

Analisi bioinformatica nell'era NGS: biotools a supporto della tipizzazione HLA.

Veronica Scaglia

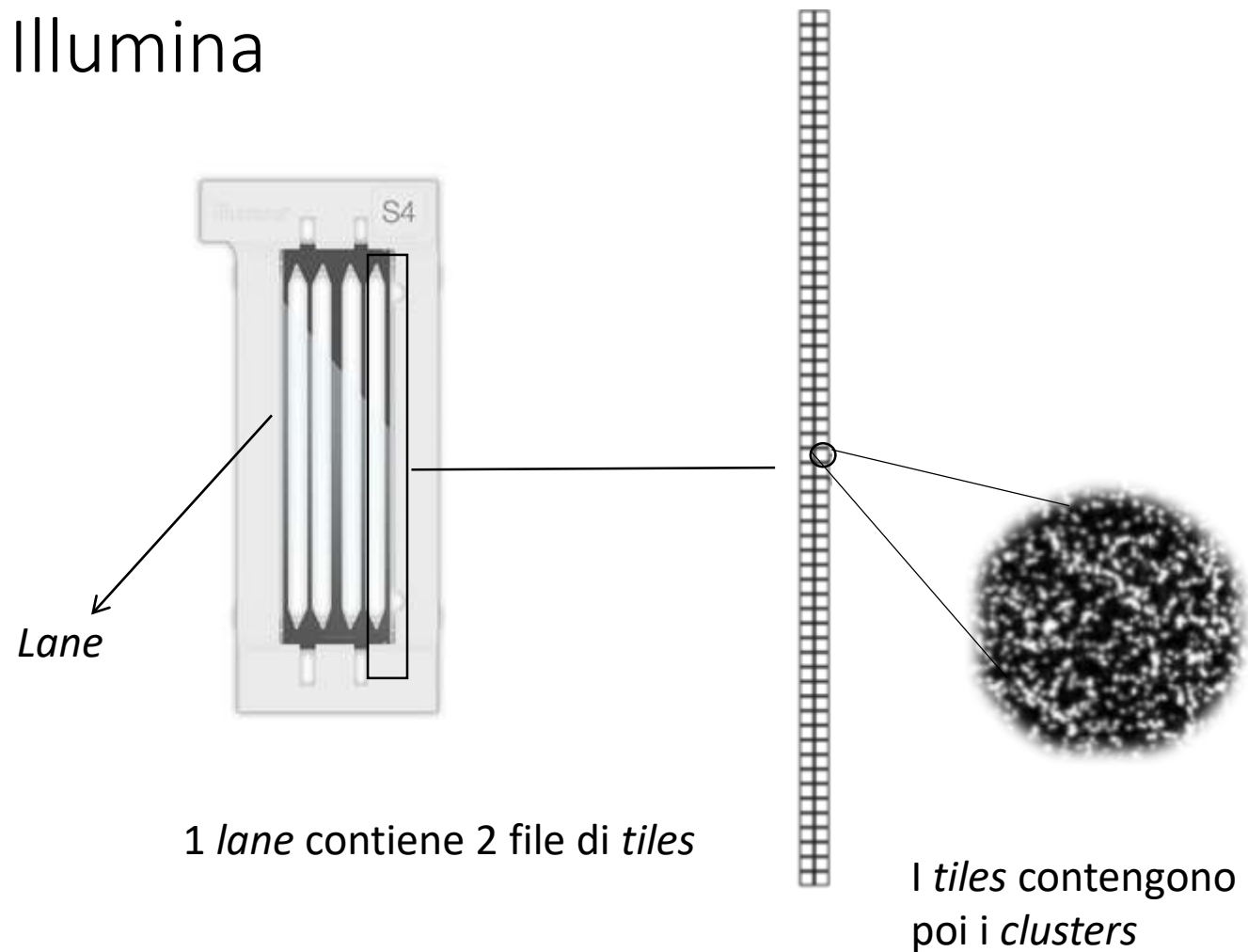
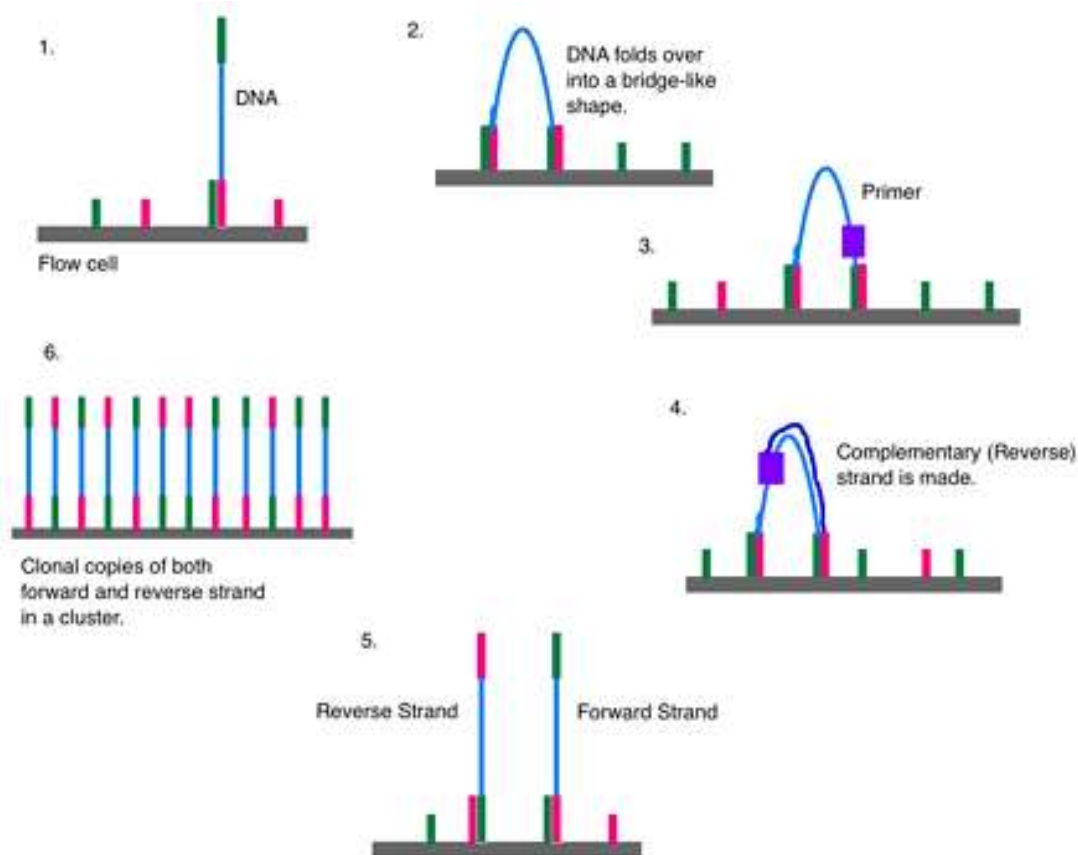
UO Biologia dei Trapianti, Diagnostica Molecolare e Manipolazione CSE

Dipartimento di Patologia Clinica, AUSL Piacenza

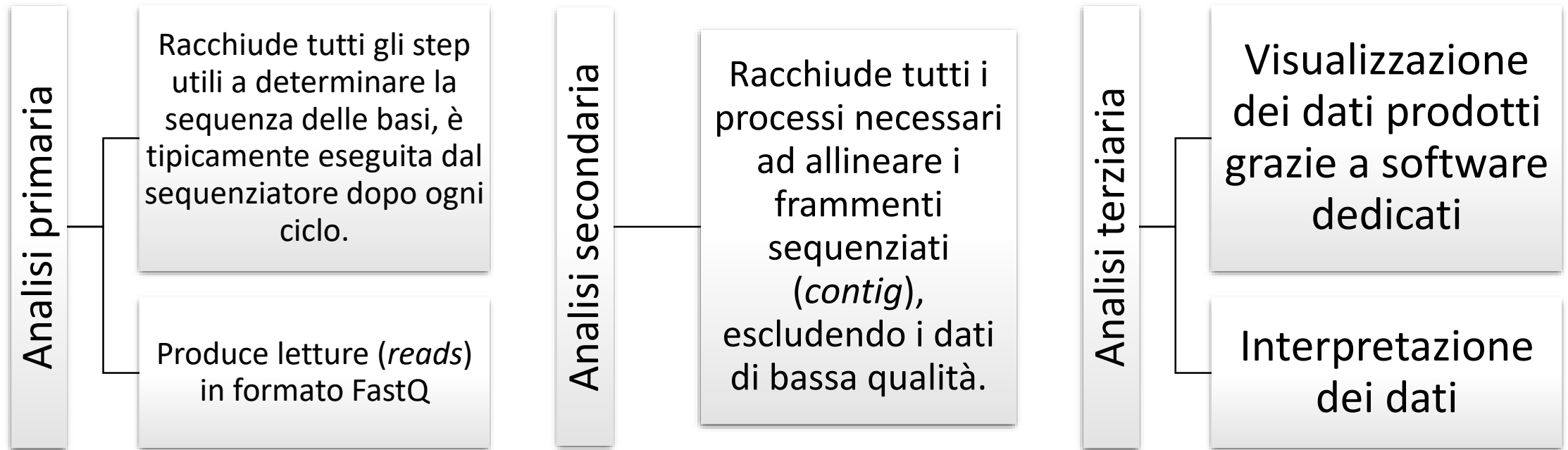


SERVIZIO SANITARIO REGIONALE
EMILIA ROMAGNA
Azienda Unità Sanitaria Locale di Piacenza

Sequenziamento NGS piattaforma Illumina



Analisi dei dati NGS: *workflow*



A. RTA Logs folder— Contains log files that describe each step performed by RTA for each Read.

B. InterOp folder— Contains binary files used by Sequencing Analysis Viewer (SAV) to summarize various primary analysis metrics.

C. Logs folder— Contains log files that describe every step performed by the instrument for each cycle.

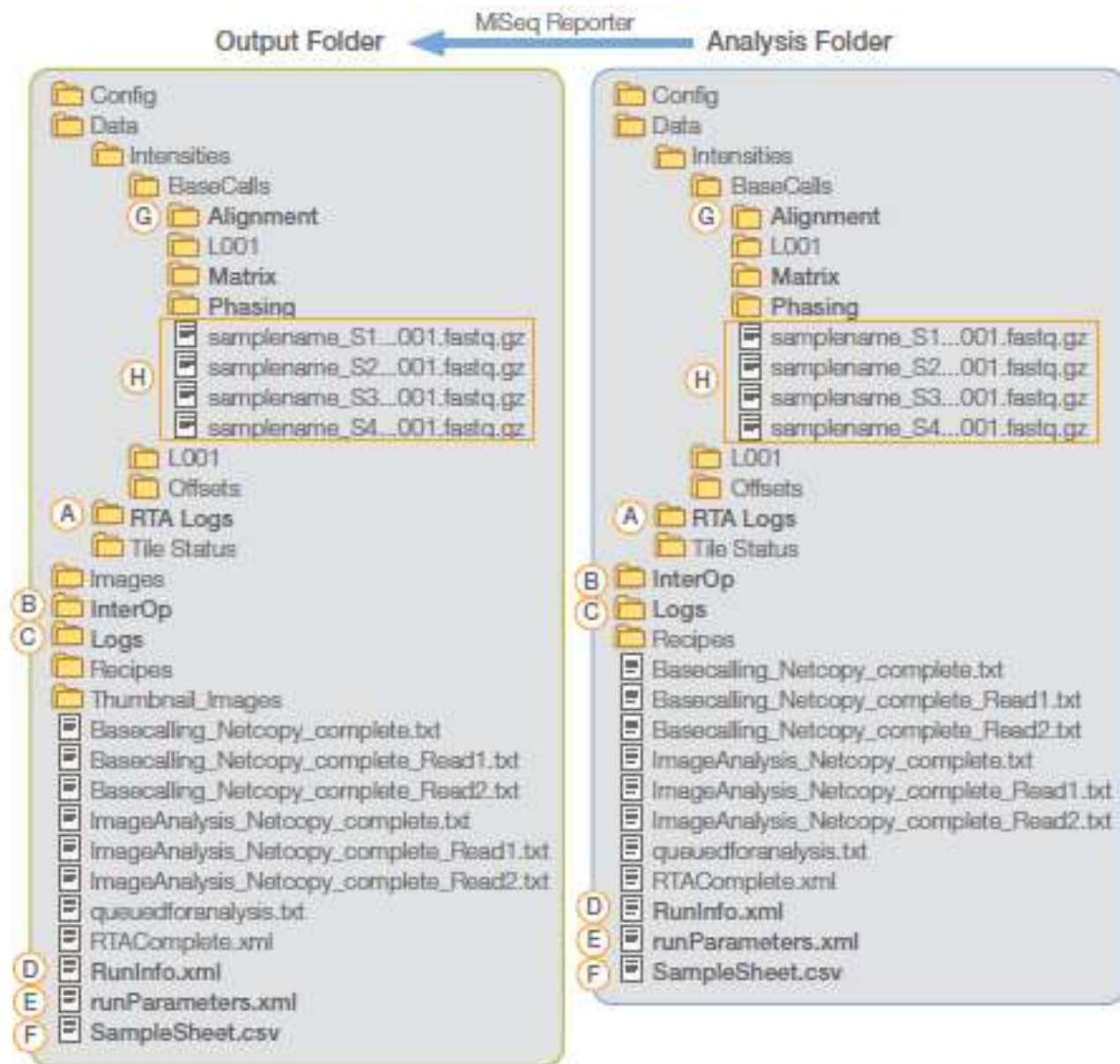
D. RunInfo.xml— Contains high-level run information, such as the number of Reads and cycles in the sequencing run.

E. runParameters.xml— Contains a summary of run parameters and information about run components.

F. SampleSheet.csv— Provides parameters for the run and subsequent analysis.

G. Alignment folder— Contains alignment (*.bam) and variant call (*.vcf) files. This folder is not created if you set up your run with the GenerateFASTQ workflow.

H. FASTQ files— Provide intermediate analysis data for downstream third-party analysis. Files are located in Data\Intensities\BaseCalls.



Sequencing Analysis Viewer

Analysis Tab

- **Flow Cell Chart** mostra i parametri di qualità delle sezioni sulla flowcell in codice colore
- **Data by Cycle** mostra un grafico che permette di seguire la progressione dei parametri di qualità durante tutta la corsa
- **Data by Lane** mostra i parametri di qualità per ogni lane
- **Q Score Distribution Plot** mostra in un grafico il numero di reads che superano il filtro di qualità
- **Q score Heatmap** mostra il Q score per ogni ciclo di sequenziamento.



Sequencing Analysis Viewer

Imaging Tab

- Mostra una lista dettagliata dei parametri e dei dati della corsa.
- Ogni lane, ciclo, sezione può essere selezionata per valutarne le caratteristiche specifiche.

Sequencing Analysis Viewer

Run Folder: Y:\101029_P22_0759_BFC805GRAB Browse Refresh

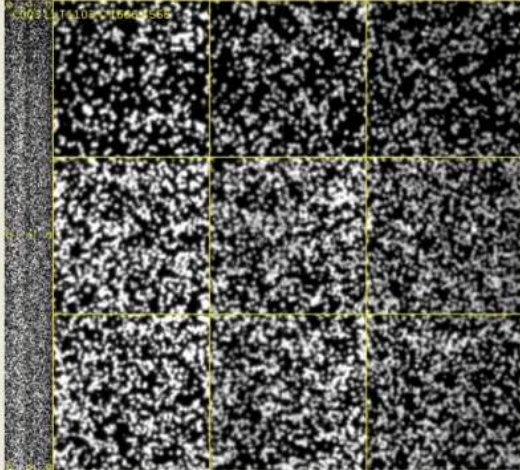
Analysis | **Imaging** | Summary | Tile Status | Controls

Cycle 1 | Lane 1 | Surface Top | Swath All | Section 3

☐ A ☐ C ☐ G ☐ T

Index	Lane	Tile	Section	Cycle	Surface	Swath	Time	P90 A	P90 C	P90 G	P90 T
1297	1	1103	3	1	Top	1	10/29/201...	3326	3944	1729	41
1298	1	1103	3	2	Top	1	10/29/201...	3306	3906	1741	40
1299	1	1103	3	3	Top	1	10/29/201...	3259	3853	1699	40
1300	1	1103	3	4	Top	1	10/29/201...	3211	3806	1720	40
1301	1	1103	3	5	Top	1	10/30/201...	3087	3703	1691	37
1302	1	1103	3	6	Top	1	10/30/201...	3066	3677	1610	35
1303	1	1103	3	7	Top	1	10/30/201...	2959	3554	1577	35
1304	1	1103	3	8	Top	1	10/30/201...	2945	3547	1558	35
1305	1	1103	3	9	Top	1	10/30/201...	2918	3518	1548	35
1306	1	1103	3	10	Top	1	10/30/201...	2906	3491	1536	35
1307	1	1103	3	11	Top	1	10/30/201...	2881	3461	1512	34
1308	1	1103	3	12	Top	1	10/30/201...	2869	3452	1498	34
1309	1	1103	3	13	Top	1	10/30/201...	2835	3423	1510	35
1310	1	1103	3	14	Top	1	10/30/201...	2854	3443	1493	34
1311	1	1103	3	15	Top	1	10/30/201...	2857	3445	1482	34
1312	1	1103	3	16	Top	1	10/30/201...	2811	3392	1452	33
1313	1	1103	3	17	Top	1	10/30/201...	2779	3356	1471	34
1314	1	1103	3	18	Top	1	10/30/201...	2756	3338	1435	33
1315	1	1103	3	19	Top	1	10/30/201...	2778	3364	1419	33
1316	1	1103	3	20	Top	1	10/30/201...	2749	3335	1435	33
1317	1	1103	3	21	Top	1	10/30/201...	2740	3324	1434	33
1318	1	1103	3	22	Top	1	10/30/201...	2692	3271	1441	33
1319	1	1103	3	23	Top	1	10/30/201...	2721	3308	1423	33
1320	1	1103	3	24	Top	1	10/30/201...	2687	3270	1420	33
1321	1	1103	3	25	Top	1	10/30/201...	2656	3229	1398	33

Rows=20736 Disp=162 Sel=1 Filter



Sequencing Analysis Viewer

Summary Tab

- Mostra i dati di qualità grezzi in una tabella riassunti per lane e per read. Le statistiche mostrano medie e deviazioni standard delle varie lane.
- Cluster Density: 600-1200 K/mm²
- Cluster Passing Filter: >80%
- % reads \geq Q30: >90%

Sequencing Analysis Viewer

Run Folder: D:\Illumina\MiSeqAnalysis\160816_M70109_0001_000000000-AJTY0

Browse

Refresh

Analysis

Imaging

Summary

Indexing

Read 3 (I)	0.13	0.13	0.00	0.00	60	97.44
Read 4	2.25	2.25	1.43	0.53	117	94.14
Non-Indexed Total	4.50	4.50	1.47	0.40	126	95.45
Total	4.76	4.76	1.47	0.40	232	95.48

Read 1

Lane	Tiles	Density (K/mm2)	Cluster PF (%)	Phas/Prephas (%)	Reads (M)	Reads PF (M)	% >= Q30	Yield (G)	Cycles Err Rated	Aligned (%)	Error Rate (%)	Error Rate 35 cycle (%)	Error Rate 75 cycle (%)
1	28	1078 +/- 21	91.22 +/- 1.29	0.070 / 0.134	20.54	18.74	96.77	2.25	120	1.51 +/- 0.04	0.27 +/- 0.03	0.13 +/- 0.01	0.17 +/- 0.02

Read 2 (I)

Lane	Tiles	Density (K/mm2)	Cluster PF (%)	Phas/Prephas (%)	Reads (M)	Reads PF (M)	% >= Q30	Yield (G)	Cycles Err Rated	Aligned (%)	Error Rate (%)	Error Rate 35 cycle (%)	Error Rate 75 cycle (%)
1	28	1078 +/- 21	91.22 +/- 1.29	0.000 / 0.000	20.54	18.74	94.61	0.13	0	0.00 +/- 0.00	0.00 +/- 0.00	0.00 +/- 0.00	0.00 +/- 0.00

Read 3 (I)

Lane	Tiles	Density (K/mm2)	Cluster PF (%)	Phas/Prephas (%)	Reads (M)	Reads PF (M)	% >= Q30	Yield (G)	Cycles Err Rated	Aligned (%)	Error Rate (%)	Error Rate 35 cycle (%)	Error Rate 75 cycle (%)
1	28	1078 +/- 21	91.22 +/- 1.29	0.000 / 0.000	20.54	18.74	97.44	0.13	0	0.00 +/- 0.00	0.00 +/- 0.00	0.00 +/- 0.00	0.00 +/- 0.00

Read 4

Lane	Tiles	Density (K/mm2)	Cluster PF (%)	Phas/Prephas (%)	Reads (M)	Reads PF (M)	% >= Q30	Yield (G)	Cycles Err Rated	Aligned (%)	Error Rate (%)	Error Rate 35 cycle (%)	Error Rate 75 cycle (%)
1	28	1078 +/- 21	91.22 +/- 1.29	0.076 / 0.081	20.54	18.74	94.14	2.25	120	1.43 +/- 0.03	0.53 +/- 0.08	0.37 +/- 0.23	0.44 +/- 0.11

Il file FASTQ

Nome Strumento ID corsa ID flowcell lane flowcell numero del *tile* Coordinate del *cluster* Paired-end N read ha passato il filtro, altrimenti Y indice

```

1 @M06634:435:000000000-GLH9R:1:1101:19310:1923 2:N:0:11
2 CTCAGAAATGTGAGGAGCCTTTGGTAGCGGCTCTCTCCATAGACTTTTCCAGTGGAGGAAATAGTGC
3 +
4 1>111111111B331111AAFGG10A111000ABGFHBA1D211ADGHH21212A01A//00B11221
  
```

1. Si riferisce all'identificativo della sequenza (inizia con @)
2. La seconda riga mostra la sequenza grezza
3. Contiene il simbolo "+" seguito da spazio
4. Contiene i parametri di qualità della sequenza codificati come conversione decimale del codice ASCII (**A**merican **S**tandard **C**ode for **I**nformation **I**nterchange)

I file *.SAM e *.BAM

- Il SAM (Sequence Alignment Map) è un formato di archiviazione di file di allineamento con le relative coordinate di interpretazione.
- Viene utilizzato in diversi *workflow* di sequenziamento per allineare le reads ad una sequenza di riferimento.
- BAM è un file SAM compresso, scritto in codice binario. Non decifrabile da un operatore, ma solo da software dedicati.
- Alcuni tools permettono di convertire i file *.SAM in *.BAM (es. Samtools)

Controllo qualità dei FASTQ nella tipizzazione HLA

- Permette di affiancare i software commerciali per migliorare la qualità del dato;
- Permette di identificare eventuali problemi sorti durante il sequenziamento;
- Utile in fase di validazione / monitoraggio continuo per identificare le caratteristiche intrinseche della metodica;
- Identificazione e miglioramento dell'eventuale errore introdotto durante le fasi della metodica.

FastQC

- Tool gratuito di controllo qualità dei dati NGS
- File input *.fastq, *.bam, *.sam
- Generazione di un report *.html con i parametri di qualità del file inserito (visualizzabile tramite browser)
- In caso di sequenziamenti *paired-end* vanno inseriti entrambi i *files* *.fastq generati e verranno analizzati separatamente.



Andrews, S. (2010). FastQC: A Quality Control Tool for High Throughput Sequence Data

- Configuration
- Help
- net
- org
- Templates
- uk
- cisd-jhdf5.jar
- fastqc
- fastqc_icon.ico
- htsjdk.jar
- INSTALL.txt
- jbzip2-0.9.jar
- LICENSE
- LICENSE.txt
- LICENSE_JHDF5.txt
- README.md
- README.txt
- RELEASE_NOTES.txt
- run_fastqc.bat

adapter_list.txt
contaminant_list.txt
limits.txt

1 Introduction
2 Basic Operations
3 Analysis Modules

N.B. Questi file *.txt
possono essere modificati
manualmente
dall'operatore per adattarli
alle esigenze specifiche di
analisi.



FastQC Report

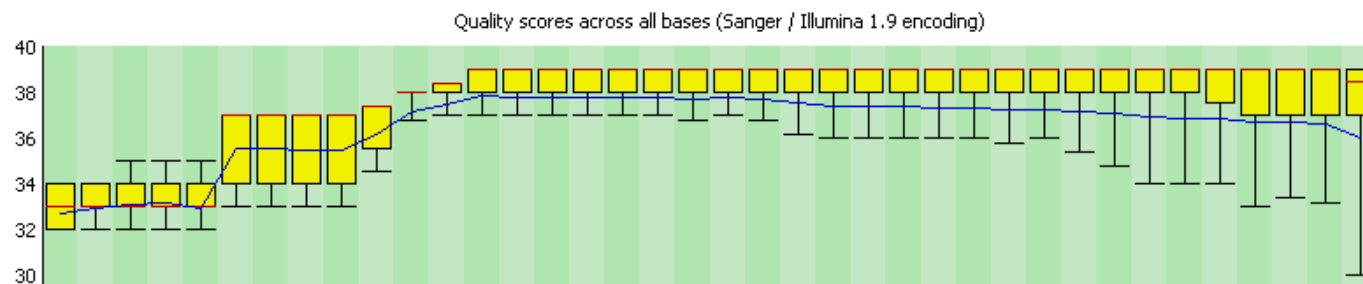
Summary

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

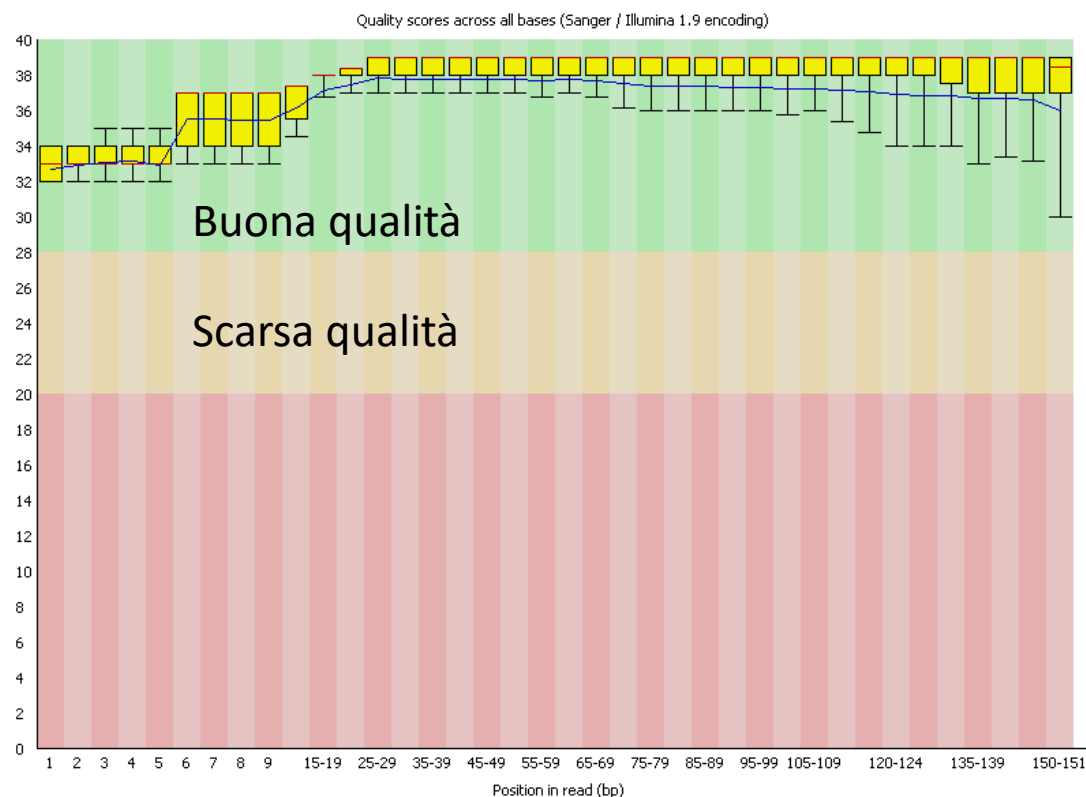
✓ Basic Statistics

Measure	Value
Filename	HA-241163-514-706_S8_L001_R1_001.fastq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	313531
Total Bases	46.6 Mbp
Sequences flagged as poor quality	0
Sequence length	35-151
%GC	50

✓ Per base sequence quality

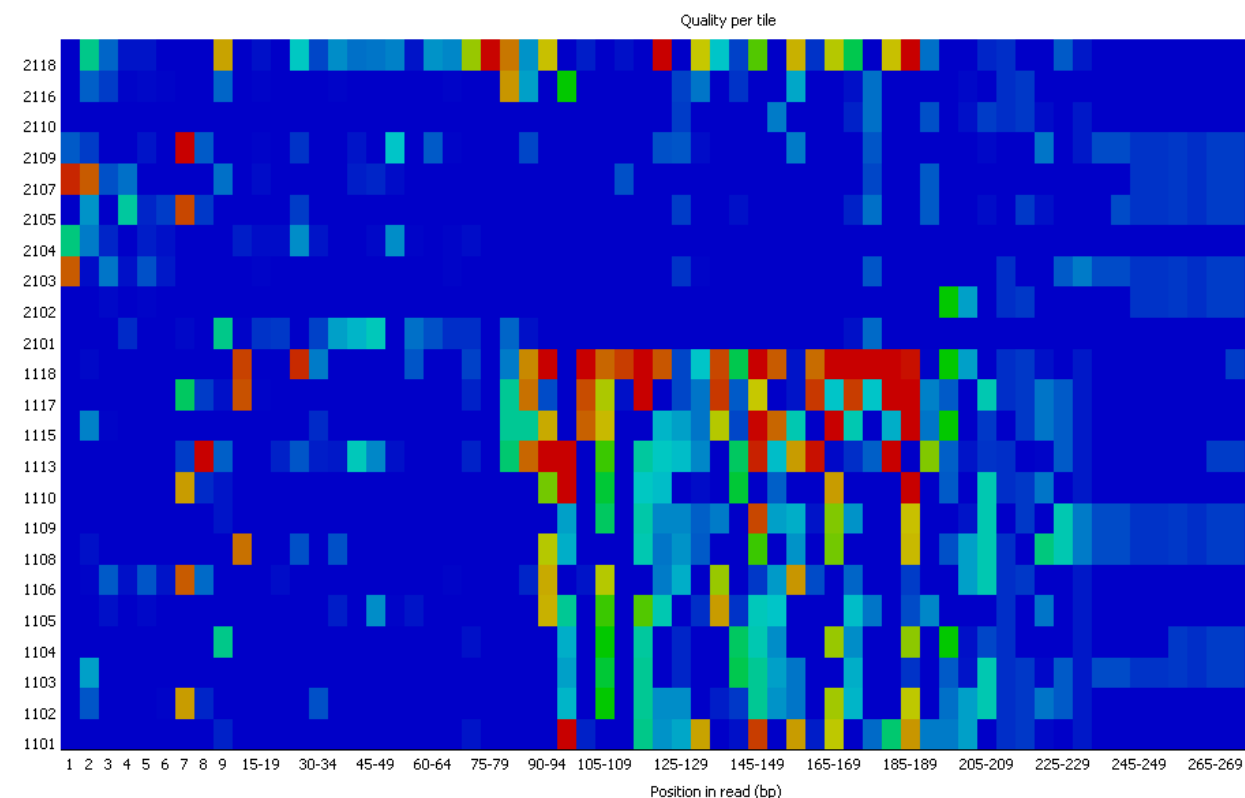


✓ **Per base sequence quality**



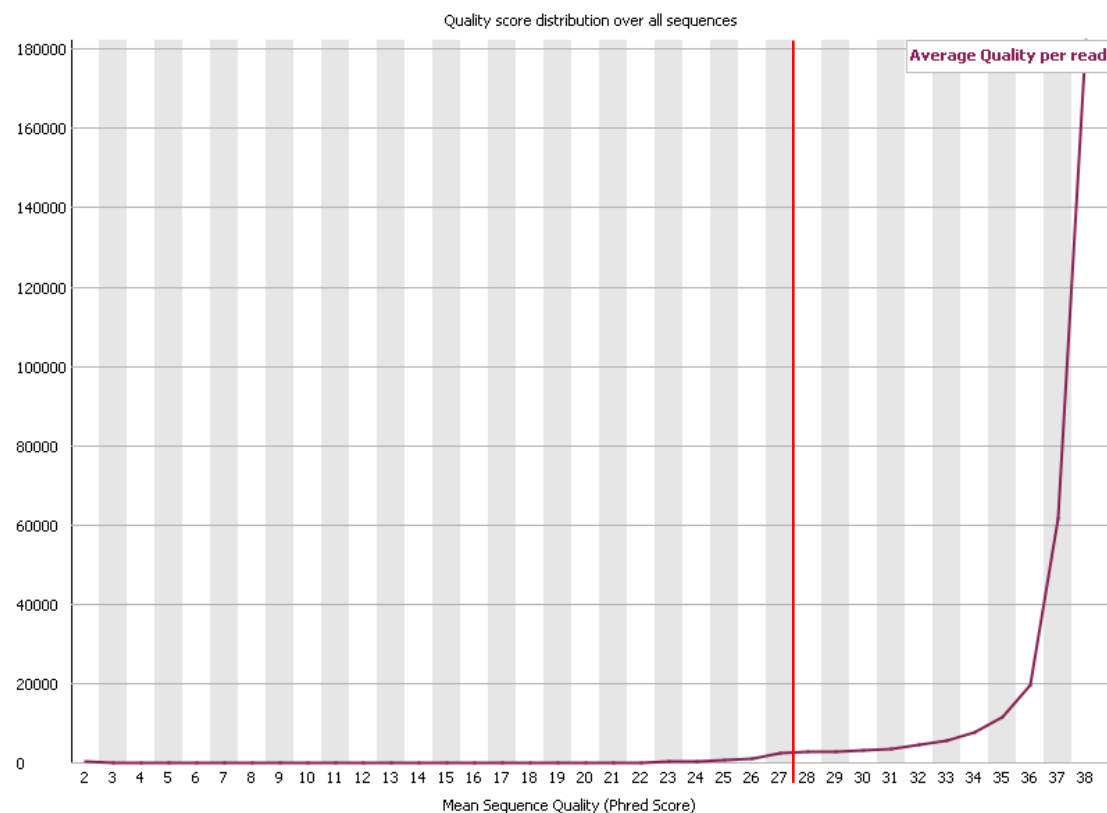
- Il box giallo rappresenta la qualità del 25-75% delle reads in quella posizione, i whiskers sotto il box giallo rappresentano il limite del 10% e 90%. Tendenzialmente la qualità tende ad abbassarsi per corse molto lunghe.

✗ **Per tile sequence quality**



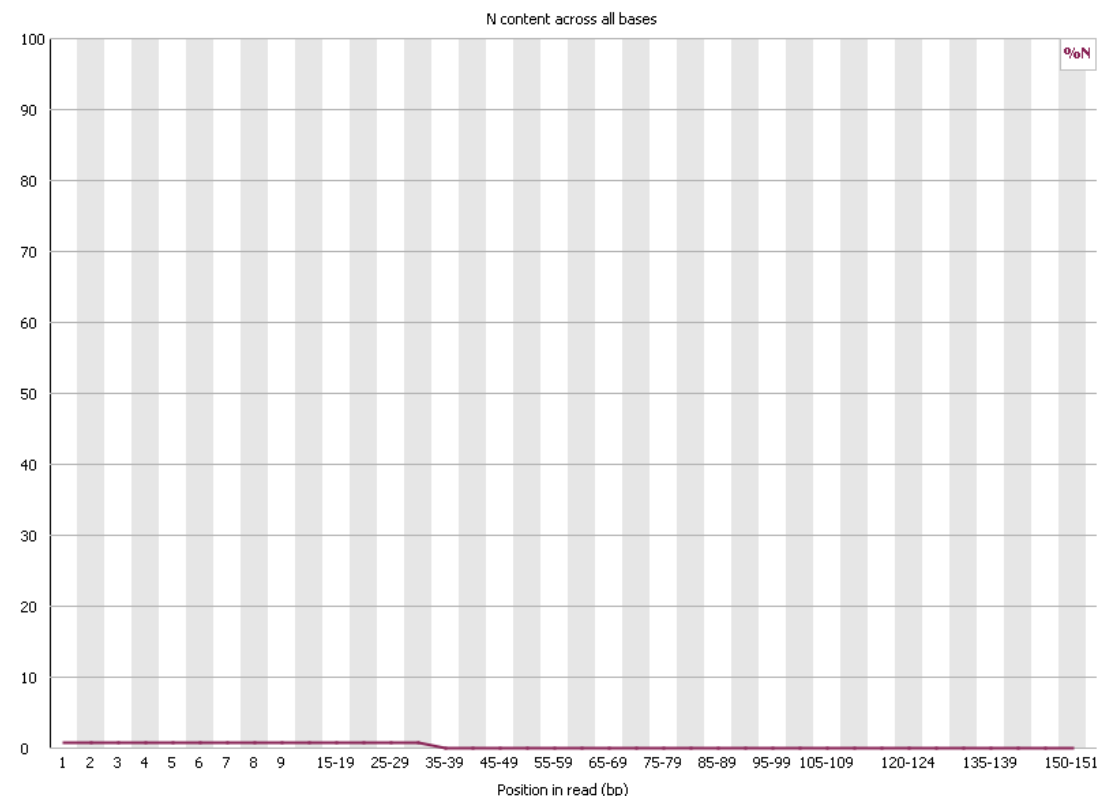
- Nelle librerie di Illumina sono presenti informazioni relative all'identificativo della sequenza, per cui è possibile ricostruire la qualità associata ai *tile* per identificare eventuali problemi in una regione della flowcell. Le sfumature di blu indicano buona qualità.

✓ Per sequence quality scores



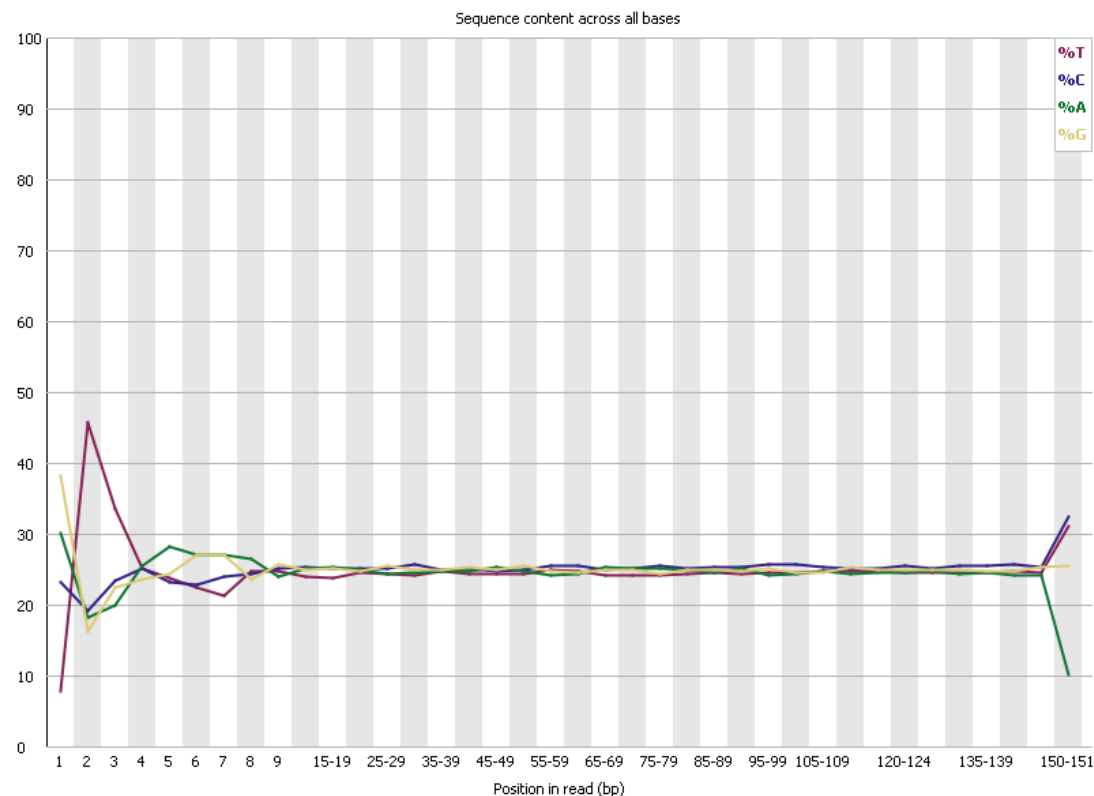
- Parametri di qualità medi per read, valori sotto 20 indicano probabilità di errore dell' 1% nell'assegnazione delle basi, il valore ottimale è sopra 27.

✓ Per base N content

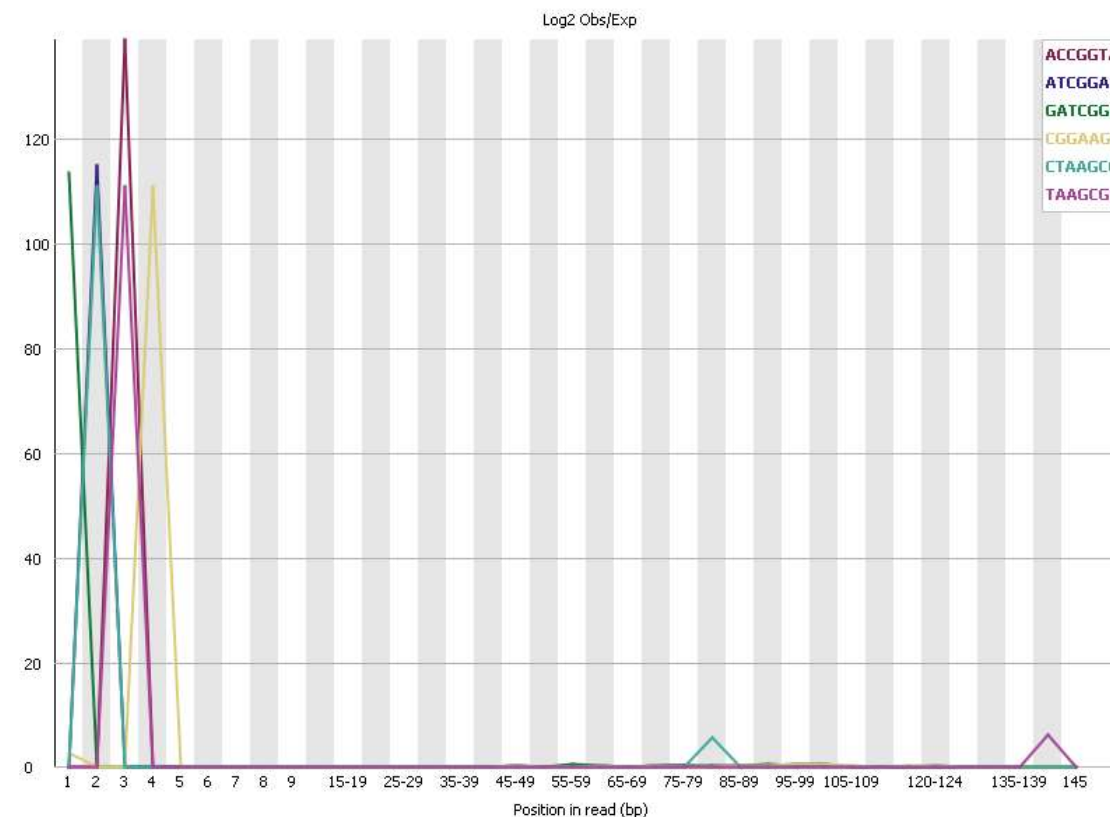


- Il valore N viene assegnato in una posizione quando il sequenziamento non è in grado di assegnare la base. Questo grafico mostra la % di basi chiamate in ogni posizione, per le quali è stata assegnata una N. Un aumento di %N può indicare che un bias della reazione rende il *basecaller* instabile.

✖ Per base sequence content



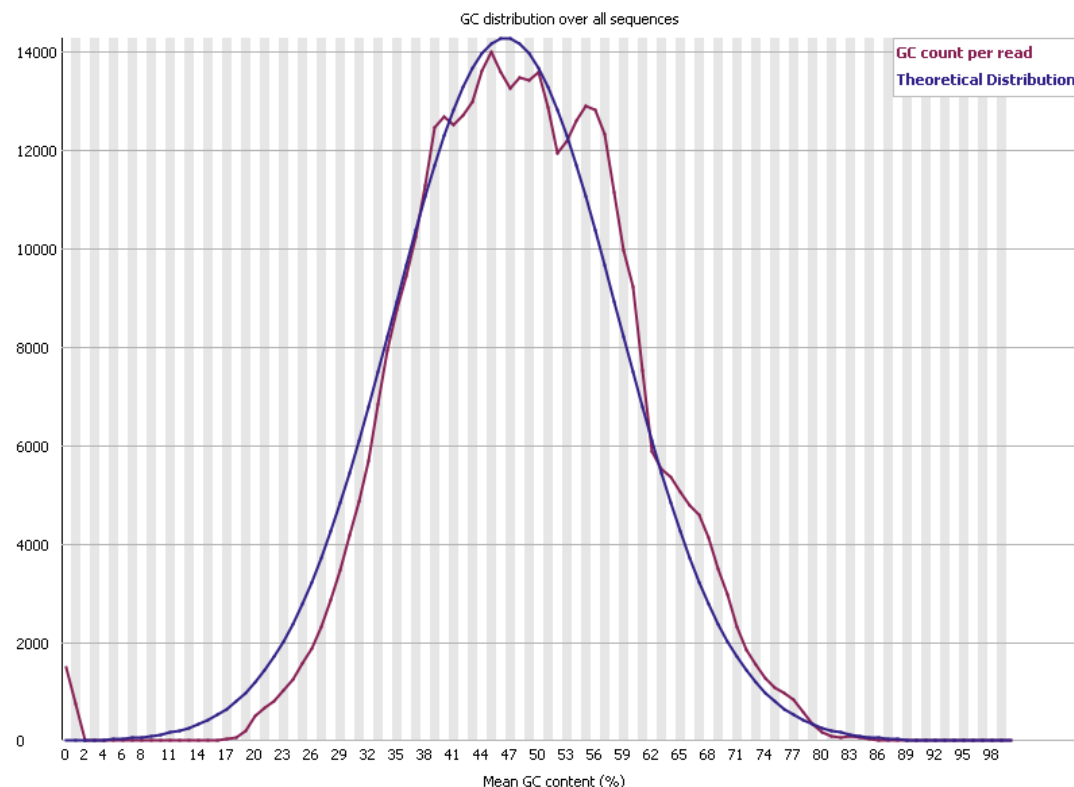
✖ Kmer Content



- Questo parametro mostra la proporzione della presenza di ogni base nel sequenziamento. In una libreria ottimale le proporzioni dovrebbero mantenersi costanti. All'inizio della sequenza si può verificare sbilanciamento per bias nella frammentazione o *random priming*; in questo caso non interferisce con la qualità del dato a valle.

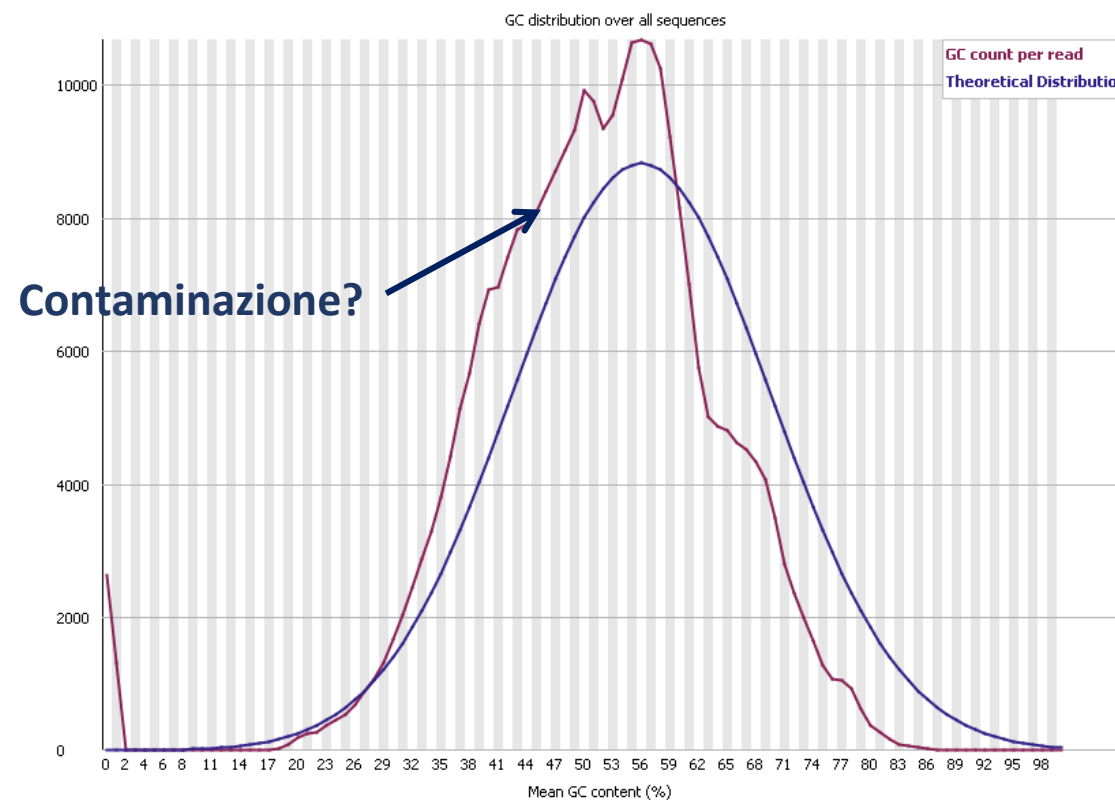
- Il Kmer content parte dal presupposto che piccoli segmenti di sequenza (*Kmeri*) dovrebbero essere equamente rappresentati all'interno di una *library*. All'inizio della sequenza questo dato può aumentare in caso di *random priming*, ed è correlato al "Per base sequence content".

✓ **Per sequence GC content**



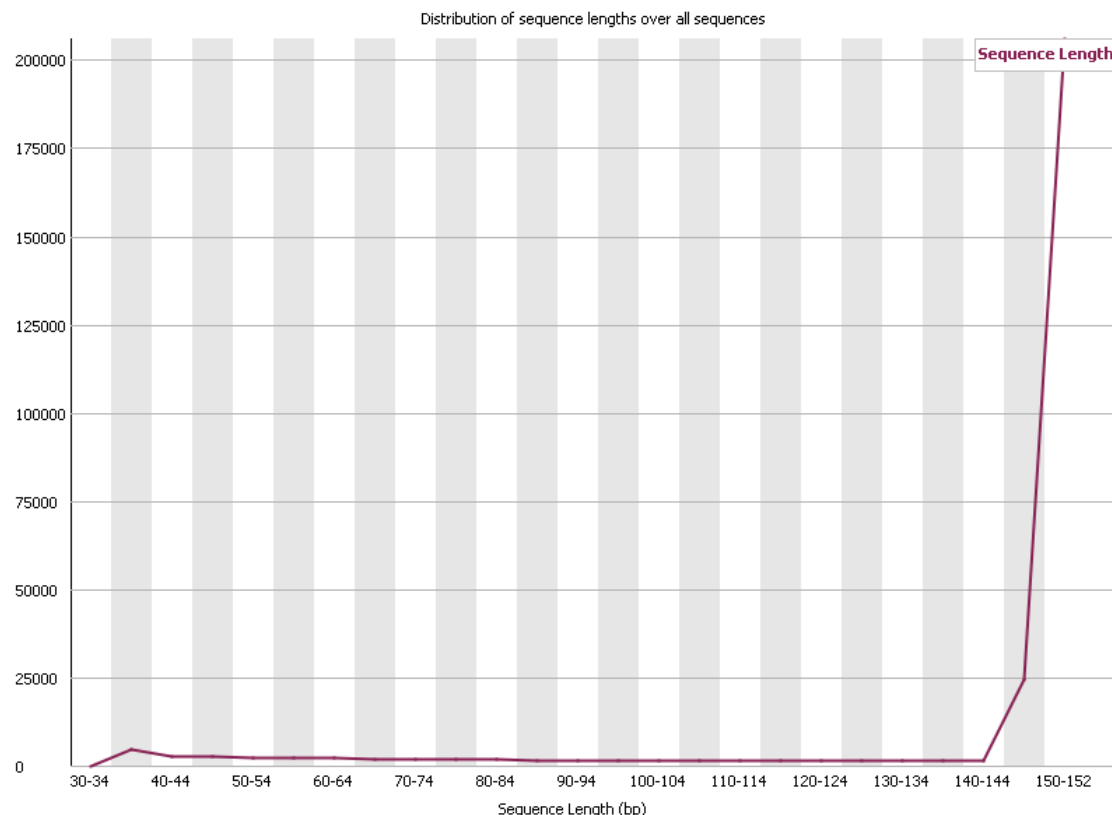
- Questa sezione mostra il *GC content* sull'intera distribuzione delle sequenze, viene comparata ad una distribuzione normale teorica del *GC content*. La curva relativa al sequenziamento deve essere il più possibile sovrapponibile alla distribuzione normale.

✗ **Per sequence GC content**



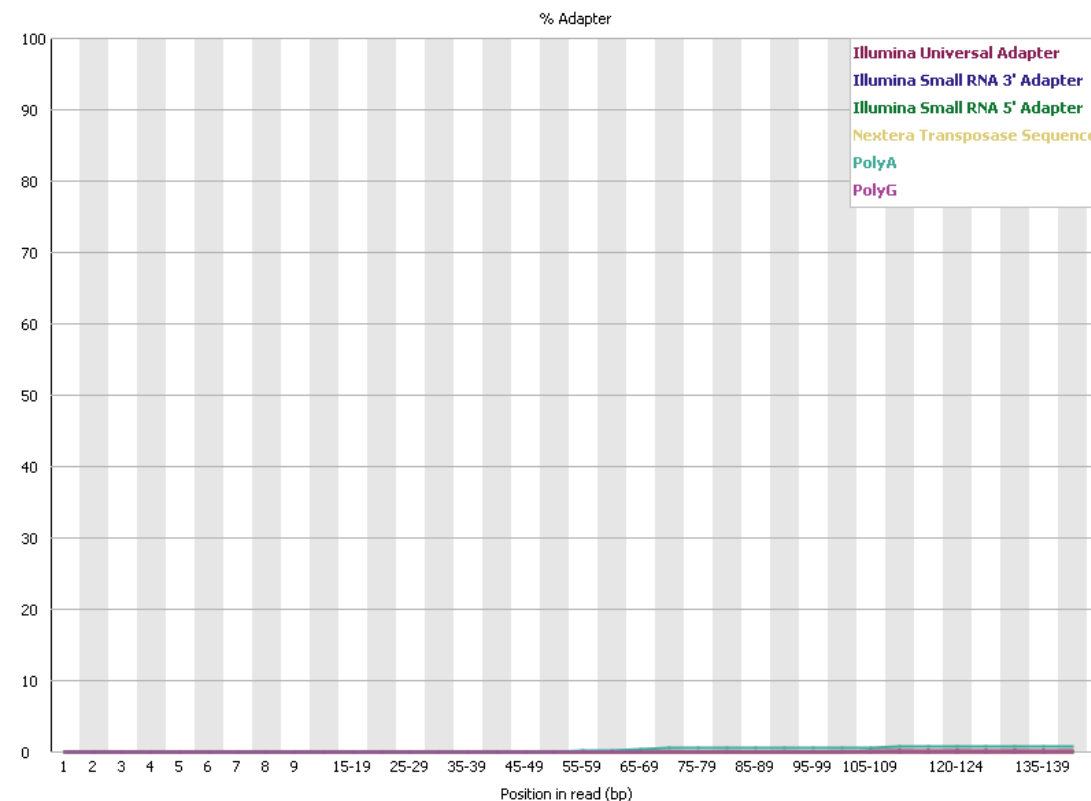
-Uno *shift* nella curva osservata indica un bias nella rappresentazione delle basi, che può essere dato da contaminazione o da problemi della *library*. Uno *shift* sistematico non viene indicato come errore, siccome può essere una caratteristica intrinseca della *library*.

! Sequence Length Distribution



- Le sequenze dovrebbero avere tutte la stessa lunghezza, che dipende dalla metodica utilizzata.

✓ Adapter Content



- Questo grafico mostra l'eventuale sovra rappresentazione di adapter nel sequenziamento. Nel software è presente una lista di adapter di default che può essere modificata manualmente dall'operatore nella sezione Configuration > adapter_list.txt

-Nella maggioranza dei sequenziamenti HLA da noi analizzati, sono presenti sequenze over rappresentate. Il software riporta la percentuale delle sequenze più abbondanti sopra lo 0,1% del totale. Nel software è presente una lista di contaminati di default identificabili, che può essere modificata manualmente dall'operatore nella cartella Configuration > contaminat list.txt

✖ Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GATCGGAAGAGCACACGTCTGAACTCCAGTCACGTCTATGAATCTCGTAT	5189	1.2673902848880378	TruSeq Adapter, Index 11 (97% over 38bp)
NN	1483	0.36221618664269806	No Hit
CTCATCTTCTGCCTCTCCATTTCTTGCTTGTATAGTTCAGTCCCTCTA	1148	0.28039391926218293	No Hit
CTCATCTTCTGCCTCTCCATTTCTTACTTTCTATACATCAGCCCCTCTA	466	0.11381843760991052	No Hit

- Nel sequenziamento HLA la sequenza più abbondante (sopra l'1%) fa spesso riferimento agli indici di Illumina
- Sono presenti altre sequenze di origine ignota che si è reso necessario indagare



Ricerca delle sequenze nel database
NCBI BLAST



UCLAHLA1114ANNO2024-241036-513-702_S8_L001_R1_001.fastq.gz UCLAHLA1114ANNO2024-241036-513-702_S8_L001_R2_001.fastq.gz

		Overrepresented sequences			
		Sequence	Count	Percentage	Possible Sou...
✔	Basic Statistics				
✔	Per base sequence quality	GATCGGAAGAGCACACGTCTGAACTCCAGTCACGACATAGTATCTCGTAT	7459	2,679	TruSeq Adapt...
		NN	2629	0,944	No Hit
✔	Per tile sequence quality	GGCAGACAGTGTGACAAAGAGGCTGGTGTAGGAGAAGAGGGATCAGGACG	677	0,243	No Hit
		TGGATACTCACGACGCGGACCCAGTTCCTCACTCCCATTTGGGTGTGGGTT	580	0,208	No Hit
✔	Per sequence quality scores	CTCATCTTCTGCCTCTCCATTTCTTGTCTGCTATAGTTCAAGTCCCTCTA	426	0,153	No Hit
		CTCATCTTCTGCCTCTCCATTTCTTGTCTGCTATAGTTCAAGTCCCTCTA	406	0,146	No Hit
✘	Per base sequence content	GCAGGTGCCTTTGCAGAAACAAAGTCAGGGTCTTCAAGTCACAAAGGGA	400	0,144	No Hit
		CTGGGGAGGAAACACAGGTTCAGCATGGGAACAGGGGTCACAGTGGACACG	294	0,106	No Hit
✘	Per sequence GC content				
✔	Per base N content				
!	Sequence Length Distribution				
✘	Sequence Duplication Levels				
✘	Overrepresented sequences				
✔	Adapter Content				

[Link To This Page | Feedback](#)

E' emerso come nella metodica in uso siano le regioni 5' e 3' UTR della prima classe e l'introne 1 dei loci DRB3/4/5 a creare contaminazione. E' possibile utilizzare questo strumento nella determinazione di un eventuale mix-up.

Conclusioni

- La metodica descritta si propone come supporto nella valutazione della qualità dei file FASTQ di sequenza, parametro che riveste particolare importanza nella routine laboratoristica.
- L'analisi in oggetto permette inoltre di caratterizzare i propri dati di sequenziamento in fase di validazione e affiancare i software commerciali nell'interpretazione del dato nella pratica lavorativa.



Grazie per l'attenzione

