# Métodos de secuenciación masiva

Biocomputación Grado en Bioquímica

## Contenido

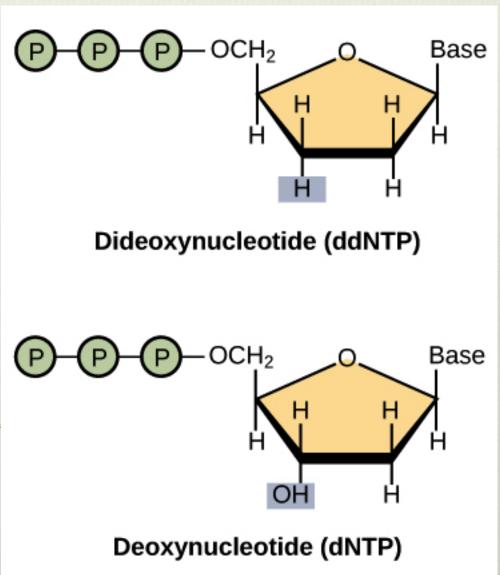
- · Principios básicos del método de Sanger
- Secuenciación masiva: sequencing by synthesis (Illumina)
- · Aplicaciones de la secuenciación masiva
- · La salida: el fichero fastq
- Control de calidad:
  - Detección del adaptador
  - Recortar bases con baja calidad

## Método de Sanger

- Este método se basa en que cuando la polimerasa incorpora un didesoxiribonuleótido se para la síntesis del ADN
- No se pueden forma enlaces fosfodiéster por la ausencia del grupo hidroxilo

https://www.khanacademy.org/science/biology/biotech-dna-technology/dna-sequencing-pcr-electrophoresis/a/dna-sequencing

Image credit: "Whole-genome sequencing: Figure 1," by OpenStax College, Biology (CC BY 4.0).



## Método de Sanger

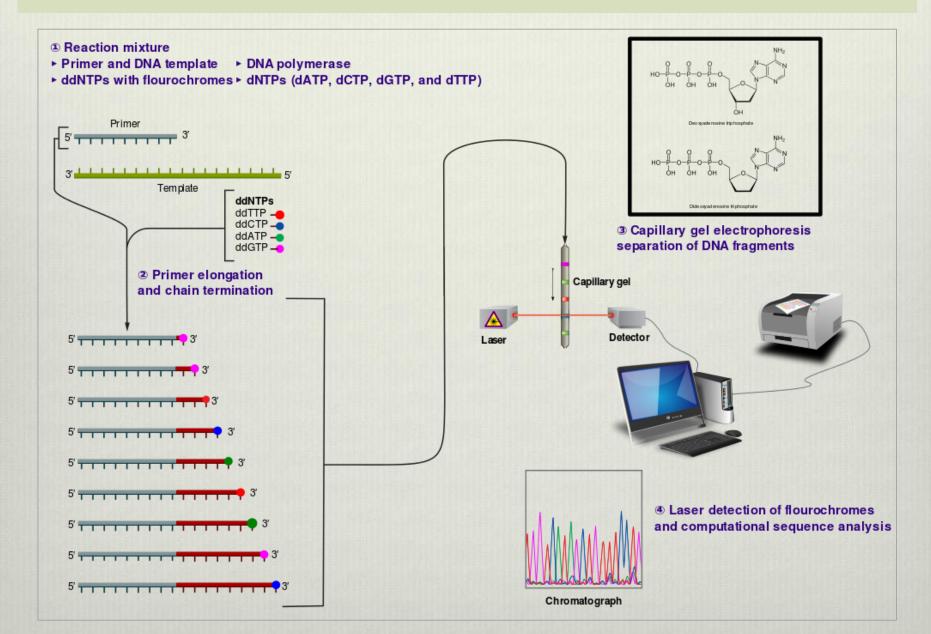
- Molde de DNA desnaturalizado (monocatenario)
- Polimerasa de DNA
- Cebador
- Nucleótidos
- Didesoxinucleótidos
- (didesoxinucleótido trifosfato)

#### Definición:

Con secuenciación masiva nos referimos a los métodos que no se basan en el método de Sanger

TGCAGGCATCAG a. **Denatured** Labelled **Template Primer** CGTCCGTAGT Add dNTPs and **Polymerase** Template/Product ddÇĢŢÇÇĢŢĄĢ **Denaturing Gel** Labelled Strands GATC b. ddA C G T C C G T A G T ddC C ddC ddC ddG

## Método de Sanger



## Comparación Sanger vs. Illumina

Podemos distinguir 3 pasos

#### Preparación de la muestra:

- Fragmentar el ADN
- Ligar los adaptadores
- Seleccionar la longitud

¡Este paso es distinto para los diferentes protocolos!

Generación de clusters y amplificación

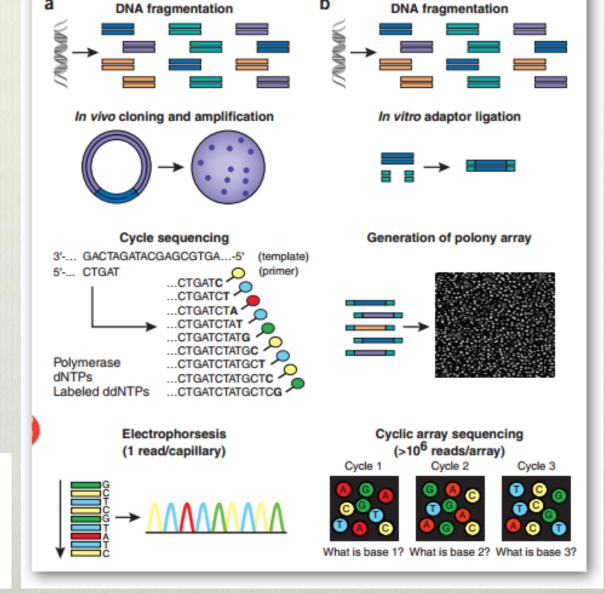
Secuenciación (sequencing by synthesis)

Review

Next-generation DNA sequencing

Jay Shendure 🏲 & Hanlee Ji 🏲

Published online: 09 October 2008



Nature Biotechnology **26**, 1135–1145 (2008) doi:10.1038/nbt1486

Download Citation

## **Clusters**

- El material de entrada es ADN de doble cadena ligado a unos adaptadores
- Los adaptadores son complementarios a oligos fijados en la superficie del flowcell
- Se desnaturaliza el ADN y este se une al *flowcell* formando puentes mediante la hibridación con los oligos

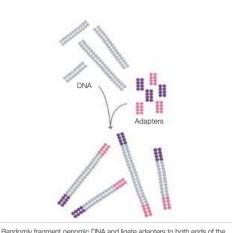
#### Figuras de Illumina:

https://www.illumina.com/documents/products/techspotlights/techspotlight sequencing.pdf

https://www.youtube.com/watc
h?v=womKfikWlxM

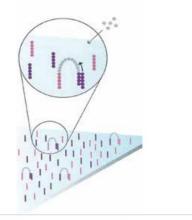
Technology Spotlight: Illumina® Sequencing

Figure 2: Prepare Genomic DNA Sample



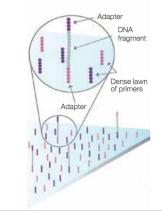
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

Figure 4: Bridge Amplification



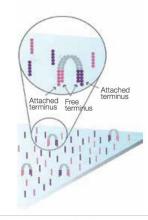
Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

Figure 3: Attach DNA to Surface



Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

Figure 5: Fragments Become Double Stranded



The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

## Método de Illumina

- Desnaturalización para generar un molde de cadena sencilla
- Este proceso se repite hasta obtener un cluster con la secuencia original amplificado
- Añadir los nucleótidos marcados con fluorocromo y con un terminador en 3', polimerasa y cebadores
- Tomar la imagen de todo el flowcell
- Eliminar el terminador y el fluorocromo y iniciar el próximo ciclo

#### Figuras de Illumina:

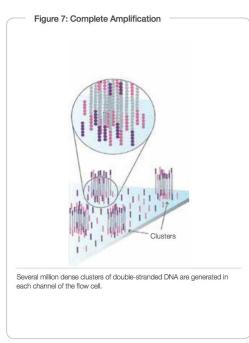
https://www.illumina.com/documents/products/techspotlights/techspotlight sequencing.pdf

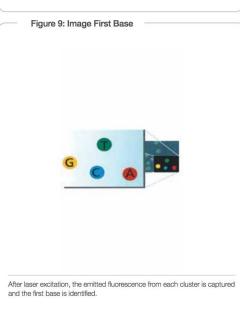
https://www.youtube.com/watc
h?v=womKfikWlxM

Figure 6: Denature the Double-Standed Molecules

Attached
Attached
Attached
Attached templates anchored to the substrate.



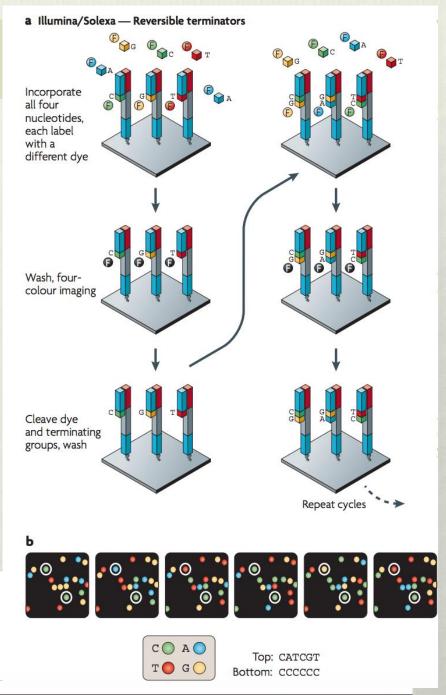




## Método de Illumina

# Sequencing technologies — the next generation

Michael L. Metzker\*\*

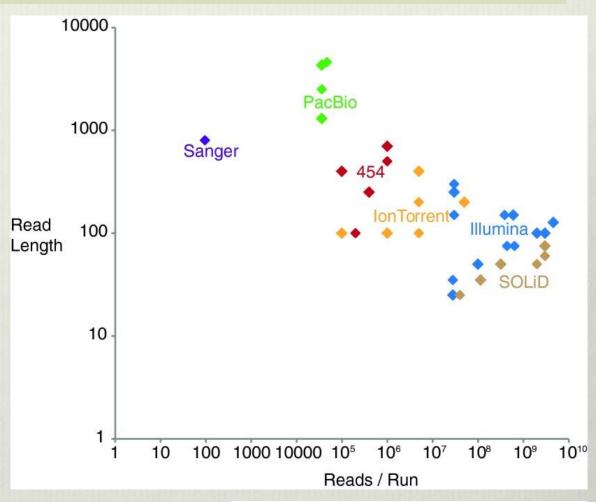


## Comparación de métodos

## Cada método en el mercado tiene sus fuertes y debilidades

#### Principales diferencias:

- La longitud de la lectura
- El número de lecturas
- El tiempo que lleva un proceso
- Los tipos de errores que podemos esperar
  - o 454 indels
  - Illumina cambios de una base







## **APLICACIONES**



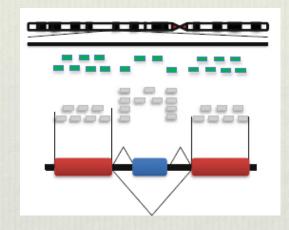


## Secuenciación del ADN genómico

# ACCCGTTACGTAAACGTTT C AGATGACGATGACCAAGGTTGACGA ACCCGTTACGTAAACGTTT G ACCCGTTACGTAAACGTTT G AGA ACCCGTTACGTAAACGTTT G AGATGAC ACCCGTTACGTAAACGTTT G AGATGAC ACCCGTTACGTAAACGTTT G AGATGACGATA ACCCGTTACGTAAACGTTT G AGATGACGATGACCA CGTTACGTAAACGTTT G AGATGACGATGACCAAGG ACGTAAACGTTT G AGATGACGATGACCAAGGTTGA TAAACGTTT G AGATGACGATGACCAAGGTTGACGA CGTTT G AGATGACGATGACCAAGGTTGACGA CGTTT G AGATGACGATGACCAAGGTTGACGA CGTTT G AGATGACGATGACCAAGGTTGACGA CGTTT G AGATGACGATGACCAAGGTTGACGA

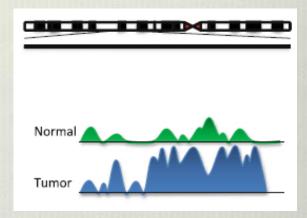
- Ensamblar un genoma
- SNVs and CNVs
- Inserciones y deleciones

## Expresión génica y su regulación



- Expresión génica
- ARNs pequeños
- TFBSs

#### Epigenómica

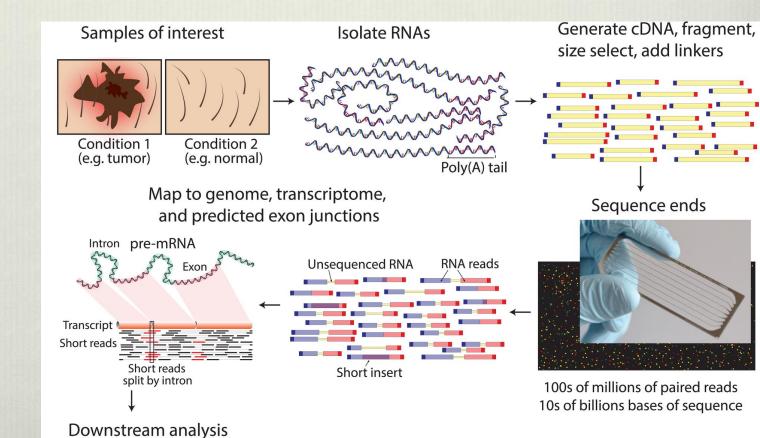


- Metilación del ADN
- Modificaciones químicas en histonas
- Cromatina

Imágenes extraídas de Schweiger et al., 2010.

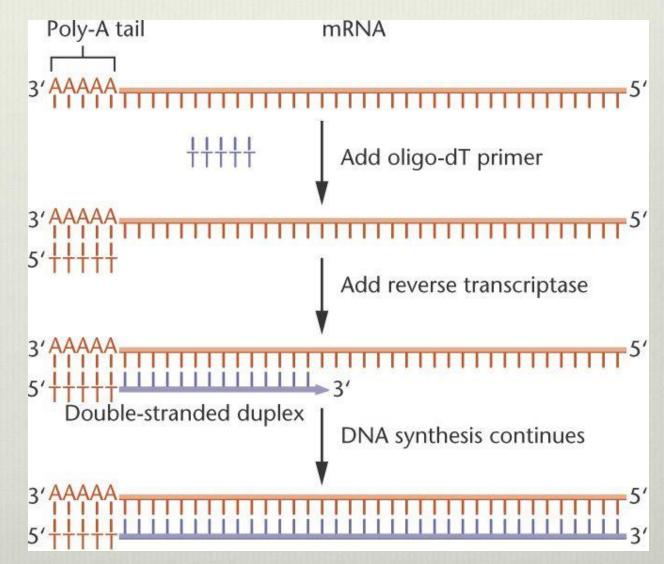
#### Permite:

- Detectar transcritos raros sin conocer el gen previamente
- Detectar splicing alternativo
- Detectar variación de secuencia



https://en.wikipedia .org/wiki/RNA-Seq

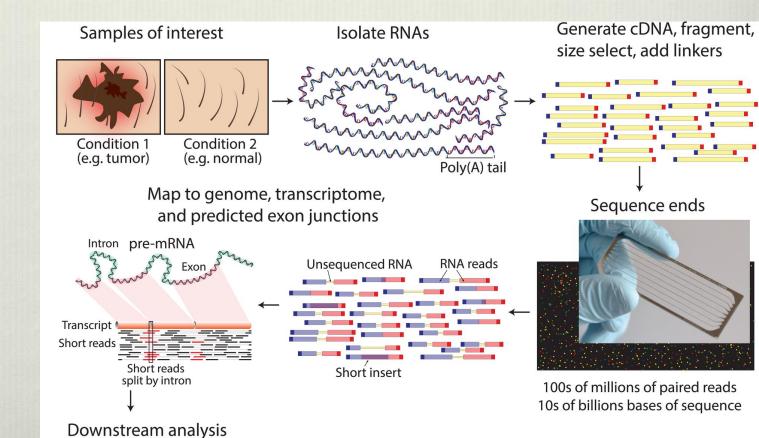
#### Transcripción inversa usando la poli-A



http://bio3400.nicer web.com/Locked/m edia/ch19/cDNA.ht ml

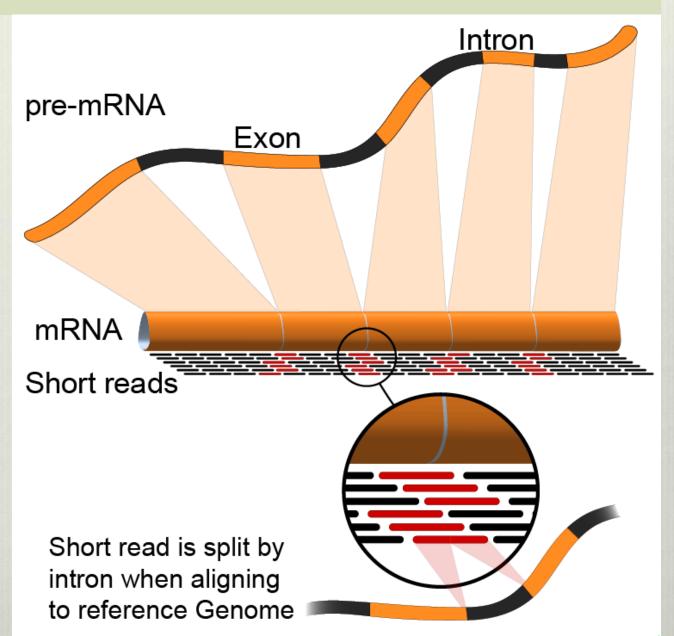
#### Permite:

- Detectar transcritos raros sin conocer el gen previamente
- Detectar splicing alternativo
- Detectar variación de secuencia

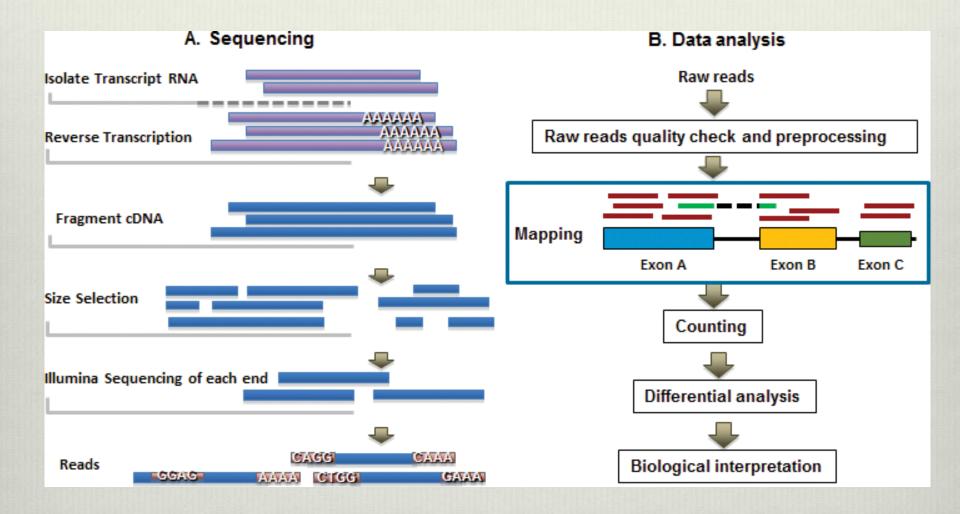


https://en.wikipedia .org/wiki/RNA-Seq

## Detectar fronteras entre exones y intrones



https://en.wikipedia.org /wiki/RNA-Seq#/media/File:RNA-Seq-alignment.png



## miRNA-seq

- Extracción del RNA total
- Ligar adaptadores
- Seleccionar por longitud
- Generar una librería de cDNA
- Secuenciar
- Análisis bioinformático

## Source: <a href="http://en.wikipedia.org/wiki/MicroRNA">http://en.wikipedia.org/wiki/MicroRNA</a> Sequencing

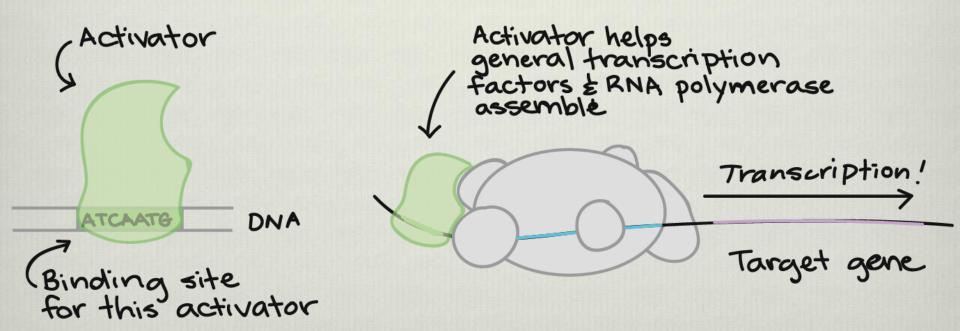
#### MIRNA-SEQ LIBRARY PREPARATION A Isolate 5ug of total RNA from your sample TOTAL RNA ISOLATION B Size fractionate total RNA using denaturing PAGE SIZE FRACTIONATION Select small RNA fraction (17-25 nt) 3' adapter ligation ADAPTOR LIGATION 5' adapter ligation Reverse transcribe RNA sequences RT & PCR PCR amplify sequences Flow Cell Attachment & Bridge Amplification TTEATRACTTACATECAPET SEQUENCING' Annealing of Sequencing Primers & Base extension Sequencing: Base Call, TAAGTGCTTCCATGTTTGAGTGT Deblock Extension,

Extension, Base Call

\*Illumina sequencing method depicted however other sequencing platforms can also be used.

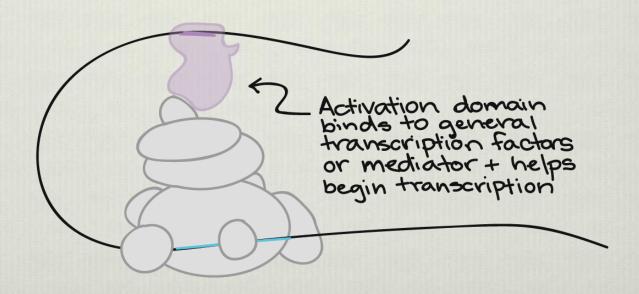
## CHIP-seq

- Un factor de transcripción (TF) es una proteína que se une a una secuencia de DNA específica
- Controlan la expresión génica
- Pueden tanto promover (activadores activator) o bloquear (represores- repressor) la unión de la polimerasa de RNA
- Estas proteínas tienen dominios específicos que les permite unirse al DNA
- Los TF se unen o al promotor o a un enhancer ("potenciador")



## CHIP-seq

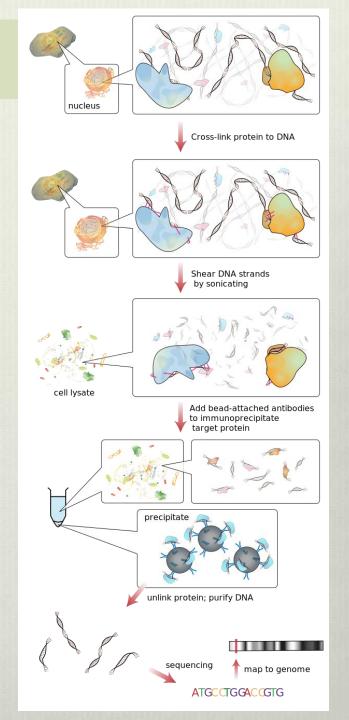
- Un factor de transcripción (TF) es una proteína que se une a una secuencia de DNA específica
- Controlan la expresión génica
- Pueden tanto promover (activadores activator) o bloquear (represores repressor) la unión de la polimerasa de RNA
- Estas proteínas tienen dominios específicos que les permite unirse al DNA
- Los TF se unen o al promotor o a un enhancer ("potenciador")



## CHIP-seq

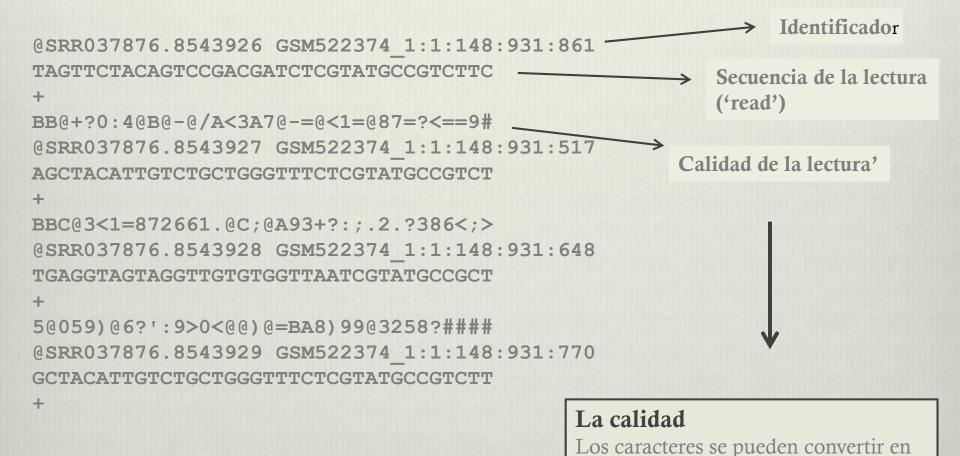
Detección experimental de los sitios de unión a factores de transcripción

- 1. Enlaces covalentes entre ADN y proteínas mediante formaldehído (reversible)
- 2. Inmunoprecipitación mediante un anticuerpo específico al factor de transcripción
- 3. Extraer /enriquecer las secuencias de ADN a las que se había unido une la proteína
- 4. Secuenciar el ADN
- 5. Localizar el ADN en el genoma



## La salida: el formato fastq

La salida del secuenciador: Los reads (lecturas) en formato fastq



un score (Q, Phred score) de calidad

## La calidad

#### Línea de calidad:

BB@+?0:4@B@-@/A<3A7

#### **ASCII** codificación

B = 66

(a) = 64

+ = 43

. . .

#### Codificación depende del fabricante

(por ejemplo codificación Sanger para la línea de calidad de de arriba)

Q(B) = 66 - 33 = 33 (primer base)

Q(@) = 64 - 33 = 31 (tercer base)

Q(+) = 43 - 33 = 10 (cuarta base)

```
Dec Hx Oct Html Chr
                      Dec Hx Oct Html Chr Dec Hx Oct Html Chr
 32 20 040 4#32; Space
                       64 40 100 a#64; 🛚
                                           96 60 140 4#96;
33 21 041 4#33;
                       65 41 101 6#65; A
                                           97 61 141 6#97;
34 22 042 6#34;
                       66 42 102 6#66; B
                                           98 62 142 4#98;
35 23 043 4#35; #
                       67 43 103 4#67; C
                                           99 63 143 4#99;
36 24 044 4#36; 6
                       68 44 104 4#68; D
                                          100 64 144 6#100; 6
37 25 045 4#37;
                       69 45 105 6#69; E
                                          101 65 145 6#101;
38 26 046 4#38; 4
                                          102 66 146 6#102;
                       70 46 106 6#70; F
39 27 047 4#39;
                       71 47 107 4#71; 6
                                          103 67 147 4#103; 9
                                          104 68 150 6#104; h
40 28 050 4#40;
                       72 48 110 6#72; H
41 29 051 6#41;
                       73 49 111 6#73; I
                                          105 69 151 6#105; 1
                       74 4A 112 6#74;
 42 2A 052 4#42;
                                          106 6A 152 6#106;
43 2B 053 6#43; +
                       75 4B 113 6#75; K
                                          |107 6B 153 4#107; k
44 2C 054 6#44;
                       76 4C 114 6#76; L
                                          108 6C 154 6#108; 1
45 2D 055 4#45;
                       77 4D 115 6#77; H
                                          109 6D 155 4#109; 1
 46 2E 056 4#46;
                       78 4E 116 6#78; N
                                          110 6E 156 6#110; n
47 2F 057 6#47;
                       79 4F 117 6#79; 0
                                          111 6F 157 6#111; 0
 48 30 060 4#48; 0
                       80 50 120 4#80; P
                                          |112 70 160 4#112; p
                       81 51 121 6#81; 0
 49 31 061 4#49; 1
                                          1113 71 161 6#113; q
 50 32 062 4#50; 2
                       82 52 122 6#82; R
                                          1114 72 162 6#114; I
                       83 53 123 4#83; $
51 33 063 4#51; 3
                                          115 73 163 4#115;
 52 34 064 4#52: 4
                       84 54 124 6#84; T
                                          116 74 164 a#116; t
                                          117 75 165 6#117: U
53 35 065 4#53; 5
                       85 55 125 6#85; U
54 36 066 4#54; 6
                       86 56 126 4#86; V
                                          118 76 166 4#118; 🔻
55 37 067 4#55: 7
                       87 57 127 6#87; 🐰
                                          119 77 167 4#119; 9
56 38 070 4#56; 8
                       88 58 130 4#88; X
                                          120 78 170 6#120; X
                       89 59 131 4#89; Y
                                          121 79 171 6#121; Y
                       90 5A 132 4#90; Z
                                          122 7A 172 6#122; 2
                       91 5B 133 6#91;
                                          123 7B 173 6#123;
61 3D 075 4#61; =
                                          125 7D 175 6#125;
```

```
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)
```

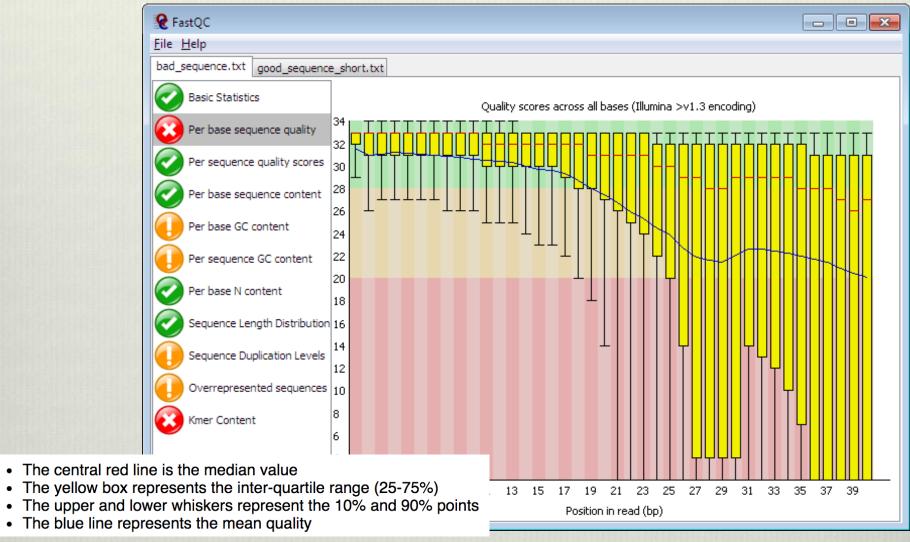
## La calidad

#### Phred Score y probabilidad de un error de secuenciación

$$Q=-10~\log_{10}P$$
 or 
$$P=10^{rac{-Q}{10}}$$

Phred quality scores are logarithmically linked to error probabilities		
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 10000	99.99 %
50	1 in 100000	99.999 %

## Control de calidad: FastQC



https://www.bioinformatics.babraham.ac.uk/pr
ojects/fastqc/

OQuality.html

https://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/3%20Analysis%20Modules/2%20Per%20Base%20Sequence%20Ouality.html

## Adaptadores y lecturas con baja calidad

## Si la molécula de RNA/DNA o el fragmento de DNA es mas corto que el número de ciclos

- → Se secuencia parte del 3' adaptador
- → Hay que encontrar esta secuencia y eliminarla de la lectura

#### La calidad de la secuenciación va bajando hacia el extremo 3'

→ Cortar la parta con menor calidad en 3'

#### **Trimmomatic:**

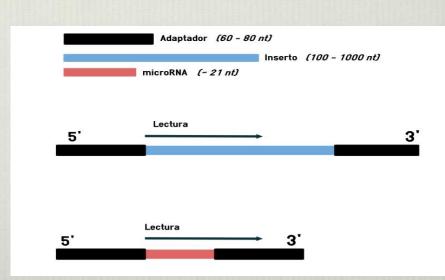
http://www.usadellab.org/cms/?page= trimmomatic

#### Cutadapt:

https://cutadapt.readthedocs.io/en/stable/

#### Trim\_galore:

https://www.bioinformatics.babraham. ac.uk/projects/trim\_galore/



## Adaptadores y lecturas con baja calidad

La idea básica es eliminar las bases del extremo 3' que están por debajo de cierto umbral de calidad -> ¡usar la información valida del extremo 5'!

Para eso, cutadapt usa la sumación parcial desde el extremo 3':

Suponemos un umbral de calidad de Q=10 y la siguiente serie de valores Q

Se resta primero el umbral de los valores Q:

Se calcula la suma parcial (empezando en el extremo 3')

$$(70)$$
,  $(38)$ ,  $8$ ,  $-8$ ,  $-25$ ,  $-23$ ,  $-20$ ,  $-21$ ,  $-15$ ,  $-7$ 

- El proceso se para en un valor > 0
- La lectura se corta en la posición que presenta el mínimo en la suma parcial
- La lectura final se quedaría en las primeras 4 nucleótidos: ACTGACTGAC → ACTG