

Genomics and Transcriptomics

Class 05 - Sequence Assembly



INSTRUCTOR:

Aureliano Bombarely

Department of Bioscience

Università degli Studi di Milano

aureliano.bombarely@unimi.it

Outline of Topics

1. Brief history about genome assembly
2. Basics about sequence assembly
3. Whole genome assembly
 - 3.1. From reads to contigs and scaffolds.
 - 3.2. From scaffolds to chromosomes.
4. Transcriptome assembly



Outline of Topics

1. Brief history about genome assembly
2. Basics about sequence assembly
3. Whole genome assembly
 - 3.1. From reads to contigs and scaffolds.
 - 3.2. From scaffolds to chromosomes.
4. Transcriptome assembly



1. Brief history about genome assembly

WHAT IS A GENOME?

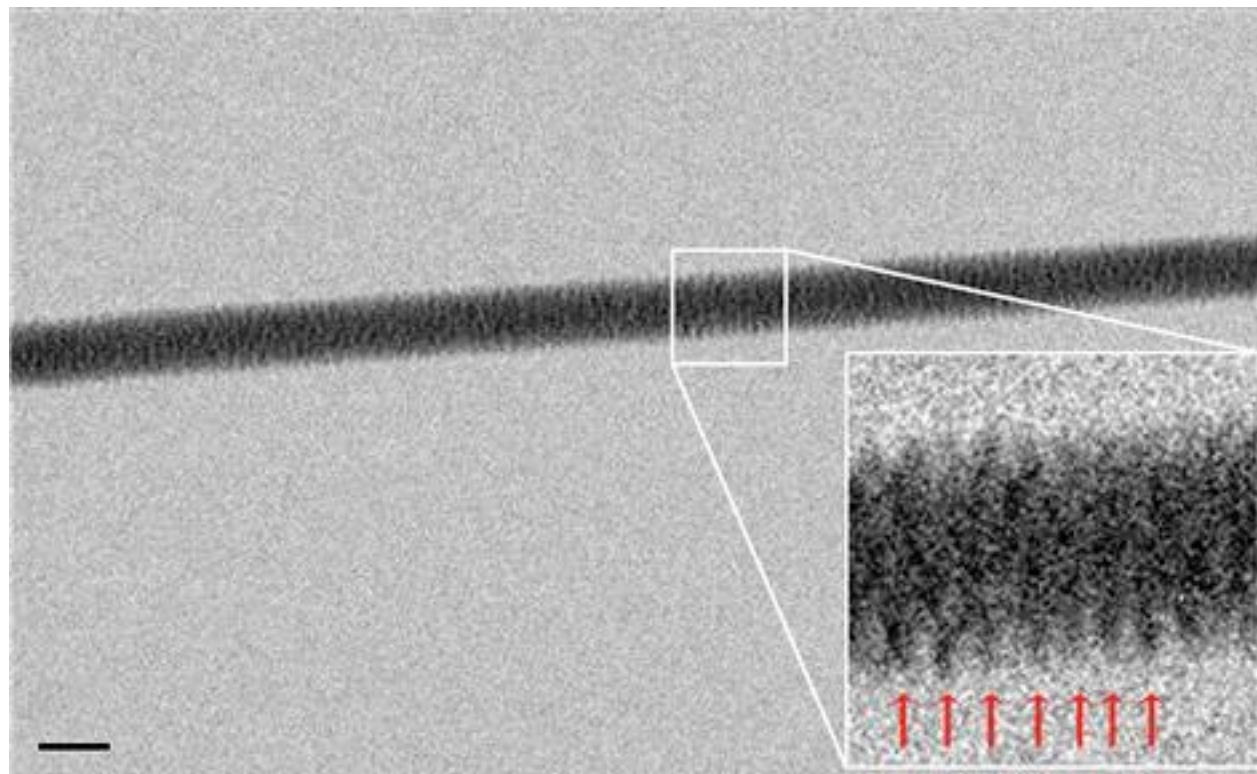
Life is specified by **genomes**. Every organism, including humans, has a genome that contains all of the biological information needed to build and maintain a living example of that organism. The biological information contained in a genome is encoded in its **deoxyribonucleic acid (DNA)** and is divided into discrete units called **genes**. Genes code for proteins that attach to the genome at the appropriate positions and switch on a series of reactions called gene expression.

In 1909, Danish botanist Wilhelm Johanssen coined the word **gene** for the hereditary unit found on a chromosome. Nearly 50 years earlier, Gregor Mendel had characterized hereditary units as **factors**— observable differences that were passed from parent to offspring. Today we know that a single gene consists of a unique sequence of DNA that provides the complete instructions to make a functional product, called a protein. Genes instruct each cell type— such as skin, brain, and liver—to make discrete sets of proteins at just the right times, and it is through this specificity that unique organisms arise.



1. Brief history about genome assembly

Genome = $N \times$ Sequence of DNA



???

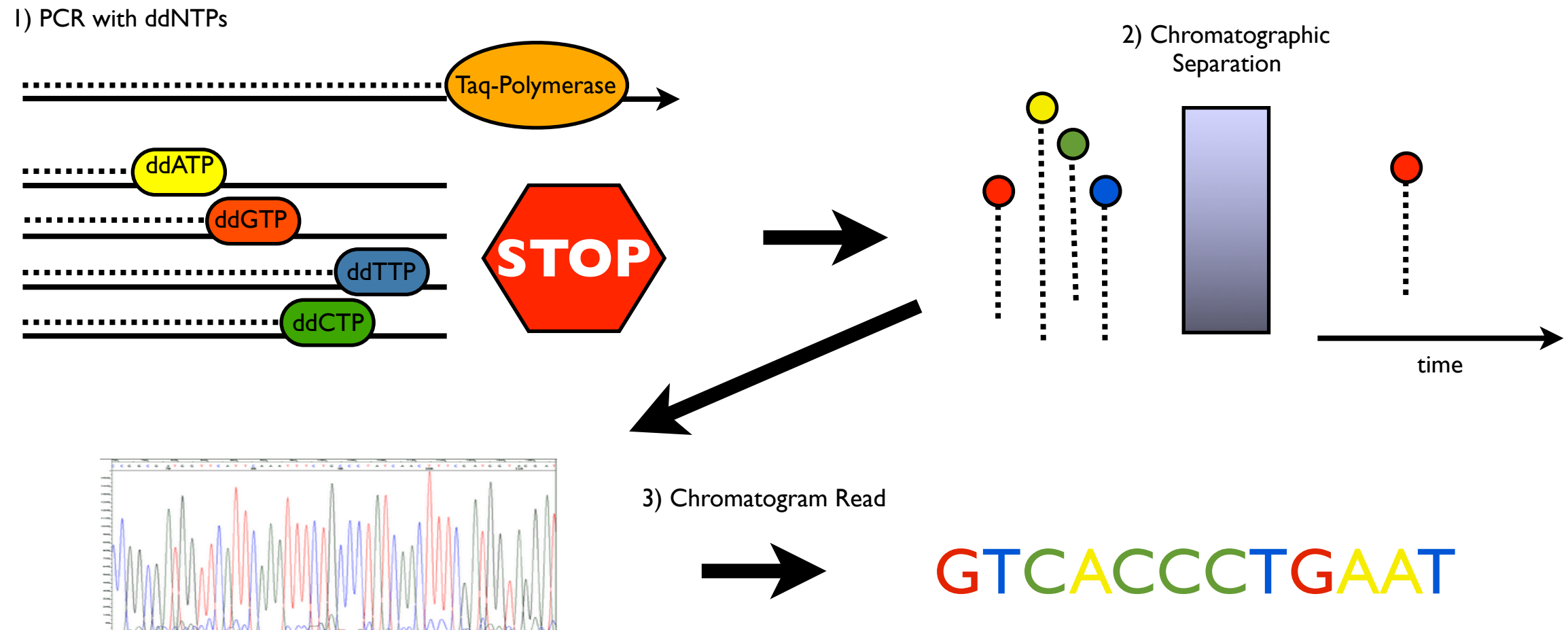
→ DNA sequencing

Gentile F. et al. *Direct Imaging of DNA Fibers: The Visage of Double Helix*
Nano Lett., 2012, 12 (12), pp 6453–6458

1. Brief history about genome assembly

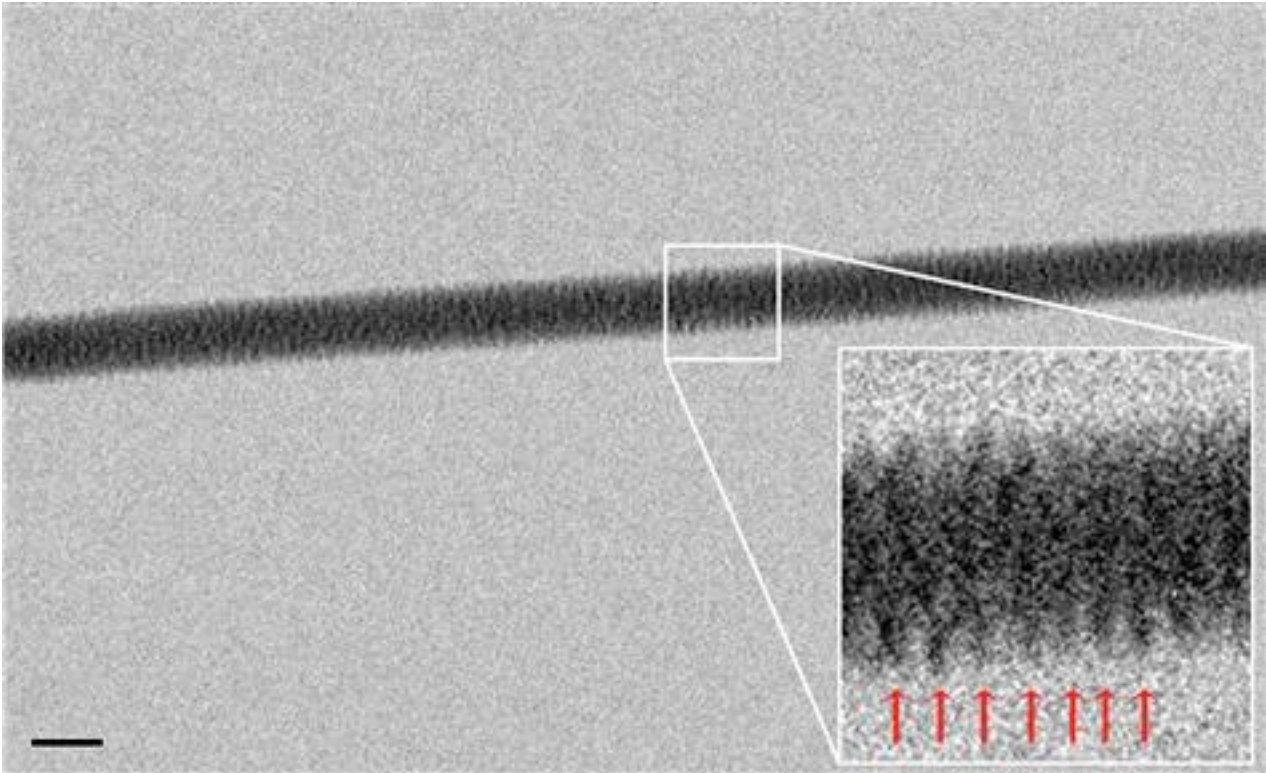
DNA Sequencing:

“Process of determining the precise order of nucleotides within a DNA molecule.”
-Wikipedia

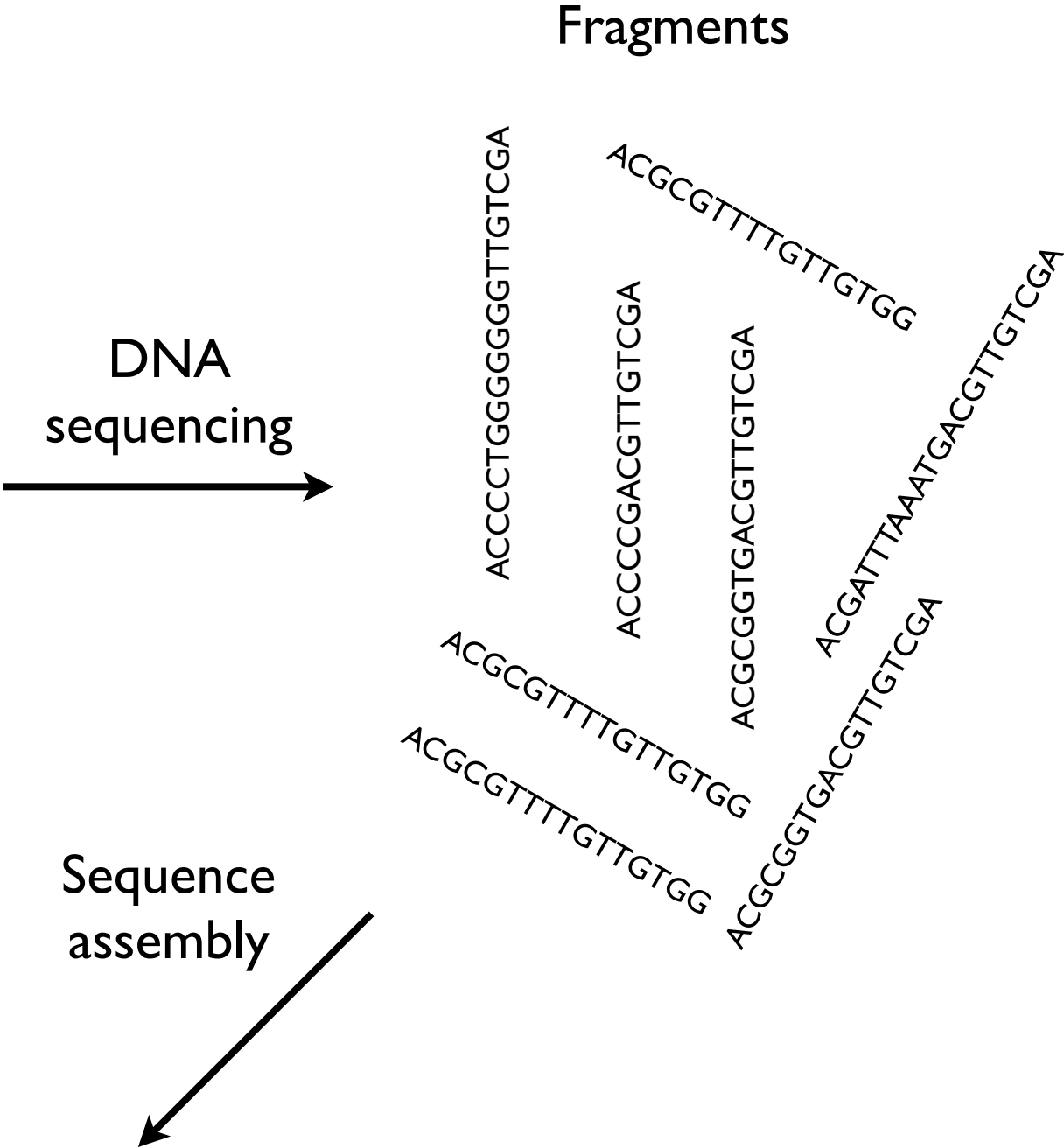


DNA Sanger Sequencing

1. Brief history about genome assembly



Gentile F. et al. *Direct Imaging of DNA Fibers: The Visage of Double Helix*
Nano Lett., 2012, 12 (12), pp 6453–6458



ACGCGTTTTGTTGTGGTGGCCACACCACGCAGTGACGGAGATAACGGCGAGAGCATGGACGGAGGATGAGGATGG



1. Brief history about genome assembly

Sequence assembly

=

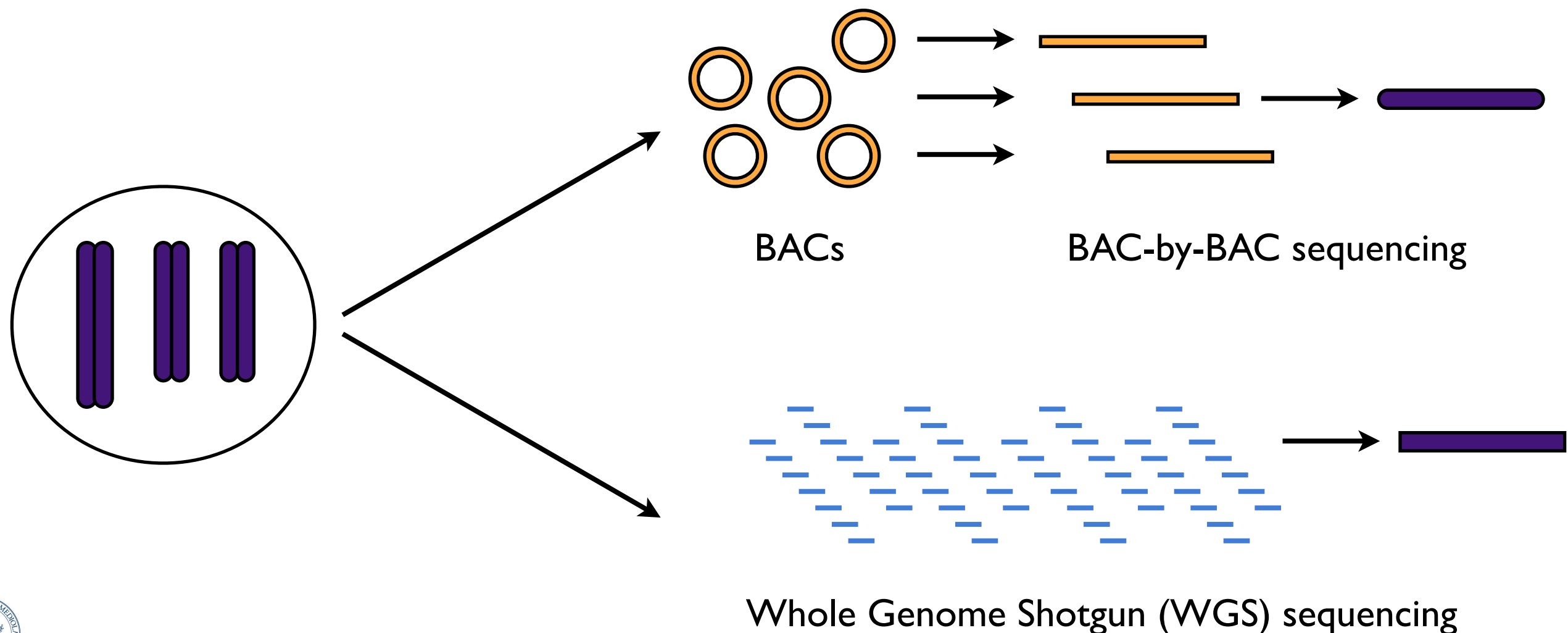
Resolve a puzzle and rebuild a DNA sequence from its pieces (fragments)



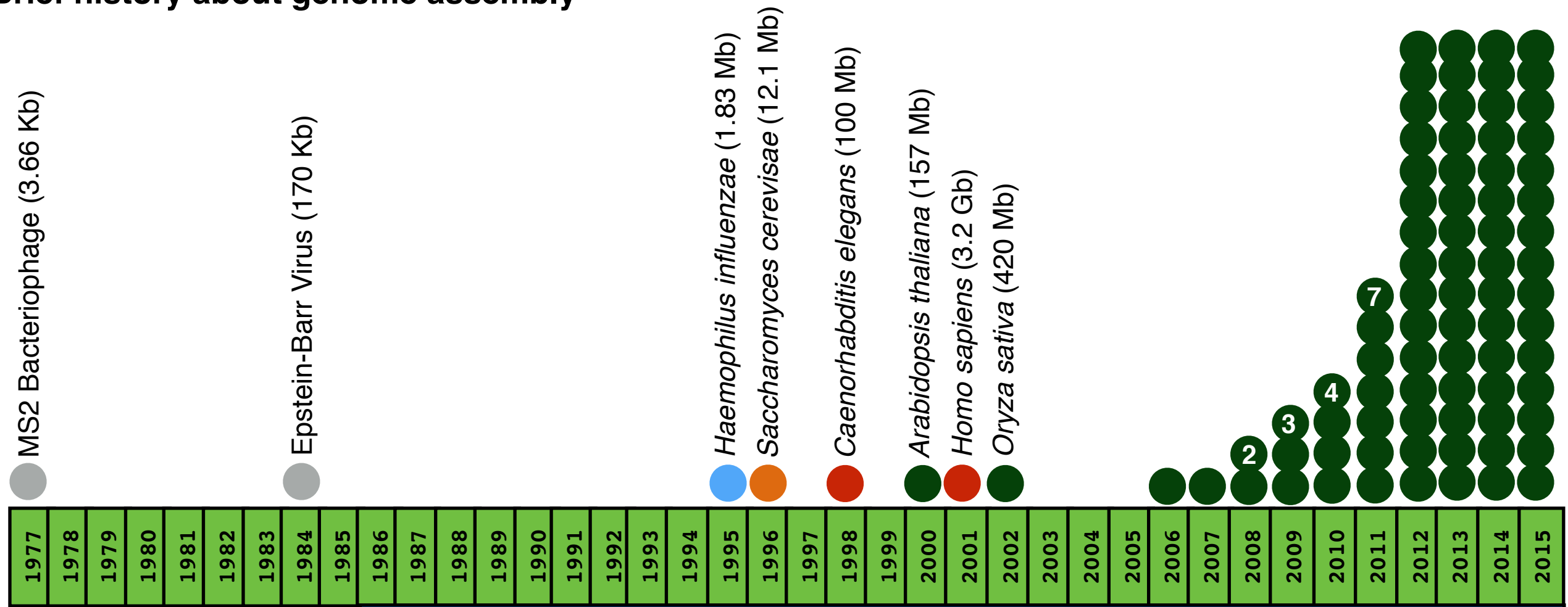
Image courtesy of iStock photo

1. Brief history about genome assembly

Resolve a puzzle and rebuild a DNA sequence from its pieces (fragments)



1. Brief history about genome assembly



Sanger

454

Solexa/Illumina

SOLiD

Ion Torrent

PacBio

OxN

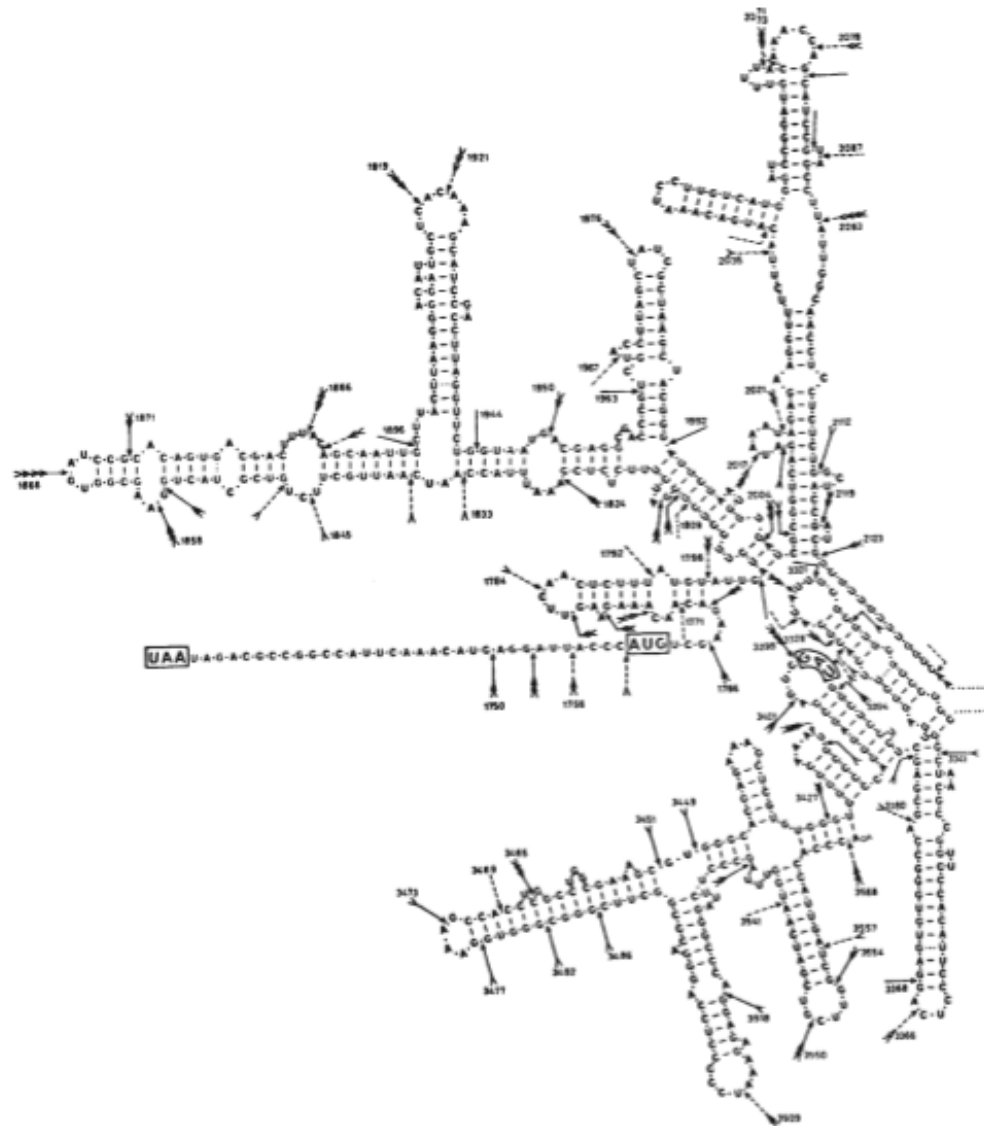
2020/03/22	Sequenced Genomes*
Plants	1,225
Animals	3,105
Fungi	5,801
Bacteria	244,944



1. Brief history about genome assembly

MS2 Bacteriophage (3.658 Kb)

1977



Complete nucleotide sequence of bacteriophage MS2 RNA: primary and secondary structure of the replicase gene

W. Fiers, R. Contreras, F. Duerinck, G. Haegeman, D. Iserentant, J. Merregaert,
W. Min Jou, F. Molemans, A. Raeymaekers, A. Van den Berghe, G. Volckaert & M. Ysebaert

Laboratory of Molecular Biology, University of Ghent, 9000 Ghent, Belgium



1. Brief history about genome assembly

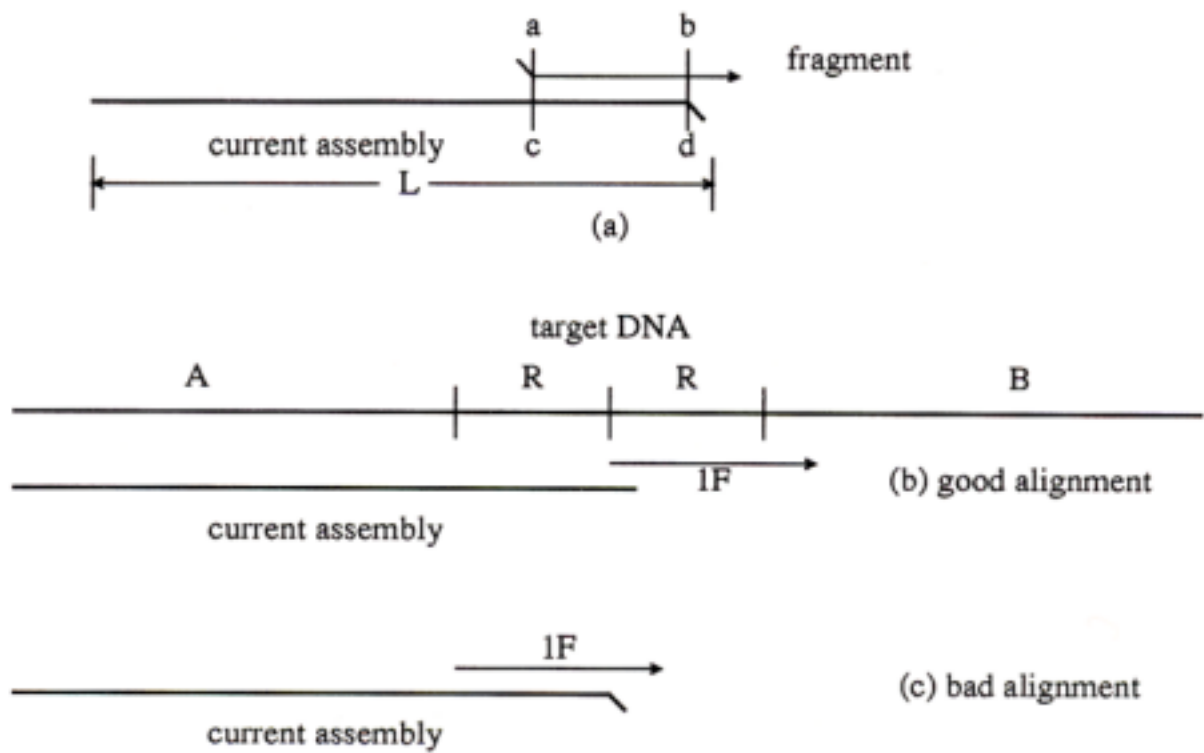
<i>Haemophilus influenzae</i> (1.83 Mb)	1995
---	------

Software:

TIGR ASSEMBLER



Smith-Waterman alignments



OVERLAPPING
METHODOLOGY

FIG. 1. **a.** A graphic representation of a Smith-Waterman alignment is shown. The overlap is shown between vertical bars, and the diagonal lines represent overhangs where the sequences do not match. The current assembly has a length of L. **b.** A tandem repeat of two copies of repeat region R is shown along with the proper alignment of fragment 1F with the current assembly. **c.** A bad alignment of the current assembly and fragment 1F is produced if the overlap is maximized without regard to the length of the overhang. The bad alignment can result in two outcomes. If the overhang is short, it will be ignored, and the two repeat regions will be compressed into a single region. If the overhang is long, the merge will not be allowed, and the current assembly will not be extended.

Sutton GG. et al.TIGR Assembler:A New Tool to Assembly Large Shotgun Sequencing Projects
Genome Science and Technology, 1995,1:9-19



1. Brief history about genome assembly

Homo sapiens (3.2 Gb)	2001
-----------------------	------

BAC-by-BAC sequencing

articles

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium*

** A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.*

International Human Genome Sequencing Consortium. Initial Sequencing and Analysis of the Human Genome. Nature. 2001. 409:860-921

Whole Genome Shotgun (WGS) sequencing

The Sequence of the Human Genome

J. Craig Venter,^{1*} Mark D. Adams,¹ Eugene W. Myers,¹ Peter W. Li,¹ Richard J. Mural,¹ Granger G. Sutton,¹ Hamilton O. Smith,¹ Mark Yandell,¹ Cheryl A. Evans,¹ Robert A. Holt,¹

Venter JC. et al. The Sequence of the Human Genome. Science. 2001. 291:1304-1351

OVERLAPPING
METHODOLOGY



1. Brief history about genome assembly

Homo sapiens (3.2 Gb)

2001

Whole Genome Shotgun (WGS) sequencing

THE HUMAN GENOME

The Sequence of the Human Genome

J. Craig Venter,^{1*} Mark D. Adams,¹ Eugene W. Myers,¹ Peter W. Li,¹ Richard J. Mural,¹
Granger G. Sutton,¹ Hamilton O. Smith,¹ Mark Yandell,¹ Cheryl A. Evans,¹ Robert A. Holt,¹

Venter JC. et al. The Sequence of the Human Genome. Science. 2001. 291:1304-1351

Software:

WGA ASSEMBLER (CABOG)

Hardware:

40 machines AlphaSMPs (4 Gb RAM/each and 4 cores/each, total=160 Gb RAM and 160 cores); 5 days.



1. Brief history about genome assembly

Ailuropoda melanoleura (2.3 Gb)

2009

Whole Genome Shotgun (WGS) sequencing

ARTICLES

The sequence and *de novo* assembly of the giant panda genome

Ruiqiang Li^{1,2*}, Wei Fan^{1*}, Geng Tian^{1,3*}, Hongmei Zhu^{1*}, Lin He^{4,5*}, Jing Cai^{3,6*}, Quanfei Huang¹, Qingle Cai^{1,7}, Bo Li¹, Yinqi Bai¹, Zhihe Zhang⁸, Yaping Zhang⁶, Wen Wang⁶, Jun Li¹, Fuwen Wei⁹, Heng Li¹⁰, Min Jian¹, Jianwen Li¹,

Li R. et al. The Sequence and the Novo Assembly of the Giant Panda Genome. Nature. 2009. 463:311-317

Software:

SOAPdenovo

Hardware:


Supercomputer with 32 cores and 512 Gb RAM.

BRUIJN GRAPHS
METHODOLOGY



1. Brief history about genome assembly

<https://www.ncbi.nlm.nih.gov/genome/>



Genome

This resource organizes information on genomes including sequences, maps, chromosomes, assemblies, and annotations.

Using Genome

[Help](#)

[Browse by Organism](#) **UPDATED**

[Download / FTP](#)

[Download FAQ](#)

[Submit a genome](#)

Custom resources

[Human Genome](#)

[Microbes](#)

[Organelles](#)

[Viruses](#)

[Prokaryotic reference genomes](#)

Other Resources

[Assembly](#)

[BioProject](#)

[BioSample](#)

[Genome Data Viewer](#) **NEW**

Genome Tools

[BLAST the Human Genome](#)

[Microbial Nucleotide BLAST](#)

Genome Annotation and Analysis

[Eukaryotic Genome Annotation](#)

[Prokaryotic Genome Annotation](#)

[PASC \(Pairwise Sequence Comparison\)](#)

External Resources

[GOLD - Genomes Online Database](#)

[Bacteria Genomes at Sanger](#)

[Ensembl](#)



1. Brief history about genome assembly

<https://www.yourgenome.org/facts/timeline-history-of-genomics>



Outline of Topics

1. Brief history about genome assembly
- 2. Basics about sequence assembly**
3. Whole genome assembly
 - 3.1. From reads to contigs and scaffolds.
 - 3.2. From scaffolds to chromosomes.
4. Transcriptome assembly



2. Basics about sequence assembly

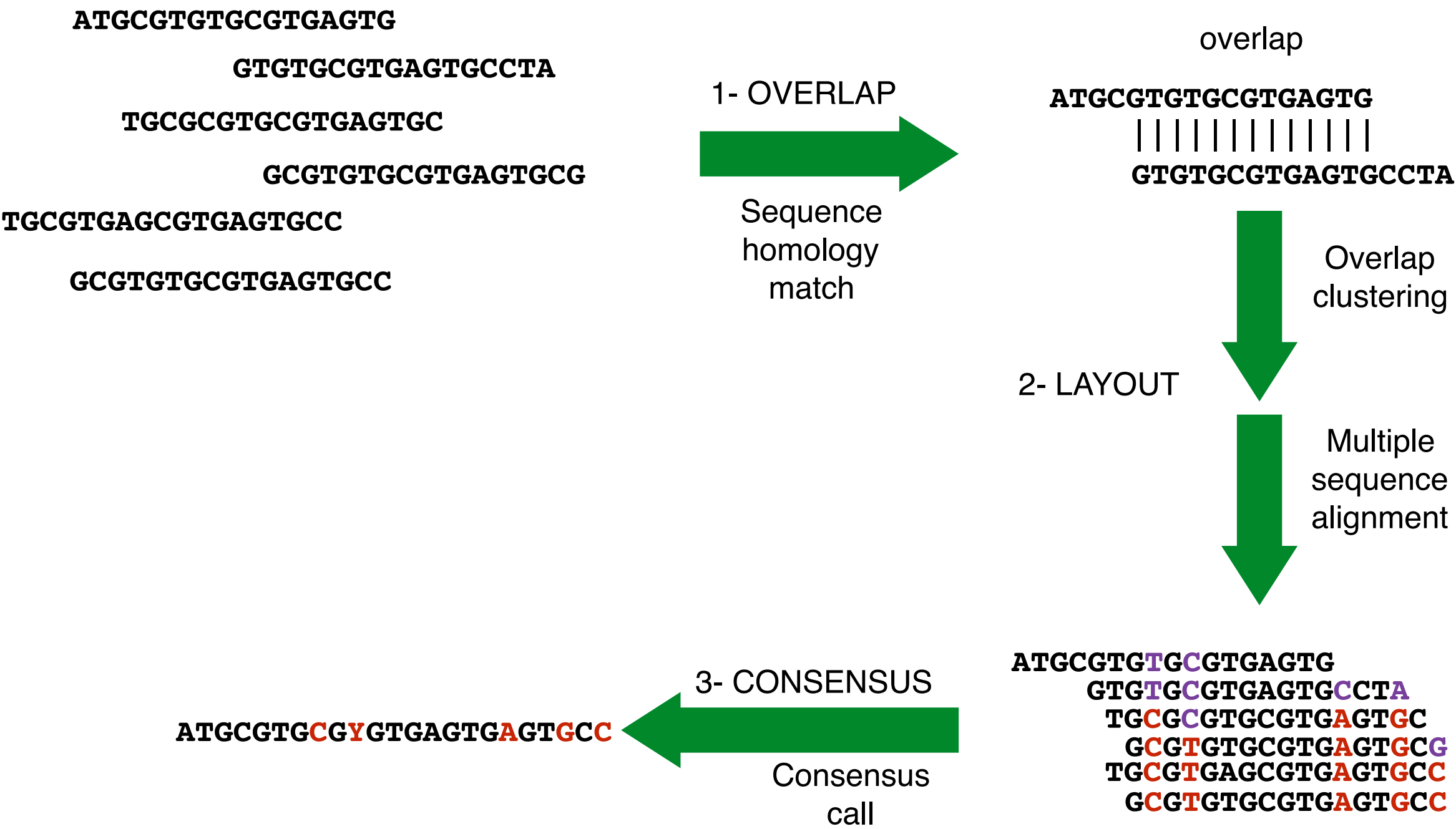
There are two sequence assembly main approaches:

- Overlap-Layout-Consensus (OLC) based in sequence alignments.
- De-Bruijn Graphs (DBG) based in Kmer decomposition and De-Bruijn graphs.



2. Basics about sequence assembly

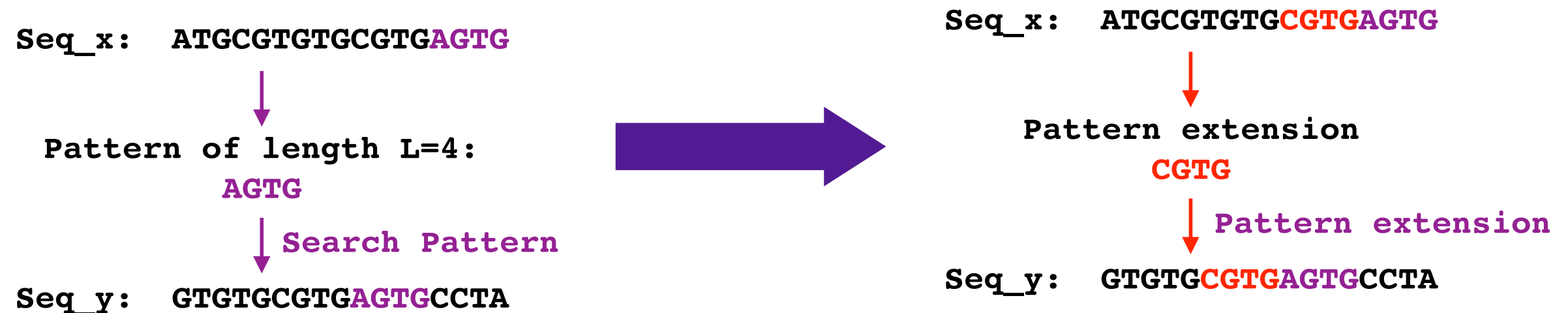
Overlap-Layout-Consensus (OLC)



2. Basics about sequence assembly

Overlap-Layout-Consensus (OLC)

1- Finding the Overlap



... and repeat with each pair

2 reads — 2 comparisons
4 reads — 16 comparisons
10 reads — 100 comparisons
1000000 reads — 1000000000000 comparisons
...



2. Basics about sequence assembly

Overlap-Layout-Consensus (OLC)

1- Finding the Overlap

1.1- Suffix tree

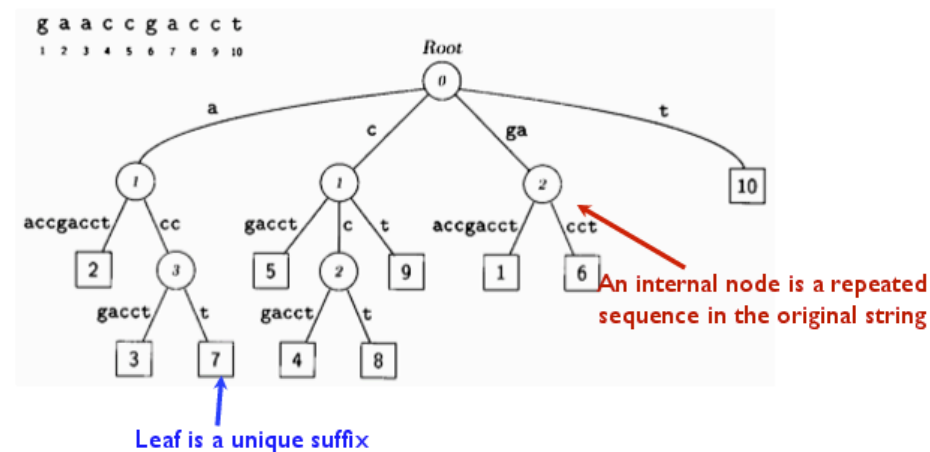
Suffix tree allows one to find, extremely efficiently, all distinct subsequences in a given sequence.

Suffix tree and suffix array for string matching

Suffix tree

- Definition: a suffix tree is a compressed trie containing all the suffixes of the given text as their keys and positions in the text as their values. See an example of suffix tree:

Suffix tree of a short sequence



2. Basics about sequence assembly

Overlap-Layout-Consensus (OLC)

1- Finding the Overlap

1.2- Dynamic Programming: Smith Waterman Algorithms

The Smith–Waterman algorithm performs local sequence alignment; that is, for determining **similar regions between two strings of nucleic acid sequences or protein sequences**. Instead of looking at the entire sequence, the Smith–Waterman algorithm compares segments of all possible lengths and optimizes the similarity measure.

The Smith–Waterman algorithm has four steps:

1. Determine the substitution matrix and the gap penalty scheme.
2. Initialize the scoring matrix.
3. Scoring.
4. Traceback.

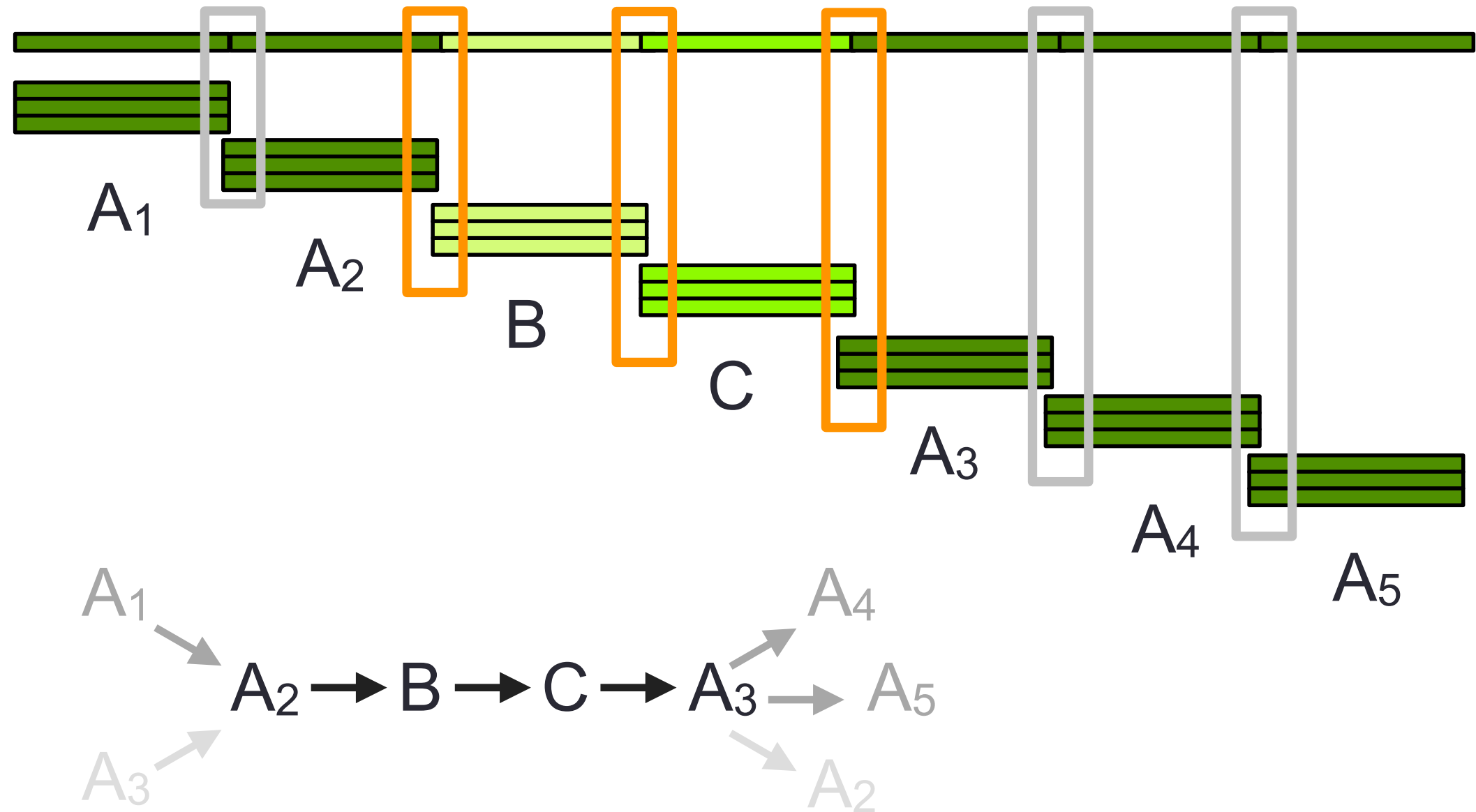
The Smith–Waterman algorithm example: <https://www.youtube.com/watch?v=QphFHG9tmOY>



2. Basics about sequence assembly

Overlap-Layout-Consensus (OLC)

2- Resolving the layout



2. Basics about sequence assembly

Overlap-Layout-Consensus (OLC)

3- Calling the consensus

ATGCGTGTCGTGAGTG
GTGTGCGTGAGTGCTA
TGCCTGTGCGTGAGTGC
TGCCTGTGCGTGAGTGCC
GCGTGTGCGTGAGTGCG
GCGTGTGCGTGAGTGCC

4 x C	3 x C	4 x A	4 x G	1 x G
2 x T	3 x T	1 x C	1 x A	2 x C

Consensus by
majority vote

ATGCGTGCGYGTGAGTGAGTGCC



2. Basics about sequence assembly

De-Bruijn Graphs (DBG)

What is a Kmer ?

Specific n-tuple or n-gram of nucleic acid or amino acid sequences.

-Wikipedia

ordered list
of elements

contiguous sequence
of n items from a given
sequence of text

ATGCGCAGTGGAGAGAGAGCGATG Sequence A with 25 nt

↓ 5 Kmers of 20-mer

ATGCGCAGTGGAGAGAGAGC

TGCGCAGTGGAGAGAGAGCG

GCGCAGTGGAGAGAGAGCGA

CGCAGTGGAGAGAGAGCGAT

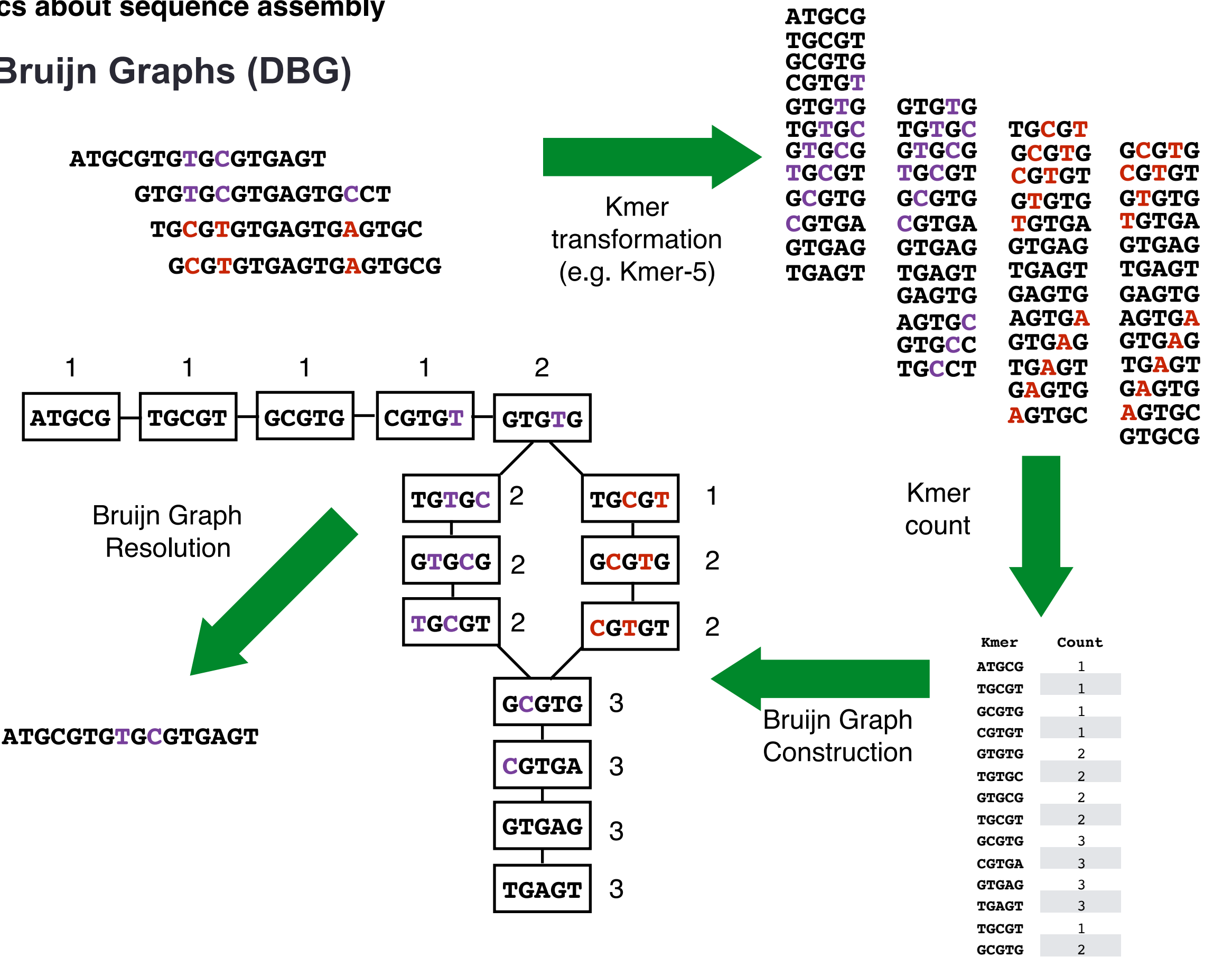
GCAGTGGAGAGAGAGCGATG

$N_kmers = L_read - Kmer_size$



2. Basics about sequence assembly

De-Bruijn Graphs (DBG)



2. Basics about sequence assembly

De-Bruijn Graphs (DBG)

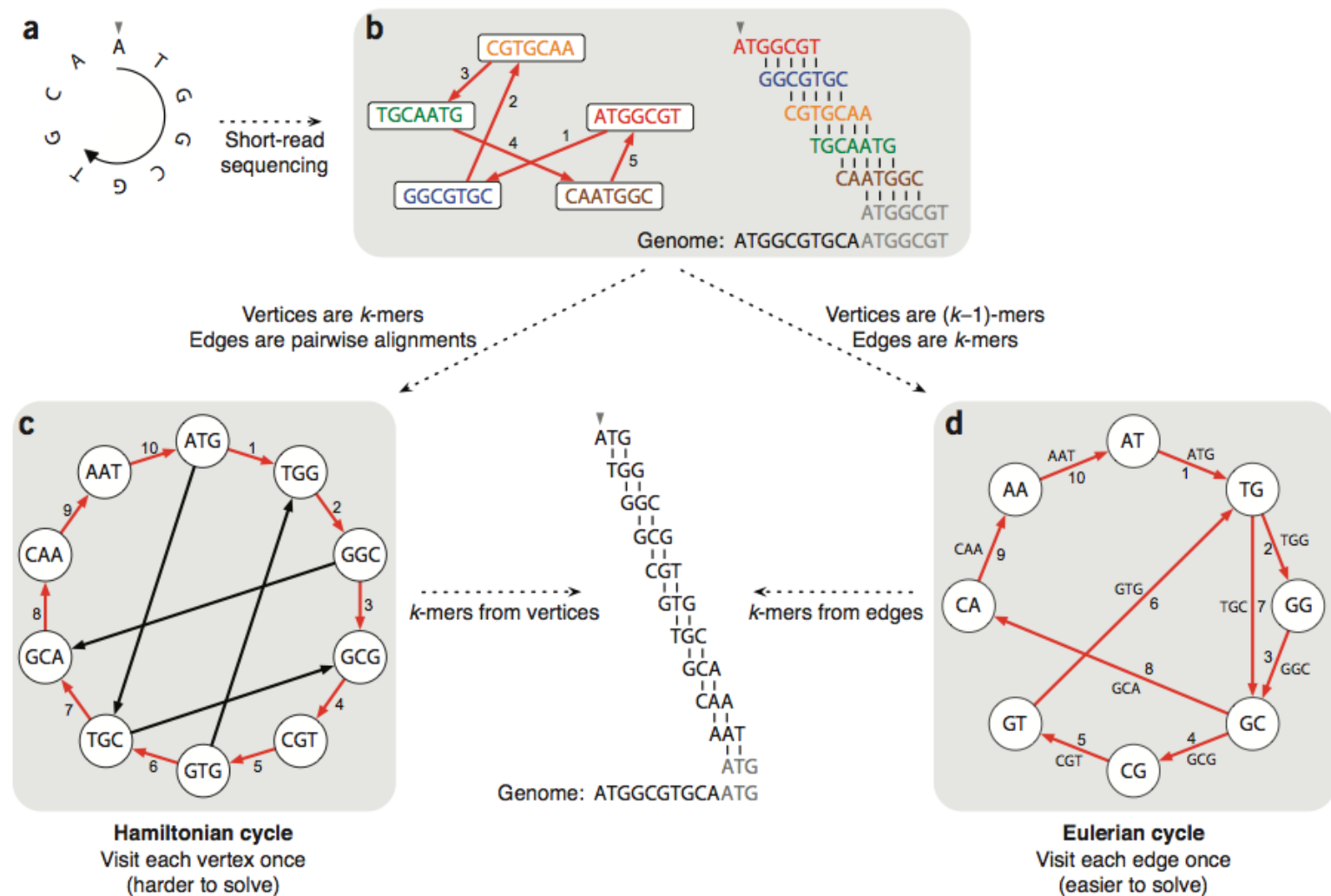


Figure 3 Two strategies for genome assembly: from Hamiltonian cycles to Eulerian cycles. **(a)** An example small circular genome. **(b)** In traditional Sanger sequencing algorithms, reads were represented as nodes in a graph, and edges represented alignments between reads. Walking along a Hamiltonian cycle by following the edges in numerical order allows one to reconstruct the circular genome by combining alignments between successive reads. At the end of the cycle, the sequence wraps around to the start of the genome. The repeated part of the sequence is grayed out in the alignment diagram. **(c)** An alternative assembly technique first splits reads into all possible k -mers: with $k = 3$, ATGGCGT comprises ATG, TGG, GGC, GCG and CGT. Following a Hamiltonian cycle (indicated by red edges) allows one to reconstruct the genome by forming an alignment in which each successive k -mer (from successive nodes) is shifted by one position. This procedure recovers the genome but does not scale well to large graphs. **(d)** Modern short-read assembly algorithms construct a de Bruijn graph by representing all k -mer prefixes and suffixes as nodes and then drawing edges that represent k -mers having a particular prefix and suffix. For example, the k -mer edge ATG has prefix AT and suffix TG. Finding an Eulerian cycle allows one to reconstruct the genome by forming an alignment in which each successive k -mer (from successive edges) is shifted by one position. This generates the same cyclic genome sequence without performing the computationally expensive task of finding a Hamiltonian cycle.



Comparison of the two major classes of assembly algorithms: overlap–layout–consensus and de-bruijn-graph

Zhenyu Li*, Yanxiang Chen*, Desheng Mu*, Jianying Yuan, Yujian Shi, Hao Zhang, Jun Gan, Nan Li, Xuesong Hu, Binghang Liu, Bicheng Yang and Wei Fan

Advance Access publication date 19 December 2011

OLC (Overlap-layout-consensus) algorithm is more suitable for the low-coverage long reads, whereas the DBG (De-Bruijn-Graph) algorithm is more suitable for high-coverage short reads and especially for large genome assembly

Key Points

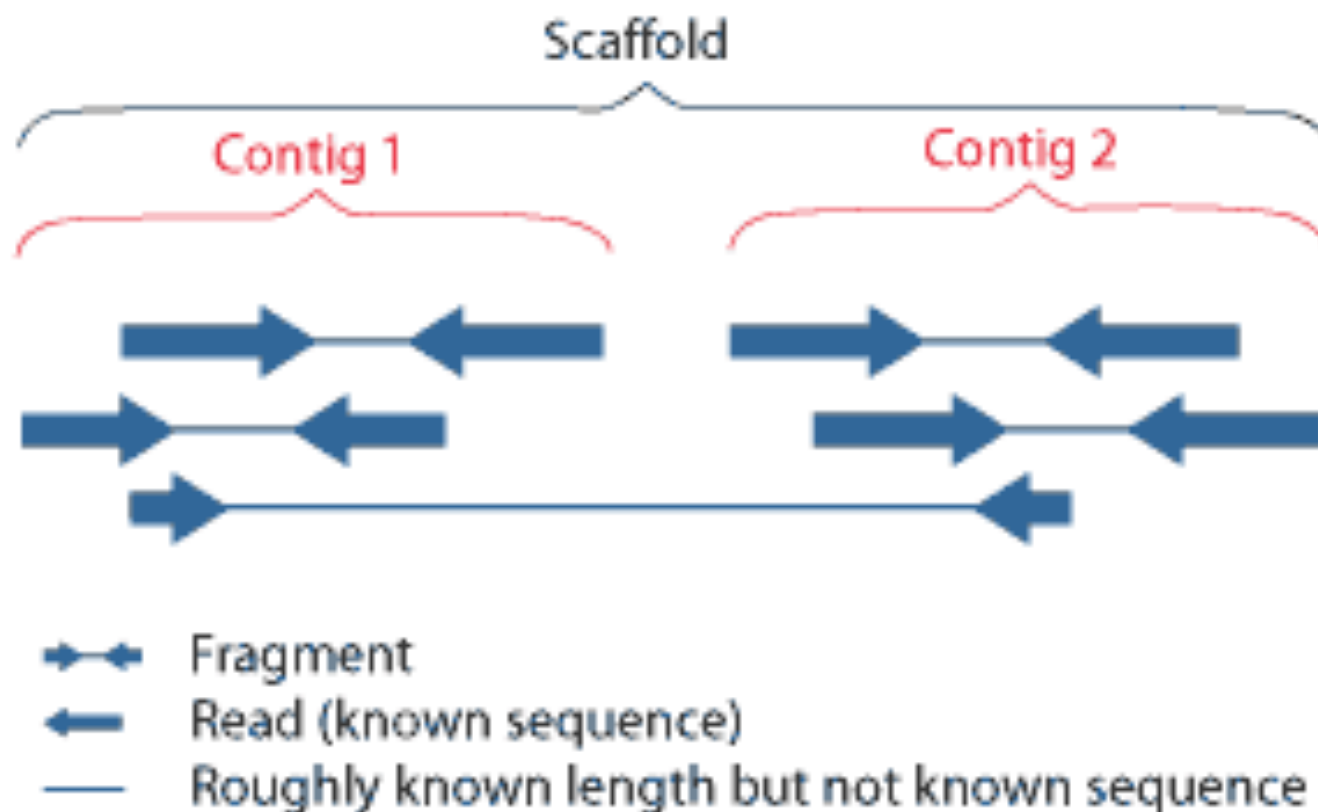
- High-quality genome sequences for many species are still strongly desired by the genomics community. With the rapid development of sequencing technologies and assembly algorithms, we have seen practical improvements and a bright future lies ahead.
- There are two major types of assembly algorithms: OLC and DBG; both of them are in accordance with Lander–Waterman model, but suit the assembly of different read lengths and sequencing depths, and have significant differences in computational efficiency.
- How well a genome can be assembled depends not only on sequencing technologies such as read length and sequencing error rate, but also on the characteristics of the genome, including repeat and the heterozygosity rate of the sequenced sample.



2. Basics about sequence assembly

What is a Scaffold?

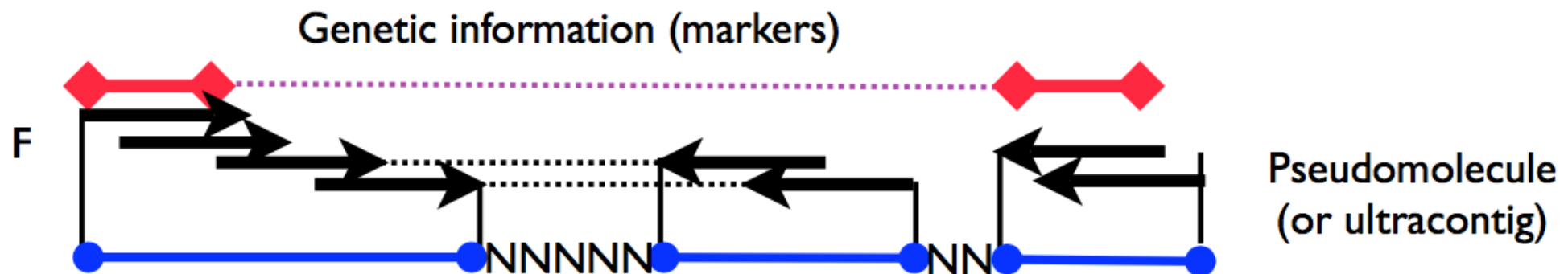
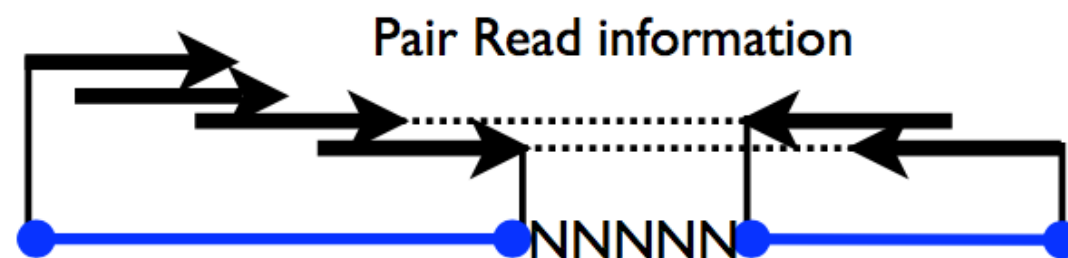
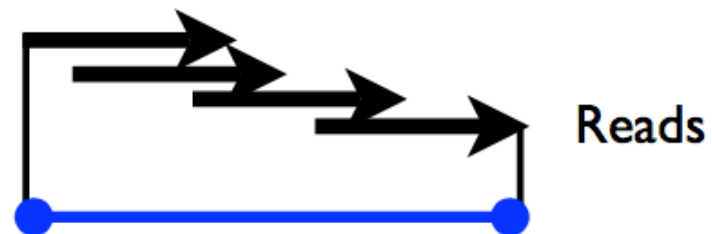
A scaffold is a portion of the genome sequence reconstructed from end-sequenced whole-genome shotgun clones. Scaffolds are composed of contigs and gaps. A contig is a contiguous length of genomic sequence in which the order of bases is known to a high confidence level. Gaps occur where reads from the two sequenced ends of at least one fragment overlap with other reads in two different contigs (as long as the arrangement is otherwise consistent with the contigs being adjacent). Since the lengths of the fragments are roughly known, the number of bases between contigs can be estimated.



2. Basics about sequence assembly

What is a Scaffold?

- Why is important the pair information ?
 - *novο* assembly:



Outline of Topics

1. Brief history about genome assembly
2. Basics about sequence assembly
- 3. Whole genome assembly**
 - 3.1. From reads to contigs and scaffolds.
 - 3.2. From scaffolds to chromosomes.
4. Transcriptome assembly



3. Whole genome assembly

Decisions during the Experimental Design

Technology

Library Preparation

Sequencing Amount

Fastq raw

Reads
processing
and
filtering

1. Low quality reads (qscore) (Q30)
2. Short reads (L50)
3. PCR duplications (Only Genomes).
4. Contaminations.
5. Corrections

Fastq Processed

Consensus

Contigs

Scaffolding

Scaffolds

Gap Filling

Scaffolds

Long distance
scaffolding

Chromosomes

Decisions during the Assembly Optimization

Software

Identity %, Kmer ...

Post-assembly Filtering



3. Whole genome assembly

Decisions during the Experimental Design

Technology

Library Preparation

Sequencing Amount

3. Whole genome assembly

Technologies

Technology	Read length (bp)	Accuracy	Reads/Run	Time/Run	Cost/Mb
Applied Bio 3730XL (Sanger)	400 - 900	99.9%	384	4 h (12 runs/day)	US\$2,400
Roche 454 GS FLX (Pyrosequencing)	700 Single/Pairs	99.9%	1,000,000	24h	US\$10
Illumina HiSeq2500 (Seq. by synthesis)	50-250 Single/Pairs	99%	4,000,000,000	24 to 120 h	\$0.05 to \$0.15
Illumina MiSeq (Seq. by synthesis)	50-300 Single/Pairs	99%	44,000,000	24 to 72 h	US\$0.17
SOLiD 4 (Seq. by ligation)	25-50 Single/Pairs	99.9%	1,400,000,000	168 h	US\$0.13
ION Torrent (Seq. by semiconductor)	170-400 Single	98%	80,000,000	2 h	US\$2
Pacific Biosciences RSII (SMRT)	10,000 Single	85% (99.9%)	750,000	4 h	US\$0.6
Oxford N. Minion (Nanopore sequencing)	10,000 Single	62% (96%)	4,400,000	48 h	US\$0.02
10X Genomics	20,000 Single	99%	1,000,000,000	NA	NA

3. Whole genome assembly

Libraries

★ Library types (orientations):

- Single reads



- Pair ends (PE) (150-800 bp insert size)



Illumina

- Mate pairs (MP) (2-40 Kb insert size)



Illumina



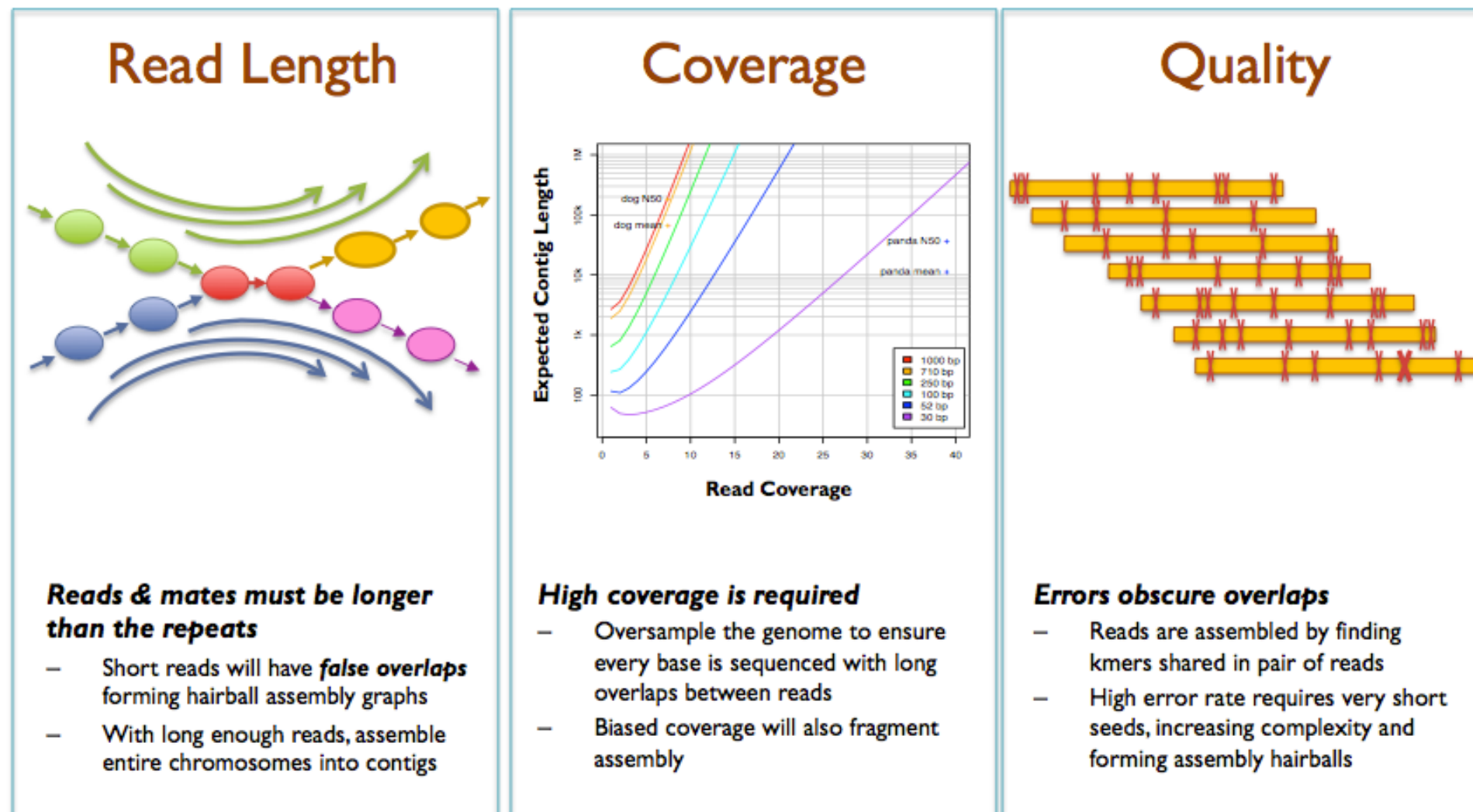
454/Roche



3. Whole genome assembly

Sequencing Amount

Ingredients for a good assembly



Current challenges in de novo plant genome sequencing and assembly

Schatz MC, Witkowski, McCombie, VWR (2012) *Genome Biology*. 12:243

3. Whole genome assembly

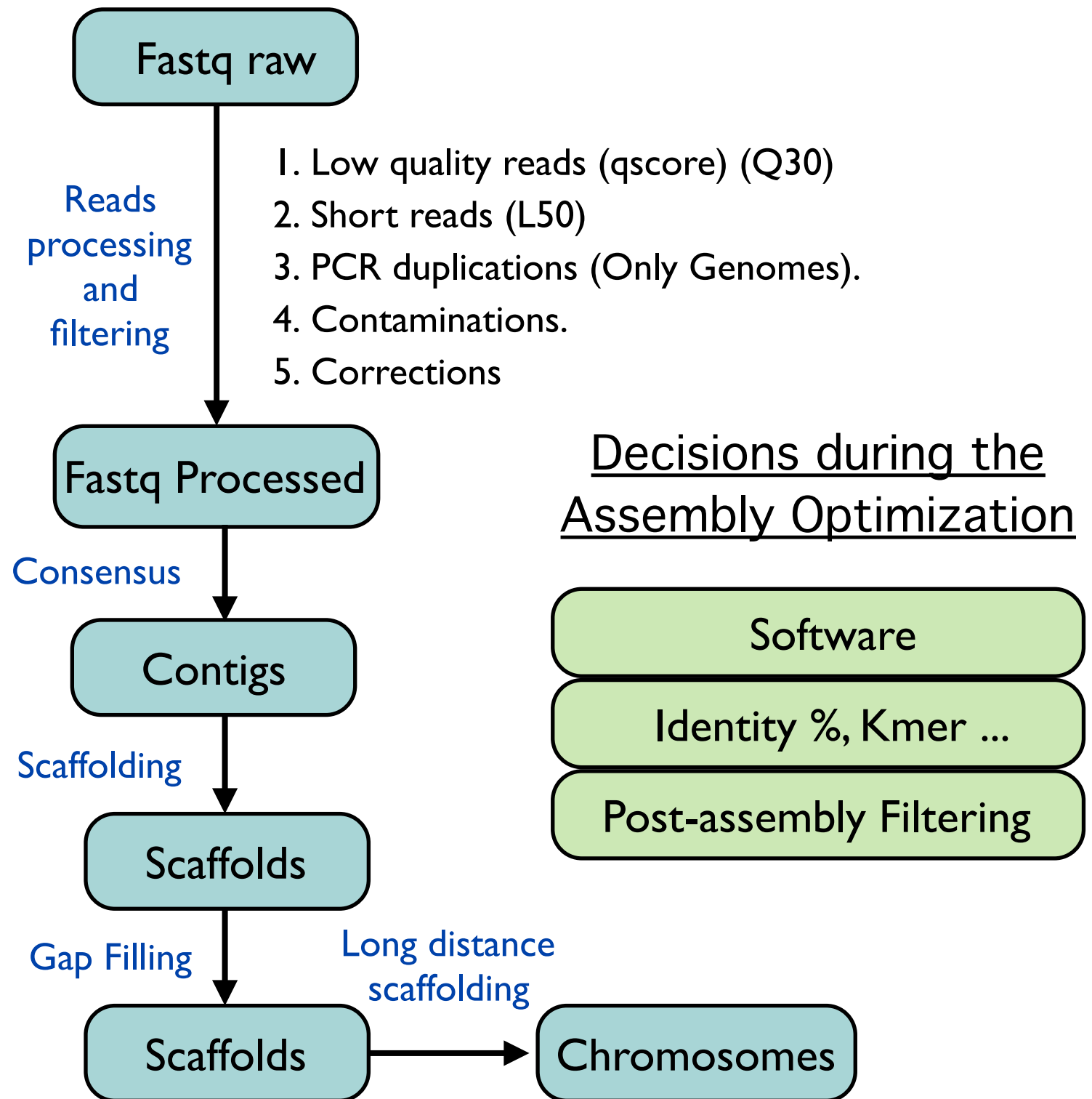
Sequencing Amount

Depending of the genome complexity, technology used and assembler:

- More is better (if you have enough computational resources).
 - Sanger > 10X (less for BACs-by-BACs approaches).
 - 454 > 20X
 - Illumina > 100X
 - PacBio > 20X (corrected by Illumina) or 50X (corrected PacBio)
 - 10X Genomics > 20X
- Polyploidy or high heterozygosity increase the amount of reads needed.
- The use of different library types (pair ends and mate pairs with different insert sizes is essential).
- Longer reads is preferable.



3. Whole genome assembly

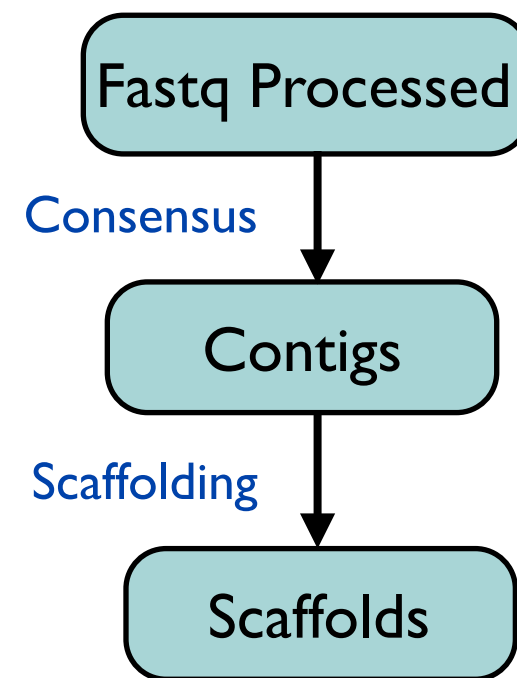


Outline of Topics

1. Brief history about genome assembly
2. Basics about sequence assembly
3. Whole genome assembly
 - 3.1. From reads to contigs and scaffolds.
 - 3.2. From scaffolds to chromosomes.
4. Transcriptome assembly



3.1. From reads to contigs and scaffolds.



Decisions during the
Assembly Optimization

Software



3.1. From reads to contigs and scaffolds.

Tools: Assemblers

	Type	Technology Used	Features	Link
Arachne	Overlap-layout-consensus	Sanger, 454	Highly configurable	http://www.broadinstitute.org/crd/wiki/index.php/Main_Page
CABOG	Overlap-layout-consensus	Sanger, 454, Illumina, PacBio	Highly configurable	http://sourceforge.net/apps/mediawiki/wgs-assembler/index.php?title=Main_Page
MIRA	Overlap-layout-consensus	Sanger, 454	Highly configurable	http://sourceforge.net/apps/mediawiki/mira-assembler
gsAssembler	Overlap-layout-consensus	Sanger, 454	Easy to use	http://454.com/products/analysis-software/index.asp
iAssembler	Overlap-layout-consensus	Sanger, 454	Improves MIRA	http://bioinfo.bti.cornell.edu/tool/iAssembler
ABYSS	Bruijn graph	454 or Illumina	Easy to use	http://www.bcgsc.ca/platform/bioinfo/software/abyss
ALLPATH-LG	Bruijn graph	454 or Illumina	Good results	http://www.broadinstitute.org/software/allpaths-lg/blog
Ray	Bruijn graph	454 or Illumina	Slow but use less memory	http://denovoassembler.sf.net/
SOAPdenovo2	Bruijn graph	454 or Illumina	Fastest	http://soap.genomics.org.cn/soapdenovo.html
Velvet	Bruijn graph	454 or Illumina or SOLiD	SOLiD	http://www.ebi.ac.uk/~zerbino/velvet/
Minia	Bloom filter + Brujn graph	Illumina	Really fast Only contigs	https://github.com/GATB/minia

3.1. From reads to contigs and scaffolds.

Tools: Assemblers

... but there are more assemblers and information... Take a look to SeqAnswers

Also highly recommendable:



Assemblathon 1: A competitive assessment of de novo short read assembly methods

Dent Earl, Keith Bradnam, John St. John, et al.

Genome Res. 2011 21: 2224-2241 originally published online September 16, 2011
Access the most recent version at doi:[10.1101/gr.126599.111](https://doi.org/10.1101/gr.126599.111)

GAGE: A critical evaluation of genome assemblies and assembly algorithms

Steven L. Salzberg, Adam M. Phillippy, Aleksey Zimin, et al.

Genome Res. 2012 22: 557-567 originally published online December 6, 2011
Access the most recent version at doi:[10.1101/gr.131383.111](https://doi.org/10.1101/gr.131383.111)



3.1. From reads to contigs and scaffolds.

Tools: Assemblers

	Type	Technology Used	Features	Link
Arachne	Overlap-layout-consensus	Sanger, 454	Highly configurable	http://www.broadinstitute.org/crd/wiki/index.php/Main_Page
CABOG	Overlap-layout-consensus	Sanger, 454, Illumina, PacBio	Highly configurable	http://sourceforge.net/apps/mediawiki/wgs-assembler/index.php?title=Main_Page
MIRA	Overlap-layout-consensus	Sanger, 454	Highly configurable	http://sourceforge.net/apps/mediawiki/mira-assembler
gsAssembler	Overlap-layout-consensus	Sanger, 454	Easy to use	http://454.com/products/analysis-software/index.asp
iAssembler	Overlap-layout-consensus	Sanger, 454	Improves MIRA	http://bioinfo.bti.cornell.edu/tool/iAssembler
ABYSS	Bruijn graph	454 or Illumina	Easy to use	http://www.bcgsc.ca/platform/bioinfo/software/abyss
ALLPATH-LG	Bruijn graph	454 or Illumina	Good results	http://www.broadinstitute.org/software/allpaths-lg/blog
Ray	Bruijn graph	454 or Illumina	Slow but use less memory	http://denovoassembler.sf.net/
SOAPdenovo2	Bruijn graph	454 or Illumina	Fastest	http://soap.genomics.org.cn/soapdenovo.html
Velvet	Bruijn graph	454 or Illumina or SOLiD	SOLiD	http://www.ebi.ac.uk/~zerbino/velvet/
Minia	Bloom filter + Brujn graph	Illumina	Really fast Only contigs	https://github.com/GATB/minia

3.1. From reads to contigs and scaffolds.

Tools: Assemblers

	Type	Technology Used	Features	Link
HGAP	Overlap-layout-consensus	PacBio	Recommended by PacBio. Multiple tools. Difficult to install	https://github.com/PacificBiosciences/Bioinformatics-Training/wiki/HGAP
Falcon	Overlap-layout-consensus	PacBio	Faster than HGAP or CABOG	https://github.com/PacificBiosciences/falcon
Sprai	Overlap-layout-consensus	PacBio	Easy to use interface for CABOG	http://zombie.cb.k.u-tokyo.ac.jp/sprai/README.html
Canu	Overlap-layout-consensus	PacBio, Oxford Nanopore	Easy to use	http://canu.readthedocs.org/
Miniasm	Overlap-layout-consensus	PacBio, Oxford Nanopore	Easy to use	https://github.com/lh3/miniasm
MECAT	Overlap-layout-consensus	PacBio, Oxford Nanopore	Fast assembler	https://github.com/xiaochuanle/MECAT
Spades	Brujn graphs	Illumina, Ion Torrent, PacBio, Oxford Nanopore	Diverse type of input	http://cab.spbu.ru/software/spades/
MaSurCa	Hybrid	PacBio + Illumina	Hybrid assembly	http://www.genome.umd.edu/masurca.html
Supernova	Hybrid	10X Genomics	Used with 10X genomics	https://support.10xgenomics.com/de-novo-assembly/software/pipelines/latest/using/running

3.1. From reads to contigs and scaffolds.

Tools: Assemblers

	Type	Technology Used	Features	Link
HGAP	Overlap-layout-consensus	PacBio	Recommended by PacBio. Multiple tools. Difficult to install	https://github.com/PacificBiosciences/Bioinformatics-Training/wiki/HGAP
Falcon	Overlap-layout-consensus	PacBio	Faster than HGAP or CABOG	https://github.com/PacificBiosciences/falcon
Sprai	Overlap-layout-consensus	PacBio	Easy to use interface for CABOG	http://zombie.cb.k.u-tokyo.ac.jp/sprai/README.html
Canu	Overlap-layout-consensus	PacBio, Oxford Nanopore	Easy to use	http://canu.readthedocs.org/
Miniasm	Overlap-layout-consensus	PacBio, Oxford Nanopore	Easy to use	https://github.com/lh3/miniasm
MECAT	Overlap-layout-consensus	PacBio, Oxford Nanopore	Fast assembler	https://github.com/xiaochuanle/MECAT
Spades	Brujn graphs	Illumina, Ion Torrent, PacBio, Oxford Nanopore	Diverse type of input	http://cab.spbu.ru/software/spades/
MaSurCa	Hybrid	PacBio + Illumina	Hybrid assembly	http://www.genome.umd.edu/masurca.html
Supernova	Hybrid	10X Genomics	Used with 10X genomics	https://support.10xgenomics.com/de-novo-assembly/software/pipelines/latest/using/running

3.1. From reads to contigs and scaffolds.

Scaffolding approaches:

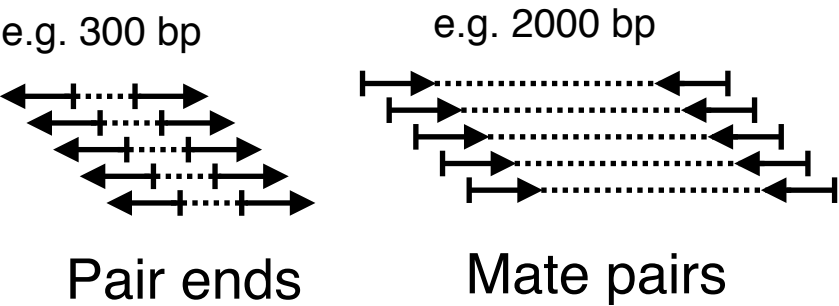
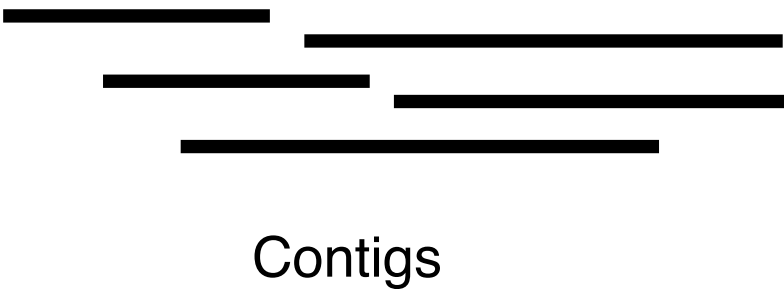
1. Pair ends/Mate pair information (Since 2008) - Short Distances
2. Optical mapping data (Since 2012) - Medium-Long Distance
3. HiC (Since 2014) - Long-Very Long Distance
4. Linkage groups anchoring (Since 2001) - Very Long Distance



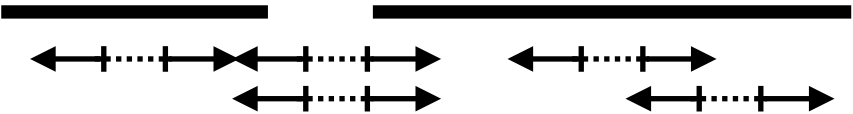
3.1. From reads to contigs and scaffolds.

I. Pair ends/Mate pair information

0. Input



1.1 Small fragments mapping



1.2. Scaffold building



2.1. Long fragments mapping



2.2. Scaffold building



3.1. From reads to contigs and scaffolds.

Tools: Scaffolders

	Technology Used	Link
SSPACE	Short reads, pairs	http://www.baseclear.com/genomics/bioinformatics/basetools/SSPACE
SSPACE-Long	Long reads	http://www.baseclear.com/genomics/bioinformatics/basetools/SSPACE-longread
Bambus2	Pairs	https://www.cbcu.edu/software/bambus2
Ragoo	Reference genome	https://github.com/malonge/RaGOO
PEP_scaffolder	Protein sequence	https://github.com/CAFS-bioinformatics/PEP_scaffolder



3.1. From reads to contigs and scaffolds.

Tools: Scaffolders

[Journal List](#) > [Genome Biol](#) > v.15(3); 2014 > [PMC4053845](#)

[this article](#)[search](#)[submit a manuscript](#)[register](#)

Genome Biol. 2014; 15(3): R42.

PMCID: PMC4053845

Published online 2014 Mar 3. doi: [10.1186/gb-2014-15-3-r42](https://doi.org/10.1186/gb-2014-15-3-r42)

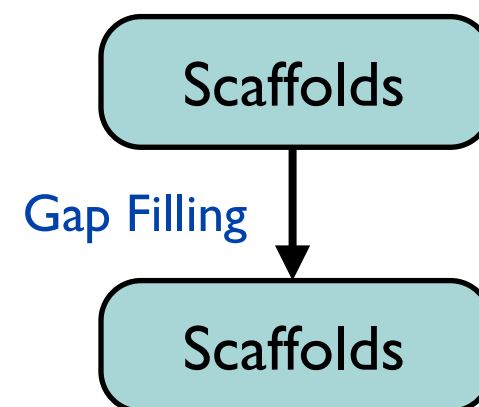
A comprehensive evaluation of assembly scaffolding tools

[Martin Hunt](#),¹ [Chris Newbold](#),^{2,1} [Matthew Berriman](#),¹ and [Thomas D Otto](#)¹

[Author information](#) ► [Article notes](#) ► [Copyright and License information](#) ►

This article has been [cited by](#) other articles in PMC.

3.1. From reads to contigs and scaffolds.



3.1. From reads to contigs and scaffolds.

Tools: Gap fillers

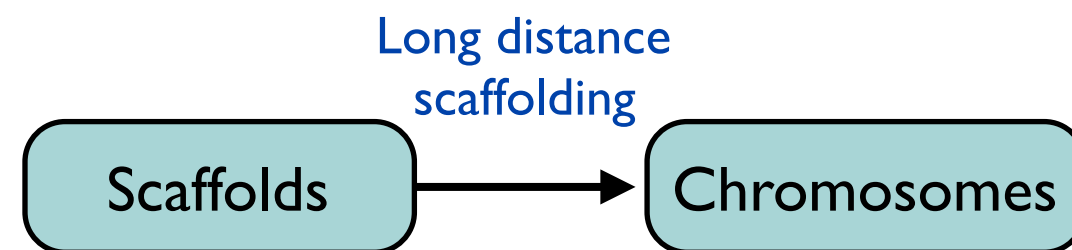
	Type	Technology Used	Features	Link
GapCloser	Overlap-layout-consensus	Illumina	C++ Free	http://soap.genomics.org.cn/soapdenovo.html
GapFiller	Overlap-layout-consensus	Any	Perl Commercial*	http://www.baseclear.com/landingpages/basetools-a-wide-range-of-bioinformatics-solutions/gapfiller/
PBSuite	Overlap-layout-consensus	454, PacBio	Python, Long reads	http://sourceforge.net/p/pb-jelly/wiki/

Outline of Topics

1. Brief history about genome assembly
2. Basics about sequence assembly
3. Whole genome assembly
 - 3.1. From reads to contigs and scaffolds.
 - 3.2. From scaffolds to chromosomes.**
4. Transcriptome assembly



3.2. From scaffolds to chromosomes.



3.2. From scaffolds to chromosomes.

Scaffolding approaches:

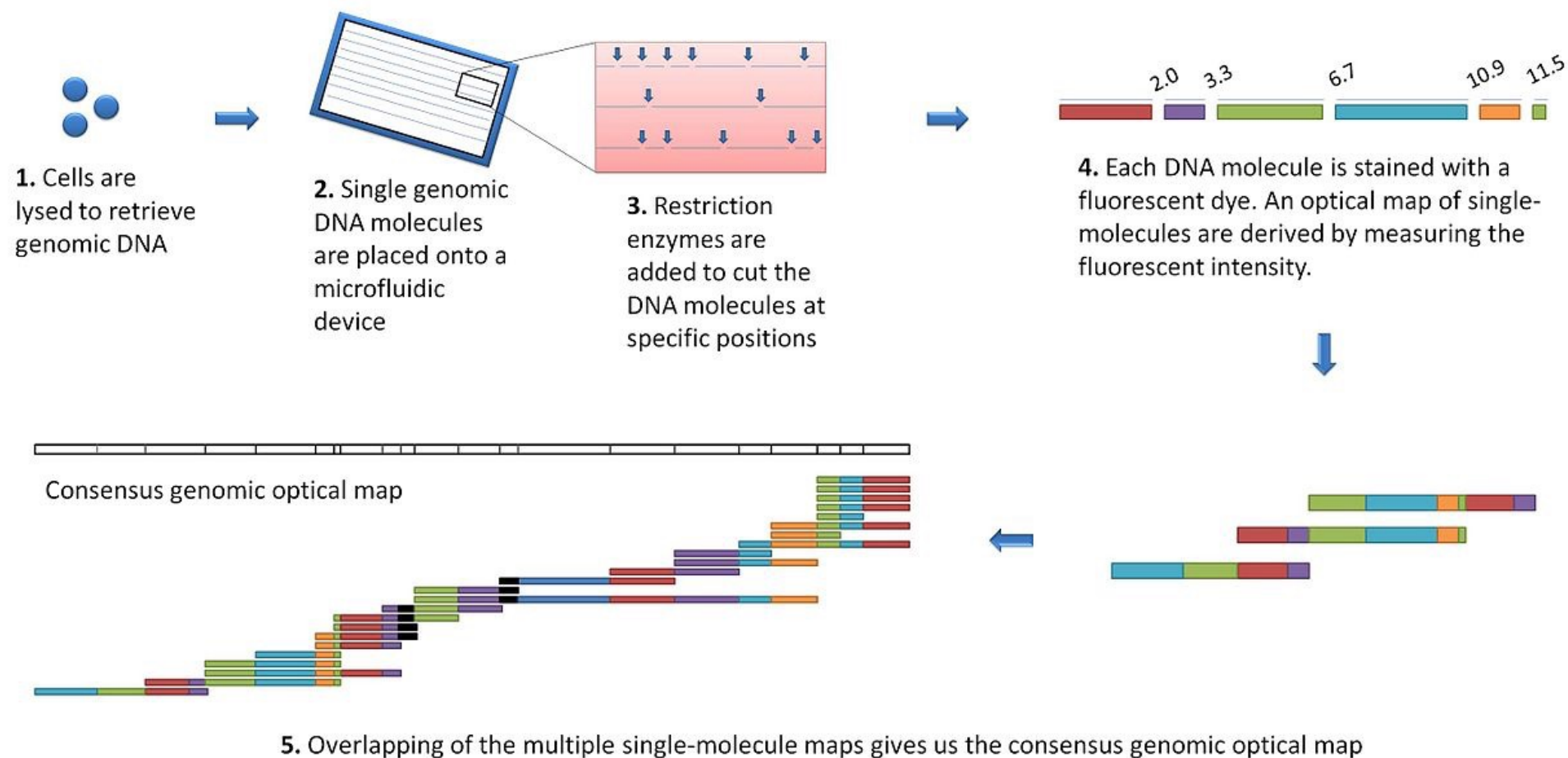
1. Pair ends/Mate pair information (Since 2008) - Short Distances
2. Optical mapping data (Since 2012) - Medium-Long Distance
3. HiC (Since 2014) - Long-Very Long Distance
4. Linkage groups anchoring (Since 2001) - Very Long Distance



3.2. From scaffolds to chromosomes.

2. Optical mapping data (Since 2012).

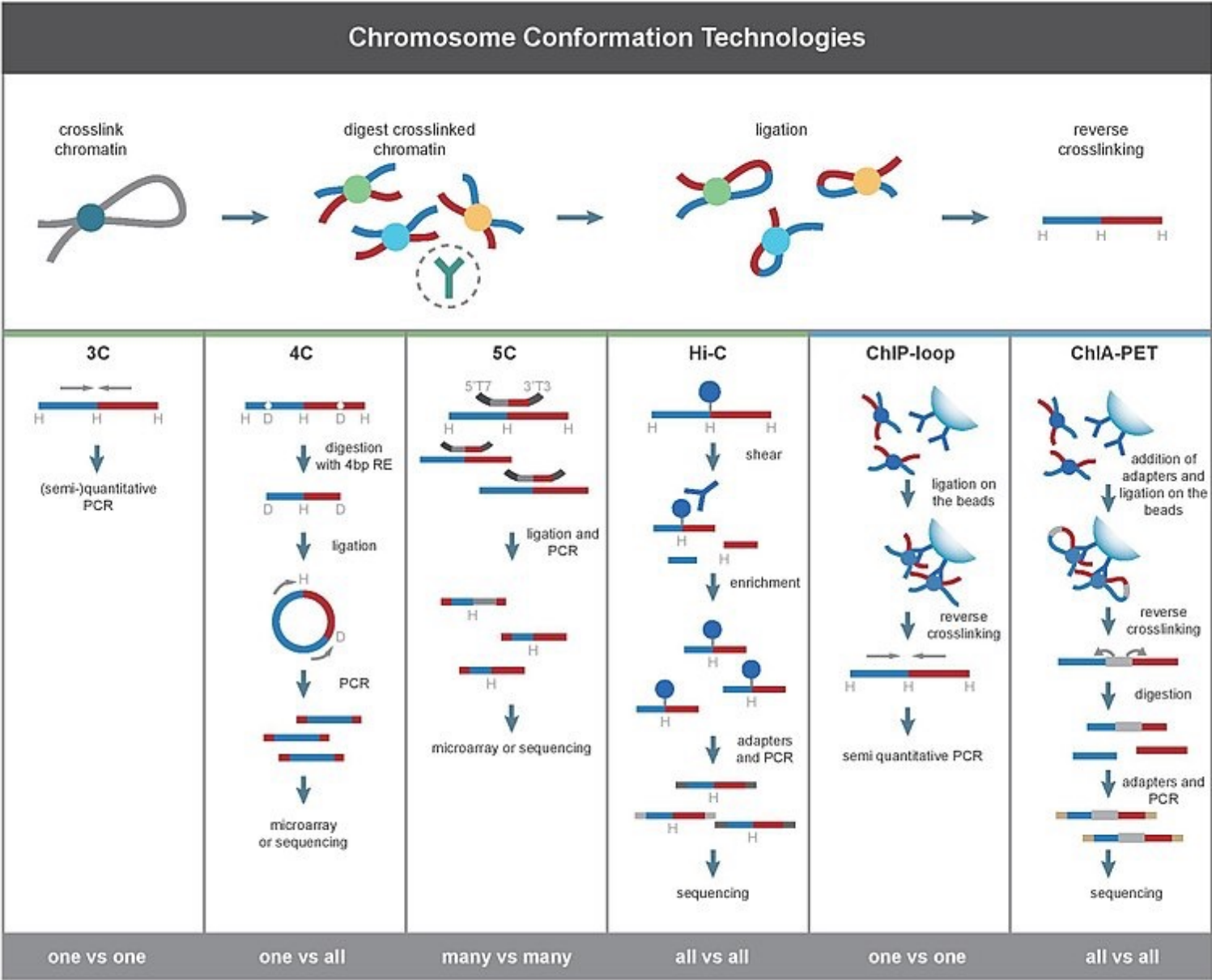
Optical mapping is a technique for constructing ordered, genome-wide, **high-resolution restriction maps from single, stained molecules of DNA**, called "optical maps". By mapping the location of restriction enzyme sites along the unknown DNA of an organism, the spectrum of resulting DNA fragments collectively serves as a unique "**fingerprint**" or "**barcode**" for that sequence



3.2. From scaffolds to chromosomes.

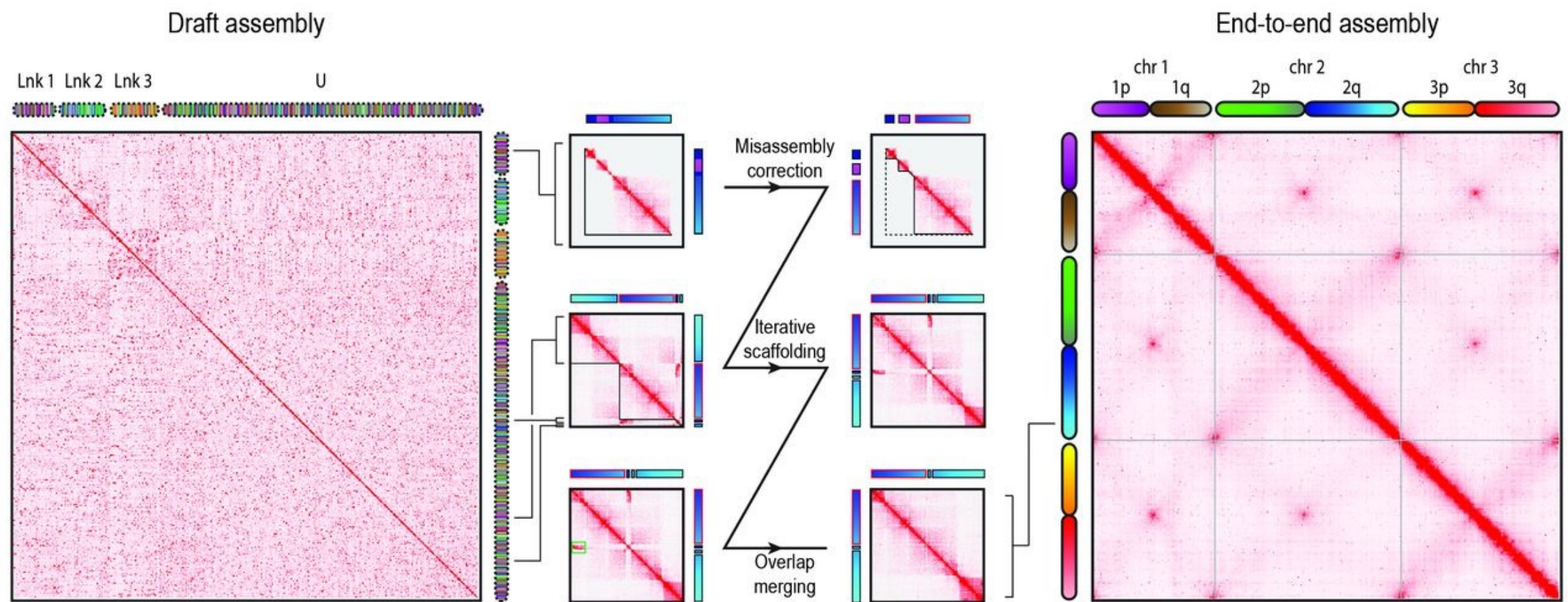
3. HiC, Chromatin Conformation Capture (Since 2014).

Chromosome conformation capture techniques (often abbreviated to 3C technologies or 3C-based methods) are a set of molecular biology methods used to analyze the **spatial organization of chromatin** in a cell. These methods quantify the number of interactions between genomic loci that are nearby in 3-D space, but may be separated by many nucleotides in the linear genome



3.2. From scaffolds to chromosomes.

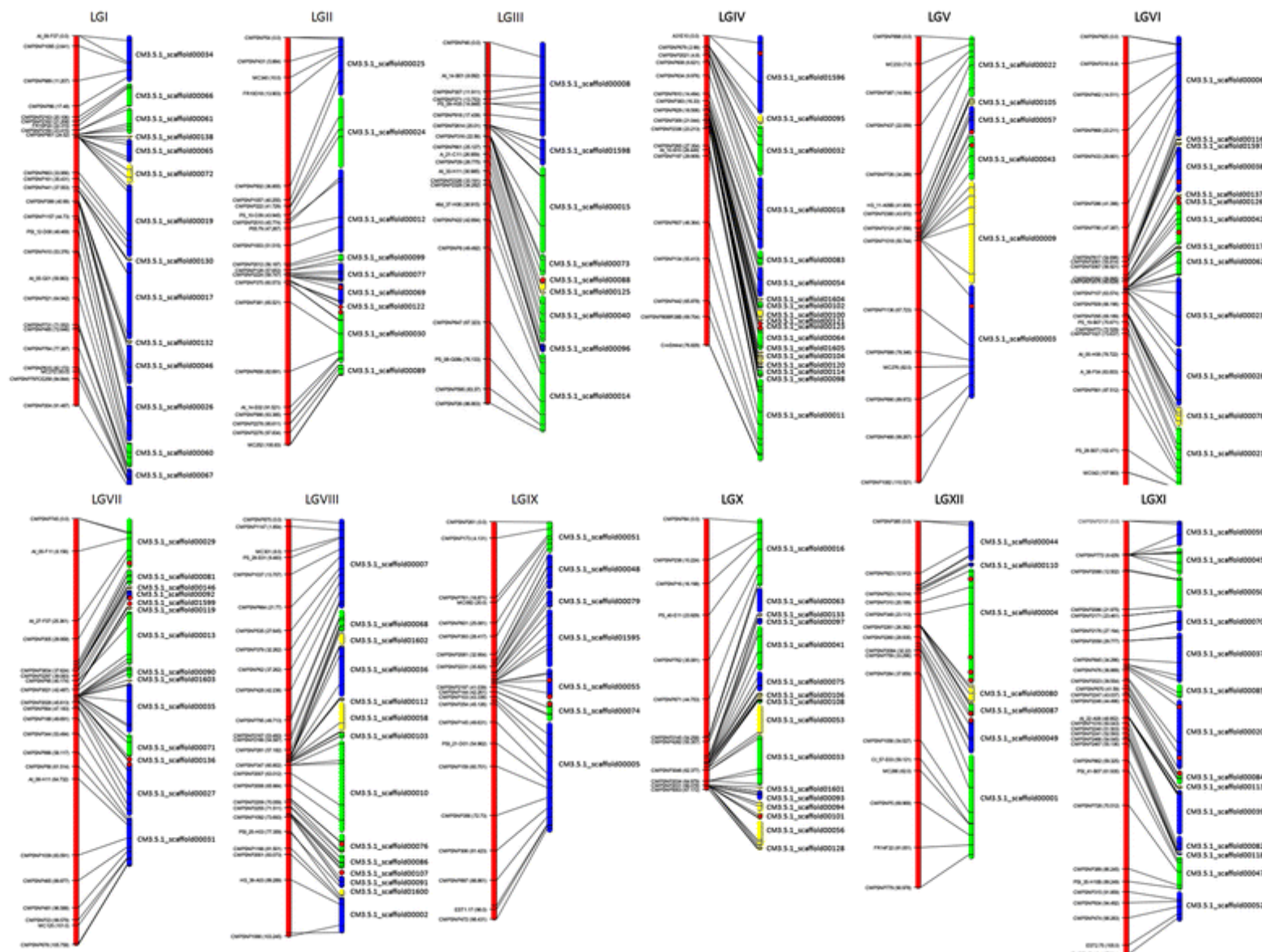
3. HiC, Chromatin Conformation Capture (Since 2014).



Dudechenko et al, Science 07 Apr 2017: Vol. 356, Issue 6333, pp. 92-95
DOI: 10.1126/science.aal3327

3.2. From scaffolds to chromosomes.

4. Linkage groups anchoring (Since 2001) - Very Long Distance



Argyris et al, BMC Genomics 2015 16:4



Outline of Topics

1. Brief history about genome assembly
2. Basics about sequence assembly
- 3. Whole genome assembly**
 - 3.1. From reads to contigs and scaffolds.
 - 3.2. From scaffolds to chromosomes.
4. Transcriptome assembly



3. Whole genome assembly

Decisions during the Assembly Optimization

Software

Identity %, Kmer ...

Post-assembly Filtering



3. Whole genome assembly

Computers

Bigger is better:

- How much do you need depends of:
 - ➡ how big is your genome ?
~ Human size (3Gb) require ~256 Gb to 1 Tb
 - ➡ how many reads do you have, ?
More reads, more memory and hard disk.
 - ➡ what software are you going to use ?
OLC uses more memory and time than DBG.
 - ➡ what parameters are you going to use ?
Bigger Kmers, more memory.

Four options:

1. Buy your own server (512 Gb, 4.5 Tb, 64 cores; ~ \$15,000).
2. Rent a server for ~ 1 month (same specs. \$1.5/h; ~\$1,000)(CBSU).
3. Use a supercomputing center associated with NSF, NIH, USDA... where they offer reduced prices (iPlant, Indiana University....).
4. Collaborate with some group with a big server.



3. Whole genome assembly

Assembly evaluation

Different methods to evaluate an assembly:

- 1. Assembly stats (total assembly size vs expected size; N50/L50...)**
- 2. Read remapping, variant calling and coverage evaluation.**
- 3. Comparison with alternative sequencing methods.**
- 4. Gene Space Completeness (RNASeq & EST mapping, BUSCO).**
- 5. Comparison with genetic information.**



3. Whole genome assembly

Assembly evaluation: Assembly stats

Rationale: Close that a sequence assembly is to the expected genome size, better it will be the genome assembly.

- Ideal case: The final assembly will have as many sequences as chromosomes have the genome.
- Real case: The assembly will be fragmented is sequences with different sizes. Lower number of fragments, better quality.



3. Whole genome assembly

Assembly evaluation: Assembly stats

During the assembly optimization will be generated several assemblies. The most used parameters to evaluate the assembly are:

1. Total Assembly Size,

How far is this value from the estimated genome size

2. Total Number of Sequences (Scaffold/Contigs)

How far is this value from the number of chromosomes.

3. Longest scaffold/contig

4. Average scaffold/contig size

5. N50/L50 (or any other N/L)

Number sequence (N) and minimum size of them (L) that represents the 50% of the assembly if the sequences are sorted by size, from bigger to smaller.



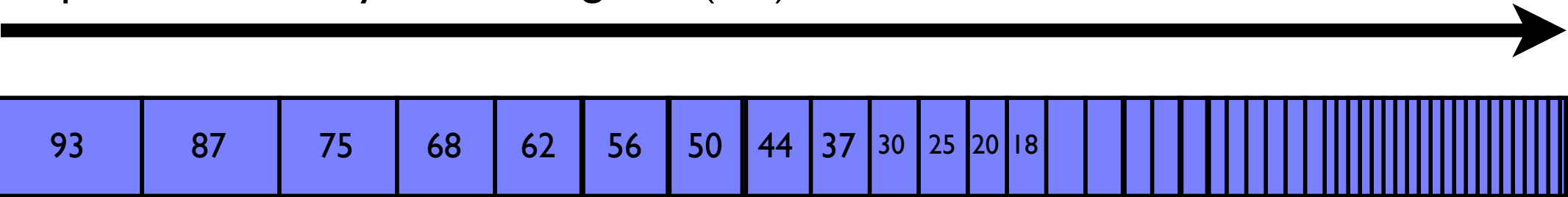
3. Whole genome assembly

Assembly evaluation: Assembly stats

N50/L50



Sequences order by descending size (Mb)



3. Whole genome assembly

Assembly evaluation: Assembly stats

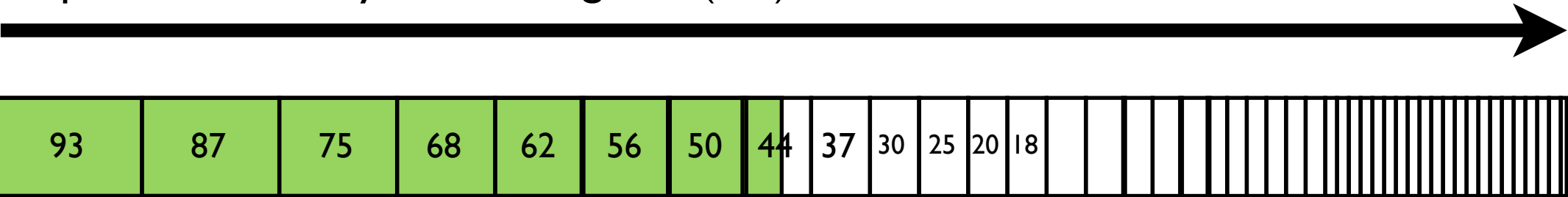
N50/L50



N50



Sequences order by descending size (Mb)



N50 = 7 sequences

L50 = 50 Mb



3. Whole genome assembly

Assembly evaluation: Assembly stats

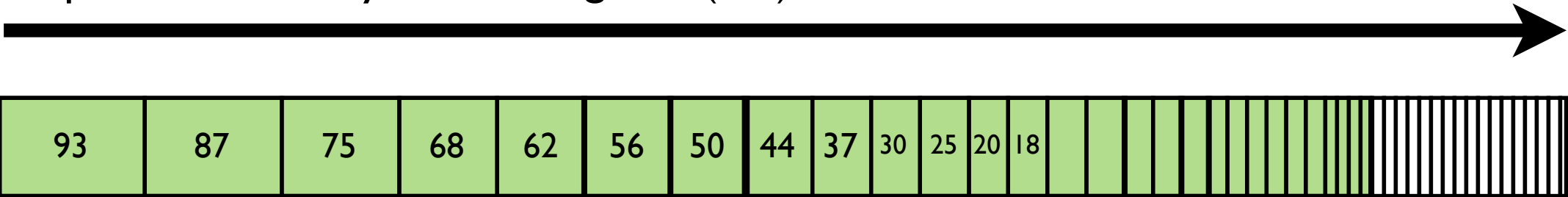
N90/L90



N90



Sequences order by descending size (Mb)



N90 = 29 sequences


L90 = 12.5 Mb



3. Whole genome assembly

Assembly evaluation: Assembly stats

Example: Sequencing of the *Petunia axillaris* genome.

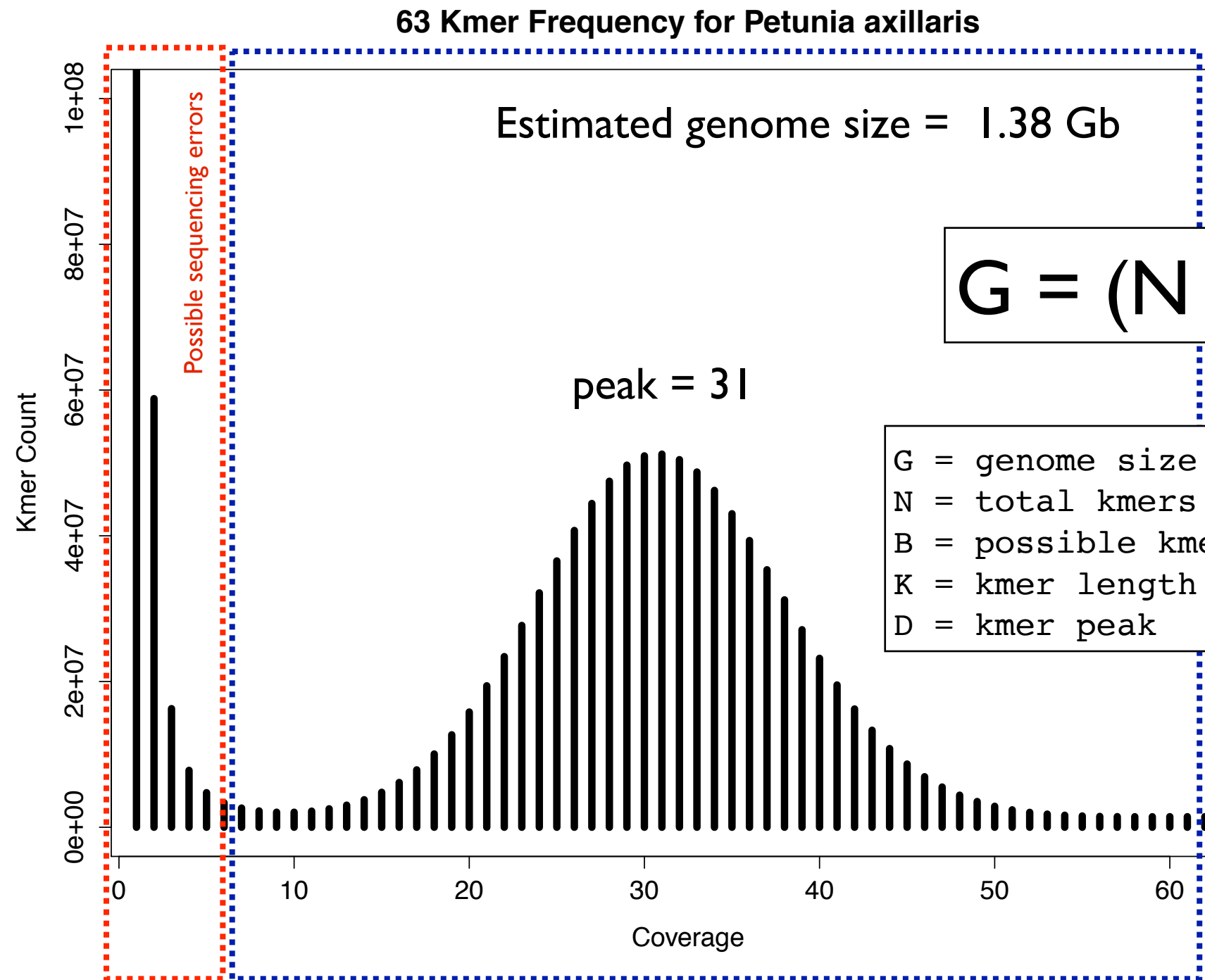
<i>P. axillaris</i>		
		
Peaxi v1.6.2		
Current Assembly		
Dataset	Contigs	Scaffolds
Total assembly size (Gb)	1.22	1.26
Total assembly sequences	109,892	83,639
Longest sequence (Mb)	0.57	8.56
Sequence length mean (Kb)	11.08	15.05
N90 (sequences)	13,481	1,051
L90 (Kb)	22.28	295.75
N50 (sequences)	3,943	309
L50 (Kb)	95.17	1,236.73



3. Whole genome assembly

Assembly evaluation: Assembly stats

Genome size (Kmer count)



$$G = (N - B) * K / D * 2$$

G = genome size
N = total kmers
B = possible kmers with errors
K = kmer length
D = kmer peak

3. Whole genome assembly

Assembly evaluation: Assembly stats

Estimated genome size for *P. axillaris*:

- Flow Cytometry (White & Rees, 1987) = 1.37 Gb
- Kmer Count* = 1.38 Gb
- Assembly size (scaffolds v1.6.2) = 1.26 Gb
- Assembly size (contigs v1.6.2) = 1.22 Gb

Estimated genome not assembled = 110 - 120 Mb



3. Whole genome assembly

Assembly evaluation: Read mapping

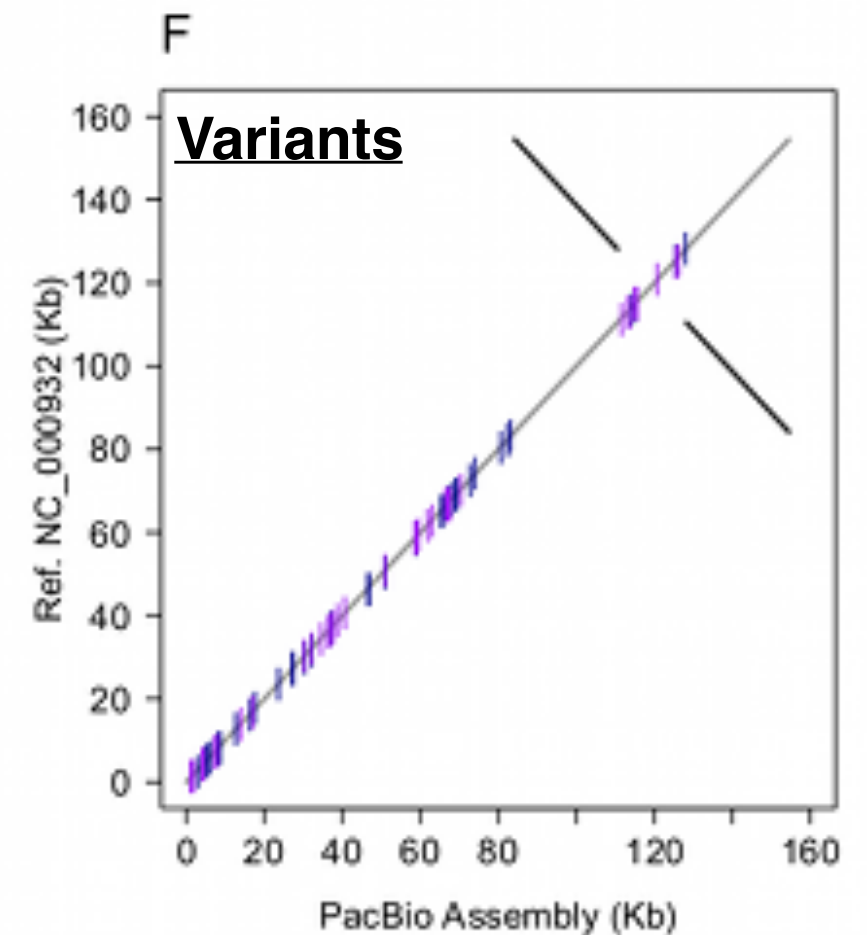
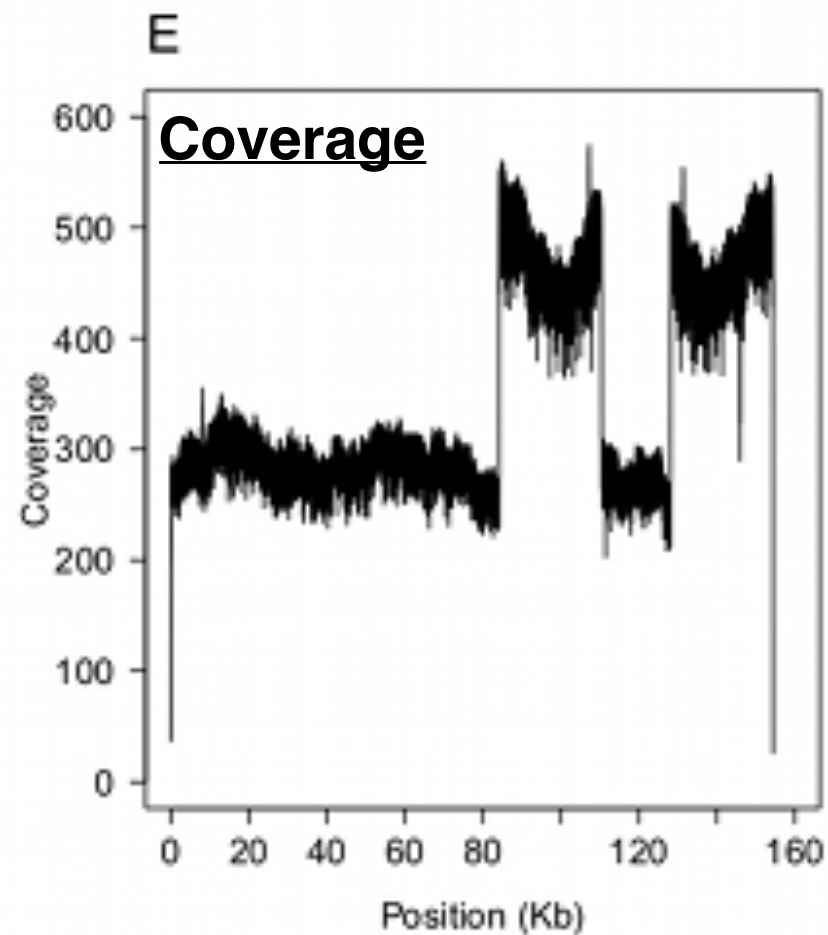
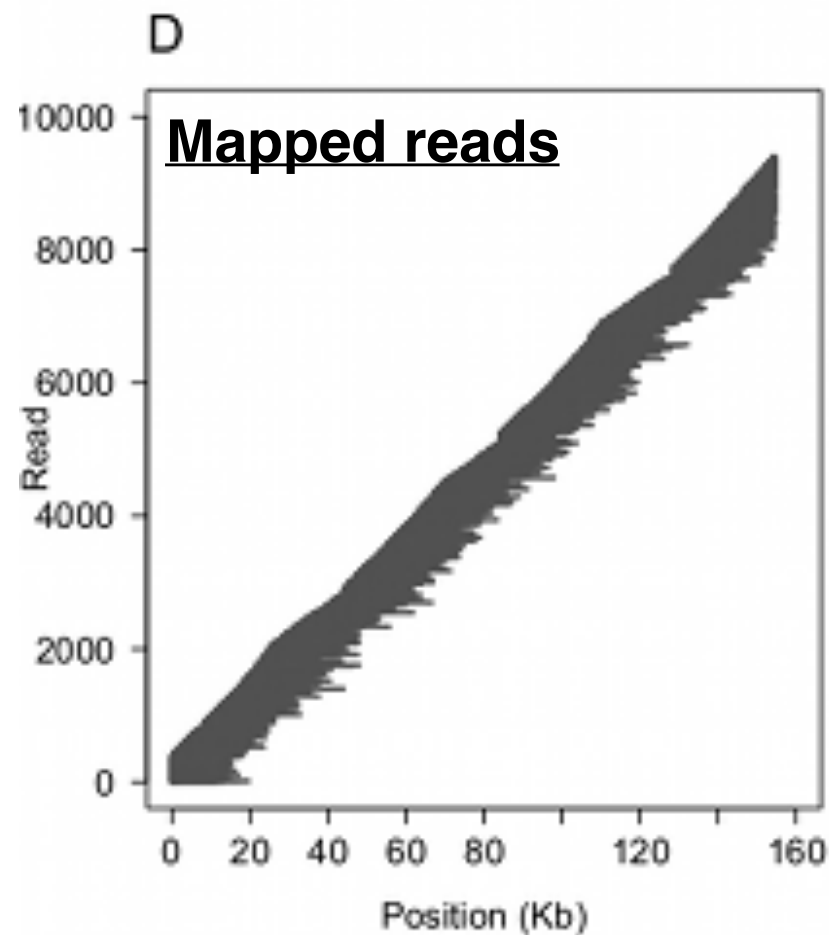
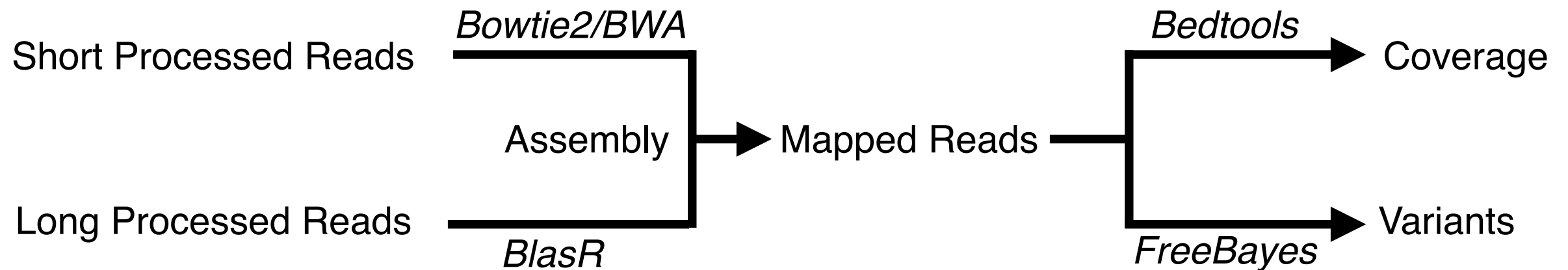
Rationale: Low quality assembled regions will have:

- Low coverage (breaking points)
- High coverage and high number of polymorphisms (due the collapsing of multiple copies).



3. Whole genome assembly

Assembly evaluation: Read mapping



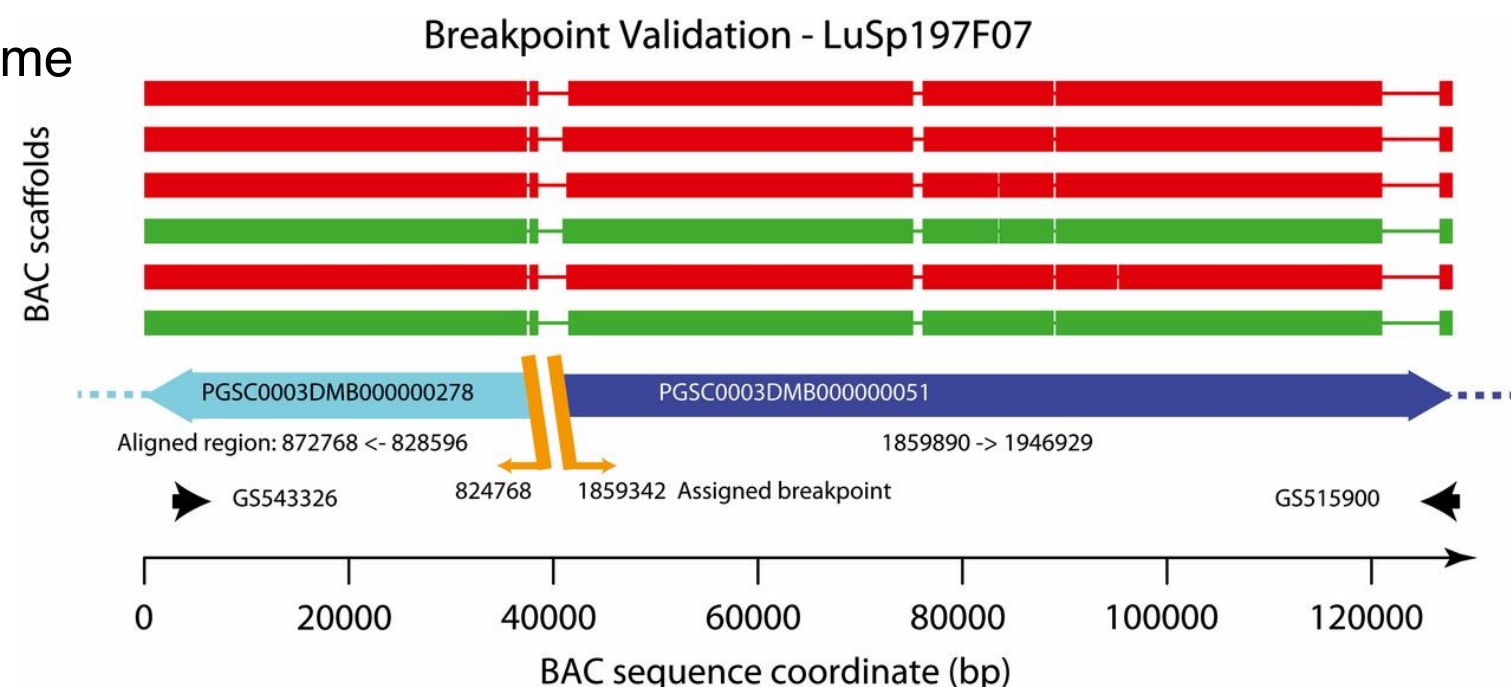
3. Whole genome assembly

Assembly evaluation: Comparative with alternative sequencing methods

Rationale: A sequence assembly should agree with alternative sequencing methods.

- Illumina based assembly → PacBio based assembly
- Short distance scaffolding (< 200 Kb) → Optical Mapping scaffolding.
- Long distance scaffolding (> 200 Kb) → BACs, YACs and Genetic maps

E.g. Potato genome



3. Whole genome assembly

Assembly evaluation: Gene space completeness

Rationale: The gene space is defined as the whole collection of genes in a genome. A complete genome should have a complete gene space. There are two ways to evaluate the completeness of a gene space.

- Check the completeness of a set of conserved genes (CEGMA/BUSCO).
 - Maybe not all the conserved genes are in the studied genome
- Check the mapping rate of a transcriptome to a genome (RNASeq).
 - The RNASeq may have contaminations of other organisms.



3. Whole genome assembly

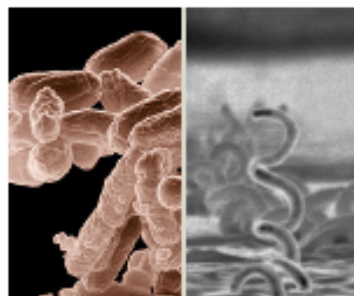
Assembly evaluation: Gene space completeness

Other way to estimate how good is an assembly, it is to run BUSCO over the assembly and evaluate how much of the gene space was captured.

<http://busco.ezlab.org/>

BUSCO  **v3**

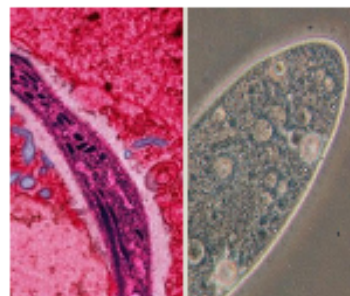
Datasets



Bacteria sets



Eukaryota sets



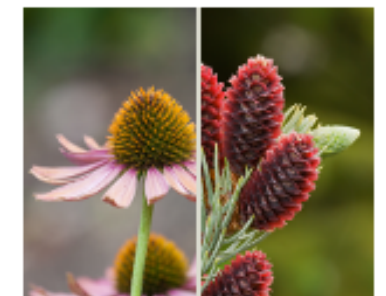
Protists sets



Metazoa sets



Fungi sets



Plants set

3. Whole genome assembly

Assembly evaluation: Gene space completeness

Other way to estimate how good is an assembly, it is to run BUSCO over the assembly and evaluate how much of the gene space was captured.

```
# BUSCO version is: 3.0.2
# The lineage dataset is: embryophyta_odb9 (Creation date: 2016-02-13, number of
species: 30, number of BUSCOs: 1440)
# To reproduce this run: python /data/software/busco/scripts/run_BUSCO.py -i Arth
a_TAIR10_Genome.fasta -o TAIR10_GENOME_BUSCO_TEST -l /data/software/busco/dataset
s/embryophyta_odb9/ -m genome -c 60 -sp arabidopsis
#
# Summarized benchmarking in BUSCO notation for file Artha_TAIR10_Genome.fasta
# BUSCO was run in mode: genome

C:98.2%[S:97.3%,D:0.9%],F:0.5%,M:1.3%,n:1440

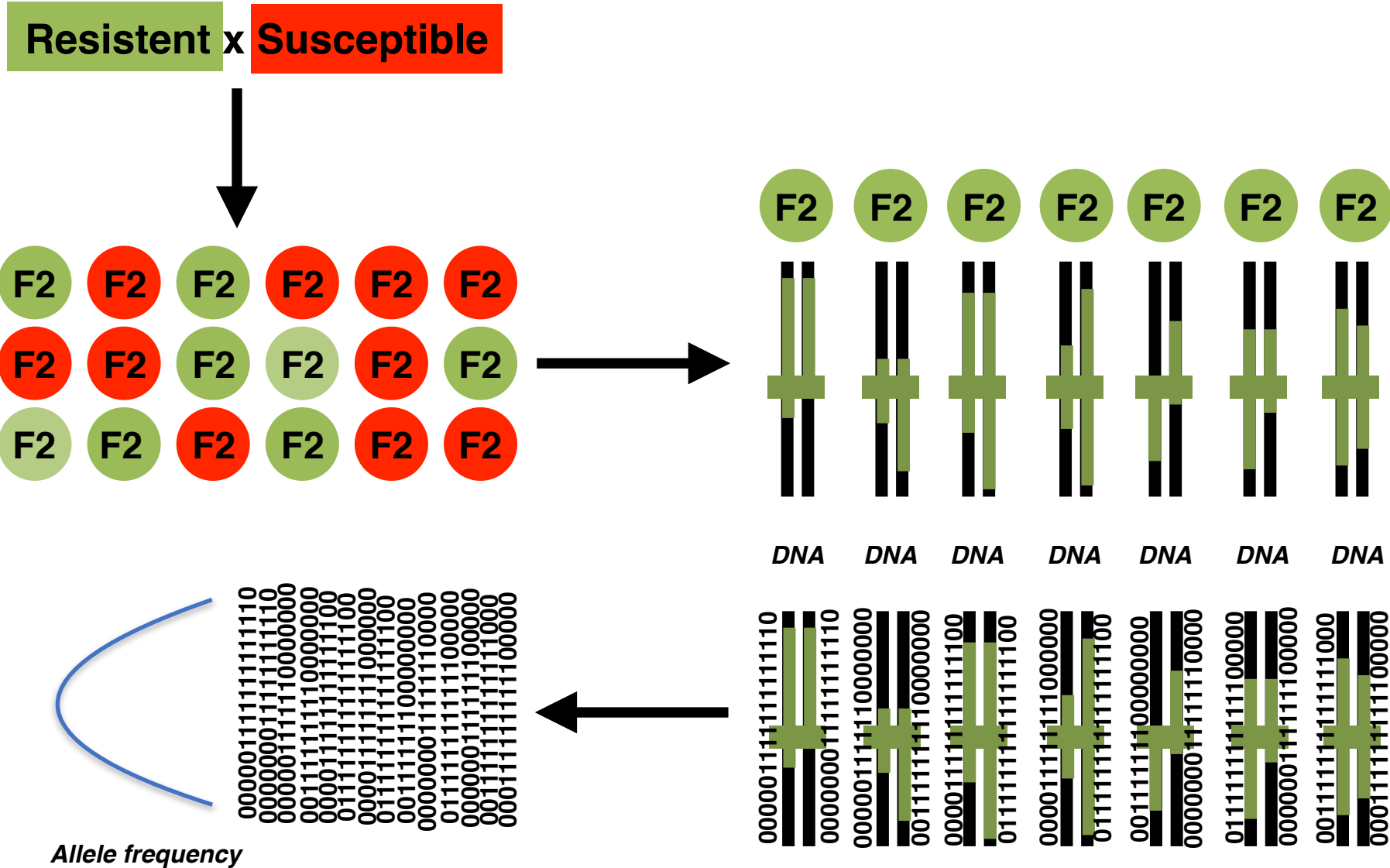
1414    Complete BUSCOs (C)
1401    Complete and single-copy BUSCOs (S)
13      Complete and duplicated BUSCOs (D)
7       Fragmented BUSCOs (F)
19      Missing BUSCOs (M)
1440    Total BUSCO groups searched
```



3. Whole genome assembly

Assembly evaluation: Comparison with Genetic Information

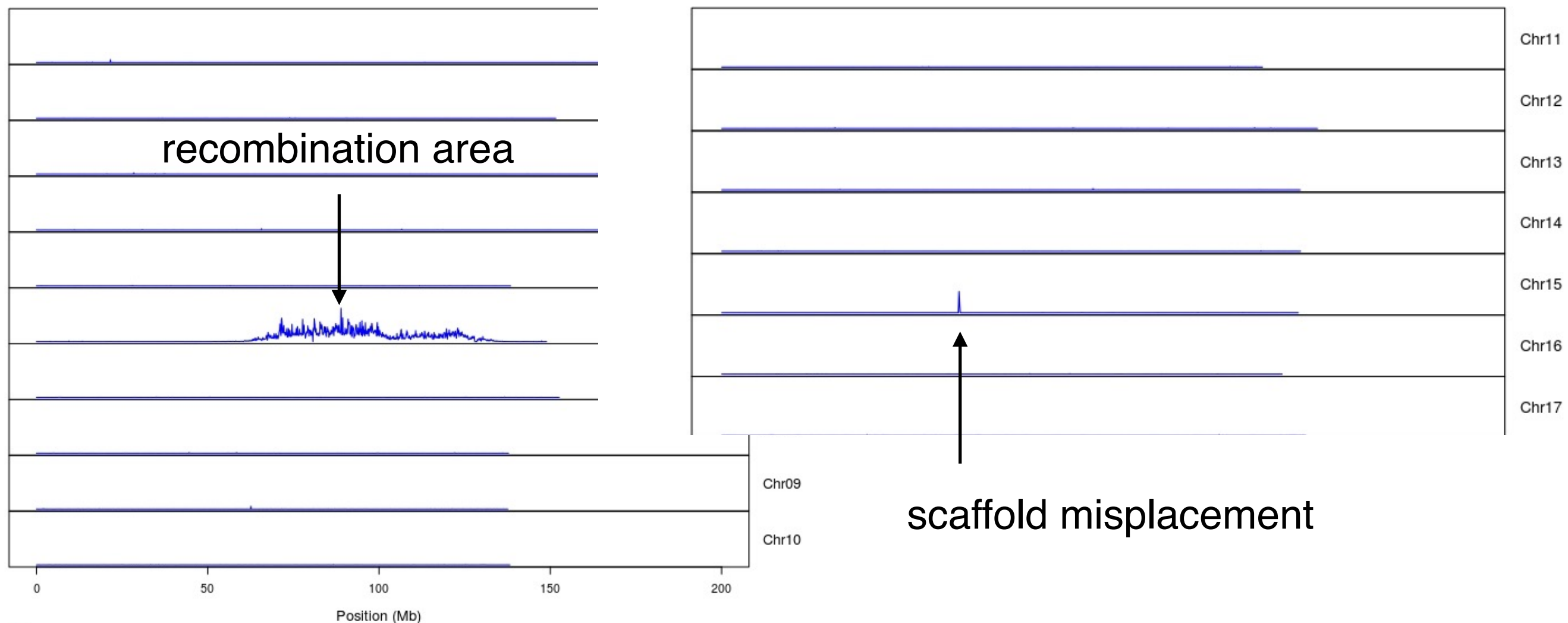
Rationale: The assembly information should be in agreement with genetic maps and recombination studies.



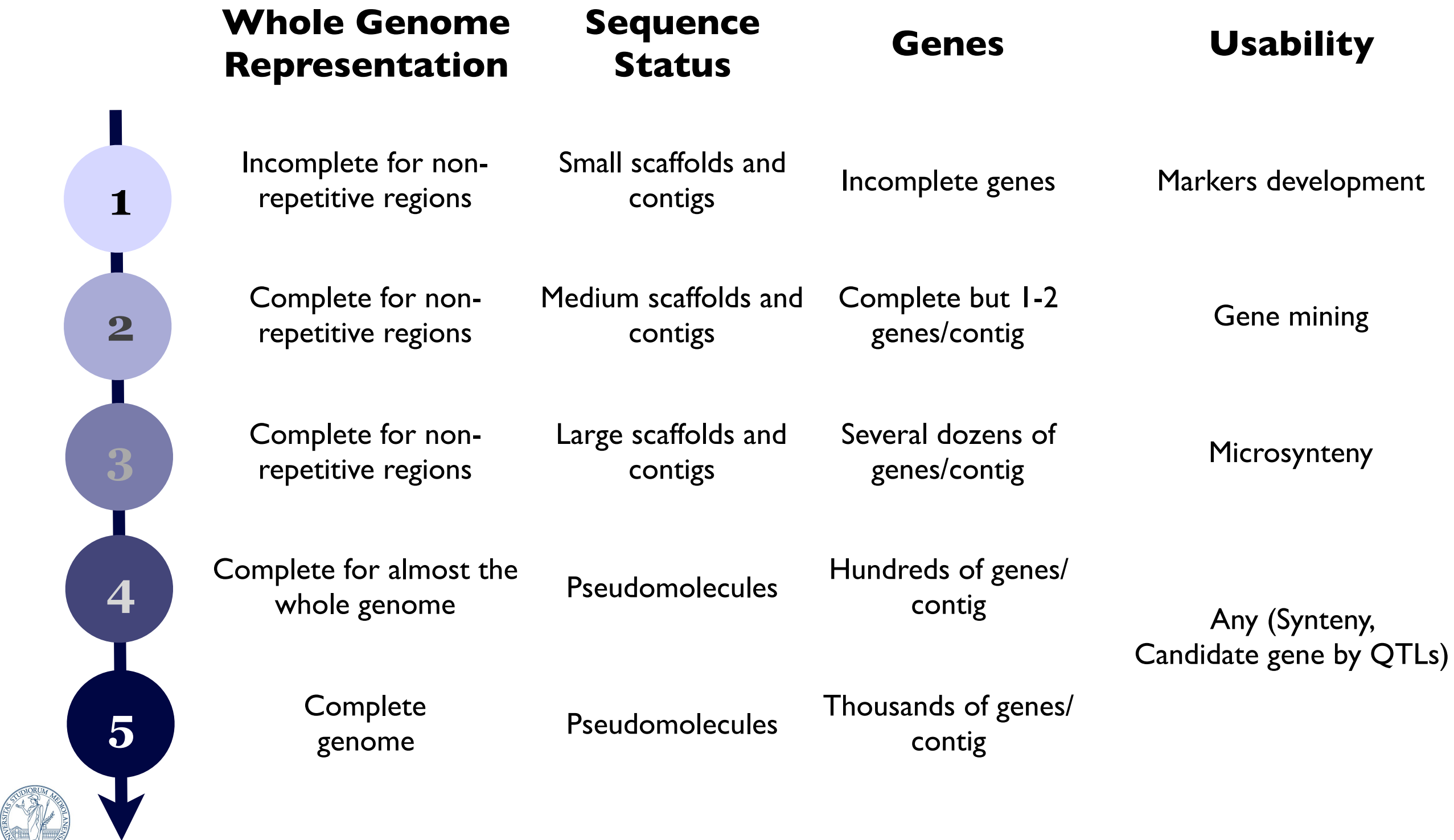
3. Whole genome assembly

Assembly evaluation: Comparison with Genetic Information

Rationale: The assembly information should be in agreement with genetic maps and recombination studies.



3. Whole genome assembly



Outline of Topics

1. Brief history about genome assembly
2. Basics about sequence assembly
3. Whole genome assembly
 - 3.1. From reads to contigs and scaffolds.
 - 3.2. From scaffolds to chromosomes.
4. Transcriptome assembly



4. Transcriptome assembly

<https://github.com/trinityrnaseq/trinityrnaseq/wiki>

Home

Brian Haas edited this page on 10 Apr 2018 · 37 revisions

RNA-Seq De novo Assembly Using Trinity



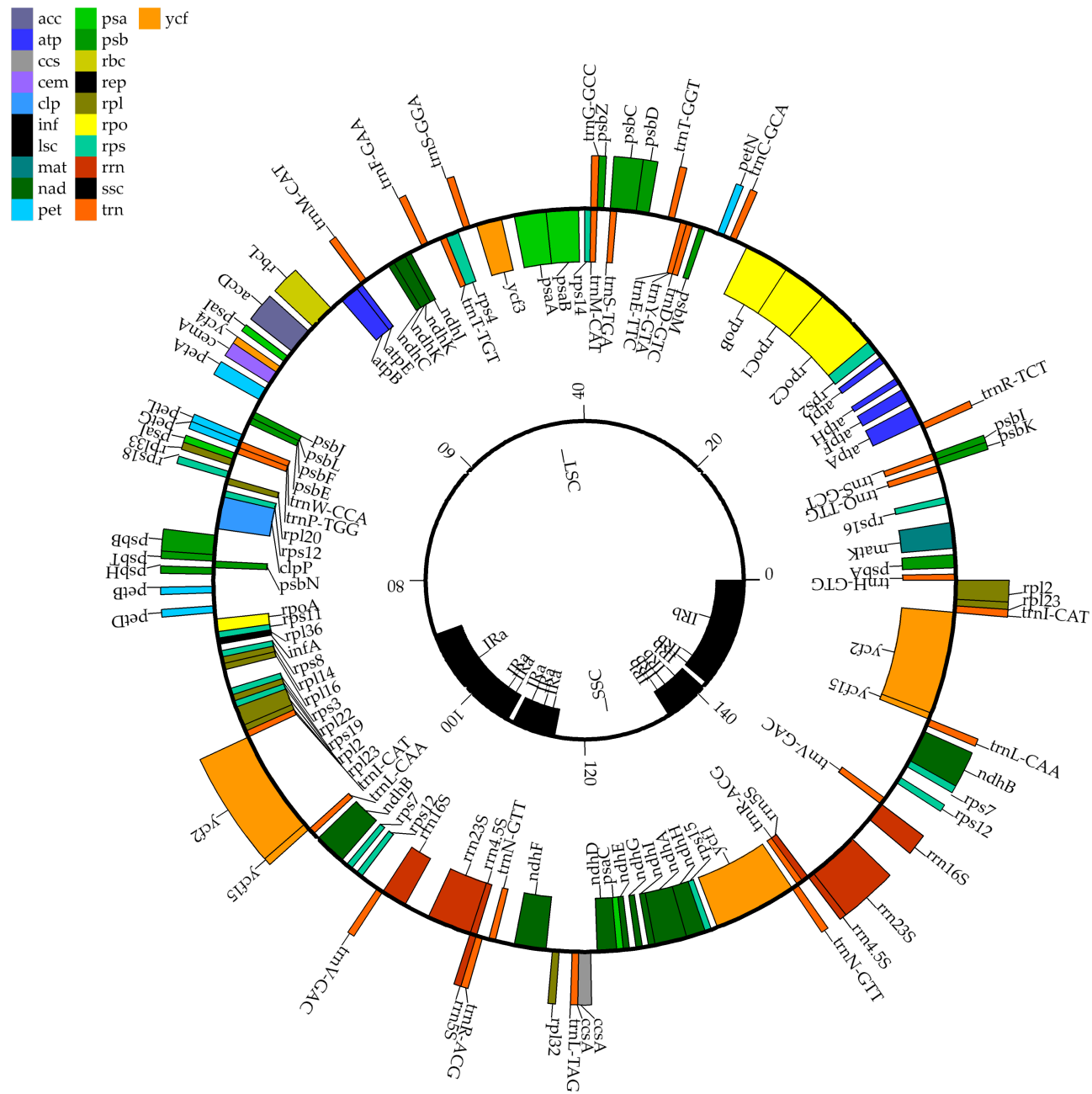
Exercises 5

Chloroplast genome assembly

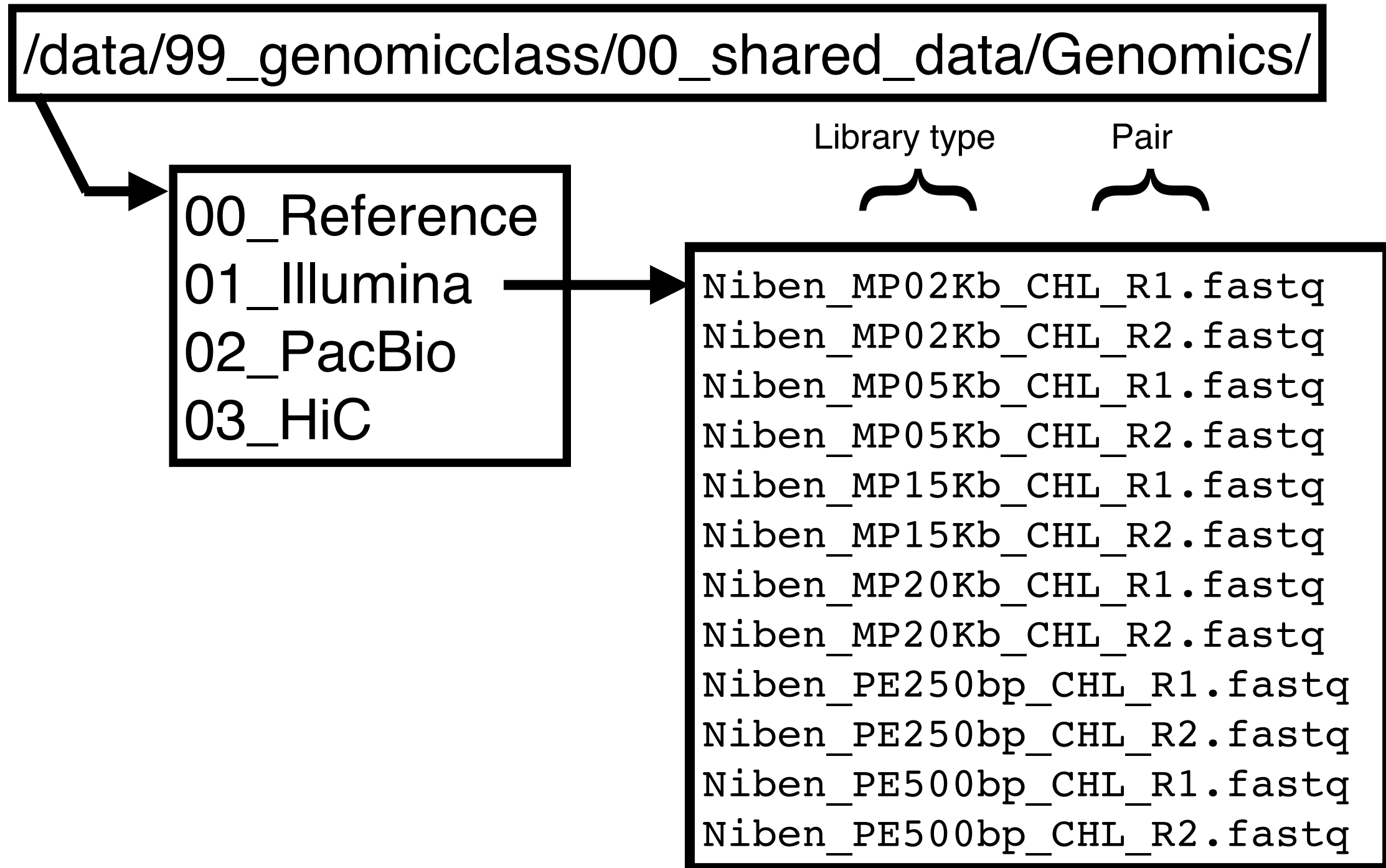


- Exercise 5.1: Get read dataset stats.
- Exercise 5.2: Short single read assembly - Coverage.
- Exercise 5.3: Short single read assembly - Kmer.
- Exercise 5.4: Short pair end read assembly - Scaffolding.
- Exercise 5.5: Long read assembly.
- Exercise 5.6: Assembly evaluation - Coverage.
- Exercise 5.7: Assembly evaluation - Variants.
- Exercise 5.8: Assembly structural annotation.

N. benthamiana chloroplast



Input data for the exercise



- **Exercise 5.1: Get read dataset stats.**
- Exercise 5.2: Short single read assembly - Coverage.
- Exercise 5.3: Short single read assembly - Kmer.
- Exercise 5.4: Short pair end read assembly - Scaffolding.
- Exercise 5.5: Long read assembly.
- Exercise 5.6: Assembly evaluation - Coverage.
- Exercise 5.7: Assembly evaluation - Variants.
- Exercise 5.8: Assembly structural annotation.

- Exercise 5.1: Get read dataset stats.

OBJECTIVE: To obtain the stats for all the read files from the 01_Illumina and 02_PacBio directories

TOOLS: Fastq-stats

COMMAND:

```
fastq-stats input_file.fq > input_file.stats.txt  
## To get all the read data  
grep read *stats.txt
```

INPUT FILES LOCATION:

```
/data/99_genomicclass/00_shared_data/Genomics/  
01_Illumina/  
/data/99_genomicclass/00_shared_data/Genomics/  
02_PacBio/
```

- Exercise 5.1: Get read dataset stats.
- **Exercise 5.2: Short single read assembly - Coverage.**
- Exercise 5.3: Short single read assembly - Kmer.
- Exercise 5.4: Short pair end read assembly - Scaffolding.
- Exercise 5.5: Long read assembly.
- Exercise 5.6: Assembly evaluation - Coverage.
- Exercise 5.7: Assembly evaluation - Variants.
- Exercise 5.8: Assembly structural annotation.

- Exercise 5.2: Short single read assembly - Coverage.

OBJECTIVE: To run different assemblies with different number of single reads and then compare the size of the contains. We will run **minia** over the file Niben_PE250bp_CHL_R1.fastq. Then we will select the 10% (0.1), 1% (0.01) and 0.1% (0.001) of the reads using the tool **seqtk sample** and we run mini again.

Name the inputs as Nb_IN100.0, Nb_IN010.0, Nb_IN001.0 and Nb_IN000.1 and the outputs as Minia_CHL100.0, Minia_CHL010.0, Minia_CHL001.0 and Minia_CHL000.1

TOOLS: Minia, Seqtk

- Exercise 5.2: Short single read assembly - Coverage.

COMMANDS:

Assembly

```
minia -in Nb_INXXX.fastq -kmer-size 63 -abundance-  
min 5 -out Minia_CHLXXX -nb-cores 10
```

Get a subset

```
seqtk sample 01_Illumina/Niben_PE250bp_CHL_R1.fastq  
XXX > Nb_INXXX.fastq
```

```
fastq-stats Nb_INXXX.fastq > Nb_INXXX.stats.txt
```

Get the stats

```
FastaSeqStats -i Minia_CHLXXX.contigs.fasta >  
Minia_CHLXXX.contigs.stats.txt
```


- Exercise 5.2: Short single read assembly - Coverage.

INPUT FILES LOCATION:

**`/data/99_genomicclass/00_shared_data/Genomics/
Niben_PE250bp_CHL_R1.fastq`**

SUGGESTIONS FOR DATA VISUALIZATION

1. Once you get the assembly stats with FastaSeqStats retrieve all the assemblies sizes running:
`grep Total
NibenSR_CHL*.contigs.stats.txt | sed -r 's/.*:\s+//' | sed -
r 's/\s+bp//' > NibenSR_CHL_ALL_TotalSize.txt`
2. Copy the file into your computer using FileZilla.
3. Load the file in R Studio
4. Use barplot to represent the different sizes.

- Exercise 5.2: Short single read assembly - Coverage.

GROUP DISCUSSION

Estimate the coverage used for each of the assemblies based in the input datasets and discuss the which dataset delivered better results.

- Exercise 5.1: Get read dataset stats.
- Exercise 5.2: Short single read assembly - Coverage.
- **Exercise 5.3: Short single read assembly - Kmer.**
- Exercise 5.4: Short pair end read assembly - Scaffolding.
- Exercise 5.5: Long read assembly.
- Exercise 5.6: Assembly evaluation - Coverage.
- Exercise 5.7: Assembly evaluation - Variants.
- Exercise 5.8: Assembly structural annotation.

- Exercise 5.3: Short single read assembly - Kmer.

OBJECTIVE: To run different assemblies with different kmers. You will run **minia** over the file Nb_IN010.0.fastq from the previous exercise. Run four different assemblies using the Kmers of 17, 31, 63 and 81.

Name the outputs as Minia_CHL_K17, Minia_CHL_K31, Minia_CHL_K63 and Minia_CHL_K81

TOOLS: Minia, Seqtk

- Exercise 5.3: Short single read assembly - Kmer.

COMMANDS:

Assembly

```
minia -in Nb_IN010.0.fastq -kmer-size 63 -abundance-  
min 5 -out Minia_CHL_KXX -nb-cores 10
```

Get the stats

```
FastaSeqStats -i Minia_CHLXXX.contigs.fasta >  
Minia_CHLXXX.contigs.stats.txt
```

- Exercise 5.3: Short single read assembly - Kmer.

GROUP DISCUSSION

Discuss the use of different Kmer sizes for the assembly. Is longer always better?

- Exercise 1: Get read dataset stats.
- Exercise 2: Short single read assembly - Coverage.
- Exercise 3: Short single read assembly - Kmer.
- **Exercise 4: Short pair end read assembly - Scaffolding.**
- Exercise 5: Long read assembly.
- Exercise 6: Assembly evaluation - Coverage.
- Exercise 7: Assembly evaluation - Variants.
- Exercise 8: Assembly structural annotation.

- Exercise 5.4: Short pair end read assembly - Scaffolding.

OBJECTIVE: To run different assemblies with different combinations of pair end and mate pair datasets. For this assembly you will be using ABYSS instead Minia because it has a scaffolding pipeline. We want to evaluate the effect of scaffolding rather than the assembly so we will need to prepare different datasets containing the same number of total reads.

The different assemblies to run are:

- ABYSS_01SR
- ABYSS_02PE
- ABYSS_03PE
- ABYSS_04PEMP
- ABYSS_05PEMP

- Exercise 5.4: Short pair end read assembly - Scaffolding.

COMMANDS

- ABYSS_01SR

```
## Select 50,000 reads
```

```
seqtk sample 01_Illumina/Niben_PE250bp_CHL_R1.fastq 50000  
> Niben_CHL50K.fq
```

```
## Run ABYSS
```

```
abyss-pe name="NibenCHL_Ab01SR" k=63 in="Niben_CHL50K.fq"
```

```
## Get the stats
```

```
FastaSeqStats -i NibenCHL_Ab01SR-3.fa  
>NibenCHL_Ab01SR-3.stats.txt
```

- Exercise 5.4: Short pair end read assembly - Scaffolding.

COMMANDS

- **ABYSS_02PE**

```
## Select 50,000 reads (25,000 per pair)
seqtk sample -s1000 01_Illumina/Niben_PE250bp_CHL_R1.fastq
25000 > Niben_PE250bp_25K_R1.fq
seqtk sample -s1000 01_Illumina/Niben_PE250bp_CHL_R2.fastq
25000 > Niben_PE250bp_25K_R2.fq
```

```
## Run ABYSS
abyss-pe name="NibenCHL_Ab02PE" k=63
in="Niben_PE250bp_25K_R1.fq Niben_PE250bp_25K_R2.fq"
```

```
## Get the stats
FastaSeqStats -i NibenCHL_Ab02PE1-scaffolds.fa >
NibenCHL_Ab02PE1-scaffolds.stats.txt
```

- Exercise 5.4: Short pair end read assembly - Scaffolding.

COMMANDS

- ABYSS_03PE

```
## Select 50,000 reads (12,500 per pair)
seqtk sample -s1000 01_Illumina/Niben_PE250bp_CHL_R1.fastq
12500 > Niben_PE250bp_12K_R1.fq
seqtk sample -s1000 01_Illumina/Niben_PE250bp_CHL_R2.fastq
12500 > Niben_PE250bp_12K_R2.fq
seqtk sample -s1000 01_Illumina/Niben_PE500bp_CHL_R1.fastq
12500 > Niben_PE500bp_12K_R1.fq
seqtk sample -s1000 01_Illumina/Niben_PE500bp_CHL_R2.fastq
12500 > Niben_PE500bp_12K_R2.fq
```

```
## Run ABYSS
```

```
abyss-pe name="NibenCHL_Ab03PE2" k=63 lib="pe250 pe500"
pe250="Niben_PE250bp_12K_R1.fq Niben_PE250bp_12K_R2.fq"
pe500="Niben_PE500bp_12K_R1.fq Niben_PE500bp_12K_R2.fq"
```

- Exercise 5.4: Short pair end read assembly - Scaffolding.

COMMANDS

- ABYSS_03PE

```
## Get the stats
```

```
FastaSeqStats -i NibenCHL_Ab03PE3-scaffolds.fa >  
NibenCHL_Ab03PE2-scaffolds.stats.txt
```

- Exercise 4: Short pair end read assembly - Scaffolding.

COMMANDS

- ABYSS_04PEMP

Note: The MP reads will not be used to generate the contigs

```
## Select 50,000 reads (12,500 per pair) and 50,000 for
the mate pairs
seqtk sample -s1000 01_Illumina/Niben_PE250bp_CHL_R1.fastq
12500 > Niben_PE250bp_12K_R1.fq
seqtk sample -s1000 01_Illumina/Niben_PE250bp_CHL_R2.fastq
12500 > Niben_PE250bp_12K_R2.fq
seqtk sample -s1000 01_Illumina/Niben_PE500bp_CHL_R1.fastq
12500 > Niben_PE500bp_12K_R1.fq
seqtk sample -s1000 01_Illumina/Niben_PE500bp_CHL_R2.fastq
12500 > Niben_PE500bp_12K_R2.fq
seqtk sample -s1000 01_Illumina/Niben_MP02Kb_CHL_R1.fastq
25000 > Niben_MP02Kb_25K_R1.fq
seqtk sample -s1000 01_Illumina/Niben_MP02Kb_CHL_R2.fastq
25000 > Niben_MP02Kb_25K_R1.fq
```

- Exercise 5.4: Short pair end read assembly - Scaffolding.

COMMANDS

- **ABYSS_04PEMP**

```
## Run ABYSS
```

```
abyss-pe name="NibenCHL_Ab04MP1" k=63 lib="pe250 pe500"  
mp="mp02" pe250="Niben_PE250bp_CHL00025_R1.fastq  
pe250="Niben_PE250bp_12K_R1.fq Niben_PE250bp_12K_R2.fq"  
pe500="Niben_PE500bp_12K_R1.fq Niben_PE500bp_12K_R2.fq"  
mp02="Niben_MP02Kb_25K_R1.fq Niben_MP02Kb_25K_R2.fq"
```

```
## Get the stats
```

```
FastaSeqStats -i NibenCHL_Ab04MP1-scaffolds.fa >  
NibenCHL_Ab04MP1-scaffolds.stats.txt
```

- Exercise 5.4: Short pair end read assembly - Scaffolding.

COMMANDS

- ABYSS_05PEMP

Note: The MP reads will not be used to generate the contigs

```
## Use the same dataset than in the previous exercise plus  
the following ones.
```

```
seqtk sample -s1000 01_Illumina/Niben_MP05Kb_CHL_R1.fastq  
25000 > Niben_MP05Kb_25K_R1.fq
```

```
seqtk sample -s1000 01_Illumina/Niben_MP05Kb_CHL_R2.fastq  
25000 > Niben_MP05Kb_25K_R1.fq
```

```
seqtk sample -s1000 01_Illumina/Niben_MP15Kb_CHL_R1.fastq  
25000 > Niben_MP15Kb_25K_R1.fq
```

```
seqtk sample -s1000 01_Illumina/Niben_MP15Kb_CHL_R2.fastq  
25000 > Niben_MP15Kb_25K_R1.fq
```

- Exercise 5.4: Short pair end read assembly - Scaffolding.

COMMANDS

- **ABYSS_05PEMP**

```
## Run ABYSS
```

```
abyss-pe name="NibenCHL_Ab05MP2" k=63 lib="pe250 pe500"  
mp="mp02 mp05 mp15" pe250="Niben_PE250bp_CHL00025_R1.fastq  
pe250="Niben_PE250bp_12K_R1.fq Niben_PE250bp_12K_R2.fq"  
pe500="Niben_PE500bp_12K_R1.fq Niben_PE500bp_12K_R2.fq"  
mp02="Niben_MP02Kb_25K_R1.fq Niben_MP02Kb_25K_R2.fq"  
mp05="Niben_MP05Kb_25K_R1.fq Niben_MP05Kb_25K_R2.fq"  
mp15="Niben_MP15Kb_25K_R1.fq Niben_MP15Kb_25K_R2.fq"
```

```
## Get the stats
```

```
FastaSeqStats -i NibenCHL_Ab05MP2-scaffolds.fa >  
NibenCHL_Ab05MP2-scaffolds.stats.txt
```


- Exercise 5.4: Short pair end read assembly - Scaffolding.

GROUP DISCUSSION

Discuss the use of different combination of pair ends and mate pair reads to generate different assemblies.

- Exercise 5.1: Get read dataset stats.
- Exercise 5.2: Short single read assembly - Coverage.
- Exercise 5.3: Short single read assembly - Kmer.
- Exercise 5.4: Short pair end read assembly - Scaffolding.
- **Exercise 5.5: Long read assembly.**
- Exercise 5.6: Assembly evaluation - Coverage.
- Exercise 5.7: Assembly evaluation - Variants.
- Exercise 5.8: Assembly structural annotation.

- Exercise 5.5: Long read assembly.

OBJECTIVE: To run different a assembly using PacBio reads and compare with the previous assembly made with Illumina reads (NibenCHL_Ab05MP2-scaffolds.fa). You will select the longest sequence of both assemblies, copy in your computer with FileZilla and Blast them.

TOOLS: Canu, FastaExtract

COMMANDS:

```
## Run a canu assembly  
canu -d NibenCHL_canu01 -p NibenCHL genomeSize=153k -  
pacbio-raw 02_PacBio/Niben_PB_rawreads.fastq
```

- Exercise 5.5: Long read assembly.

COMMANDS:

```
## Get the stats
```

```
FastaSeqStats -i NibenCHL_canu01/NibenCHL.contigs.fasta >  
NibenCHL.contigs.stats.txt
```

```
## Select the longest canu assembled contig
```

```
FastaExtract -f NibenCHL_canu01/NibenCHL.contigs.fasta -l  
100000 -o NibenCHL_canu01.longest.fa
```

```
## Select the longest ABYSS assembled contig
```

```
FastaExtract -f ../Exercise04/NibenCHL_Ab05MP2-scaffolds.fa  
-l 100000 -o NibenCHL_abyss05.longest.fa
```

- Exercise 5.5: Long read assembly.

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>

Standard Nucleotide BLAST

blastn blastp blastx tblastn tblastx

Enter Query Sequence

BLASTN programs search nucleotide databases using a nucleotide query. [more...](#) [Reset page](#) [Boc](#)

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#)

Or, upload file [Choose File](#) No file chosen [?](#)

Job Title [?](#)
Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Query subrange [?](#)

From

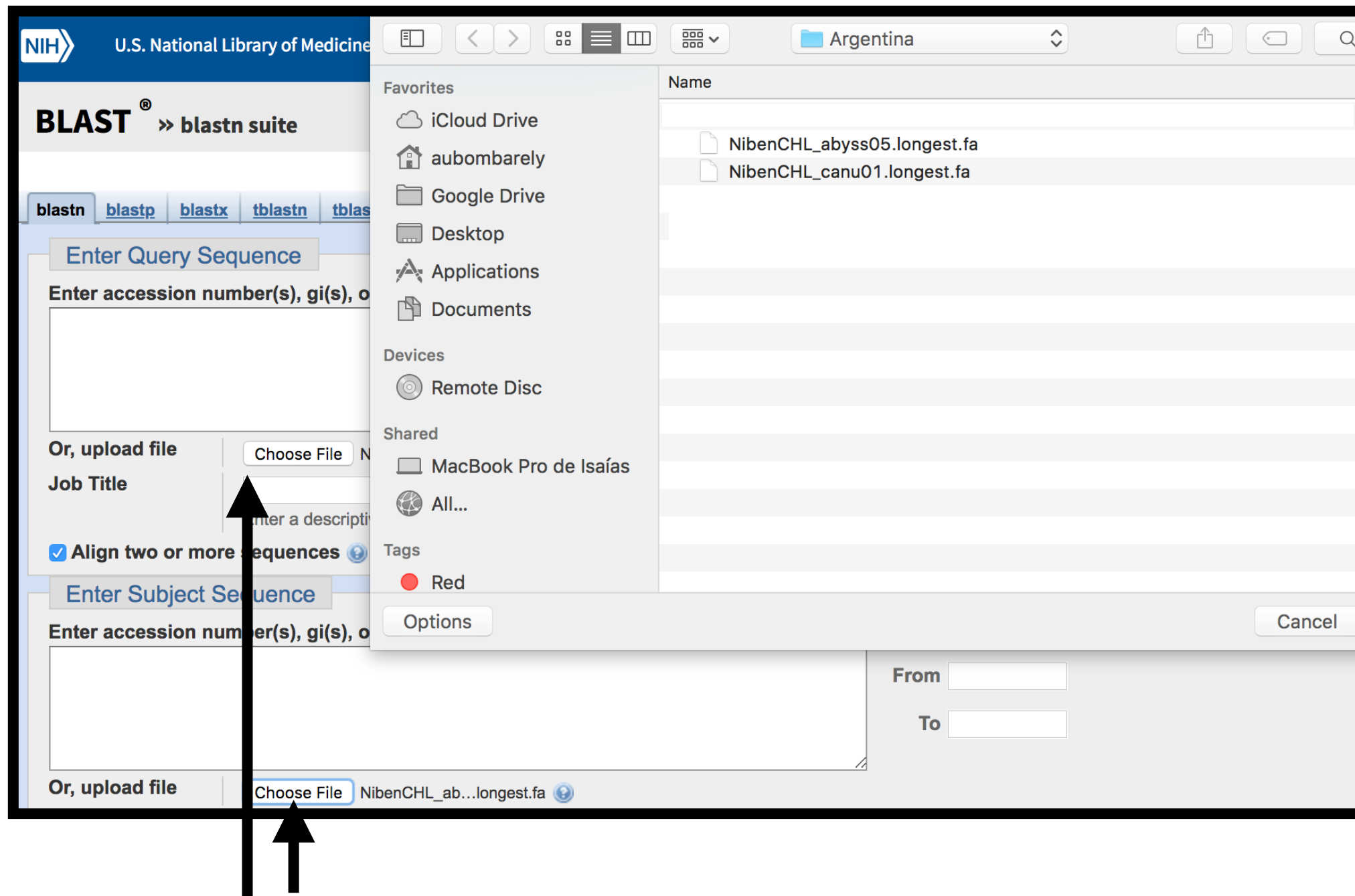
To



Click “Align two or more sequences”

- Exercise 5.5: Long read assembly.

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>



Upload both sequences and run the blast

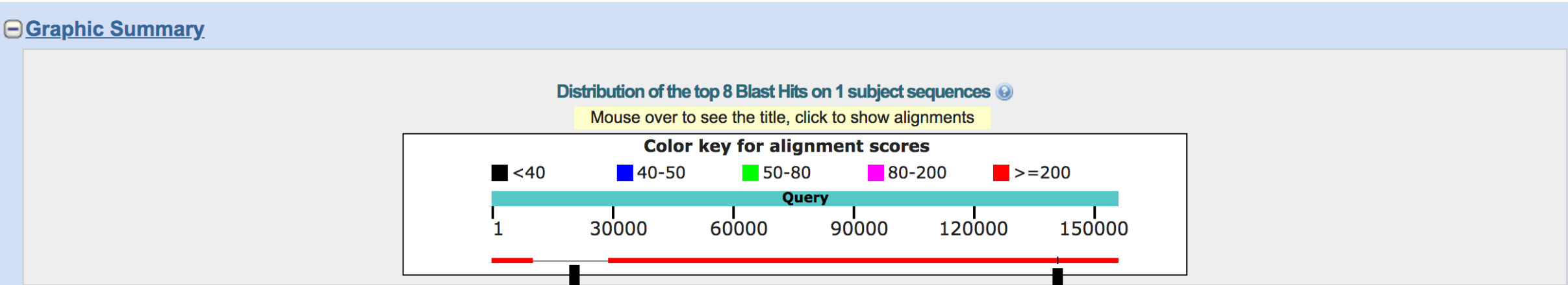
• Exercise 5.5: Long read assembly.

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>

RID TWNDX9U3114 (Expires on 09-17 21:08 pm)	
Query ID	lcl Query_45069
Description	tig00000001 len=156083 reads=3393 covStat=4464.05 gappedBases=no class=contig suggestRepeat=no suggestCircular=no
Molecule type	nucleic acid
Query Length	156083
Subject ID	lcl Query_45071
Description	4 111764 4920694 1+,0+
Molecule type	nucleic acid
Subject Length	111764
Program	BLASTN 2.8.0+ Citation

query = PacBio canu01

subject = Illumina ABYSS05

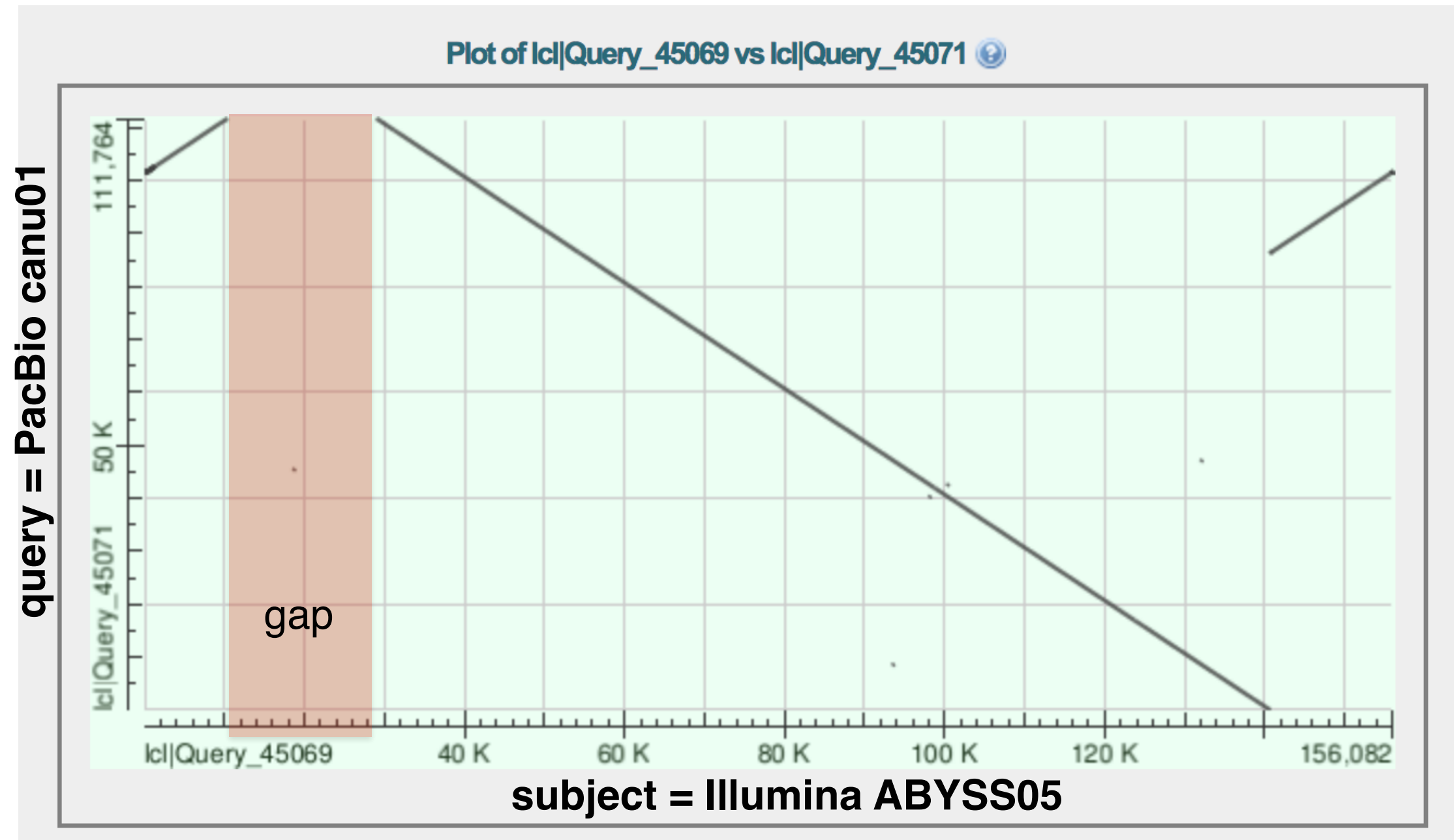


gap

break

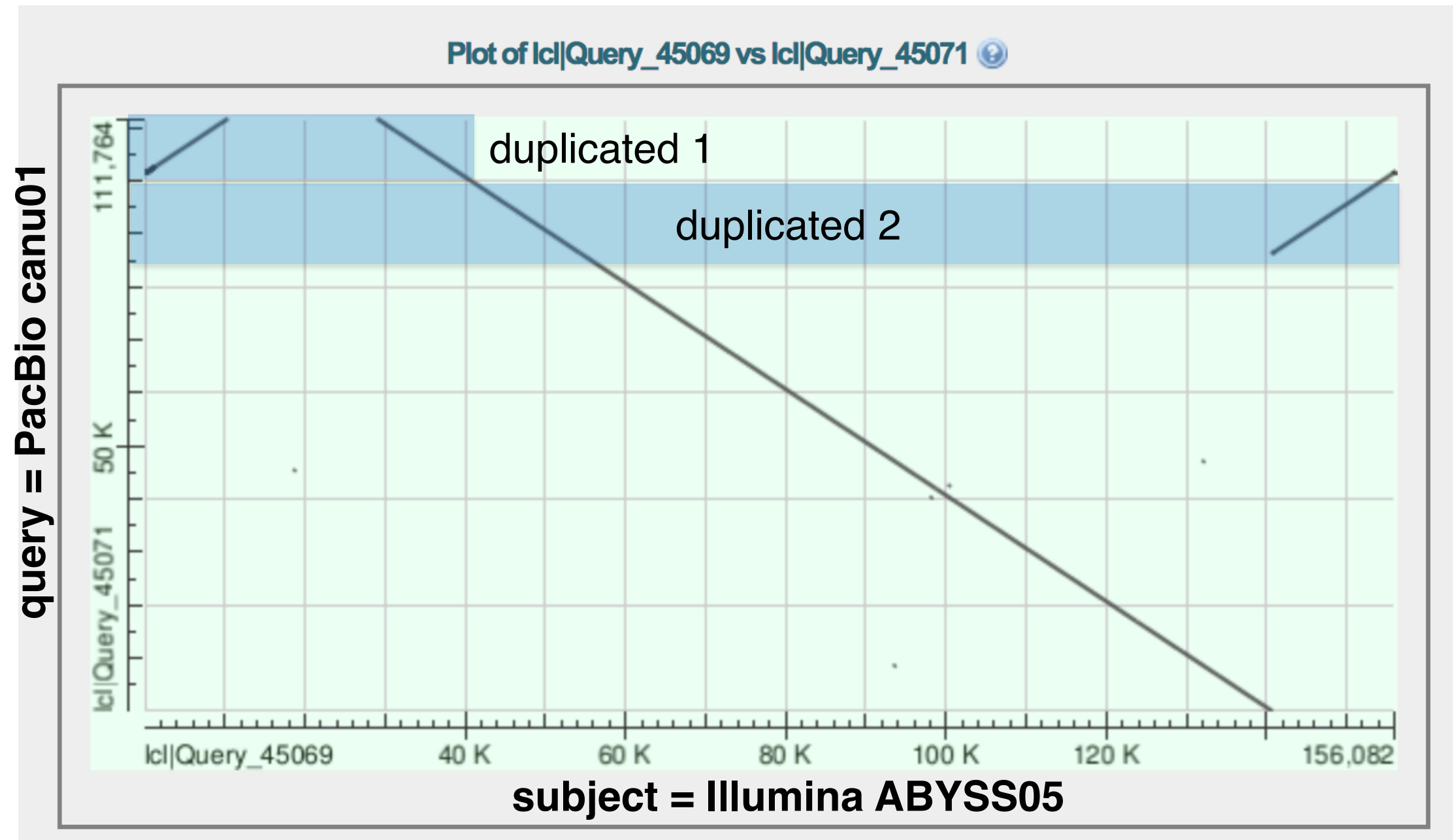
- Exercise 5.5: Long read assembly.

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>



- Exercise 5.5: Long read assembly.

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>



- Exercise 5.5: Long read assembly.

<https://blast.ncbi.nlm.nih.gov/Blast.cgi>

The screenshot shows the NCBI BLAST Standard Nucleotide BLAST interface. The browser address bar displays the URL: https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastn&PAGE_TYPE=BlastSearch&LINK_LOC=blasthome. The page header includes the NIH logo, "U.S. National Library of Medicine", and "NCBI National Center for Biotechnology Information". The main heading is "BLAST® >> blastn suite". Below this, the "Standard Nucleotide BLAST" section is active, with tabs for "blastn", "blastp", "blastx", "tblastn", and "tblastx". A blue banner states: "BLASTN programs search nucleotide databases using a nucleotide query. [more...](#)". The "Enter Query Sequence" section contains a large text input field for "Enter accession number(s), gi(s), or FASTA sequence(s)" with a "Clear" link. To the right, the "Query subrange" section has "From" and "To" input fields. Below the main input field, there is a section for "Or, upload file" with a "Choose File" button and "No file chosen" text. A "Job Title" input field is also present with the placeholder text "Enter a descriptive title for your BLAST search". At the bottom left, there is a checkbox labeled "Align two or more sequences".

Run the canu assembly against the whole Nucleotide Collection NR Genbank database

- Exercise 5.5: Long read assembly.

RID [TWP12V7V015](#) (Expires on 09-17 21:18 pm)

Query ID lcl|Query_126885

Description tig00000001 len=156083 reads=3393 covStat=4464.05 gappedBases=no
class=contig suggestRepeat=no suggestCircular=no

Molecule type nucleic acid

Query Length 156083

Database Name nr

Description Nucleotide collection (nt)

Program BLASTN 2.8.0+ [Citation](#)

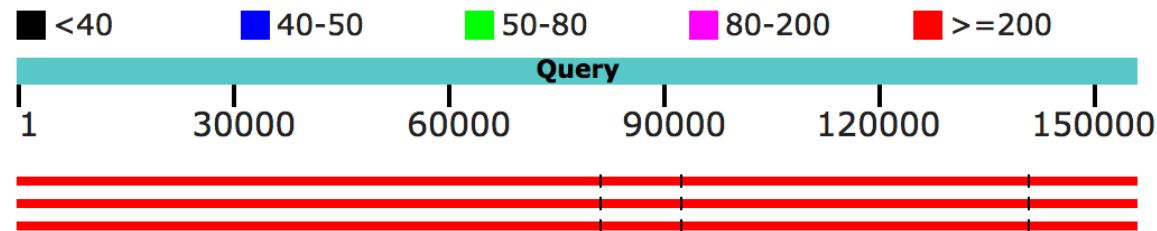
Other reports: [Search Summary](#) [Taxonomy reports](#) [Distance tree of results](#) [MSA viewer](#)

Graphic Summary

Distribution of the top 200 Blast Hits on 100 subject sequences

Mouse over to see the title, click to show alignments

Color key for alignment scores



Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

[Alignments](#) [Download](#) [GenBank](#) [Graphics](#) [Distance tree of results](#)

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	Nicotiana tabacum cultivar TN90 plastid, complete genome	1.465e+05	3.749e+05	99%	0.0	99%	KU199713.1
<input type="checkbox"/>	Nicotiana tabacum chloroplast genome DNA	1.465e+05	3.747e+05	99%	0.0	99%	Z00044.2
<input type="checkbox"/>	Nicotiana sylvestris chloroplast DNA, complete sequence						

Nicotiana tabacum cultivar TN90 plastid, complete genome
Sequence ID: [KU199713.1](#) Length: 155992 Number of Matches: 21

- Exercise 5.1: Get read dataset stats.
- Exercise 5.2: Short single read assembly - Coverage.
- Exercise 5.3: Short single read assembly - Kmer.
- Exercise 5.4: Short pair end read assembly - Scaffolding.
- Exercise 5.5: Long read assembly.
- **Exercise 5.6: Assembly evaluation - Coverage.**
- Exercise 5.7: Assembly evaluation - Variants.
- Exercise 5.8: Assembly structural annotation.

- Exercise 5.6: Assembly evaluation - Coverage.

OBJECTIVE: To map Illumina and PacBio reads back to the PacBio assembly and evaluate the coverage.

TOOLS: Bowtie2, ngmlr, Samtools, Bedtools, R

- Exercise 5.6: Assembly evaluation - Coverage.

COMMANDS:

```
##Copy the reference file
```

```
cp NibenCHL_canu01.longest.fa NibenCHL_canu01.longest_ref.fa
```

```
## Create the index using Bowtie2-Build
```

```
bowtie2-build NibenCHL_canu01.longest_ref.fa
```

```
NibenCHL_canu01.longest_ref
```

```
## Create a 200X dataset: Genome size ~ 150Kb x 200 =
```

```
30,000Kb = 30 Mb / 150 = 200,000 reads (100,000 pairs).
```

```
seqtk sample -s1000 01_Illumina/Niben_PE250bp_CHL_R1.fastq  
100000 > Niben_PE250bp_100K_R1.fq
```

```
seqtk sample -s1000 01_Illumina/Niben_PE250bp_CHL_R2.fastq  
100000 > Niben_PE250bp_100K_R2.fq
```

```
## Map the reads and filter the output with samtools
```

```
bowtie2 -p 8 -x NibenCHL_canu01.longest_ref -1
```

```
Niben_PE250bp_100K_R1.fq -2 Niben_PE250bp_100K_R2.fq |
```

```
samtools view -Sb -F 4 -o Niben_PE250bp.map.canu01.bam -
```

- Exercise 5.6: Assembly evaluation - Coverage.

COMMANDS:

```
##Sort the mapped reads by position
```

```
samtools sort -o Niben_PE250bp.map.canu01.bam
```

```
Niben_PE250bp.map.canu01.bam
```

```
##Create a fasta index file for bedtools
```

```
samtools faidx NibenCHL_canu01.longest_ref.fa
```

```
##Run bedtools and get the coverage per position
```

```
bedtools genomecov -d -ibam Niben_PE250bp.map.canu01.bam -g
```

```
NibenCHL_canu01.longest_ref.fa.fai >
```

```
Niben_PE250bp.map.canu01.covbed.txt
```

- Exercise 5.6: Assembly evaluation - Coverage.

To visualize the coverage:

- 1- Copy reference fasta and the bam into your computer and use IGV or Tablet.
- 2- Copy the bedtools coverage file in your computer and use R plot to visualize it.

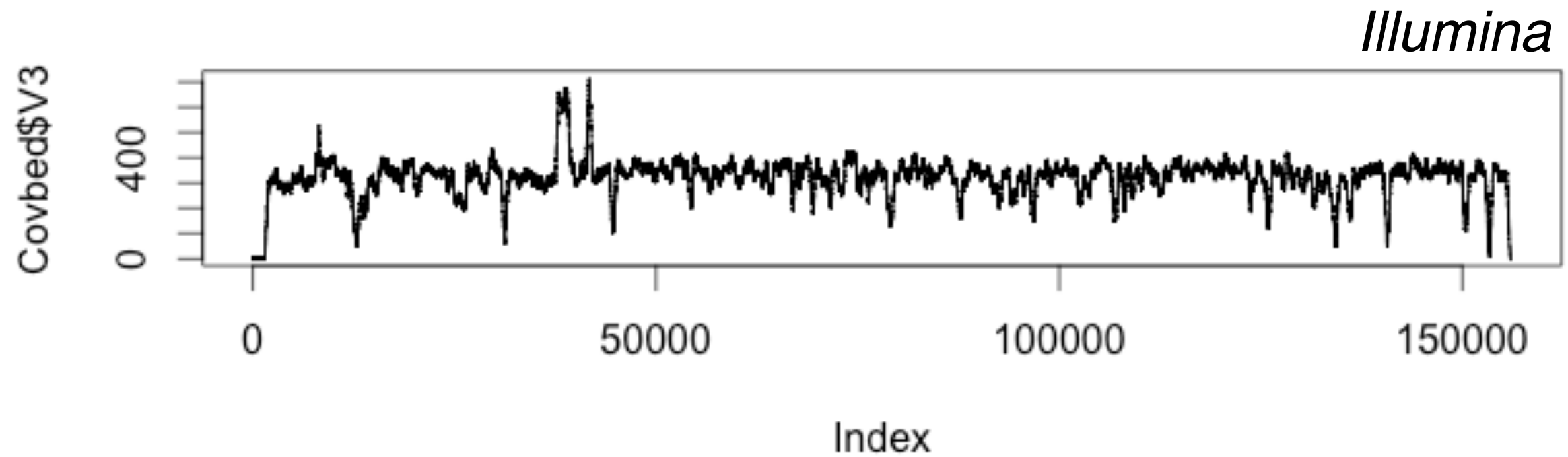
R commands

Covbed <-

```
read.delim("Niben_PE250bp.map.canu01.covbed.txt",  
header=FALSE)
```

```
plot(Covbed$V3, type="l")
```


- Exercise 5.6: Assembly evaluation - Coverage.



- Exercise 5.6: Assembly evaluation - Coverage.

COMMANDS:

```
## Select 200X of PacBio reads. The average read size is  
6.6Kb. Genome size ~ 150Kb x 200 = 30,000Kb / 6.6Kb = 4,545  
seqtk sample -s1000 ../../00_data/02_PacBio/  
Niben_PB_rawreads.fastq 4545 > Niben_PB_5K_R1.fq
```

```
##Map the PacBio reads with ngmlr  
ngmlr -t 4 -r NibenCHL_canu01.longest_ref.fa -q  
Niben_PB_5K_R1.fq | samtools view -Sb -F 4 -o  
Niben_PB.map.canu01.bam -
```

```
##Sort the reads  
samtools sort -o Niben_PB.map.canu01.bam  
Niben_PB.map.canu01.bam
```

```
##Run bedtools and get the coverage per position  
bedtools genomecov -d -ibam Niben_PB.map.canu01.bam -g  
NibenCHL_canu01.longest_ref.fa.fai >  
Niben_PB.map.canu01.covbed.txt
```

- Exercise 5.6: Assembly evaluation - Coverage.

To visualize the coverage:

- 1- Copy reference fasta and the bam into your computer and use IGV or Tablet.
- 2- Copy the bedtools coverage file in your computer and use R plot to visualize it.

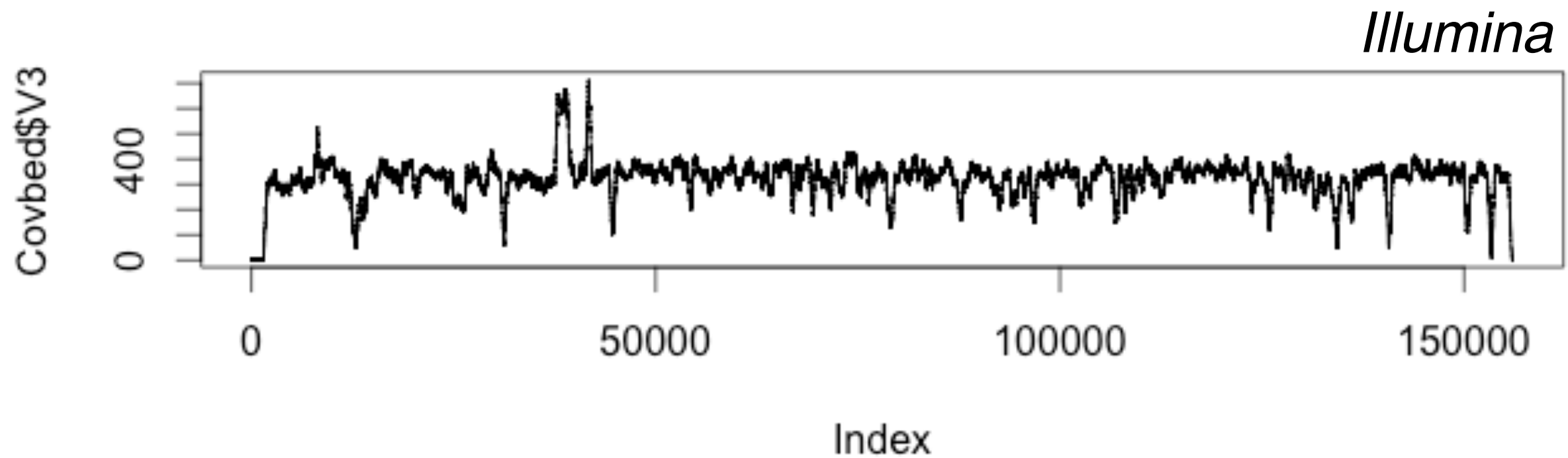
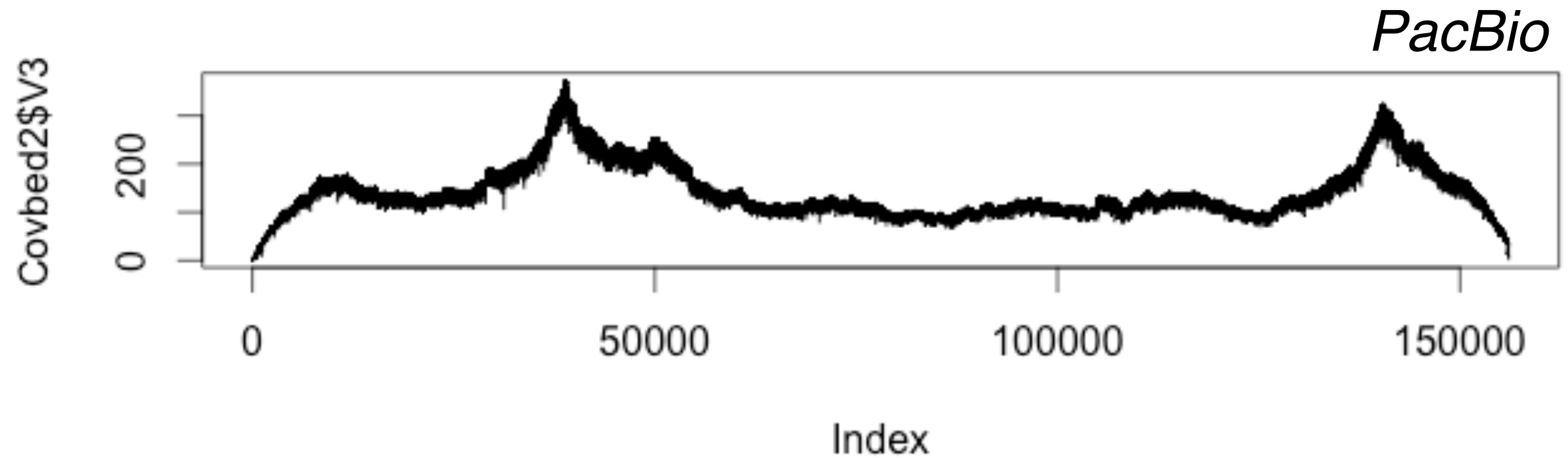
```
## R commands
```

```
Covbed2 <-
```

```
read.delim("Niben_PB.map.canu01.covbed.txt",  
header=FALSE)
```

```
plot(Covbed2$V3, type="l")
```

- Exercise 5.6: Assembly evaluation - Coverage.



- Exercise 5.1: Get read dataset stats.
- Exercise 5.2: Short single read assembly - Coverage.
- Exercise 5.3: Short single read assembly - Kmer.
- Exercise 5.4: Short pair end read assembly - Scaffolding.
- Exercise 5.5: Long read assembly.
- Exercise 5.6: Assembly evaluation - Coverage.
- **Exercise 5.7: Assembly evaluation - Variants.**
- Exercise 5.8: Assembly structural annotation.

- Exercise 5.7: Assembly evaluation - Variants

OBJECTIVE: Call the variants over the previously mapped Illumina reads.

TOOLS: Freebayes, R

- Exercise 5.7: Assembly evaluation - Variants

COMMANDS:

```
##Index the Bam file
```

```
samtools index Niben_PE250bp.map.canu01.bam
```

```
## Call the variants over the Illumina BAM file
```

```
freebayes -v Niben_PE250bp.map.canu01.vcf -b
```

```
Niben_PE250bp.map.canu01.bam -f
```

```
NibenCHL_canu01.longest_ref.fa
```

```
## Get some stats from the VCF
```

```
## 1- Count total variants
```

```
grep -v "#" Niben_PE250bp.map.canu01.vcf | wc -l
```

```
## 2- Count variants by type
```

```
grep -v "#" Niben_PE250bp.map.canu01.vcf | cut -f8 | sed -r
```

```
's/.;TYPE=/' | sort | uniq -c
```

```
## 3- Count non homozygous variants
```

```
grep -v "#" Niben_PE250bp.map.canu01.vcf | cut -f8,10 | sed -
```

```
r 's/.;TYPE=/' | sed -r 's/:.+/'
```

- Exercise 5.1: Get read dataset stats.
- Exercise 5.2: Short single read assembly - Coverage.
- Exercise 5.3: Short single read assembly - Kmer.
- Exercise 5.4: Short pair end read assembly - Scaffolding.
- Exercise 5.5: Long read assembly.
- Exercise 5.6: Assembly evaluation - Coverage.
- Exercise 5.7: Assembly evaluation - Variants.
- **Exercise 5.8: Assembly structural annotation.**

- Exercise 5.8: Assembly structural annotation.

<https://chlorobox.mpimp-golm.mpg.de/geseq.html>

[Nucleic Acids Res.](#) 2017 Jul 3; 45(Web Server issue): W6–W11.

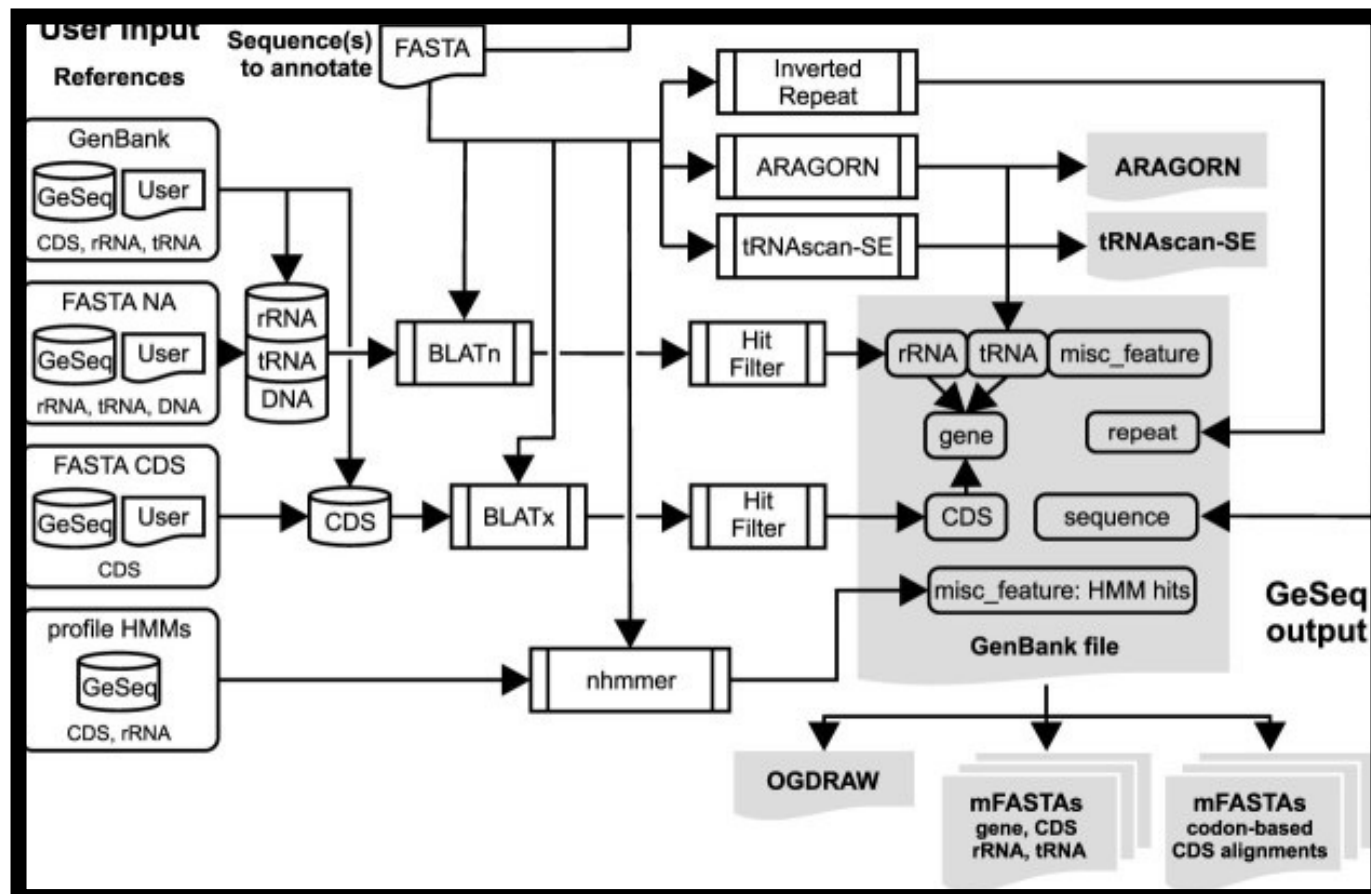
Published online 2017 May 9. doi: [10.1093/nar/gkx391](https://doi.org/10.1093/nar/gkx391)

PMCID: PMC5570176

PMID: [28486635](https://pubmed.ncbi.nlm.nih.gov/28486635/)

GeSeq – versatile and accurate annotation of organelle genomes

[Michael Tillich](#),¹ [Pascal Lehwark](#),² [Tommaso Pellizzer](#),¹ [Elena S. Ulbricht-Jones](#),¹ [Axel Fischer](#),¹ [Ralph Bock](#),¹ and [Stephan Greiner](#)¹



- Exercise 5.8: Assembly structural annotation.

<https://chlorobox.mpimp-golm.mpg.de/geseq.html>

Hoover over buttons for tips. To load an example sequence, click "Example" in the "Actions" field.

FASTA file(s) to annotate

1 file in list [+ Add Files](#)

NibenCHL_canu01.long

☐ Circular ☒ Linear

Sequence source

☒ Plastid ☐ Mitochondrial

Options

☐ Generate multi-GenBank

☐ Generate multi-FASTA

☐ Generate codon-based alignments

BLAT Reference Sequences

Server References

[Select NCBI RefSeq\(s\)](#)

no RefSeq selected

☐

MPI-MP chloroplast references (Embryophyta CDS + rRNA)

Custom References

GenBank

Actions

[Submit](#) [Reset](#)

[Example](#)

☒ I have read and accept the [Disclaimer](#)

Results

GeSeqJob-20180900-93339

ID gs_job_5b9e5bc62c619

Status: busy

2- Upload your canu01 assembly

1- Set up the annotation boxes

Annotation

☐ Annotate plastid IR

BLAT protein search

Identity

25

BLAT rRNA, tRNA, DNA search

Identity

85

☒ Find short matches (12 nucleotides)

HMMER profile search

☒ Embryophyta chloroplast (CDS + rRNA)


3rd Party tRNA annotators



☒ ARAGORN v1.2.38

- Exercise 5.8: Assembly structural annotation.

<https://chlorobox.mpimp-golm.mpg.de/geseq.html>

3- Download the results

Results 

GeSeqJob-20180900-93339  
ID gs_job_5b9e5bc62c619
Status: finished
tig00000001-len=156083-
reads=3393-covStat=4464.05-
gappedBases=no-
class=contig-
suggestRepeat=no-
suggestCircular=no

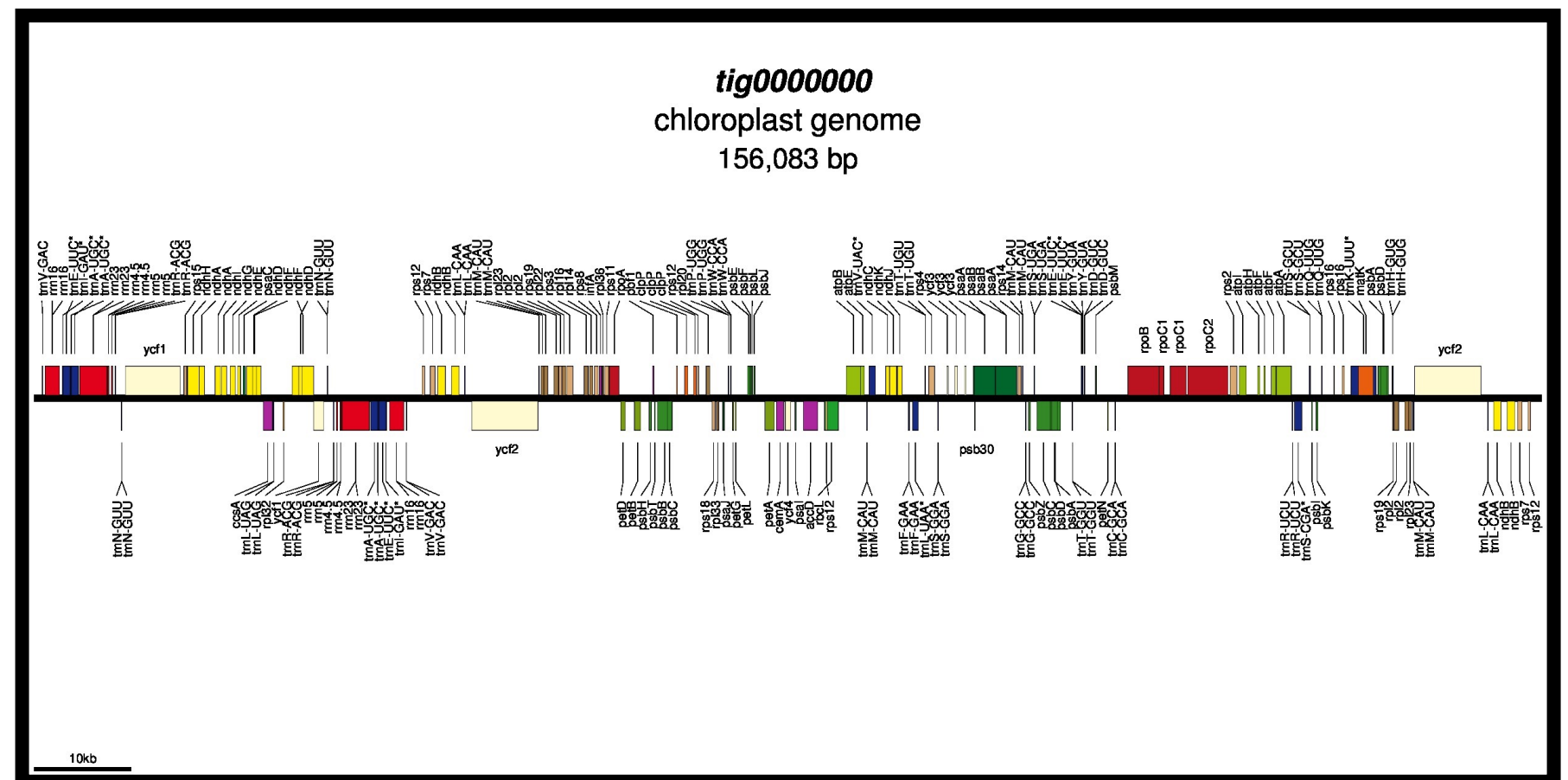
hmmer

ARAGORN v1.2.38

tRNAscan-SE v2.0

OGDRAW **GenBank**


Name	Date Modified	Size	Kind
See the location of the current folder			
GeSeqJob-20180900-9...stCircular=no_ARAGORN v1.2.38.txt	Today, 3:38 PM	2 KB	Plain Text
GeSeqJob-20180900-9...stCircular=no_GenBank.gb	Today, 3:38 PM	260 KB	Genba...text file
GeSeqJob-20180900-9...stCircular=no_hmmer.txt	Today, 3:38 PM	30 KB	Plain Text
GeSeqJob-20180900-9...stCircular=no_OGDRAW.jpg	Today, 3:38 PM	347 KB	JPEG image
GeSeqJob-20180900-9...stCircular=no_tRNAscan-SE v2.0.txt	Today, 3:38 PM	6 KB	Plain Text





- Exercise 5.8: Assembly structural annotation.

<https://chlorobox.mpimp-golm.mpg.de/geseq.html>

3- Download the results

Results 

GeSeqJob-20180900-93339  
ID gs_job_5b9e5bc62c619
Status: finished
tig00000001-len=156083-
reads=3393-covStat=4464.05-
gappedBases=no-
class=contig-
suggestRepeat=no-
suggestCircular=no

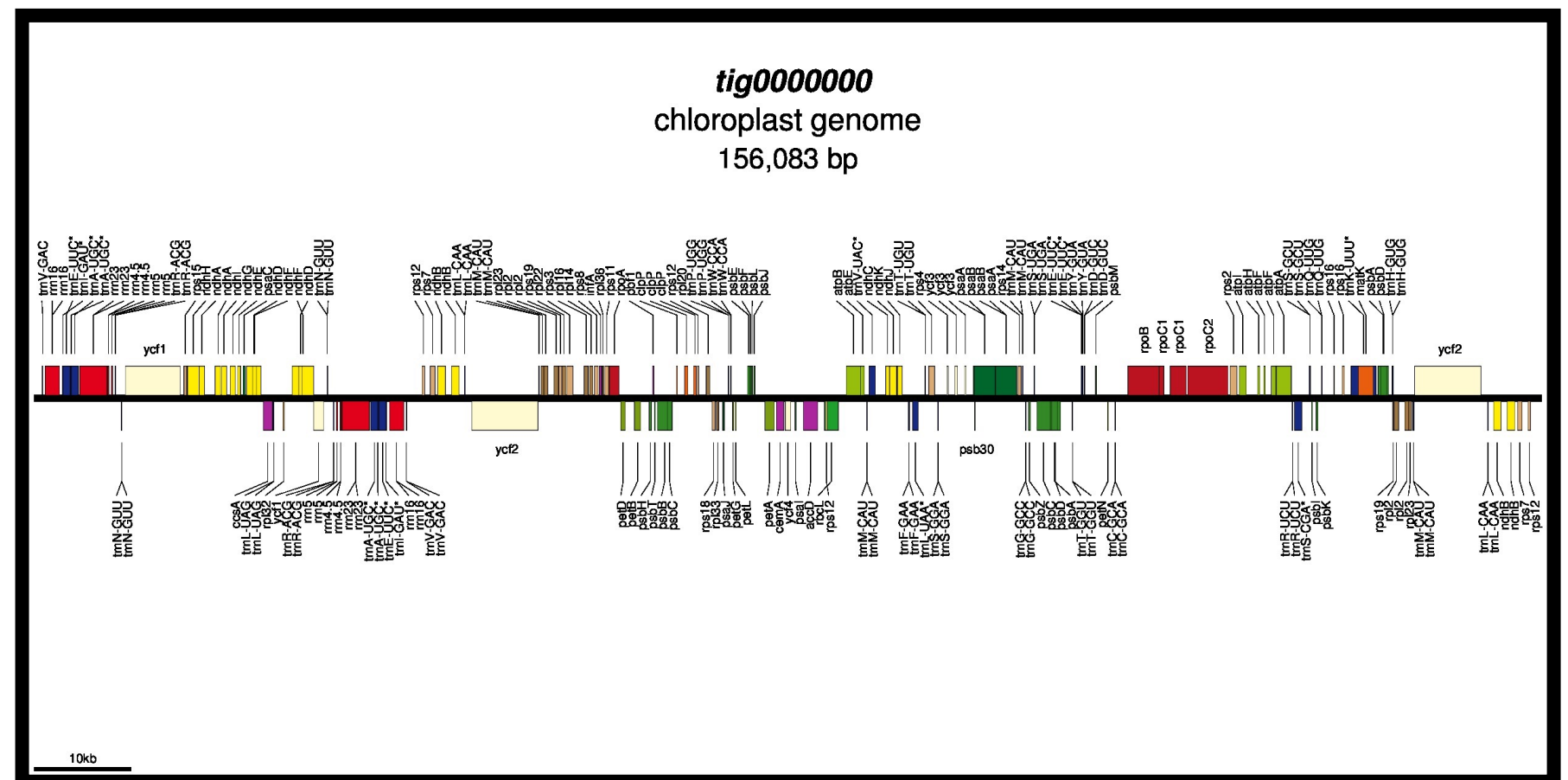
hmmer

ARAGORN v1.2.38

tRNAscan-SE v2.0

OGDRAW **GenBank**

Name	Date Modified	Size	Kind
See the location of the current folder			
GeSeqJob-20180900-9...stCircular=no_ARAGORN v1.2.38.txt	Today, 3:38 PM	2 KB	Plain Text
GeSeqJob-20180900-9...stCircular=no_GenBank.gb	Today, 3:38 PM	260 KB	Genba...text file
GeSeqJob-20180900-9...stCircular=no_hmmer.txt	Today, 3:38 PM	30 KB	Plain Text
GeSeqJob-20180900-9...stCircular=no_OGDRAW.jpg	Today, 3:38 PM	347 KB	JPEG image
GeSeqJob-20180900-9...stCircular=no_tRNAscan-SE v2.0.txt	Today, 3:38 PM	6 KB	Plain Text



Notes about sending the results

Please, send the results as two files:

- Discussion_ExerciseX.X_MyLastName.txt
- Script_ExerciseX.X_MyLastName.sh

In the subject of the email write “GT2020 - Solution Exercise X.X” to facilitate the use of search tools in the email inbox.

Deadlines:

- Exercises 5.1 to 5.3 - March 31, 2020
- Exercises 5.4 to 5.6 - April 3, 2020
- Exercises 5.7 and 5.8 - April 6, 2020