



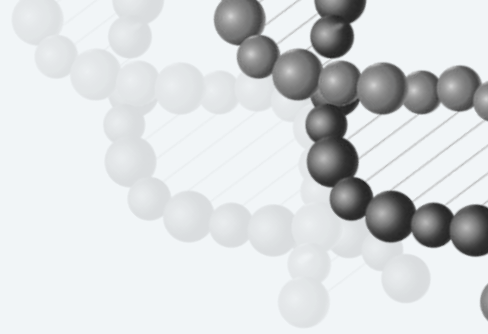
# BIOINFORMATICA

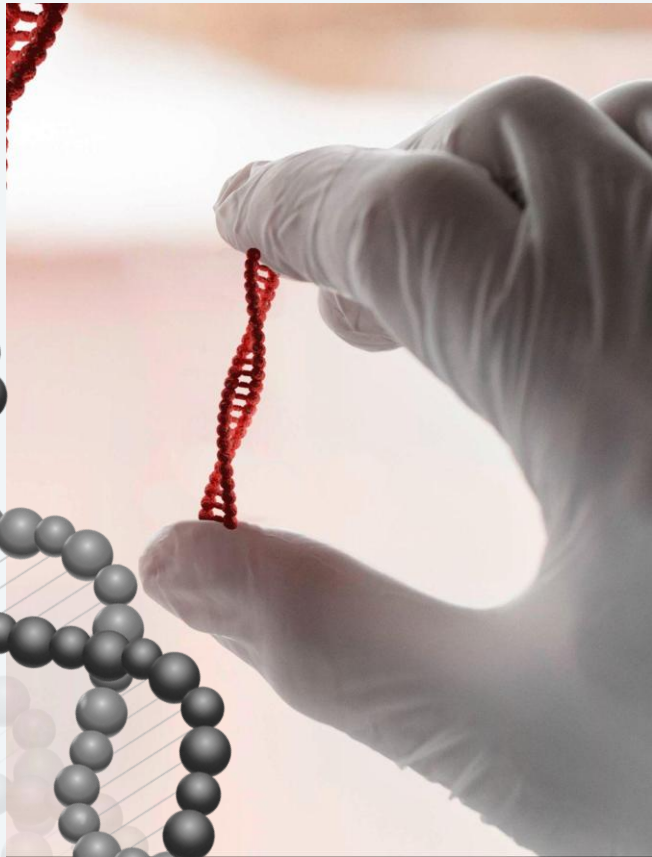
## II

### *TRIMMING*

# ARGOMENTI

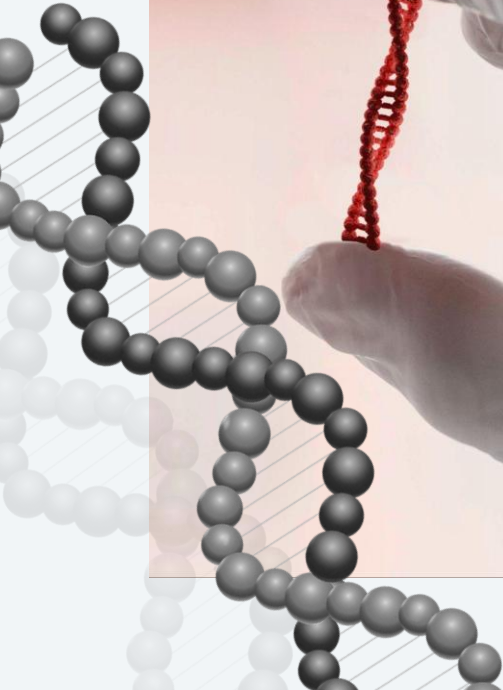
## 01 TRIMMING





01

**TRIMMING**



# TRIMMING: Definizione


Il *trimming* è il processo di rimozione delle sequenze di bassa qualità o delle porzioni di una read che non sono rilevanti per l'analisi, come gli adattatori o le basi di qualità scarsa alle estremità.

Il trimming migliora la qualità dei dati che sono usati per analisi successive, come il mapping al genoma di riferimento o l'assemblaggio *de novo*.

# TRIMMING: sequenze da eliminare



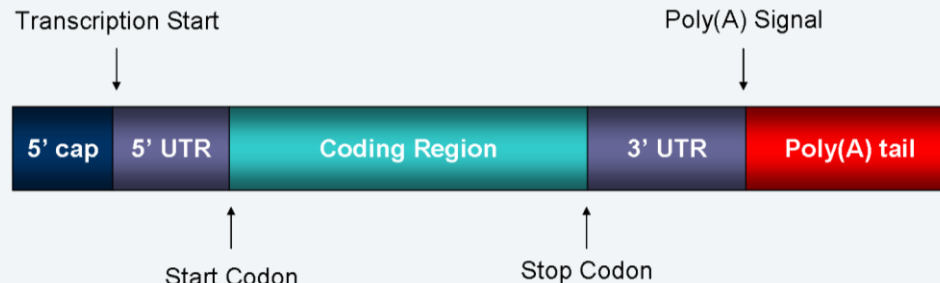
Le reads grezze prodotte dal sequenziamento possono contenere **sequenze indesiderate**, come:

- **Adattatori:** sequenze artificiali usate durante la preparazione della libreria.
  - **Basi a bassa qualità:** nucleotidi con alta probabilità di errore di chiamata.
  - **Poli-A o Poli-T tails:** le code omopolimeriche, che possono confondere l'analisi successiva.
- 

# TRIMMING: approfondimento PolyA/PolyT

Le poli-A e le poli-T sono sequenze di basi nucleotidiche specifiche che si trovano nel DNA o nell'RNA :

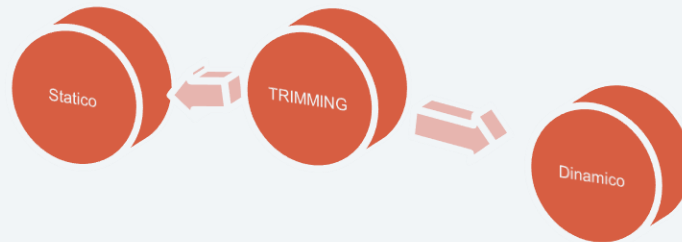
- **Poli-A (Poly-A):**  
Una sequenza costituita da molte basi adenina (A) consecutive. In RNA messaggero (mRNA) degli eucarioti, una coda poli-A viene aggiunta all'estremità 3' durante la maturazione dell'mRNA.
- **Poli-T (Poly-T):**  
Una sequenza costituita da molte basi timina (T) consecutive. Nel DNA, le sequenze di poli-T possono apparire come complementari alle regioni poli-A su un filamento opposto. Sono spesso usate in laboratorio, ad esempio: Nelle colonne di purificazione di mRNA, dove oligonucleotidi di poli-T legano la coda poli-A dell'mRNA per separarlo dagli altri tipi di RNA.



# TRIMMING: Tipologia di elaborazione

Rimuovere queste sequenze migliora l'accuratezza delle analisi bioinformatiche, riducendo gli errori durante il mapping e l'interpretazione dei dati.


Il trimming può essere statico o dinamico



# TRIMMING: Statico

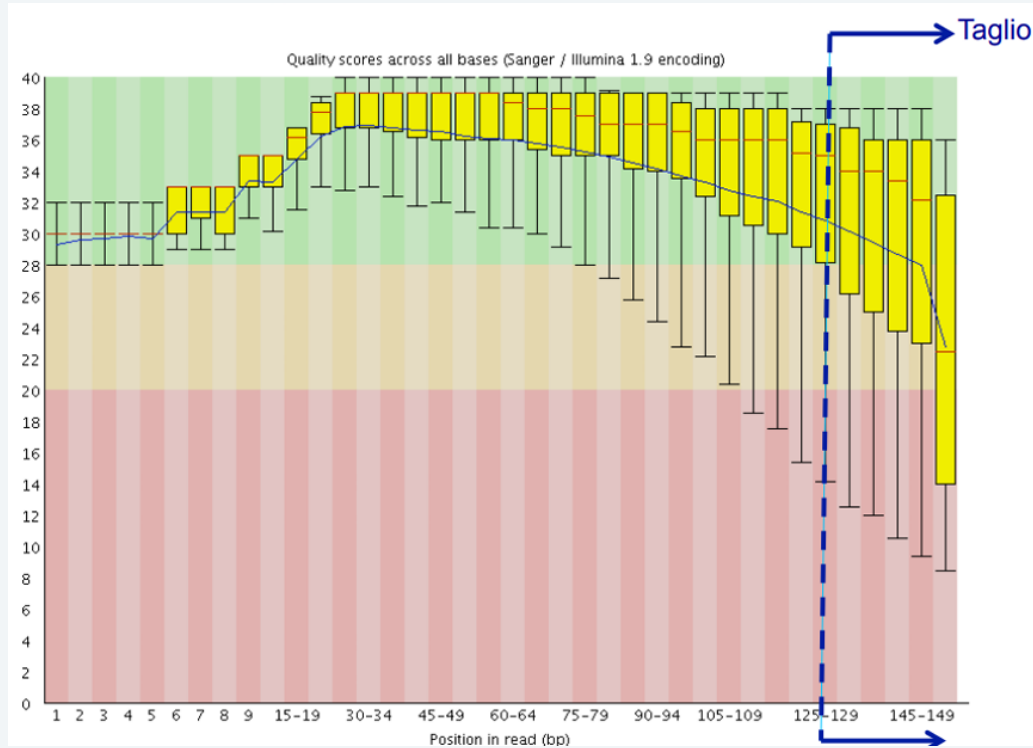


## Trimming statico:

- Nel trimming statico, le reads sono tagliate seguendo regole predefinite. Questo significa che sono rimosse sempre le stesse posizioni della sequenza, indipendentemente dalla qualità specifica di quella porzione.  
*Ad esempio, si può decidere di rimuovere sempre i primi 10 nucleotidi e gli ultimi 5, oppure eliminare una porzione fissa identificata come adattatore.*
  - È poco flessibile perché non considera la variabilità della qualità della sequenza e potrebbe rimuovere informazioni utili o lasciare parti di bassa qualità se queste non si trovano nelle posizioni preimpostate.
- 




# TRIMMING: statico (esempio)



# TRIMMING: Dinamico



## Trimming dinamico:

- Il trimming dinamico adatta il taglio delle reads in relazione alla qualità effettiva di ciascuna lettura. Generalmente, un algoritmo valuta ogni posizione della sequenza e taglia la read quando rileva che la qualità scende sotto una certa soglia. *Ad esempio, può eliminare le basi dalla fine di una read fino a quando la qualità non raggiunge un livello accettabile, oppure rimuovere adattatori rilevati sulla base di algoritmi di allineamento specifici.*
  - È più preciso e adattabile, perché analizza ogni sequenza individualmente, garantendo che solo le porzioni effettivamente di bassa qualità vengano rimosse, senza perdere informazioni utili.
- 

# TRIMMING: Dinamico (esempio)



Adattatore

READS  
pessima qualità

READS  
scarsa qualità

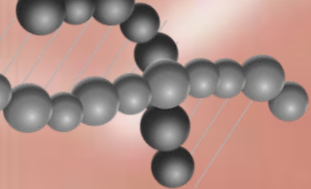
Adattatore



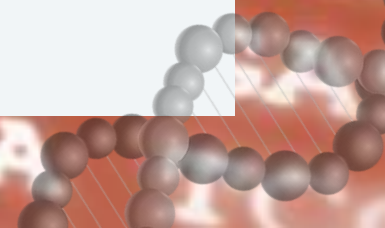
# TRIMMING: confronto

Caratteristica	Trimming Statico	Trimming Dinamico
Flessibilità	Fisso, indipendente dalla qualità	Adattivo, dipende dalla qualità della read
Precisione	Potenzialmente meno accurato	Maggiore precisione
Efficienza computazionale	Più rapido e meno dispendioso	Più complesso e dispendioso
Rischio di perdita dati	Rischio di rimuovere basi utili	Mantiene solo le porzioni necessarie
Uso comune	Rimozione di adattatori specifici	Controllo qualità per rimuovere basi degradate

In sintesi, il trimming statico è utile quando si conoscono esattamente le porzioni da tagliare e si vuole un approccio rapido e standard, mentre il trimming dinamico è preferibile quando si desidera ottenere la massima qualità possibile dalle reads, poiché permette di adattarsi alle variazioni intrinseche di ciascuna sequenza.



# BIOINFORMATICA



# TRIMMING

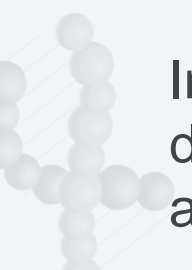


Il "trimming" in bioinformatica si riferisce al processo di rimozione di basi indesiderate o di bassa qualità dalle estremità delle sequenze nei dati di sequenziamento, come quelli nel formato FASTQ.

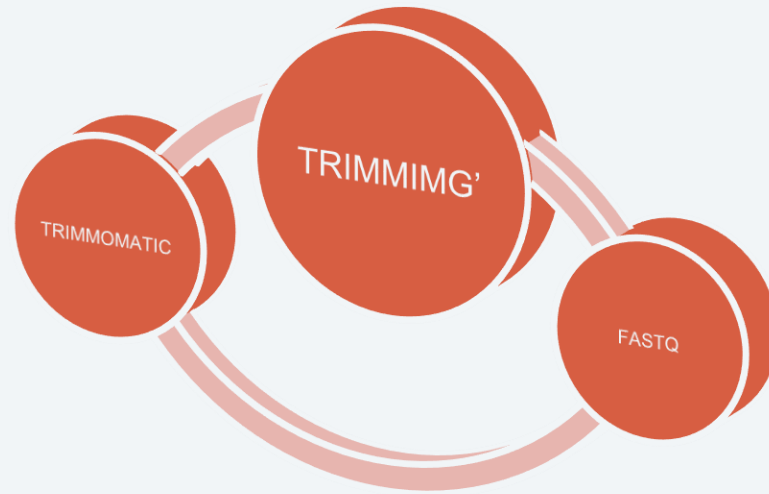
Può essere personalizzato in relazione alle specifiche esigenze dell'analisi.

*Ad esempio, è possibile impostare soglie di qualità per definire quali basi devono essere tagliate in base ai loro punteggi di qualità Phred.*

In generale, il trimming è una fase cruciale nel pre-processing dei dati di sequenziamento per garantire che le sequenze siano pulite, accurate e pronte per le successive analisi di bioinformatica



# TRIMMING



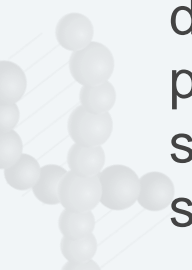
# TRIMMOMATIC



**Trimmomatic** è uno software di bioinformatica utilizzato per il trimming delle sequenze nei dati di sequenziamento.

È progettato per rimuovere basi di bassa qualità, adattatori e sequenze di scarsa qualità dalle estremità delle read nei file FASTQ.

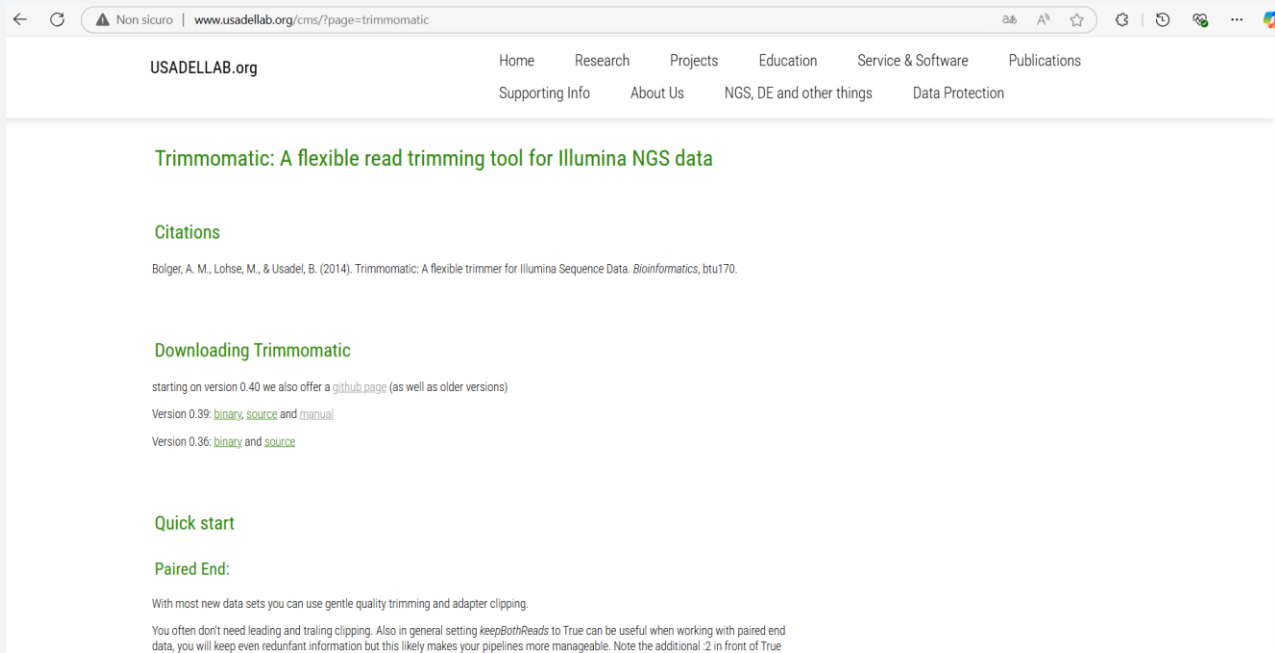
L'uso di Trimmomatic è una fase comune nella pipeline di analisi dei dati di sequenziamento, poiché contribuisce a migliorare la qualità delle sequenze e a ridurre il rumore nei dati. La sua applicazione è particolarmente importante per garantire risultati accurati nelle successive fasi di assemblaggio, allineamento e analisi dei dati di sequenziamento.





# TRIMMOMATIC: Sito

<http://www.usadellab.org/cms/?page=trimmomatic>

A screenshot of a web browser displaying the Trimmomatic website. The browser's address bar shows the URL 'www.usadellab.org/cms/?page=trimmomatic'. The website's navigation menu includes 'Home', 'Research', 'Projects', 'Education', 'Service & Software', 'Publications', 'Supporting Info', 'About Us', 'NGS, DE and other things', and 'Data Protection'. The main content area features a green heading 'Trimmomatic: A flexible read trimming tool for Illumina NGS data', followed by a 'Citations' section with a reference to Bolger et al. (2014), a 'Downloading Trimmomatic' section with links to GitHub pages and binary/source files for versions 0.39 and 0.36, a 'Quick start' section, and a 'Paired End:' section with instructions on quality trimming and adapter clipping.

USADELLAB.org

Home Research Projects Education Service & Software Publications  
Supporting Info About Us NGS, DE and other things Data Protection

## Trimmomatic: A flexible read trimming tool for Illumina NGS data

### Citations

Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible trimmer for Illumina Sequence Data. *Bioinformatics*, *btu170*.

### Downloading Trimmomatic

starting on version 0.40 we also offer a [github page](#) (as well as older versions)

Version 0.39: [binary source](#) and [manual](#)

Version 0.36: [binary](#) and [source](#)

### Quick start

#### Paired End:

With most new data sets you can use gentle quality trimming and adapter clipping.

You often don't need leading and trailing clipping. Also in general setting `keepBothReads` to True can be useful when working with paired end data, you will keep even redundant information but this likely makes your pipelines more manageable. Note the additional `.2` in front of True

# TRIMMOMATIC: Installazione

È sufficiente il download e la decompressione del file

<http://www.usadellab.org/cms/uploads/supplementary/Trimmomatic/Trimmomatic-0.39.zip>

Dipendenze:

Richiede la Java Virtual Machine

# TRIMMOMATIC: Comando (Single-End)

Single End (SE):

```
java -jar trimmomatic-0.39.jar SE -phredtype file.fastq file_trimmed.fastq  
ILLUMINACLIP:adapters.fa:val1:val2:val3 LEADING:val TRAILING:val  
SLIDINGWINDOW:val1:val2 MINLEN:val1
```

# TRIMMOMATIC: Comando (esempio)

```
java -jar trimmomatic-0.39.jar SE -phred33 input.fastq output_trimmed.fastq  
ILLUMINACLIP:adapters.fa:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15  
MINLEN:36
```

OPZIONE	Descrizione
SE	modalità single-end (per dati derivanti da un solo file FASTQ)
-phred33	specifica il formato della qualità (es. -phred33 o -phred64)
input.fastq	file di input delle sequenze non elaborate
output_trimmed.fastq	file di output delle sequenze pulite
ILLUMINACLIP adapters.fa:2:30:10	adapters.fa: file contenente le sequenze degli adattatori. 2: numero massimo di mismatch tra adattatore e sequenza. 30: soglia di punteggio per decidere il taglio. 10: lunghezza minima della sequenza di adattatore da considerare.
LEADING:3	rimuove le basi con qualità inferiore a 3 dall'inizio
TRAILING:3	rimuove le basi con qualità inferiore a 3 dalla fine
SLIDINGWINDOW:4:15	analizza le finestre mobili di 4 basi e rimuove se la qualità media è inferiore a 15
MINLEN:36	rimuove sequenze più corte di 36 basi

# TRIMMOMATIC: Comando (Paired-End)

Paired End (SE):

```
java -jar trimmomatic-0.39.jar PE -phredType input_forward.fastq  
input_reverse.fastq  
output_forward_paired.fastq output_forward_unpaired.fastq  
output_reverse_paired.fastq output_reverse_unpaired.fastq  
ILLUMINACLIP:adapters.fa:val1:val2:val3  
LEADING:val1  
TRAILING:val1  
SLIDINGWINDOW:val1:val2  
MINLEN:val1
```

# TRIMMOMATIC: Comando (esempio)

```
java -jar trimmomatic-0.39.jar PE -phred33 input_forward.fastq input_reverse.fastq
output_forward_paired.fastq output_forward_unpaired.fastq
output_reverse_paired.fastq output_reverse_unpaired.fastq
ILLUMINACLIP:adapters.fa:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36
```

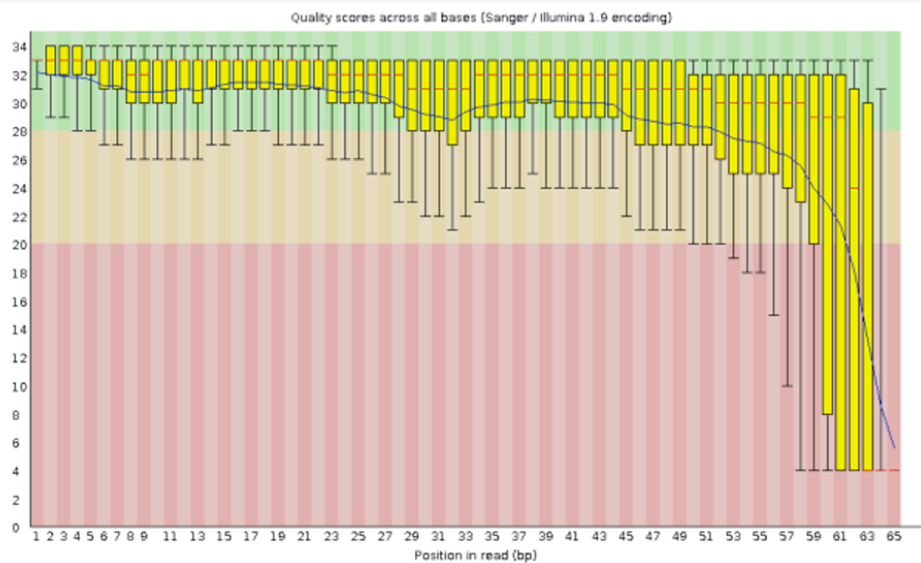
Parametro	Descrizione
SE	modalità paired-end (dati provenienti da due file FASTQ, uno per ciascun lato della lettura)
-phred33	specifica il formato della qualità (es. -phred33 o -phred64)
input_forward.fastq input_reverse.fastq	file di input per le letture forward e reverse
output_forward_paired.fastq, output_reverse_paired.fastq	file di output per le letture forward e reverse che rimangono appaiate dopo il trimming
output_forward_unpaired.fastq, output_reverse_unpaired.fastq	file di output per le letture forward e reverse che sono state scartate perché non appaiate
ILLUMINACLIP adapters.fa:2:30:10	adapters.fa: file contenente le sequenze degli adattatori. 2: numero massimo di mismatch tra adattatore e sequenza. 30: soglia di punteggio per decidere il taglio. 10: lunghezza minima della sequenza di adattatore da considerare.
LEADING:3	rimuove le basi con qualità inferiore a 3 dall'inizio
TRAILING:3	rimuove le basi con qualità inferiore a 3 dalla fine
SLIDINGWINDOW:4:15	analizza le finestre mobili di 4 basi e rimuove se la qualità media è inferiore a 15
MINLEN:36	rimuove sequenze più corte di 36 basi

# TRIMMOMATIC: Altri parametri

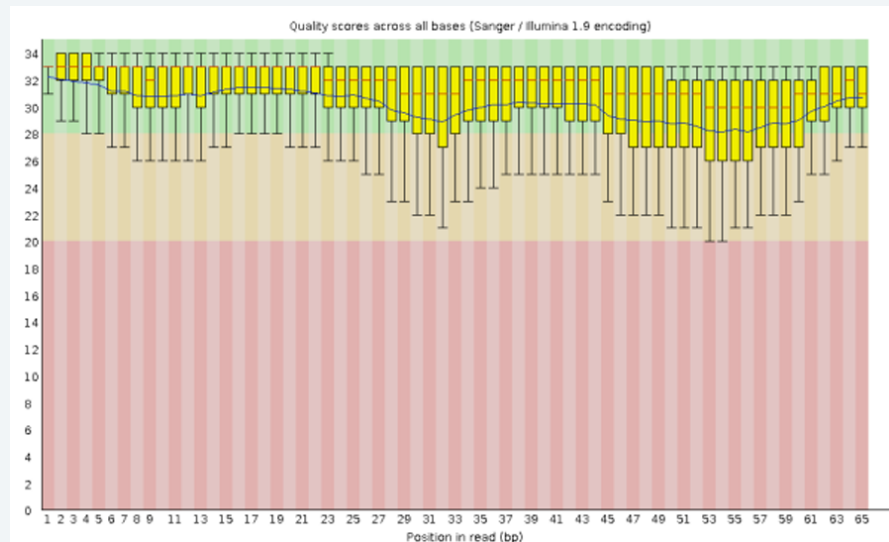
Parametro	Significato
HEADCROP	Rimuove un numero specificato di basi dall'inizio di ogni lettura
CROP	Taglia la lettura a una lunghezza specificata rimuovendo le basi in eccesso dalla fine
TOPHRED33/TOPHRED64	Converte la qualità dei punteggi dal formato Phred+64 al formato Phred+33 e viceversa
MAXINFO	Opzione di trimming adattivo che bilancia la lunghezza della lettura e la qualità per determinare la migliore lunghezza di trimming.
AVGQUAL	Rimuove la lettura se la qualità media è al di sotto di una soglia specificata

# TRIMMOMATIC

## PRIMA DEL TRIMMING



## DOPO IL TRIMMING






# FASTQ



**FASTQ** è un software di bioinformatica ampiamente utilizzato per il trimming delle sequenze nei dati di sequenziamento


Opera effettuando una profilazione completa della qualità sia prima sia dopo il filtraggio dei dati (curve di qualità, contenuti di base, KMER, Q20/Q30, GC Ratio, duplicazione, contenuti dell'adattatore...)



# FASTQ



## FASTQ opera nella seguente maniera:

- filtra le reads errate: qualità troppo bassa, troppo breve o troppi N...
  - tagliare basi di bassa qualità valutando la qualità media considerando una finestra scorrevole (come Trimmomatic).
  - Taglia tutte le reads davanti e in coda e taglia gli adattatori.
  - Corregge le coppie di basi non corrispondenti nelle regioni sovrapposte delle letture finali accoppiate, se una base è di alta qualità mentre l'altra è di qualità ultra bassa
  - Toglie le trim polyG in 3' ends, (spesso nei dati ricavati da NovaSeq/NextSeq)
  - Taglio le polyX in estremità 3'
  - riporta il risultato in formato JSON per svolgere analisi di approfondimento
  - visualizza i risultati del controllo qualità e del filtraggio su una singola pagina HTML (come FASTQC).
  - Consente l'elaborazione parallela.
- 

# FASTQ: Sito

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

Prodotto ▾ Soluzioni ▾ Risorse ▾ Open Source ▾ Azienda ▾ Prezzi

Cerca o vai a...

Accedi Registrati

OpenGene / veloce Pubblico

Notifiche Forchetta 334 Stella 1.9mila

<> Codice Problemi 341 Richieste pull 14 Azioni Progetti Sicurezza Intuizioni

padrone numero arabo Rami 42 Tag Vai al file Codice

sfchen correggi il bug di fixMGI 4F273F1 - last week 459 impegni

.github/flussi di lavoro	Correggere un errore	last year
Src	correggi il bug di fixMGI	last week
dati di prova	Supporta l'output di letture non riuscite	5 years ago
.gitignore	ignora la configurazione vscode; Aggiungi output di test	3 years ago
LICENZA	Cambia data	4 years ago
Makefile	Aggiungi l'opzione -pthread al linker per risolvere il problem...	3 years ago
README.md	Aggiorna README.md	last month

LEGGIMI Licenza MIT

Circa

Un preprocessore FASTQ all-in-one ultraveloce (QC/adattatori/rifilatura/filtraggio/divisione/fusione...)

adattatore qualità Bioinformatica per controllo qualità filtro NGS sequenziamento sovrapposizione emicranico duplicazione umi taglio Preelaborazione filtraggio illumina Veloce Q fusione Qc polig

Leggimi Licenza MIT Attività Proprietà personalizzate 1.9k stelle

# FASTQ: Installazione

È sufficiente il download del file eseguibile:

```
wget http://opengene.org/fastp/fastp  
chmod a+x ./fastp
```

Se si vuole una versione specifica:

```
wget http://opengene.org/fastp/fastp.0.23.4  
mv fastp.0.23.4 fastp  
chmod a+x ./fastp
```

Richiede:

*libdeflate* e *libisal* presenti nel compilatore GCC

# FASTQ: Comando

Per le reads Single-end

```
fastp -i in.fq -o out.fq
```

Per le reads Paired-end

```
fastp -i in.R1.fq -I in.R2.fq -o out.R1.fq -O out.R2.fq
```

Opzione	Significato
-i	File FASTQ da trimmare
-o	Nome di uscita del file trimmato

# FASTQ: Esempi di report

## fastp report

### Summary

#### General

fastp version:	0.17.0 ( <a href="https://github.com/OpenGene/fastp">https://github.com/OpenGene/fastp</a> )
sequencing:	paired end (151 cycles + 151 cycles)
mean length before filtering:	108bp, 108bp
mean length after filtering:	107bp, 107bp
duplication rate:	30.641418%
Insert size peak:	95

#### Before filtering

total reads:	16.763944 M
total bases:	1.818801 G
Q20 bases:	1.716550 G (94.378124%)
Q30 bases:	1.672955 G (91.981195%)
GC content:	47.006320%

#### After filtering

total reads:	16.034314 M
total bases:	1.722358 G
Q20 bases:	1.659759 G (96.365462%)
Q30 bases:	1.622287 G (94.189832%)
GC content:	46.794079%

#### Filtering result

reads passed filters:	16.034314 M (95.647623%)
reads with low quality:	695.022000 K (4.145934%)
reads with too many N:	14.740000 K (0.087927%)
reads too short:	19.868000 K (0.118516%)

# FASTQ: Esempi di report

## Adapters

### Adapter or bad ligation of read1

The input has little adapter percentage (~0.039878%), probably it's trimmed before.

Sequence	Occurrences
N	8579
AN	7450
AGN	6492
AGAN	4593
AGNN	374
other adapter sequences	5774

### Adapter or bad ligation of read2

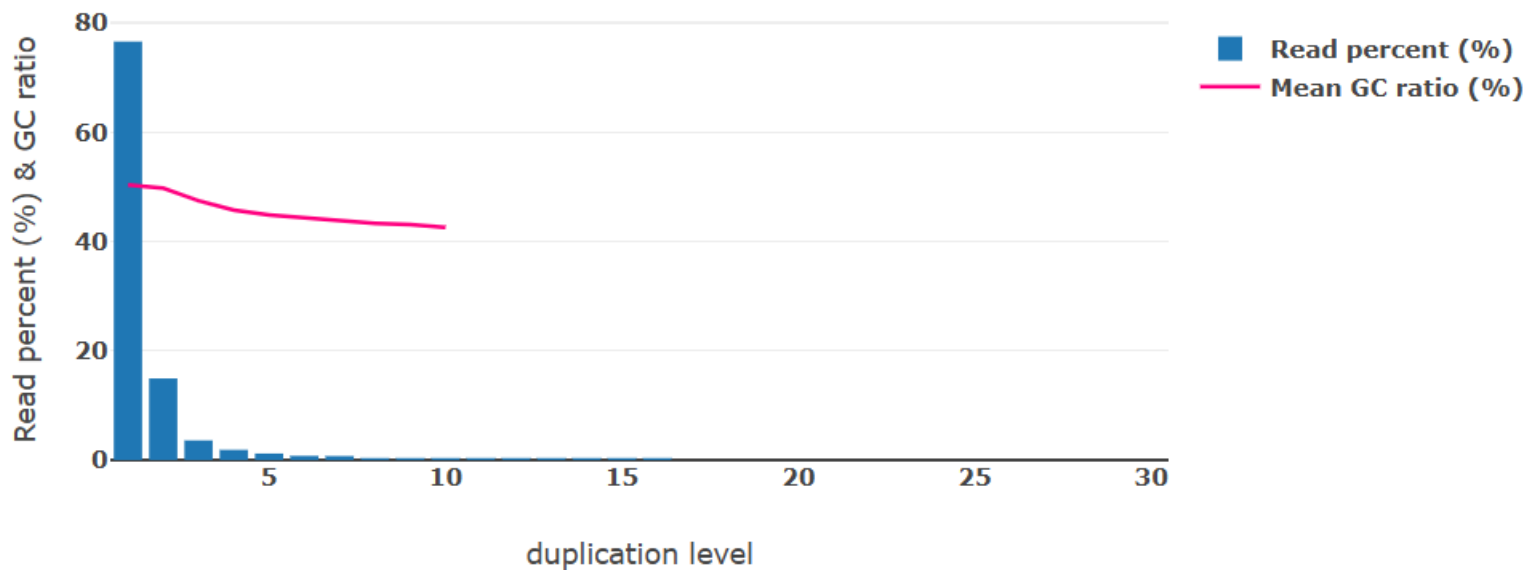
The input has little adapter percentage (~0.226118%), probably it's trimmed before.

Sequence	Occurrences
N	8529
AN	7402
AGA	1480
AGN	6429
AGAN	4663
other adapter sequences	45383

# FASTQ: Esempi di report

## Duplication

duplication rate (30.641418%)



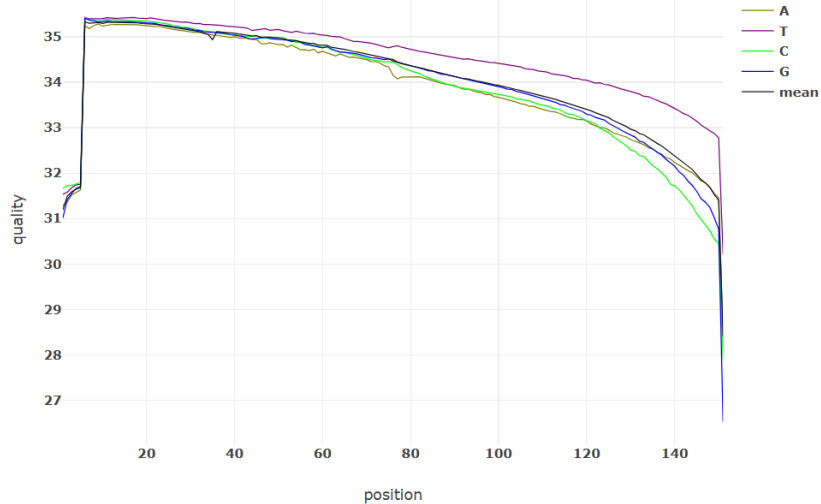


# FASTQ: Esempi di report

## Before filtering

Before filtering: read1: quality

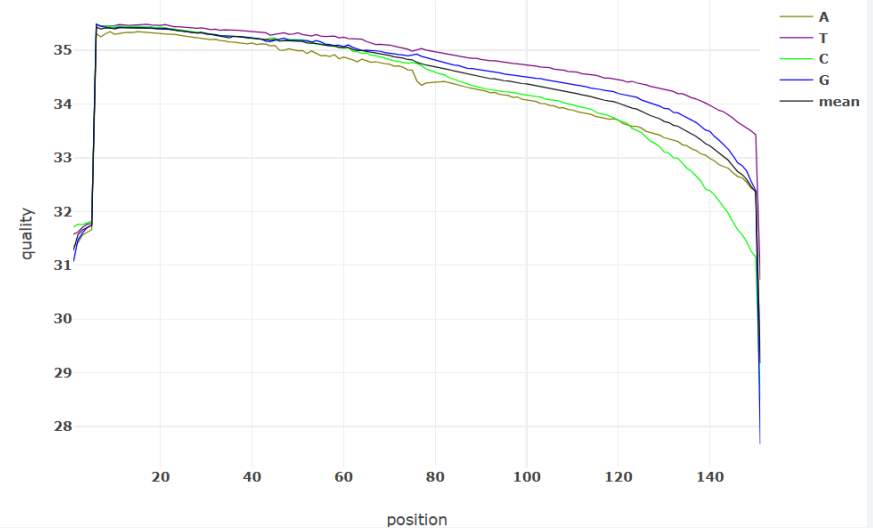
Value of each position will be shown on mouse over.



## After filtering

After filtering: read1: quality

Value of each position will be shown on mouse over.





# Grazie!

## Domande?

franco.liberati@unitus.it

deb.scienceontheweb.com

