



INSTITUTO
DE ECOLOGÍA
UNAM



INTRODUCCIÓN A LA BIOINFORMÁTICA

JAZMÍN SÁNCHEZ PÉREZ
LICENCIATURA EN CIENCIAS GENÓMICAS

DRA. VALERIA SOUZA

INSTITUTO DE ECOLOGÍA
UNIVERSIDAD NACIONAL AUTÓNOMA DE
MÉXICO

CEMCARTO EDIT MASTER TITLE STYLE

- Unidad IV – Genómica
 - Introducción al estudio de genomas
 - Ensamble de genomas
 - Click to edit the outline text format
 - Calidad de genomas
 - Second Outline Level
 - Anotación
 - Third Outline Level
 - Filogenias de genomas
 - Fourth Outline Level
- Unidad V – Metagenómica
 - Fifth Outline Level
 - Sixth Outline Level
 - Introducción al estudio de metagenómica
- Seventh Outline LevelClick to edit Master text styles
 - Metabarcoding (Anotación taxonómica)
 - Second level
 - Shotgun metagenomics (Anotación taxonómica y funcional)
 - Third level
 - Índices de diversidad
 - Fourth level
 - R como herramienta
 - Fifth level

- -OMICS es un sufijo que describe algo grande, y se refiere al campo de estudio de las ciencias de la vida que se enfoca en información a grande escala.
- “The other twist to “omics” may be associated with the “Om” (pronounced “Aum”), an ancient Sanskrit intonation, which, like music, transcends the barriers of age, race, culture, and even species. ”

¿QUÉ ES GENÓMICAS?

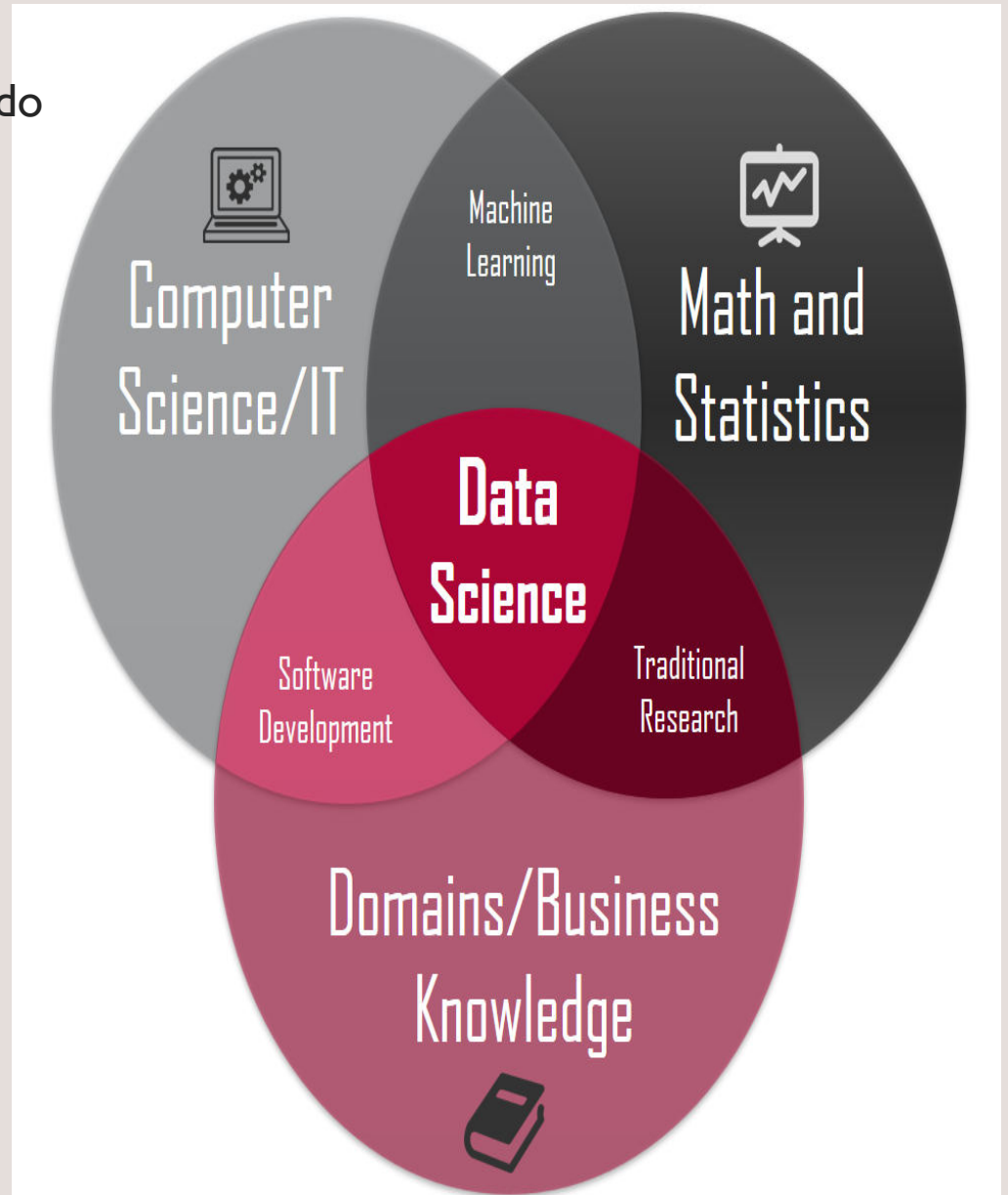
En un McDonald's Raw Bar, con una cerveza en la mano y unos amigos...

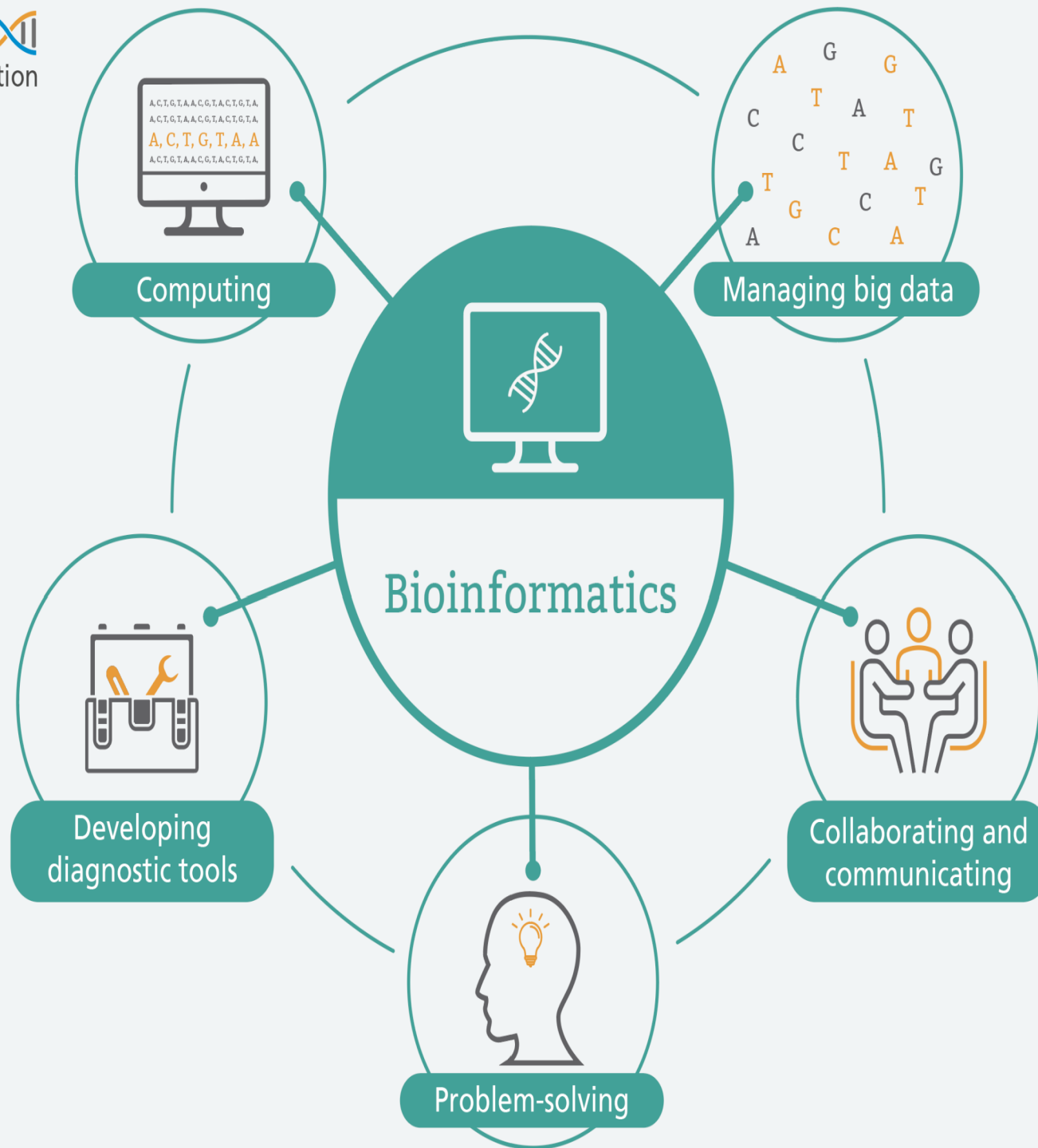
Dr. Thomas H. Roderick, en 1986, acuña el término *genomics*, como nombre para una revista que tratará de temas del genoma humano.

(Beer, Bethesda, and biology: how "genomics" came into being.)

1976 – Viral RNA-genome bacteriophage MS2
1995 – El primer genoma eucarionte secuenciado
 S.cerevisiae
1996 – Ciencias de Datos
2016 – Secuencia del Genoma Humano

“In 2012, when Harvard Business Review called it "The Sexiest Job of the 21st Century", the term "data science" became a buzzword. ”





A thick, wavy yellow line runs vertically along the left side of the slide, starting from the top and extending to the bottom. It has a slightly irregular, hand-drawn appearance.

GENOMICA

TOOLS & GENOME ASSEMBLE

CDATA
BASE64
ENCODING
XML
TEXT
TYPE
VALUE
STYLE

Full Report ▾

Send to: ▾

ASM190086v1

Organism name: [Borrelia afzelii \(spirochetes\)](#)**Intraspecific name:** Strain: BO23**BioSample:** [SAMN05941949](#)**BioProject:** [PRJNA350557](#)**Submitter:** NIH/NIAID**Date:** 2016/12/09**Assembly level:** Chromosome**Genome representation:** full**GenBank assembly accession:** GCA_001900865.1 (latest)**RefSeq assembly accession:** GCF_001900865.1 (latest)**RefSeq assembly and GenBank assembly identical:** yes**Assembly method:** SPAdes v. 3.7.1**Expected final version:** yes**Genome coverage:** 2000.0x**Sequencing technology:** Illumina NextSeq 500

IDs: 914131 [UID] 3786338 [GenBank] 3861958 [RefSeq]

History ([Show revision history](#))

Comment

Annotation was added by the NCBI Prokaryotic Genome Annotation Pipeline (released 2013). Information about the Pipeline can be found here: https://www.ncbi.nlm.nih.gov/genome/annotation_prok/

Global statistics

See [Genome](#) Information
for **Borrelia afzelii**

There are 11 assemblies
for this organism

[See more](#)

Access the data

[Download the RefSeq assembly](#)[Download the GenBank assembly](#)[Download the full sequence report](#)[Download the statistics report](#)

Assembly Information

[Assembly Help](#)[Assembly Basics](#)[NCBI Assembly Data Model](#)

Related Information

[BioProject](#)[BioSample](#)[Genome](#)[Nucleotide INSDC](#)[Nucleotide RefSeq](#)[Taxonomy](#)

Full Report ▾

Send to: ▾

Global statistics

AS

Org

Infr

Bio

Bio

Sub

Dat

Ass

Gen

Gen

Ref

Ref

Ass

Exp

Gen

Sec

IDs

His

Co

Anr

be found here: https://www.ncbi.nlm.nih.gov/genome/annotation_prok/

Global statistics

Total sequence length	1,259,498
Total assembly gap length	0
Gaps between scaffolds	0
Number of scaffolds	32
Scaffold N50	905,394
Scaffold L50	1
Number of contigs	32
Contig N50	905,394
Contig L50	1
Total number of chromosomes and plasmids	6
Number of component sequences (WGS or clone)	32

Access the data

[Download the RefSeq assembly](#)[Download the GenBank assembly](#)[Download the full sequence report](#)[Download the statistics report](#)

Assembly Information

[Assembly Help](#)[Assembly Basics](#)[ICBI Assembly Data Model](#)

Related Information

[BioProject](#)[BioSample](#)[Genome](#)[Nucleotide INSDC](#)[Nucleotide RefSeq](#)[Taxonomy](#)

CLICK TO EDIT MASTER TITLE

STYLE

chemotaxis protein CheR [Borrelia afzelii]

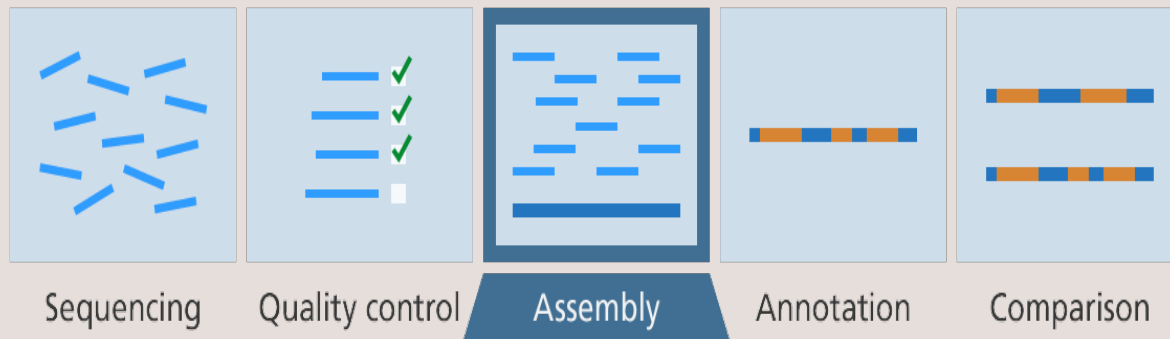
GenBank: APJ09116.1

[GenPept](#) [Identical Proteins](#) [Graphics](#)

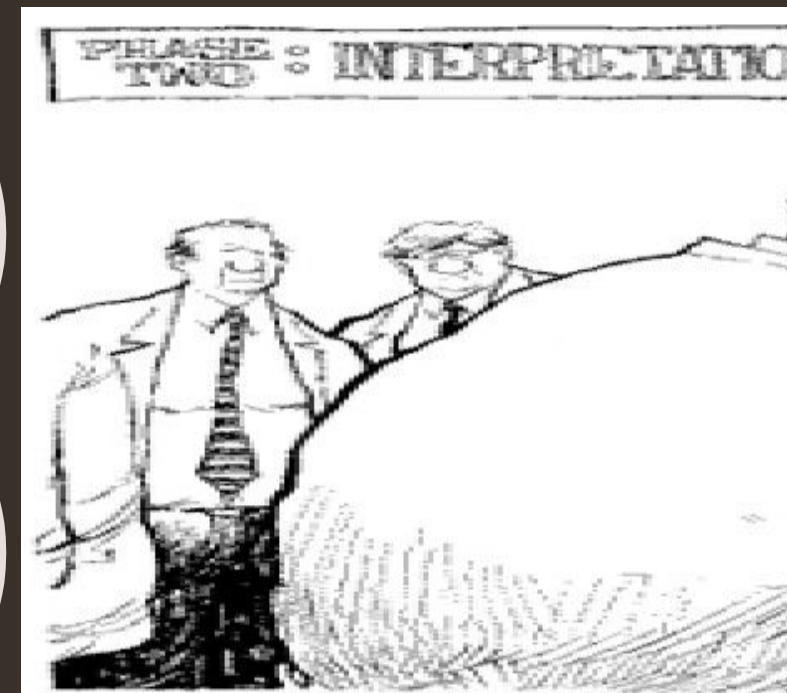
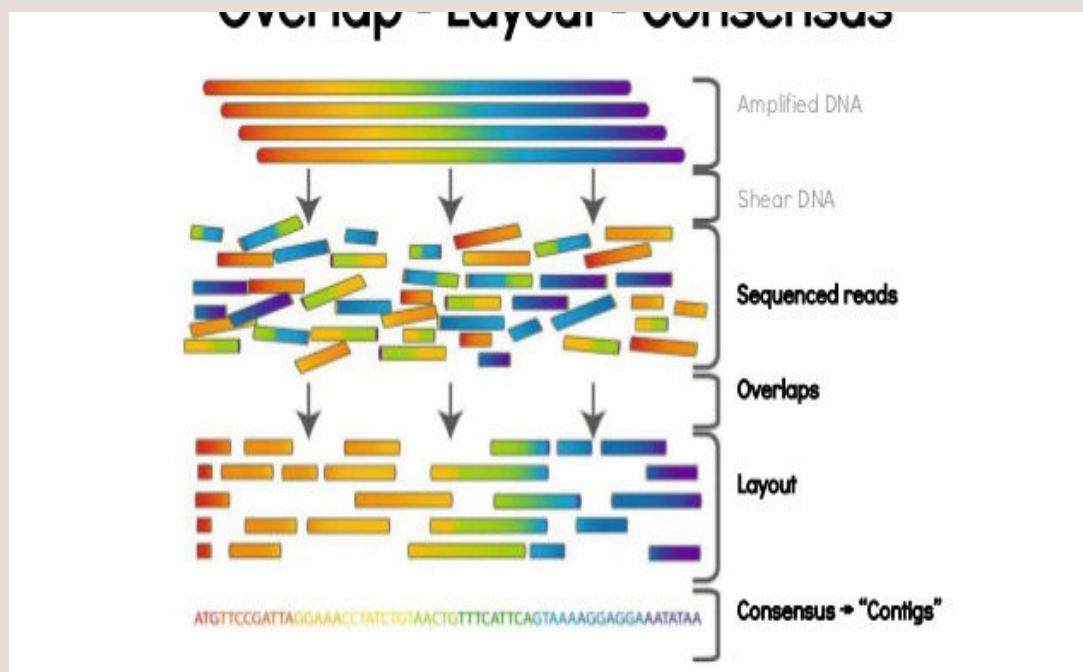
```
>APJ09116.1 chemotaxis protein CheR [Borrelia afzelii]
MNQNKFNLNINITKDELFRLIKIVYNNFGINLSEKKKMLIESRLSSLLKVKGFKNFTEYIDFLEKSAGNL
QLIELVDKISTNHTYFFRESKHFDFLNNKILPKLAEKILKSENSEIRIWSAGCSSGEEPYTAMILKEYM
ENNKVNFVKILATDISISVLHEAYEGIYPEDRTINLPKYLKTKYLNQLKDNKFQVKEILKKMVYFKKLN
LMDEKFPFSKKFDLIFCRNVMYIFDEKTRNDLANKFNYYLKNDSYLLIGHSETIRGNKNLKYIMPATYKK
N
```

FASTQ FORMAT

```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!' '*((( (**+))%%%+)(%%%).1***-+*' '))**55CCF>>>>>CCCCCCC65
```



ENSEMBLE DE GENOMAS



ENSAMBLE DE GENOMAS

CLICK TO EDIT MASTER TITLE

STYLE

Ensamble *de novo* incluye dos partes:

- El ensamble de la secuencia de nucleótidos
- La anotación. La descripción de los elementos funcionales y estructurales dentro de la secuencia

Desventajas:

- Second Outline Level
- Third Outline Level
- Se necesita mucho tiempo: para un genoma de eucarionte el ensamble puede tardar entre 1 y 5 años. Un programa de ensamble o de anotación puede correr bastante por semanas.
- Fourth Outline Level
- Fifth Outline Level
- Sixth Outline Level
- Seventh Outline Level
- Se necesita mucho poder computacional
- Third level
- Fourth level
- Fifth level
- La pregunta importante es: "¿Realmente necesito un genoma ensamblado y anotado para mi resolver pregunta de investigación?"

CHECKLIST PARA ANTES DE DISEÑAR TU PROYECTO DE ENSAMBLE DE GENOMA:

CLICK TO EDIT MASTER TITLE STYLE

- ☐ Para la extracción de ADN, seleccionar un individuo que sea un **buen representante** de la especie y que pueda **proveer suficiente ADN**
- ☐ **Extraer más ADN de lo que crees necesitar.**
 - Second Outline Level
- ☐ Recordar **extraer ARN y secuenciar ARN**, esto es muy funcional a la hora de ensamblar transcriptoma y para la anotación de funciones.
 - Third Outline Level
 - Fourth Outline Level
 - Fifth Outline Level
 - Sixth Outline Level
- ☐ **Decidir** de manera temprana **qué tecnología de secuenciación** se va a usar y que **herramientas de ensamblaje**. no quieres terminar con información que no puedes analizar con cualquier programa.
 - Seventh Outline Level
 - Second level
 - Third level
 - Fourth level
 - Fifth level

BISEÑO EXPERIMENTAL

CLICK TO EDIT MASTER TITLE

NO.1 - PROPIEDADES DEL GENOMA

STYLE

- Tamaño del Genoma
 - El número de lecturas (*reads*) depende del tamaño del genoma.
- Repeticiones
 - La cantidad y distribución de las repeticiones modifican la calidad del ensamble. Entre más repeticiones es mejor usar lecturas largas.
- Heterocigosidad
 - Los programas de ensamble colapsan alelos, si hay alta heterocigosidad las regiones homologas serán detectadas como regiones diferentes.
- Niveles de ploidia
 - Mejor usar tejidos haploides
- Contenido de GC
 - Niveles bajos o altos de GC en ciertas regiones causan problemas en ciertos secuenciadores como Illumina.

Sequoia sempervirens



II cromosomas
Hexaploide

Opuntia ficus-indica



II cromosomas
Octoploide

DESIGN EXPERIMENTAL TITLE

NO.1 - PROPIEDADES DEL GENOMA

STYLE

- Tamaño del Genoma
 - El número de lecturas (*reads*) depende del tamaño del genoma.
- Repeticiones
 - La cantidad y distribución de las repeticiones modifican la calidad del ensamble. Entre más repeticiones es mejor usar lecturas largas.
- Heterocigosidad
 - Los programas de ensamble colapsan alelos, si hay alta heterocigosidad las regiones homologas serán detectadas como regiones diferentes.
- Niveles de ploidia
 - Mejor usar tejidos haploides
- Contenido de GC
 - Niveles bajos o altos de GC en ciertas regiones causan problemas en ciertos secuenciadores como Illumina.

OTRAS COSAS A TOMAR EN CUENTA...

STYLE

- No intentar hacer *pool* de individuos.
 - Esto causa heterozigosidad en las secuencias.
- Whole Genome Amplification.
 - Click to edit the outline text format
 - Se pueden crear secuencias quiméricas al fusionar secuencias que no están relacionadas.
 - Second Outline Level
 - Third Outline Level
 - Fourth Outline Level
 - Fifth Outline Level
 - Sixth Outline Level
- Presencia de otros organismos.
 - Puede haber contaminación de otros organismos que se introducen en el laboratorio, ADN humano, etc.
- Seventh Outline LevelClick to edit Master text styles
- ADN de organelos.
 - Second level
 - En ciertas células hay una mayor cantidad de ADN mitocondria o de cloroplastos que el ADN nuclear.
 - Third level
 - Fourth level
 - Fifth level

BISEÑO EXPERIMENTAL

CLICK TO EDIT MASTER TITLE

NO.1 - PROPIEDADES DEL GENOMA

STYLE

- Tamaño del Genoma
 - El número de lecturas (*reads*) depende del tamaño del genoma.
- Repeticiones
 - La cantidad y distribución de las repeticiones modifican la calidad del ensamble. Entre más repeticiones es mejor usar lecturas largas.
- Heterocigosidad
 - Los programas de ensamble colapsan alelos, si hay alta heterocigosidad las regiones homologas serán detectadas como regiones diferentes.
- Niveles de ploidia
 - Mejor usar tejidos haploides
- Contenido de GC
 - Niveles bajos o altos de GC en ciertas regiones causan problemas en ciertos secuenciadores como Illumina.

DESIGN & EDIT MASTER TITLE

NO.3 – TECNOLOGÍAS DE SECUENCIACIÓN

SEVENTH OUTLINE
LEVELCLICK TO EDIT
MASTER TEXT STYLES

- Click to edit the outline text format
- Son más baratos
 - Second Outline Level
- Mayor cobertura
 - Third Outline Level
- Problemas con regiones repetidas y con contenido de GC
 - Fourth Outline Level
- Lecturas de ~36 bp hasta ~500 bp
 - Fifth Outline Level
 - Sixth Outline Level
- Seventh Outline LevelClick to edit Master text styles

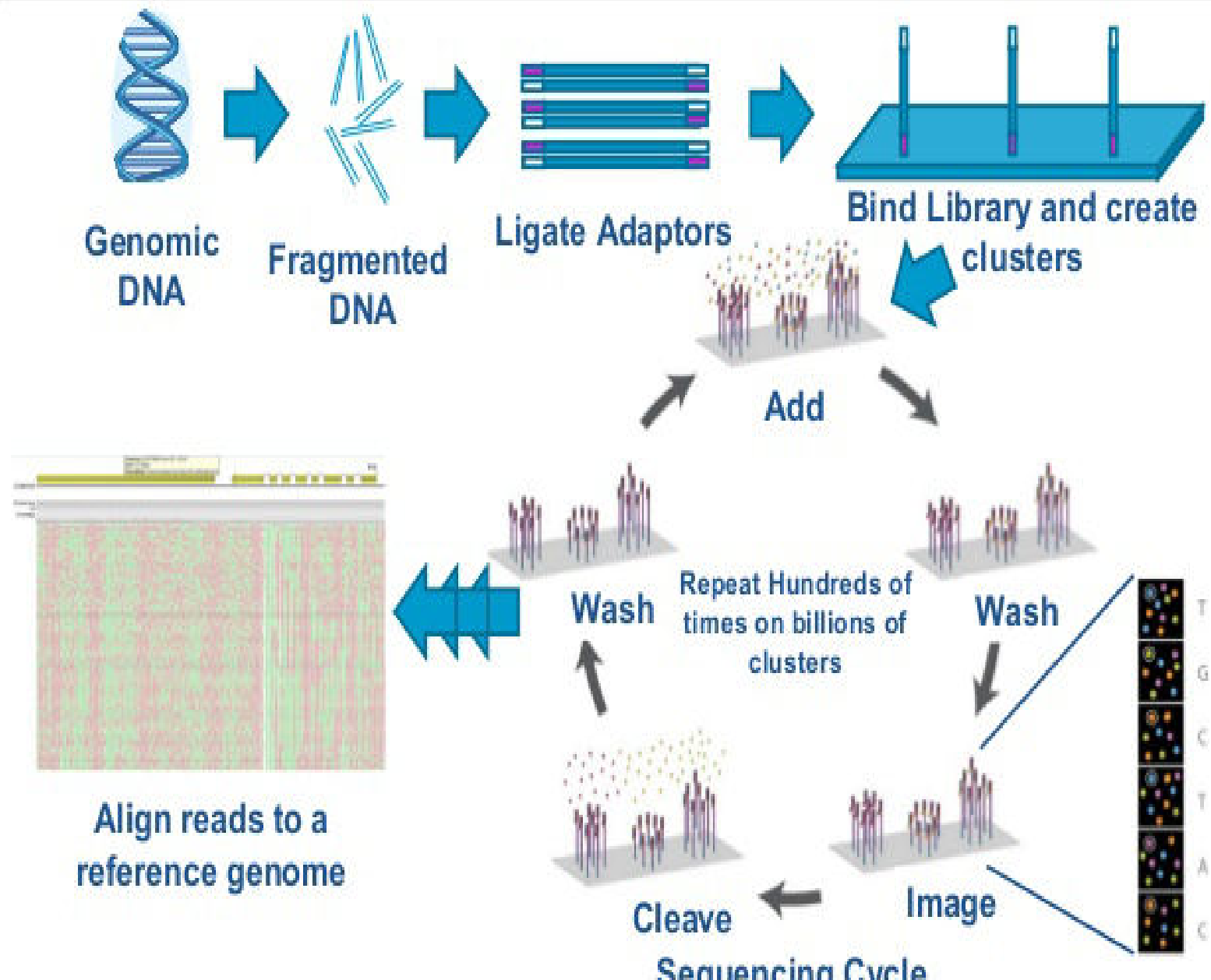
Second level

SEVENTH OUTLINE
LEVELCLICK TO EDIT
MASTER TEXT STYLES

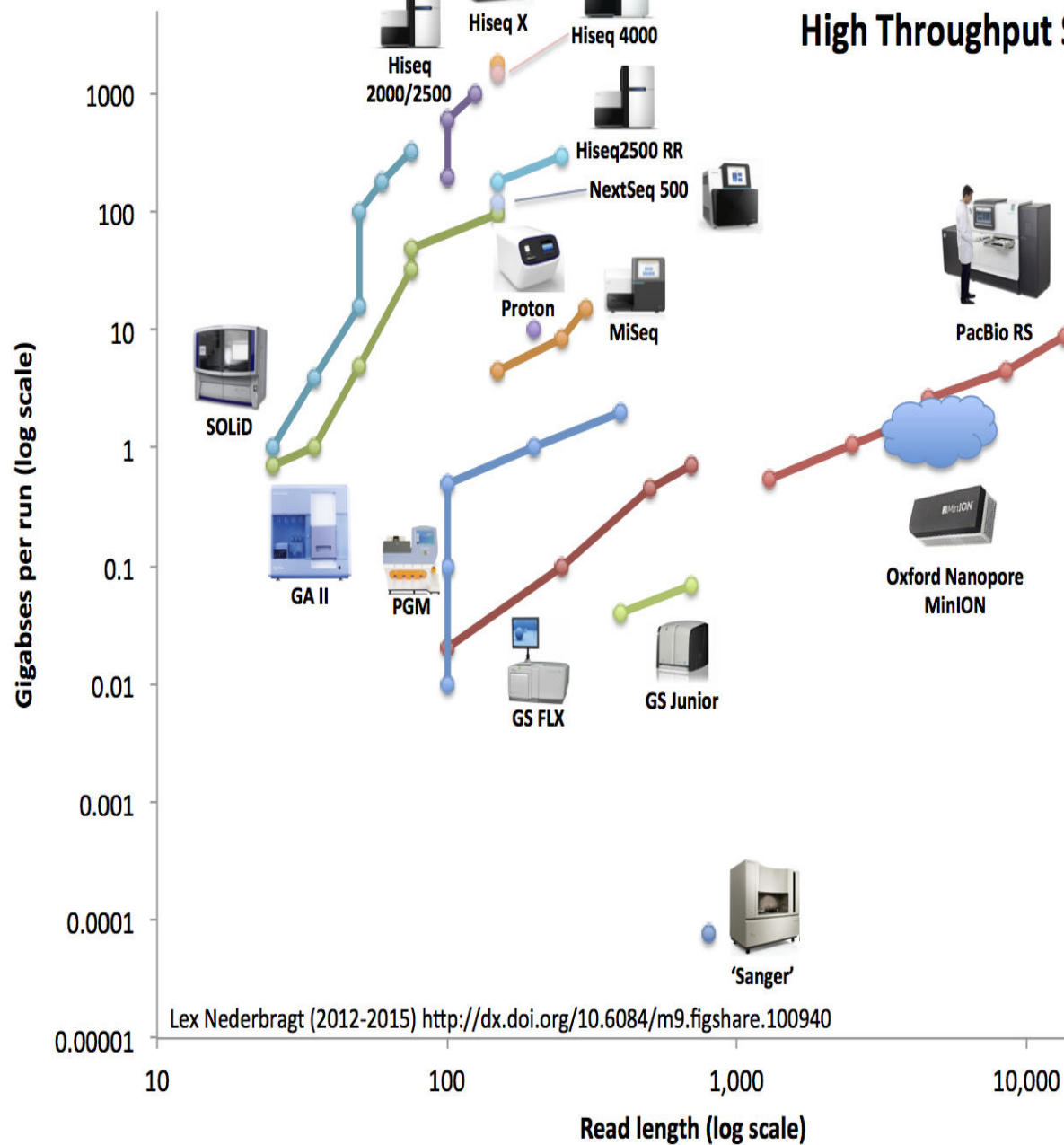
- Click to edit the outline text format
- Son más caros
 - Second Outline Level
- Lecturas de mayor tamaño, 10,00 bp hasta 100,000 bp
 - Third Outline Level
- Ensamblan mejor zonas repetidas (genomas de eucariontes)
 - Fifth Outline Level
- No son accesibles en todos los países
 - Sixth Outline Level
- Mayores requerimientos de la
- Seventh Outline LevelClick to edit Master text styles

Second level

2nd Generation Sequencing Overview



Developments in High Throughput Sequencing



BISEÑO EXPERIMENTAL

CLICK TO EDIT MASTER TITLE

NO.1 - PROPIEDADES DEL GENOMA

STYLE

- Tamaño del Genoma
 - El número de lecturas (*reads*) depende del tamaño del genoma.
- Repeticiones
 - La cantidad y distribución de las repeticiones modifican la calidad del ensamble. Entre más repeticiones es mejor usar lecturas largas.
- Heterocigosidad
 - Los programas de ensamble colapsan alelos, si hay alta heterocigosidad las regiones homologas serán detectadas como regiones diferentes.
- Niveles de ploidia
 - Mejor usar tejidos haploides
- Contenido de GC
 - Niveles bajos o altos de GC en ciertas regiones causan problemas en ciertos secuenciadores como Illumina.

ENIGMA BOULE DE QUINCE MASTER TITLE STYLE

MISIÓN: formar secuencias lo más largas posible con el menor número de gaps.

- Second Outline Level
 - Third Outline Level
 - Fourth Outline Level
 - Fifth Outline Level
 - Sixth Outline Level
- Seventh Outline LevelClick to edit Master text styles
 - Second level
 - Third level
 - Fourth level
 - Fifth level

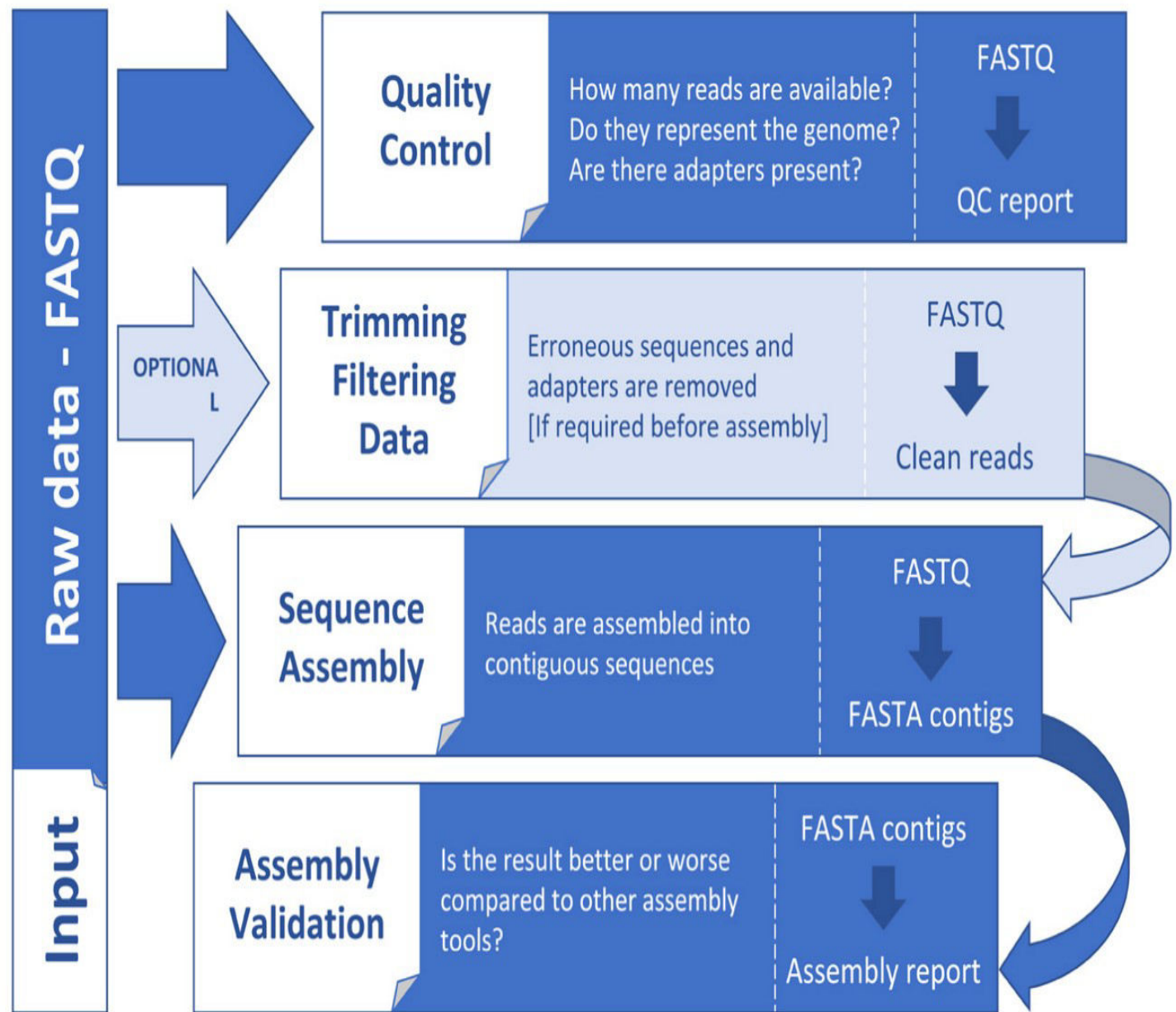


Figure 2. General steps in a genome assembly workflow. Input and output data are indicated for each step.

CLICK TO EDIT MASTER TITLE

ENSEMBLE STYLE

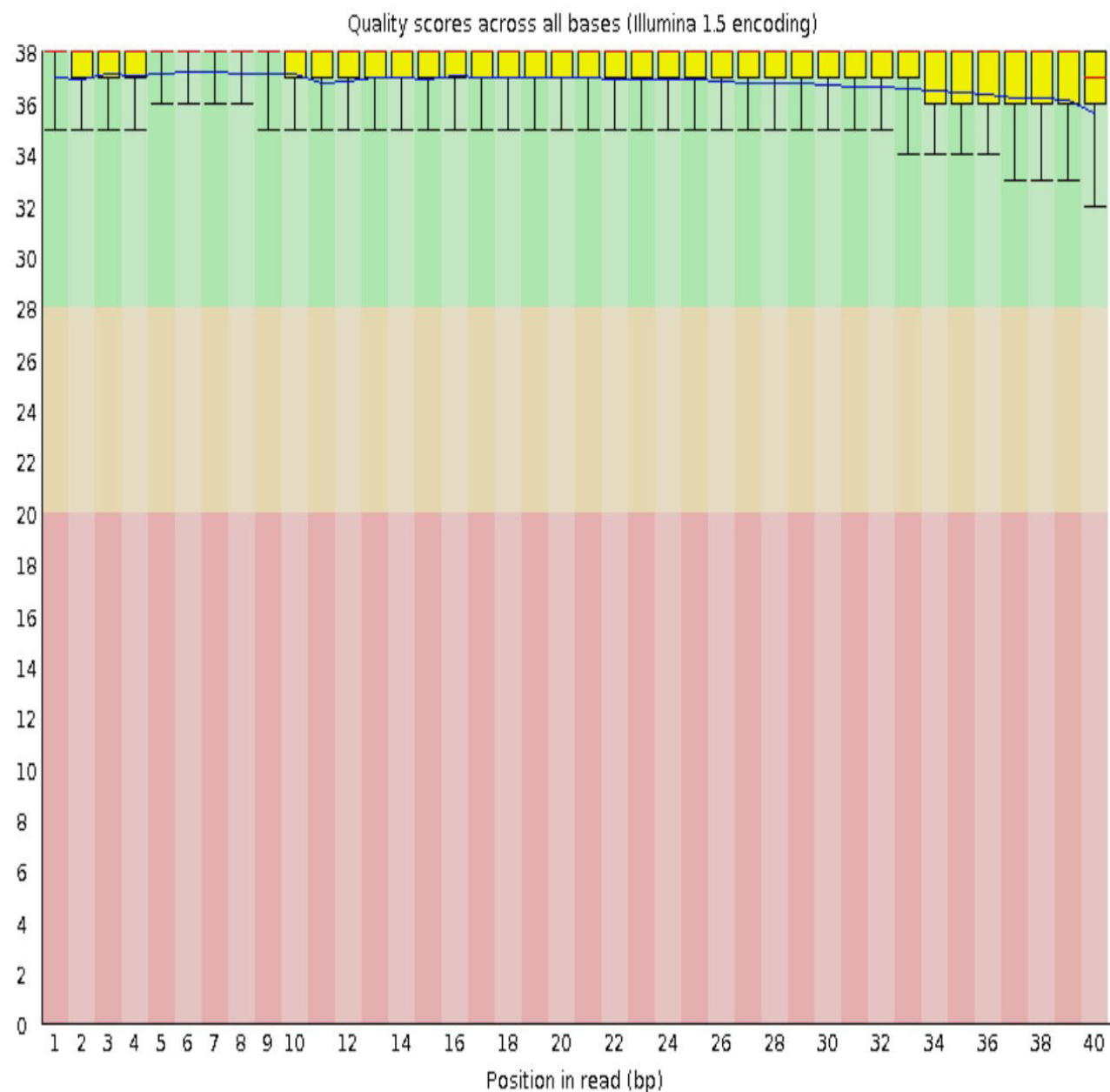
1. CONTROLES DE CALIDAD

- Basic alignment and outline text format
- Trimming: eliminar adaptadores y eliminar zonas con poca calidad
 - Third Outline Level
- FASTQC: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
 - Fourth Outline Level
 - Fifth Outline Level
 - Sixth Outline Level
- Seventh Outline LevelClick to edit Master text styles
 - Second level
 - Third level
 - Fourth level
 - Fifth level

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](#)

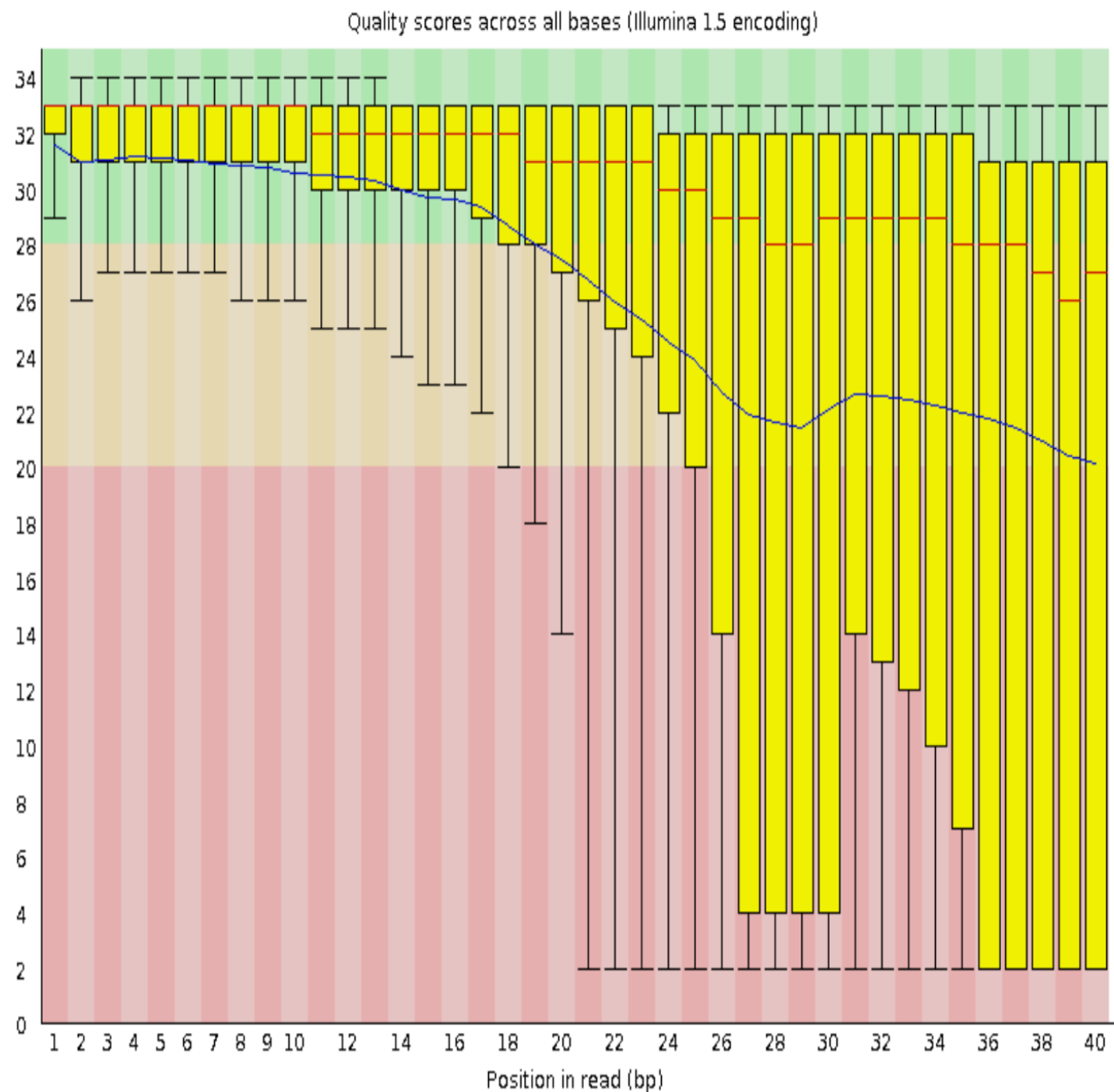
✓ Per base sequence quality



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](#)

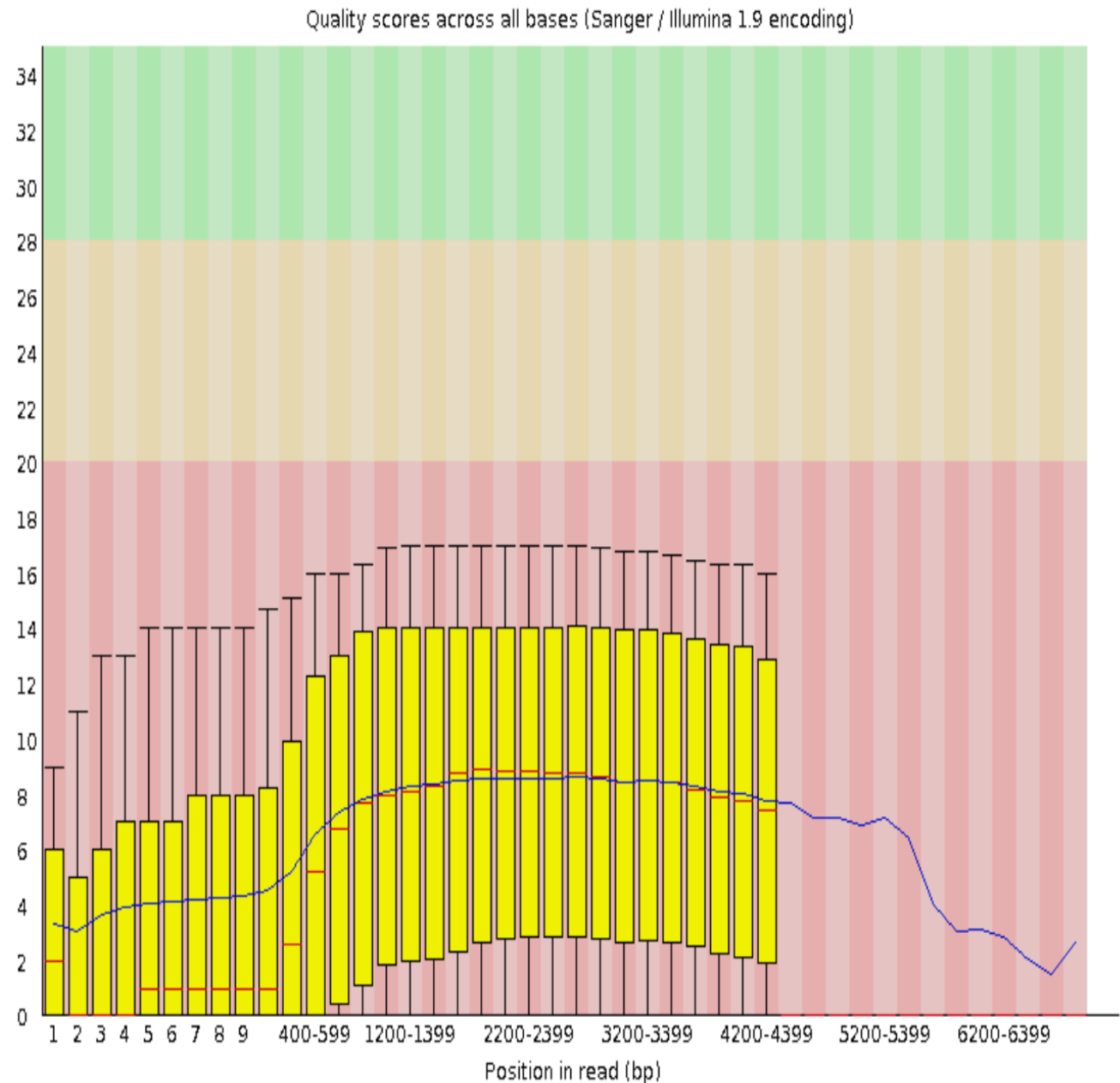
✗ Per base sequence quality



Summary

- ✓ [Basic Statistics](#)
- ✗ [Per base 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](#)

✗ Per base sequence quality



ENSEMBLE

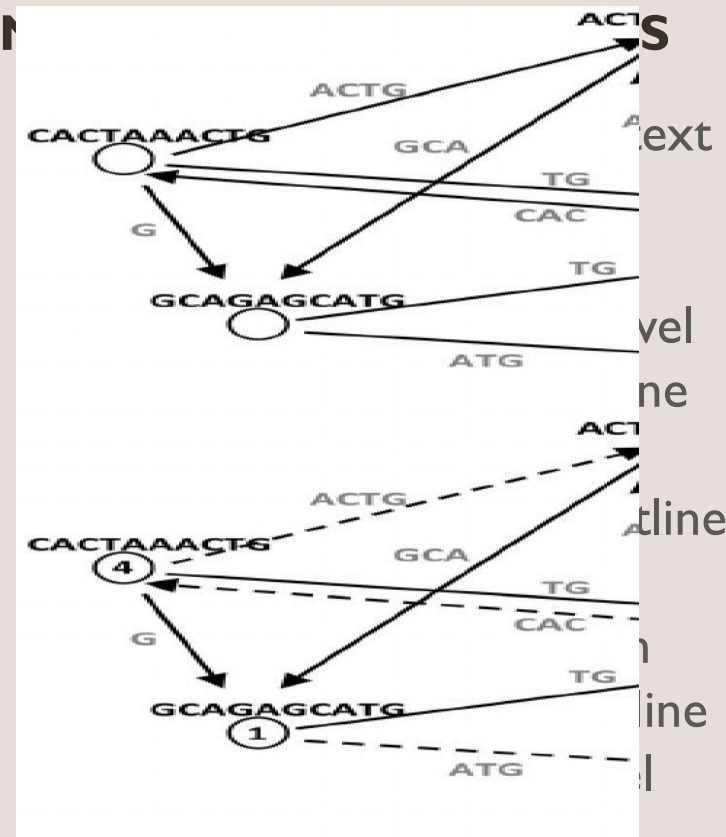
2. ENSAMBLE DE LA SECUENCIA

SEVENTH OUTLINE
LEVELCLICK TO EDIT
MASTER TEXT STYLES

- Click to edit the text
- Second Outline Level
- Third Outline Level
- Trabaja con k-meros
- Fourth Outline Level
- Menos tiempo computacional
- Fifth Outline Level
- Sixth Outline Level
- Seventh Outline LevelClick to edit Master text styles

Second level

SEVENTH OUTLINE
LEVELCLICK TO EDIT
MASTER TEXT STYLES



- Seventh Outline LevelClick to edit Master text styles

Second level

ENSEMBLE
STYLE

2. ENSAMBLE DE LA SECUENCIA

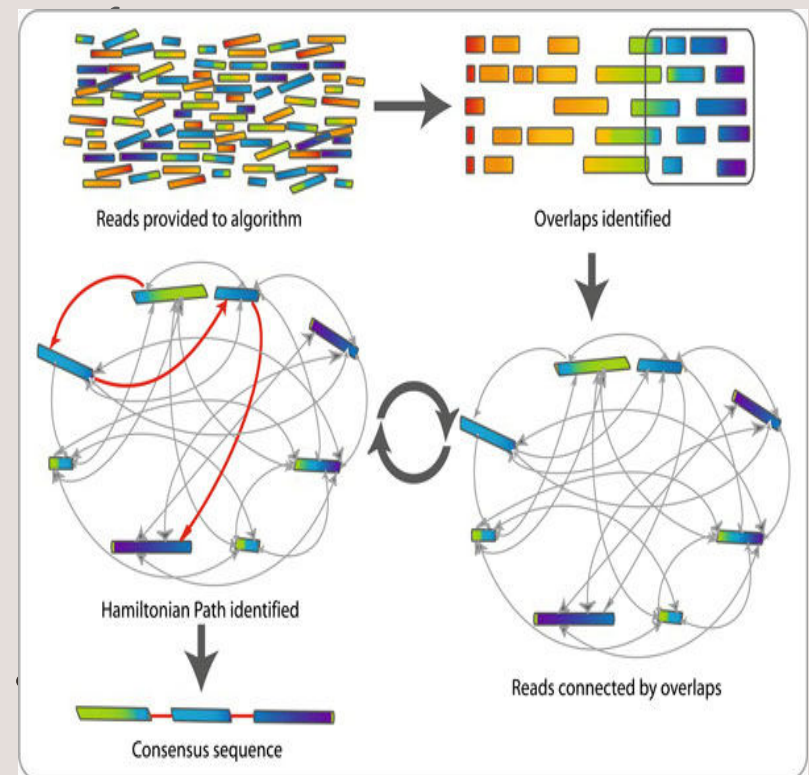
SEVENTH OUTLINE
LEVELCLICK TO EDIT
MASTER TEXT STYLES

- Click to edit the outline text format
 - Second Outline Level
 - Third Outline Level
 - Fourth Outline Level
 - Fifth Outline Level
 - Sixth Outline Level
- Seventh Outline LevelClick to edit Master text styles

Second level

SEVENTH OUTLINE
OVERLAP IDENTIFICATION
MASTER TEXT STYLES

- Click to separate the outline text



Second level

ENSEMBLE

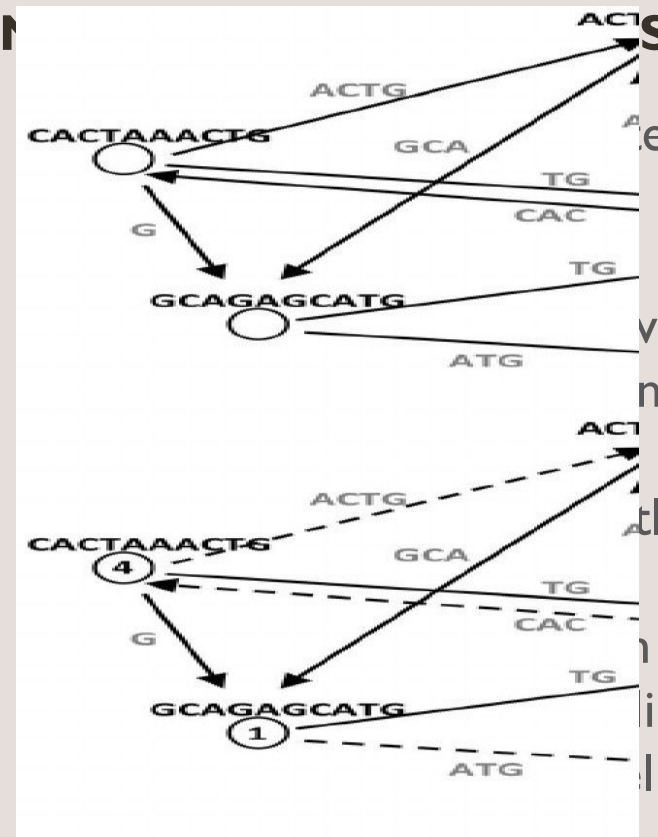
2. ENSAMBLE DE LA SECUENCIA

SEVENTH OUTLINE
LEVELCLICK TO EDIT
MASTER TEXT STYLES

- Click to edit the text
- Second Outline Level
- Third Outline Level
- Trabaja con k-meros
- Fourth Outline Level
- Menos tiempo computacional
- Fifth Outline Level
- Sixth Outline Level
- Seventh Outline LevelClick to edit Master text styles

Second level

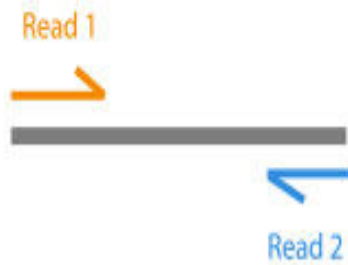
SEVENTH OUTLINE
LEVELCLICK TO EDIT



- Seventh Outline LevelClick to edit Master text styles

Second level

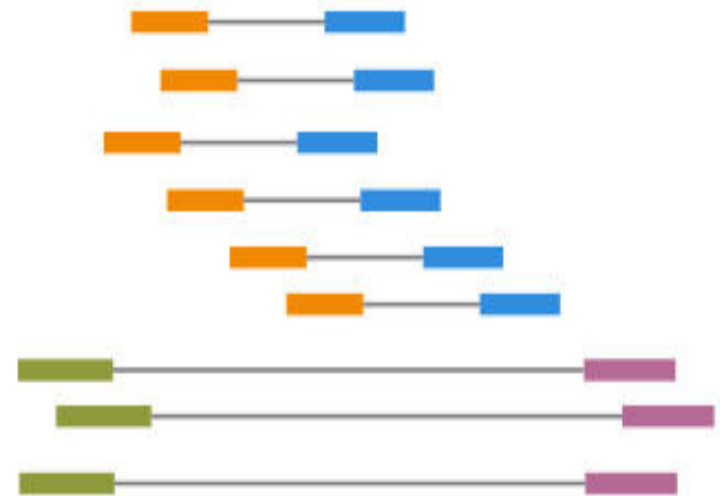
Short-Insert Paired End Reads



Long-Insert Paired End Reads (Mate Pair)



De Novo Assembly



CLICK TO EDIT MASTER TITLE

ENSEMBLE STYLE

3. DETERMINAR QUE ESTA LISTO - CALIDAD

- Click to edit the outline text format
 - Second Outline Level
 - Third Outline Level
 - Es el tamaño del scaffold o contig más pequeño que, sumado con los más largos que él, deben cubrir el 50% del genoma
 - Fourth Outline Level
 - NO IMPLICA QUE ESTE BIEN ENSAMBLADO
 - Fifth Outline Level
 - Sixth Outline Level
- Usar varios ensambladores y compararlos, correr el mismo ensamblador con distintos parámetros
 - Second level
 - Quast para comparar
 - Third level
- Determinar cuántos genes de proteínas tiene el ensamble
 - herramienta BUSCO
 - Fourth level
 - Fifth level

ANOTACIÓN ESTRUCTURAL

1. ANOTACIÓN ESTRUCTURAL

- ¿Dónde están los genes y a qué se parecen? Sobre todo genes que codifican para proteínas

- Second Outline Level

• INTRÍNSECO

- Se centra en la información extraída de la secuencia
 - Fifth Outline Level
 - Sixth Outline Level
- Es intensivo ya que se necesitan crear modelos y los software tienen que entrenarse y optimizarse.
 - Second level
 - Third level
- Son específicos para cada genoma
 - Fourth level
 - Fifth level

EXTRÍNSECO

- Utiliza información adicional a la secuencia como lo son:
 - Zonas que codifican para proteínas
 - Sitios de splicing
 - Transcritos
- Es universal
- Se utilizan bases de datos, usualmente de polipéptidos
- No da información estructural del gen

ANOTACIÓN ESTRUCTURAL

1. ANOTACIÓN ESTRUCTURAL

- ¿Dónde están los genes y a qué se parecen? Sobre todo genes que codifican para proteínas

- Second Outline Level

• INTRÍNSECO

- Se centra en la información extraída de la secuencia
 - Fifth Outline Level
 - Sixth Outline Level
- Es intensivo ya que se necesitan crear modelos y los software tienen que entrenarse y optimizarse.
 - Second level
 - Third level
- Son específicos para cada genoma
 - Fourth level
 - Fifth level

EXTRÍNSECO

- Utiliza información adicional a la secuencia como lo son:
 - Zonas que codifican para proteínas
 - Sitios de splicing
 - Transcritos
- Es universal
- Se utilizan bases de datos, usualmente de polipéptidos
- No da información estructural del gen

CLICK TO EDIT MASTER TITLE

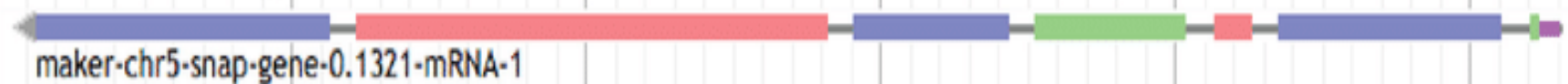
1. ANOTACIÓN FUNCIONAL

- Asignar información biológica relevante como predicción de polipeptidos y
- Click to edit the outline text format
- Second Outline Level
- Dos objetivos principales
 - Third Outline Level
 - Asignar elementos funcionales a genes
 - Fourth Outline Level
 - Hacer una revisión de calidad para los genes precedidos
 - Fifth Outline Level
 - Sixth Outline Level
- Seventh Outline LevelClick to edit Master text styles
 - Second level
 - Third level
 - Fourth level
 - Fifth level

Current annotation



MAKER prediction (homology, RNA-seq and peptides)



Gioti *et al* annotation



RNA-seq coverage



Peptide evidence

