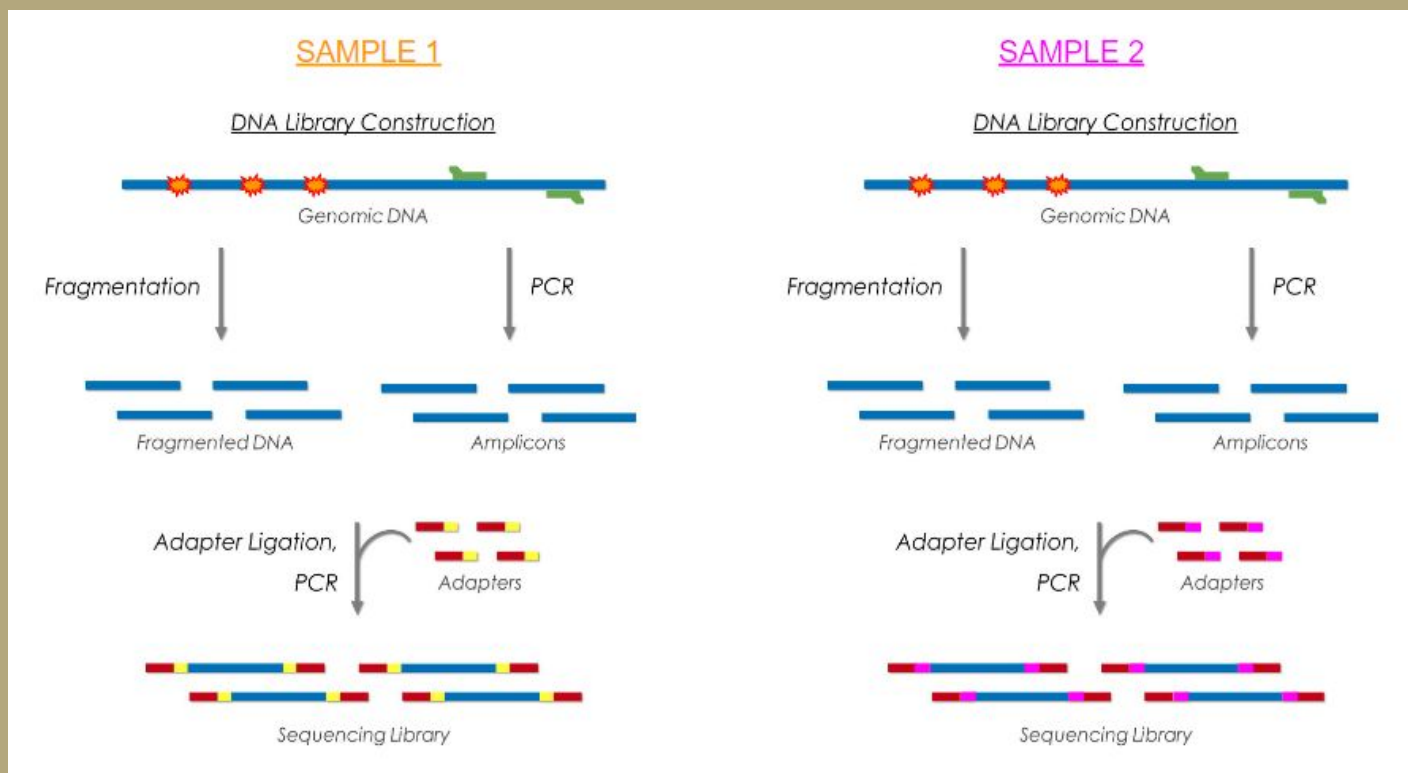


sample 3

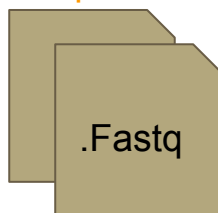


Data analyses

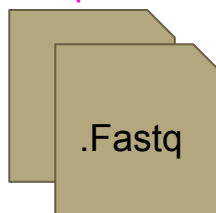
Extracting reads, Demultiplexing



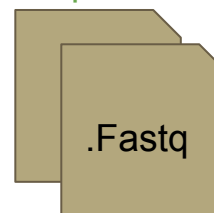
sample 1



sample 2



sample 3



Fastq files (Paired-end)

2 files : R1, R2

Reads1.fq

```
@ERR229776.100000840
CTAGGAAGCGTAGTCTGGGGTCATCTCTCTATTAATACTGTTGGGAATGTTTAGTA
+
BAEEAGEED96EHFE@BF><>EAAC;EBH<K<6:HJGFFHBC>DDIKG4AIHFFD@0/=
@ERR229776.100020365
CATTTATTTTCATAGTAGCCAAAAAGTGAAACAGTCAAAATATCCGTCAGTGAATTGACC
+
1.*/././&((&3=;B@F860C>@51(3:).6GG-68C*:CG)#B4/=HDJ6;79)<@C/
@ERR229776.100104918
TATTTCTGGAATTTTCCATTTAATATTTTCAGACTGCAGTTGACTGCGGGTAACTGAAA
+
CEEEEFEDAEAGGGFDHGFHGHIIHHHIIIGKHBKJJIGHFHKILJKLEJLJJJFJMJK
```

Reads2.fq

```
@ERR229776.100000840
TTCTGGTCAGTAAGACCTCAAAAGTTAAATACTAGCGATTTACACACCTTAAATGATT
+
CFIEEG@FFFGKFJHJ>HHKLLJIIJILLJIIILJHKAKJKKJJJJJLKMJKJJJKJ
@ERR229776.100020365
CCTAAAATGGTGTGTTTTTCGTATATTCACAATGCTGTGGAACCATCACCACTATCTGAT
+
4B@EDFF= (/CHBHEHCE6@ED8E@@I6HJB6E:6%@C46FFIBGCIGKD, DN=CBBE@
@ERR229776.100104918
TCTTTCTTTTGTTTTTTTTTCTGAGATGTCTTTTGTTTTTGTCTGAGGTCTTGTATG
+
CFIGGGKHHHFFHFIJIIJIKLIIHJIIIKLJKKIJKLLKJFJJMHJJLJFJMJIKKJJJ
```

1 interleaved paired file

Reads.fq

```
@SRR531199.1 ILLUMINA_0130:3:1101:1249:1993 length=101
TTTTTCAGAGTAGTTGGTACCCAAATTGGAAGATGTGACCCACTTCGATACCGCGCTTGAG
+
dffffffffdfeffdadffffeeefdeffeffeffffffffffddeeYdfefefe[e
@SRR531199.1 ILLUMINA_0130:3:1101:1249:1993 length=99
ANNNNNNCTTCGGTATNAACTGGGNNNGATGTTGAACTGGGTAAAGTCGAAGATCTG
+
BBBBBBSZTUVWO]YB_[cbabbWBBBBSVVUUgggadcdedbedcddffdegeggef
@SRR531199.2 ILLUMINA_0130:3:1101:1463:1964 length=101
NTGAGTAGCTCAATGCGCTGACGCCAATAGCTATACCAACGACTGGCCAGATTATGTTT
+
BXSSRU[X[Wcc_cccccccccccc_cccccccccccccccccccccccccccccccc
@SRR531199.2 ILLUMINA_0130:3:1101:1463:1964 length=99
AAGTGACCATCGCGATAAAGTCTGCGCAGTAAANAGCANCTGTTNGATGCTGGCTTA
+
ggggggggggggggggggfgfgggggggggggg^BbbbaBbbaZ]BZ[ccccfggggg
@SRR531199.3 ILLUMINA_0130:3:1101:1366:1970 length=101
NAAGTCGCGGCGACCCCTATCGTGGCTTTCGGCGTACGCCATTTCAATGCGCCGCCCG
+
B[[X[YY[YVcc_ccc_____][[V[^^^V[[]SXWUX[\\]]Z^^^B
@SRR531199.3 ILLUMINA_0130:3:1101:1366:1970 length=99
TGGTCAATACAAGCCGCAATACCTGCATCATGCGGNGGAANAATTTGCGCGCGTTTTCT
+
ggfegggggggdeggggfgcgggagggggggg^Bb`^]B[Y[[Zffffh_afeefe
```


Fastq file format

READ

1. Identifier

2. Sequence

4. Quality scores (as ASCII chars)

```
@SRR062641.6751359
CGCCCGGCCAATCATTGTGGTTTTAAGTCACTAAGTTTGAGGCTATTTTTGTTTTACAGCAAAGCTAACTGATGCAGACAGGGACAAGTCAGTCTCATCT
+
CBLNPGJQQQJPPQPPQPQRGPPPPRRRQORPSPGRQQQLRRRMEPQQPMJHQEQEHKMMFIIRH?SIIHKNJIKRLJJKIHEABHIFGCGGEFCGDGDCE
```

```
@SRR062634.16249693
CTAAGTTTGAGGCTATTTTTGTTTTACAGCAAAGCTAACTGATGCAGACAGGGACAAGTCAGTCTCATCTCTGTGCACCCAGCATTGCCCAGAACAGGGC
+
ALKMOOOOPPQJQOPPPPPQPPPPPPRJRQROQQORPQPRQQPFQSQQPRLLIMHKSJRQORMFELRPQNQRQJQRRPQQQLIRKDMKQJPN8CFDGDCCCB
```

```
@SRR062634.20060465
CTCCCAGCTTCCAACAGACCCTGTCCCAGCTCCCTCCAAGCTGAGTGTGGCCTGATACCTACCAGTGGAGCGAGGGGAACCCGAGGACTGCCAAGGGCA
+
D?KMPQEPGCPQQNPQIQIGR@DPERQHEKBED=HCHG8EHFD6<329@<:69A<6, ;<967>;=C:>AA8BBED#####
```

ASCII table:

Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char
0	0	[NULL]	32	20	[SPACE]	64	40	@	96	60	`
1	1	[START OF HEADING]	33	21	!	65	41	A	97	61	a
2	2	[START OF TEXT]	34	22	"	66	42	B	98	62	b
3	3	[END OF TEXT]	35	23	#	67	43	C	99	63	c
4	4	[END OF TRANSMISSION]	36	24	\$	68	44	D	100	64	d
5	5	[ENQUIRY]	37	25	%	69	45	E	101	65	e
6	6	[ACKNOWLEDGE]	38	26	&	70	46	F	102	66	f
7	7	[BELL]	39	27	'	71	47	G	103	67	g
8	8	[BACKSPACE]	40	28	(72	48	H	104	68	h
9	9	[HORIZONTAL TAB]	41	29)	73	49	I	105	69	i
10	A	[LINE FEED]	42	2A	*	74	4A	J	106	6A	j
11	B	[VERTICAL TAB]	43	2B	+	75	4B	K	107	6B	k
12	C	[FORM FEED]	44	2C	,	76	4C	L	108	6C	l
13	D	[CARRIAGE RETURN]	45	2D	-	77	4D	M	109	6D	m
14	E	[SHIFT OUT]	46	2E	.	78	4E	N	110	6E	n
15	F	[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

Module 1/6: Analyses ADN

- NGS Introduction
- Reads Quality Control
- Reads Cleaning
- Aligning reads on reference → *Hélène Touzet*
- Alignment parameters → *Hélène Touzet*
- Reads duplicates
- Assembly → *Hélène Touzet*

Reads quality

- Errors when reading bases
- Depends on sequencing technologie
- Error rate increases with read size

⇒ For each position in the read

- One base (A/T/C/G)
- One error probability

Phred Quality Score (for a base)

Phred quality scores Q : logarithmically related to the base-calling error probabilities P

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

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%

Quality score encoding

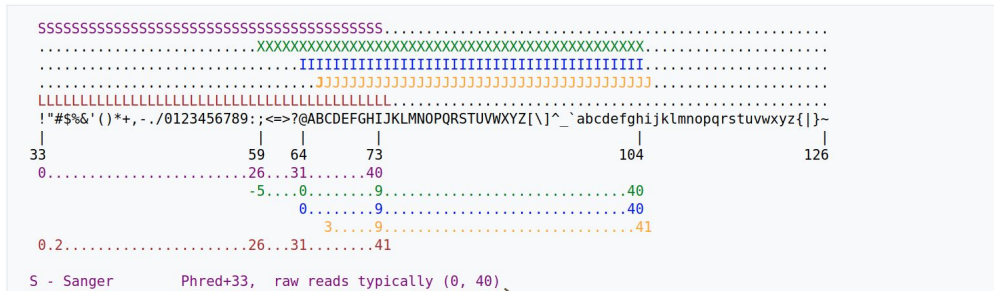
For history reasons, more than one coding convention



Source : <https://fr.wikipedia.org/wiki/FASTQ>

Galaxy : Always uses Sanger coding
=> conversion tool (*groomer*)

Example for score interpretation using sanger encoding



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%

S - Sanger Phred+33

Bad	:	0-19
Correct	:	20-29
Good	:	30-40

```
@SEQ:ID
ACTGTACGATCGATCGCATGATCAGTACGTCGTACCAGAT
+
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHI
|
0.....1.....2.....3.....4
01234567890123456789012345678901234567890
```

Goal: read cleaning

```
@SRR062641.6751359
CGCCCGGCCAATCATTGTGGTTTTAAGTCACTAAGTTTGAGGCTATTTTGTTTTACAGCAAAAGCTAACTGATGCAGACAGGGACAAGTCAGTCTCATCT
+
CBLNPGJQQQJPPQPPQPQRGPPPPRRQQRPS PGRQQQLRRRMEPQQPMJHQEHEKMMFIIRH?SIIHKNJIKRLJJIKEABHIFGCGGEFCGDGDCE
@SRR062634.16249693
CTAAGTTTGAGGCTATTTTGTTTTACAGCAAAAGCTAACTGATGCAGACAGGGACAAGTCAGTCTCATCTCTGTGCACCCAGCATTGCCCAGAACAGGGC
+
ALKMOOOOPPJQOPPPPPQPPPPPRJQRQQQQRPQPRQQPFQSQQPRLIMHKS NRJQORMFELRPQNQRQJQRRPQQLIRKDMKQJRFDFGCCCCB
@SRR062634.20060465
CTCCCAGCTTCCAACAGACCCTGTCCCAGCTCCCTCCAAGCTGAGTGTGGCCTGATACCTACCAGTGGAGCGAGGGGAACCCGAGGACTGCCAAGGGCA
+
D?KMPQEPGCPQONPQIQIGR@DPERQHEKBEHCHG8EHFD6<329@<:69A<6, ;<967>;=C:>AA8BBED#####
@SRR062635.15516129
AAAAAAAAAAAAAAAAAAAAAAAAAAGGGGGCCCCCTTTCCCCCGGGGGGGGACAGGGGGGTGTTCCGGCCCCGCGCCGCTTGACCACGG
+
EKLMPPPPPPQQQQQQQQQQQQQK#####
```

RAW



```
@SRR062641.6751359
CGCCCGGCCAATCATTGTGGTTTTAAGTCACTAAGTTTGAGGCTATTTTGTTTTACAGCAAAAGCTAACTGATGCAGACAGGGACAAGTCAGTCTCATCT
+
CBLNPGJQQQJPPQPPQPQRGPPPPRRQQRPS PGRQQQLRRRMEPQQPMJHQEHEKMMFIIRH?SIIHKNJIKRLJJIKEABHIFGCGGEFCGDGDCE
@SRR062634.16249693
CTAAGTTTGAGGCTATTTTGTTTTACAGCAAAAGCTAACTGATGCAGACAGGGACAAGTCAGTCTCATCTCTGTGCACCCAGCATTGCCCAGAACAGGGC
+
ALKMOOOOPPJQOPPPPPQPPPPPRJQRQQQQRPQPRQQPFQSQQPRLIMHKS NRJQORMFELRPQNQRQJQRRPQQLIRKDMKQJRFDFGCCCCB
@SRR062634.20060465
CTCCCAGCTTCCAACAGACCCTGTCCCAGCTCCCTCCAAGCTGAG
+
D?KMPQEPGCPQONPQIQIGR@DPERQHEKBEHCHG8EHFD
```

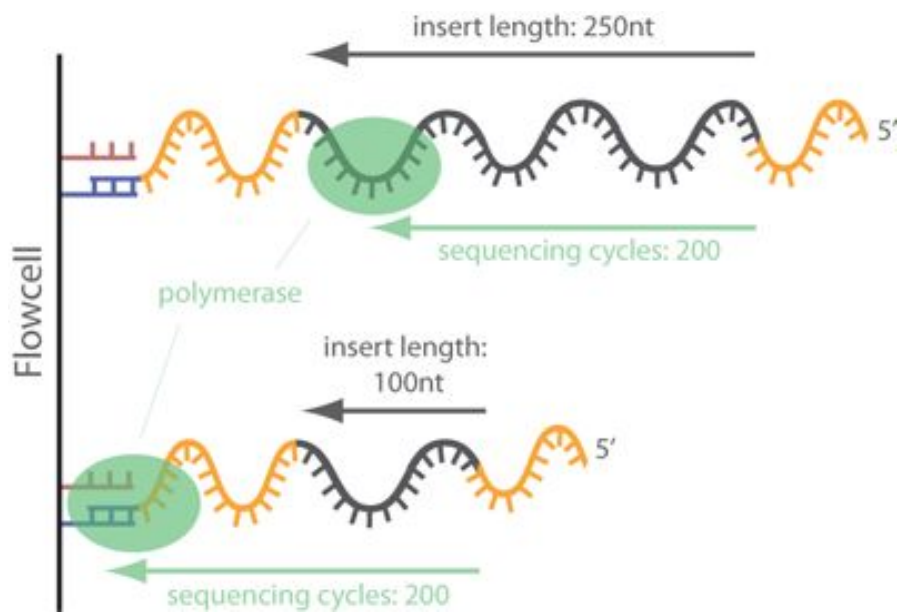
Clean

Module 1/6: Analyses ADN

- NGS Introduction
- Reads Quality Control
- Reads Cleaning
- Aligning reads on reference → *Hélène Touzet*
- Alignment parameters → *Hélène Touzet*
- Reads duplicates
- Assembly → *Hélène Touzet*

Reads cleaning

- Cut adaptators at read ends

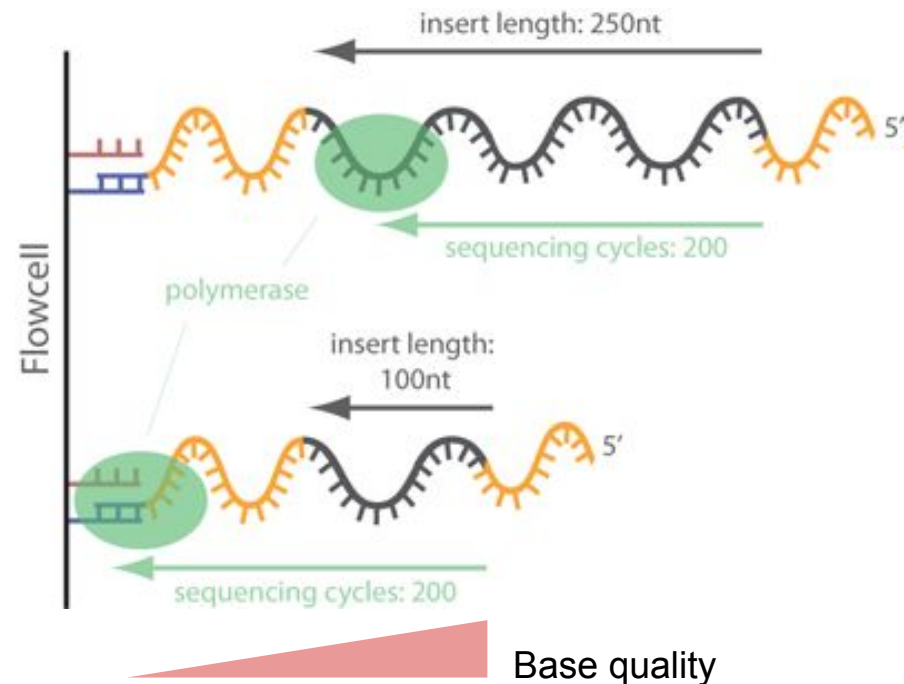


Reads cleaning

- Cut adaptators at read ends
- Trimming : cut read ends (5' ou 3')
 - Fixed number of bases
 - Individual base quality
 - Mean quality of bases in a sliding window
- Filtering : remove read
 - Size criteria (example $< 60\text{bp}$)
 - Mean base quality for all bases criteria (example < 25)

Reads cleaning

- Cut adaptators at read ends
- Trimming : cut read ends (5' ou 3')
 - Fixed number of bases
 - Individual base quality
 - Mean quality of bases in a sliding window
- Filtering : remove read
 - Size criteria (example $< 60\text{bp}$)
 - Mean base quality for all bases criteria (example < 25)



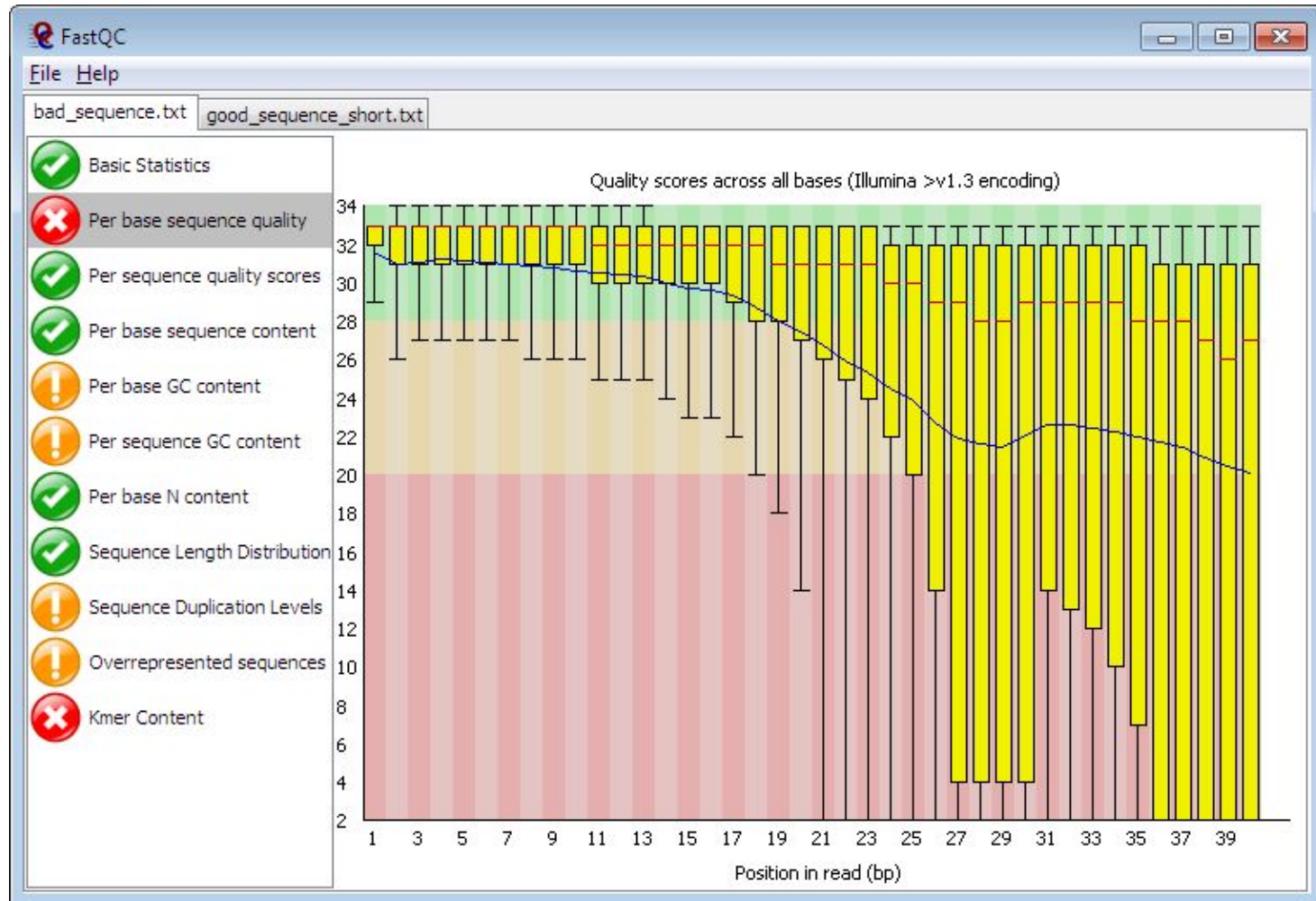
Reads cleaning example

Tool: Trimmomatic



Reads quality control (FastQC)

FastQC aims to provide a simple way to do some quality control checks on raw sequence data coming from high throughput sequencing



Workflow

