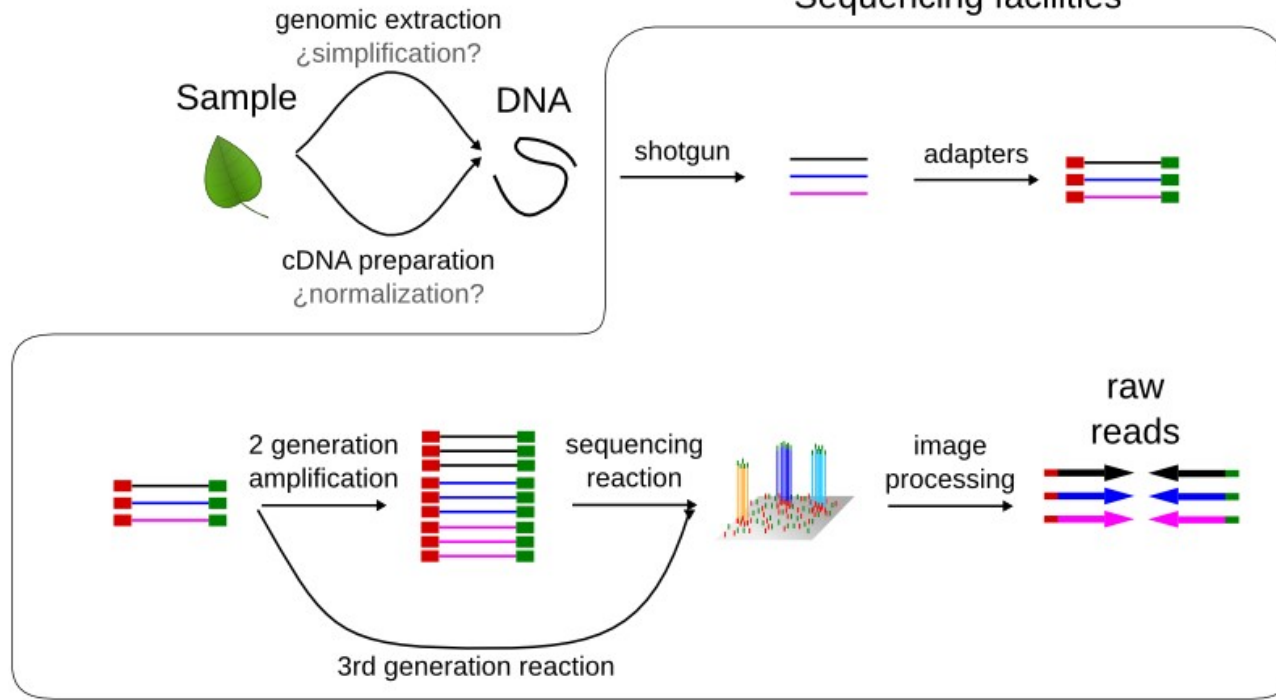


# Sequence assembly

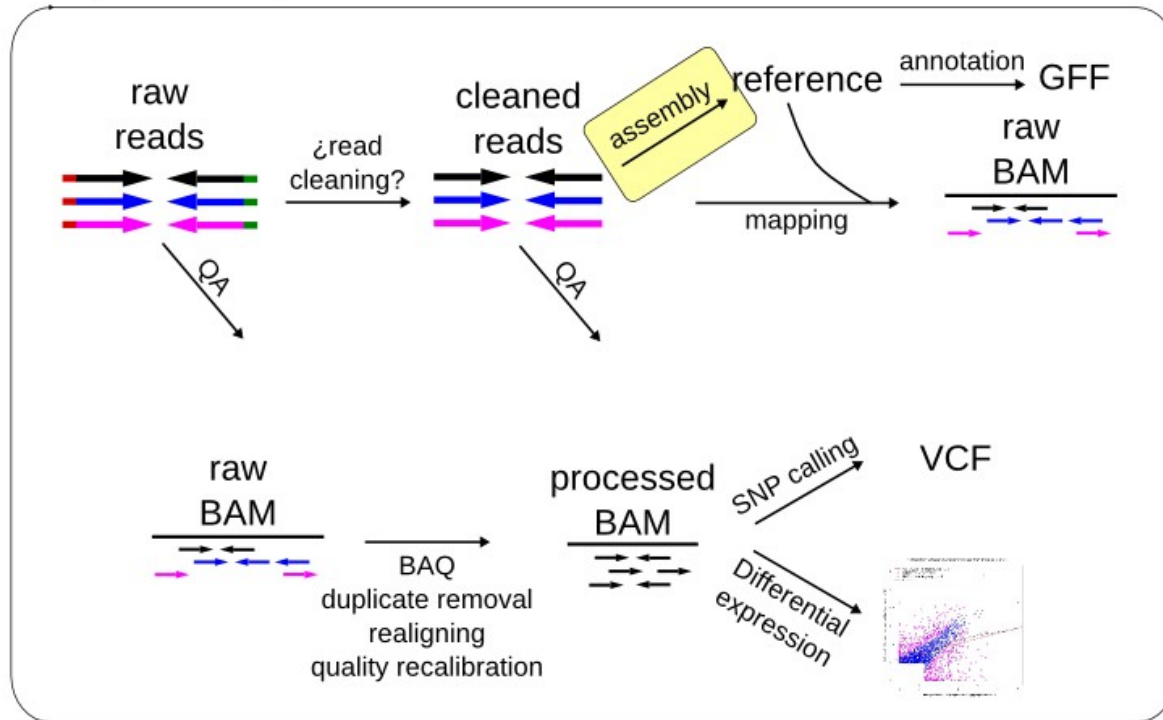
Jose Blanca  
COMAV institute  
[bioinf.comav.upv.es](http://bioinf.comav.upv.es)



## Sequencing facilities



## Sequence analysis



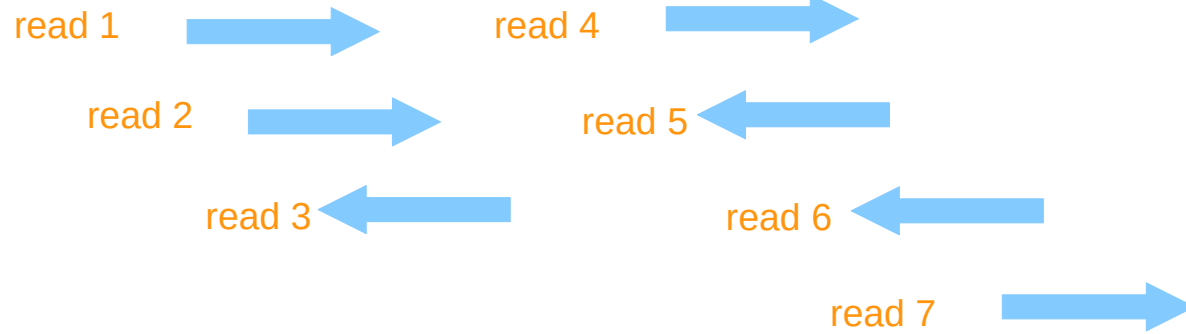
# Assembly project

---

Unknown sequence



experimental  
evidence

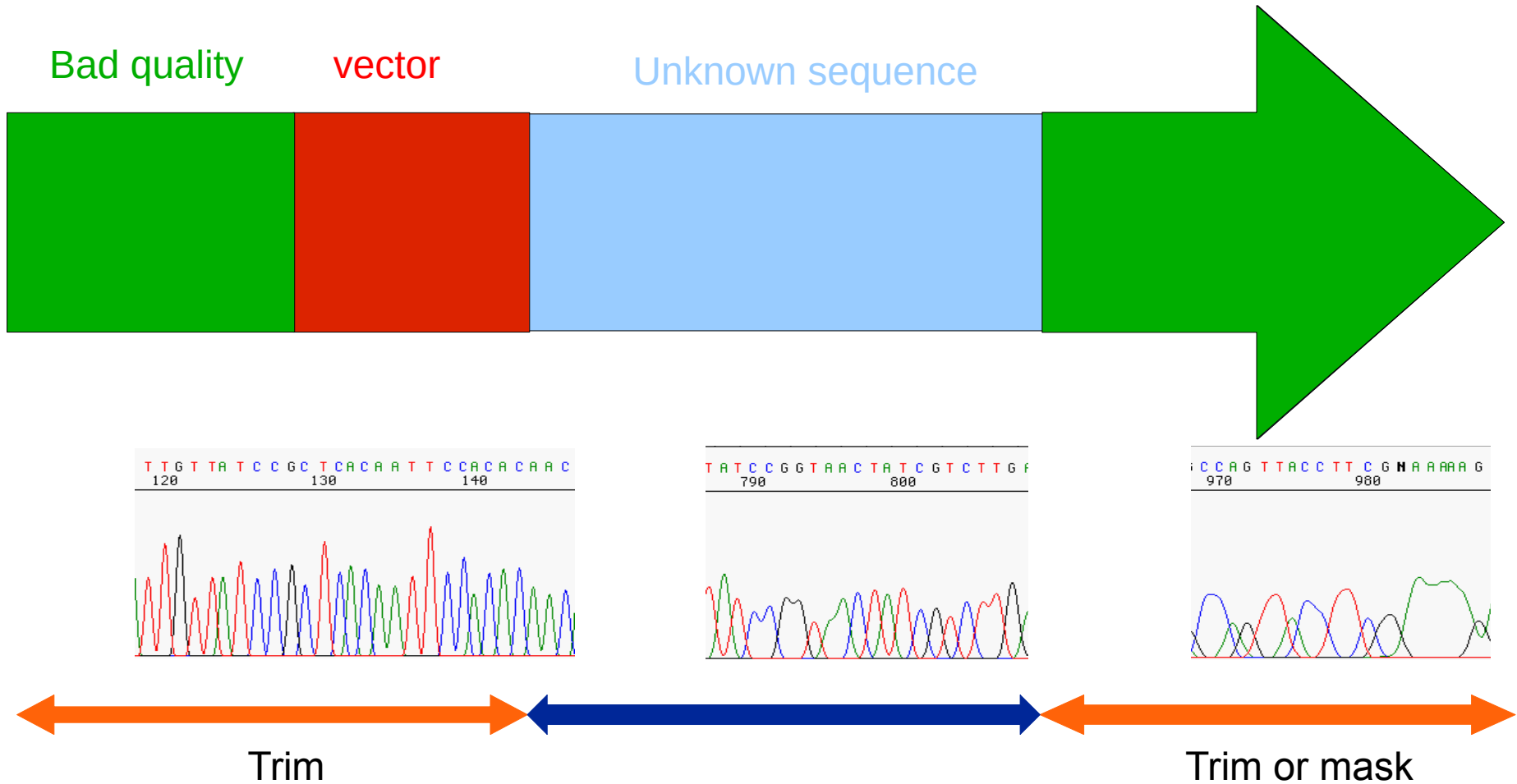


result



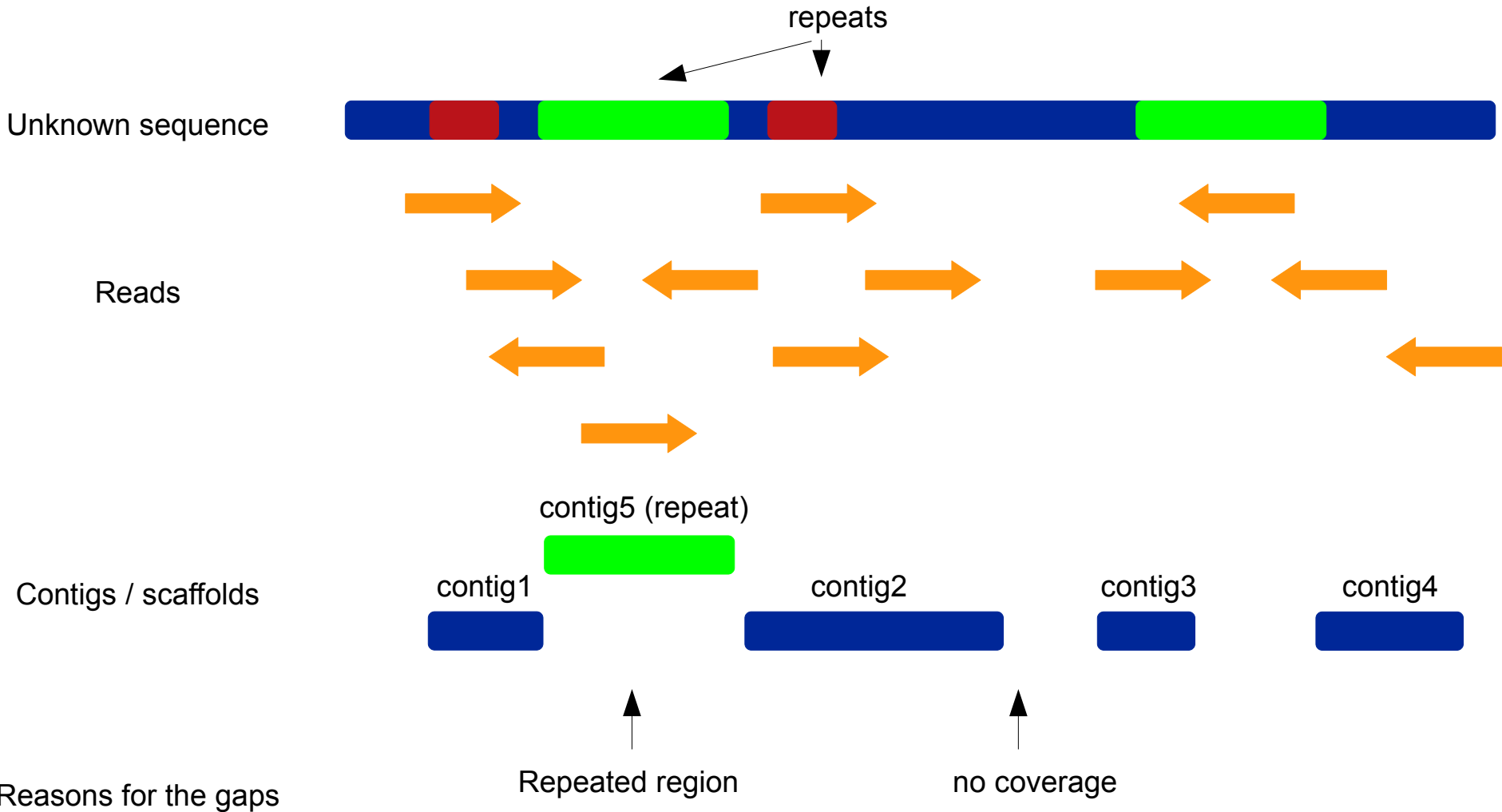
# Read cleaning

Some assemblers choke with: bad quality stretches, adapters, low complexity regions, contamination.



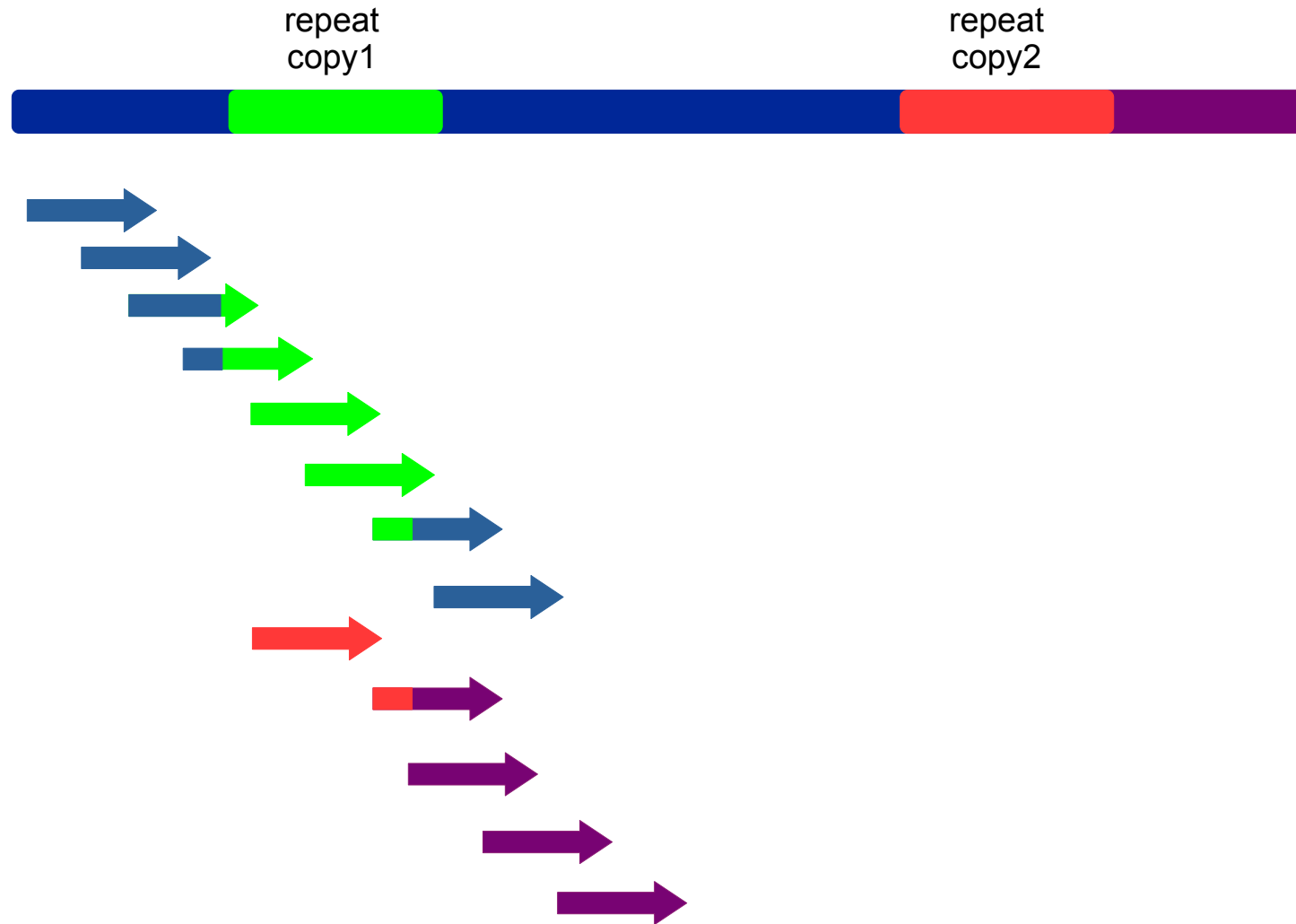
# The repeat assembly problem

Only one read is usually not capable of producing the complete sequence.  
Even if the problem sequence is short one read might have sequencing errors.



# The repetitive problem is unsurmountable

---



# Long reads

# Read size influence

---



In the ideal case each piece would be a chromosome or a transcript



# Mate pairs and paired-ends

---

Read length is critical, but constrained by sequencing technology, we can sequence molecule ends.



Useful for dealing with repetitive genomic DNA and with complex transcriptome structure.

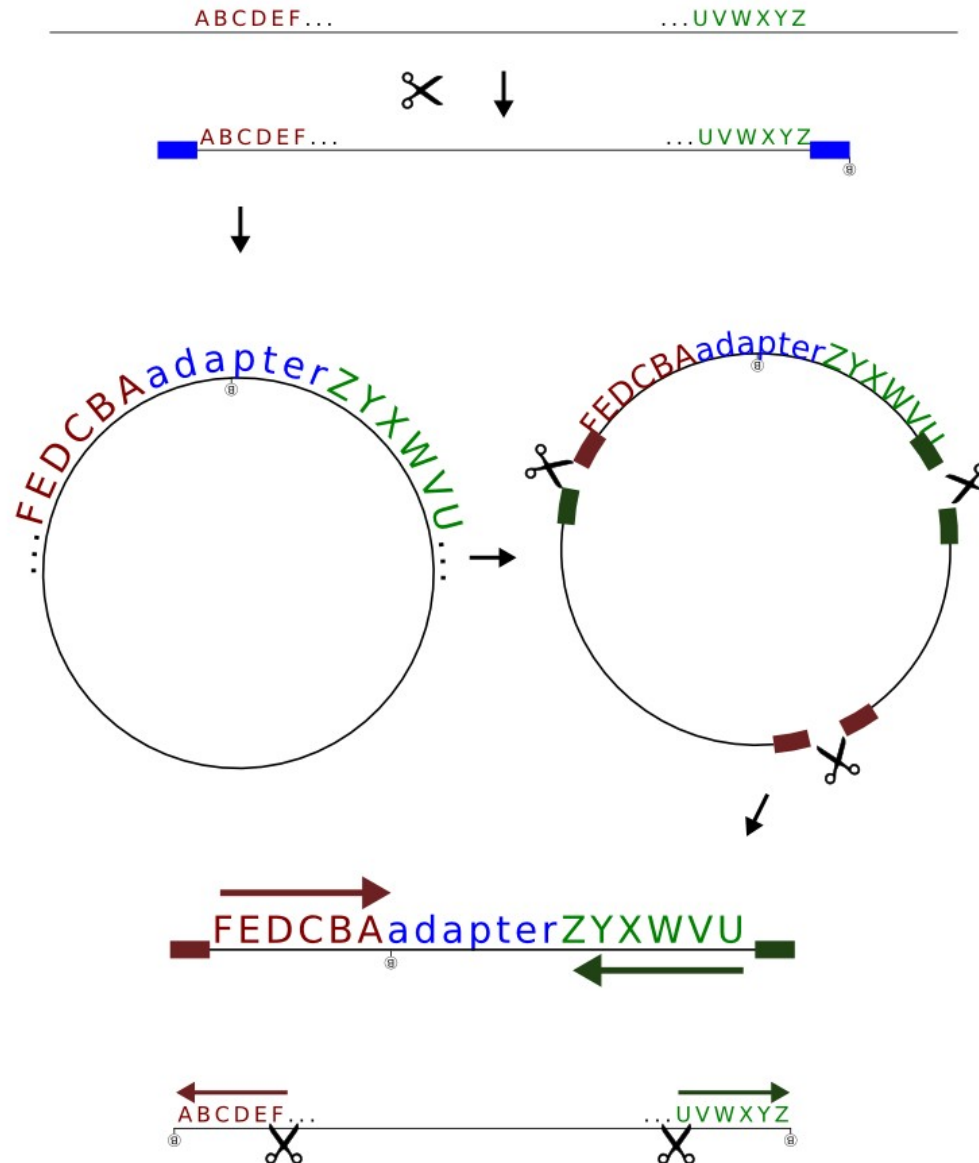
**Paired-ends:** Illumina can sequence from both ends of the molecules. (150-500 bp)

**Mate-pairs:** Can be generated from libraries with different lengths (2-20 Kb).

**BAC-ends:** Usually sequenced with Sanger.

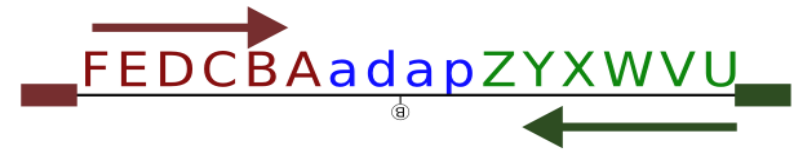
# Mate-pairs

## Nextera mate-pair (Illumina)



# Illumina mate-pair chimeras

nextera mate-pairs



chimeric



# Long reads

---

Pacbio:

- Standard
- HiFi

Nanopore:

- Long
- Ultralong

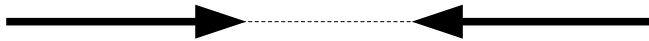
# Library types

---

Single reads



Illumina Pair Ends (150-500 pb)



Mate Pairs (2-10 kb)

Illumina



# Long reads

---

	<i>Illumina TruSeq Synthetic Long Reads</i>	<i>Pacific Biosciences</i>	<i>Oxford Nanopore Technologies</i>
Technology	Barcoded & Amplified Synthetic long reads	Single Molecule Real Time Sequencing	Nanopore Sequencing
Mean Length	3-5kbp	10-15kbp	5-10kbp
Raw Error Rate	0.1%	10-15%	10-30%
Costs / GB	~\$2500*	~\$500†	~\$1000†

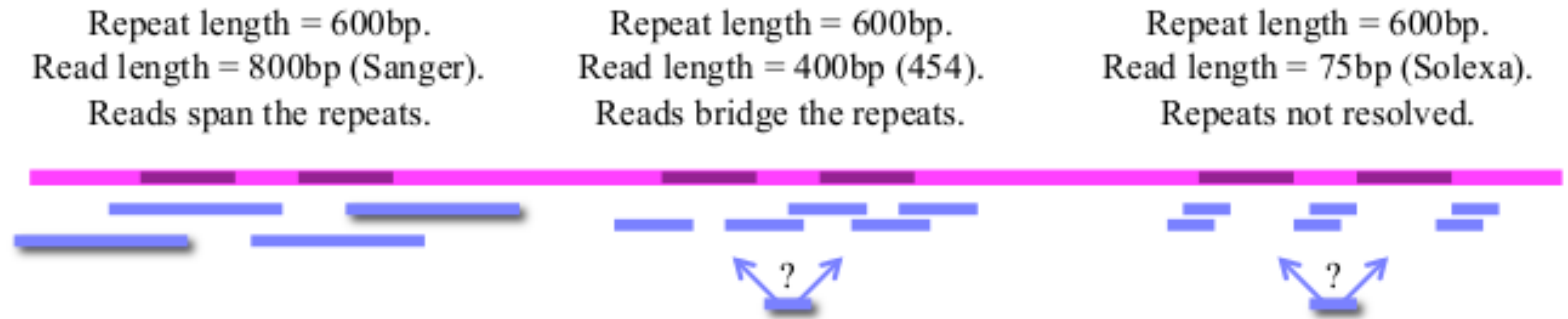
Third-generation sequencing and the future of genomics. Lee et al.

# Short reads are harder to assemble

1. **Overlap Effect:** For same number of sequenced bases, shorter reads require more coverage to achieve comparable N50.



2. **Repeat Effect:** Shorter reads resolve fewer repeats.



# The need for paired reads

---

## 1. Variety of insert sizes will span variety of repeats.

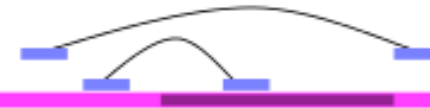
Inserts that span the repeat will enable scaffolds.



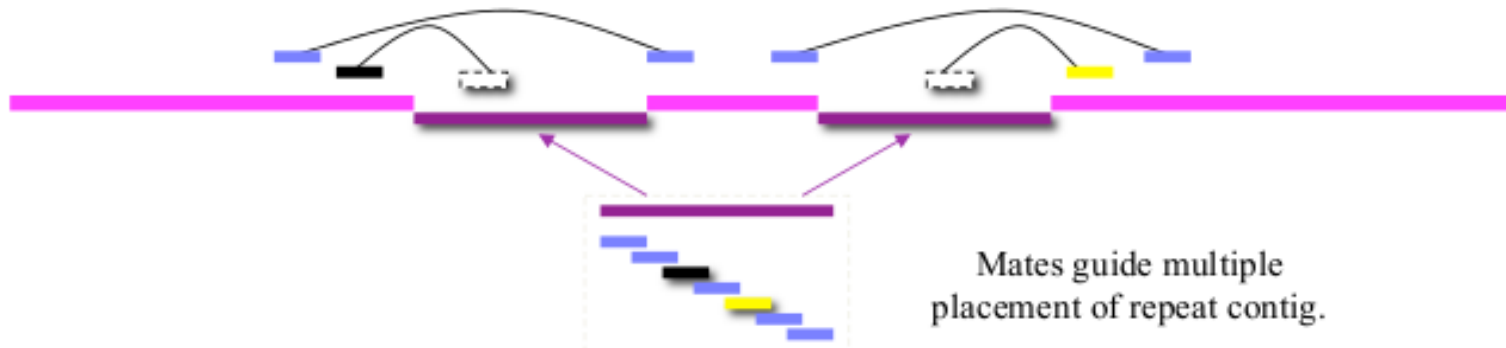
High coverage in mates will tile the repeat.



Larger repeats require larger insert sizes.



## 2. Mates can resolve repeats even if not possible to tile with reads.





# **Algorithm I**

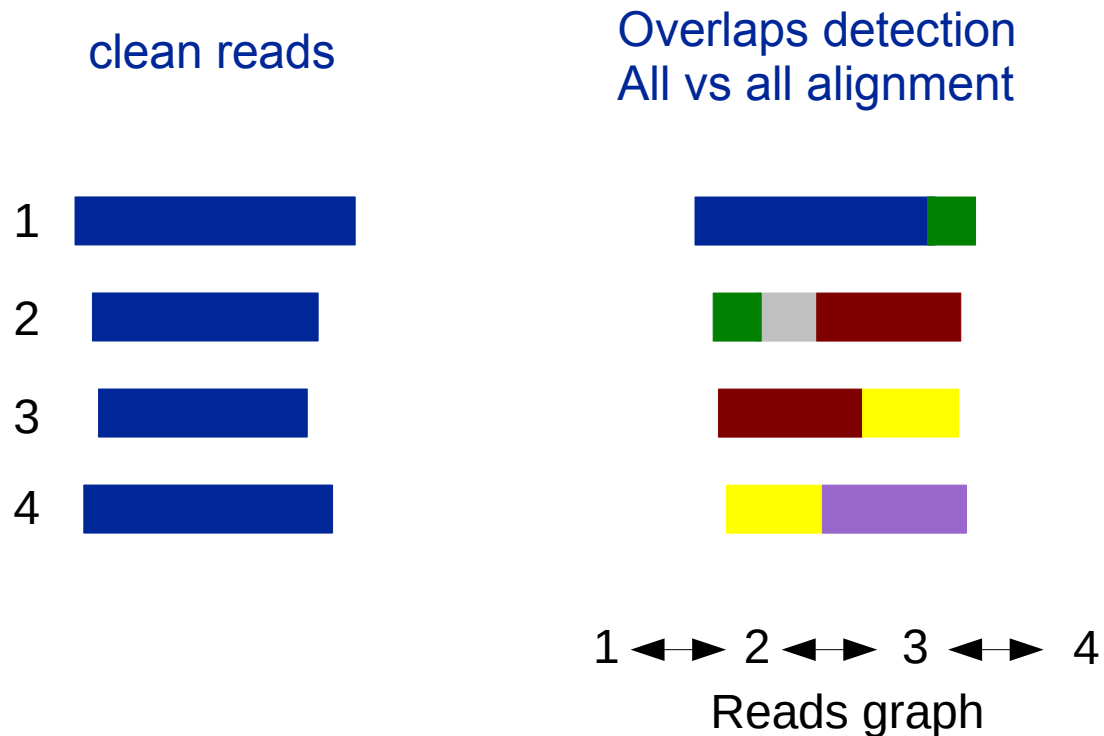
## **Overlap – Layout - Consensus**

# Overlap – Layout - Consensus

---

Algorithm:

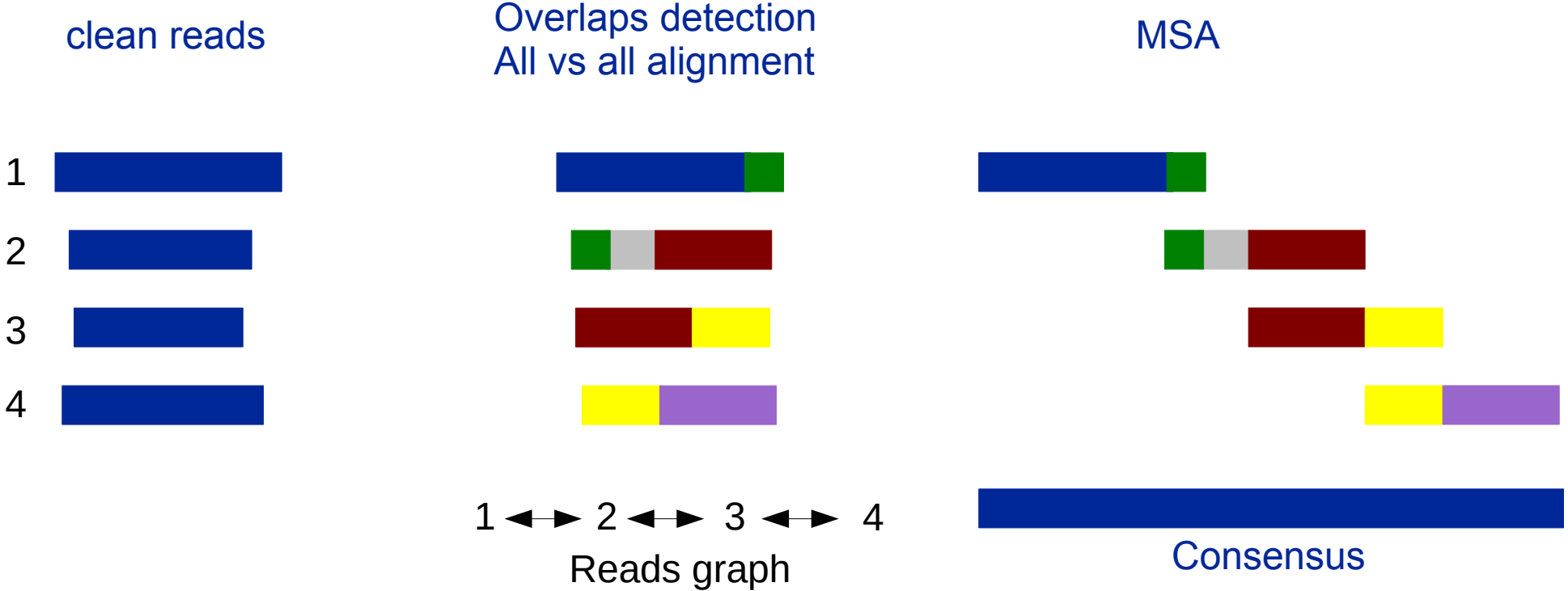
- Overlap: All-against-all, pair-wise read comparison.
- Layout: construction of an overlap graph with approximate read layout.



# Overlap – Layout - Consensus

Algorithm:

- All-against-all, pair-wise read comparison.
- Construction of an overlap graph with approximate read layout.
- Consensus: Multiple sequence alignment (MSA) determines the consensus sequence.

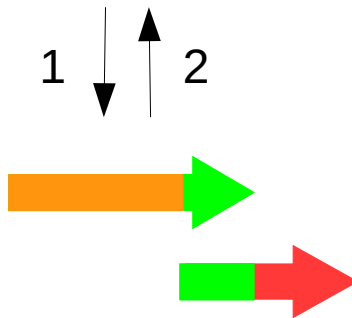


# Overlap assumption

---

We are assuming:

- 1) Two reads originated from the same region will overlap
- 2) Two reads that overlap come from the same region



# Overlap problems

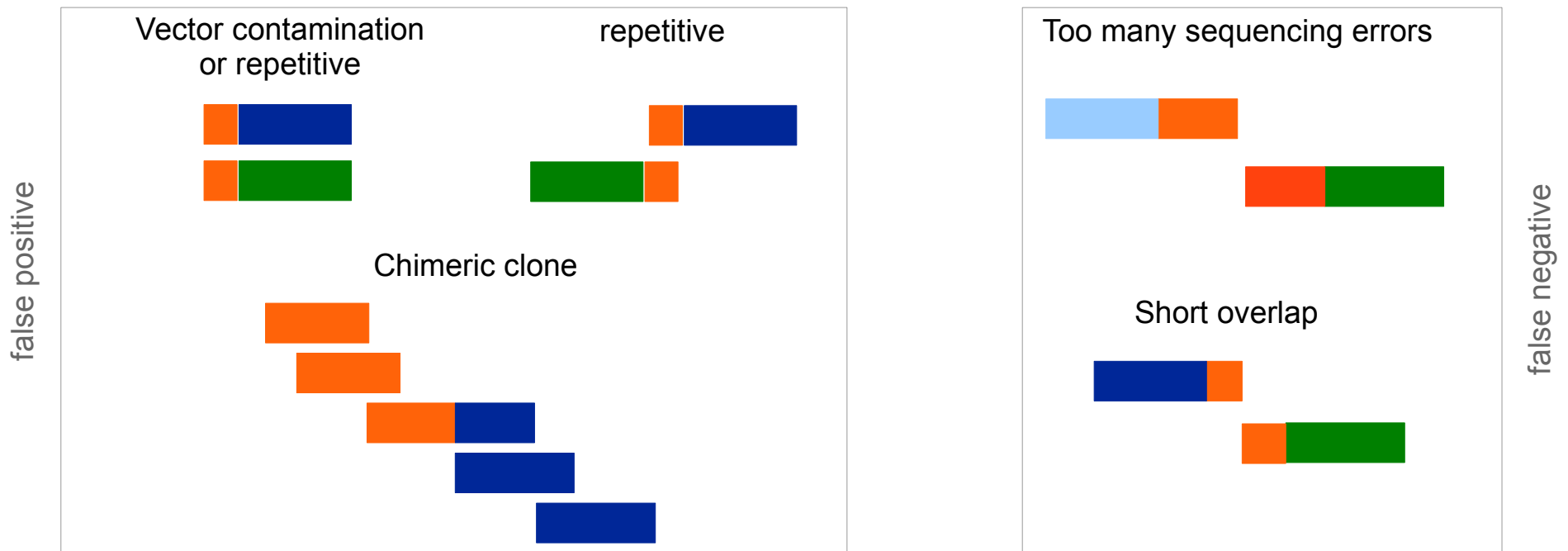
Two reads are similar if they have a good overlap.

We assume that overlap implies common origin in the genome.

This goodness depends on the overlap quality and similarity.



But several things might go wrong



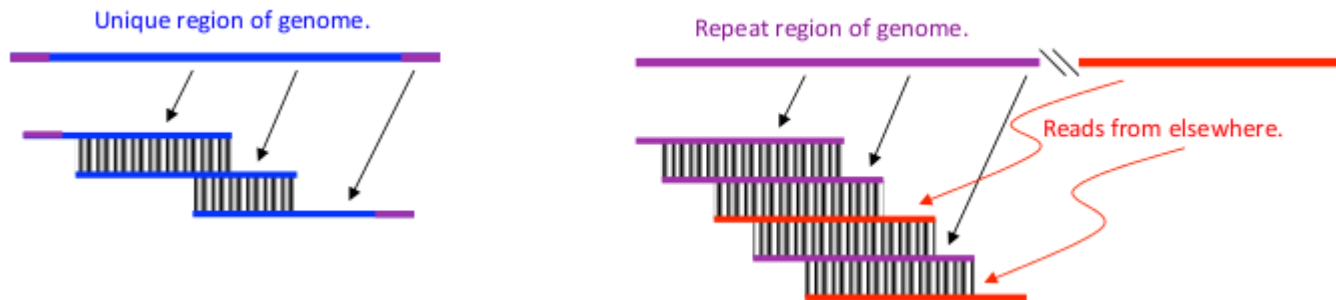
# Overlap problems

---

False positives will be induced by chance and repeats

Avoid false positives with stringent criteria:

- Overlaps must be long enough
- Sequence similarity must be high (identity threshold)
- Overlaps must reach the ends of both reads
- Ignore high-frequency overlaps (repetitive regions in genomic?)



But stringency will induce false negatives.

# Assembly terminology

Consensus

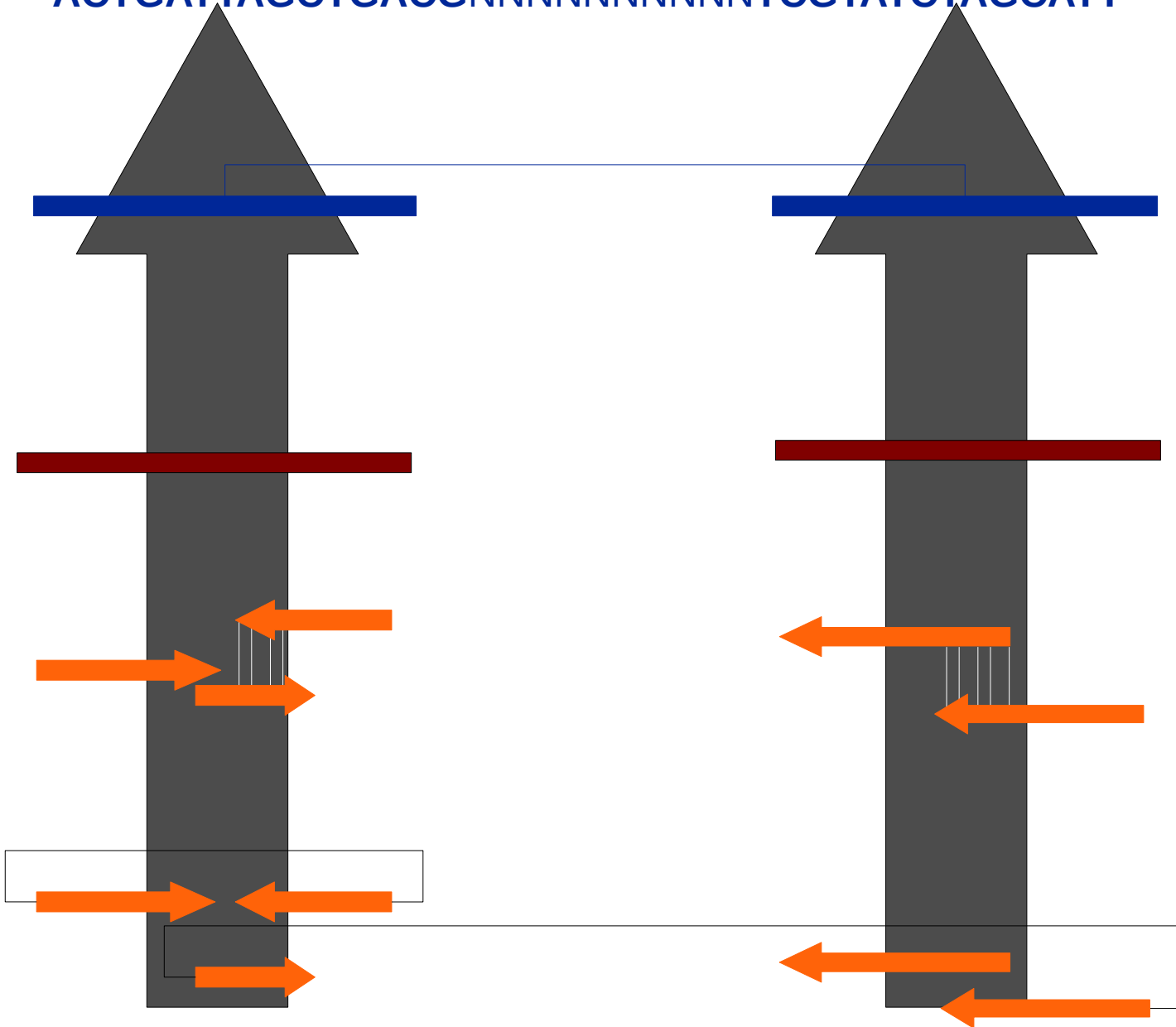
ACTGATTAGCTGACGNNNNNNNNNNNTCGTATCTAGCATT

Scaffolds = Contigs +  
Pairings  
(gap lengths derived from pairings)

Contigs = Reads +  
Overlaps

Overlaps

Reads & Pairing

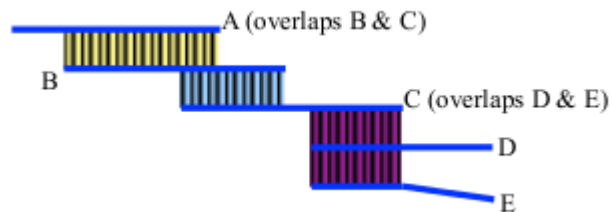


# Contigs

---

## Contigs:

- High-confidence overlaps
- Maximal contigs with no contradiction (or almost no contradiction) in the data
- Ungapped
- Contigs capture the unique stretches of the genome.
- Usually where unitigs end, repeats begin



Contig example. The ideal Contig is A, B, C.  
C, D and E have a repetitive sequence.



# Scaffolds

---

## Build with mate pairs

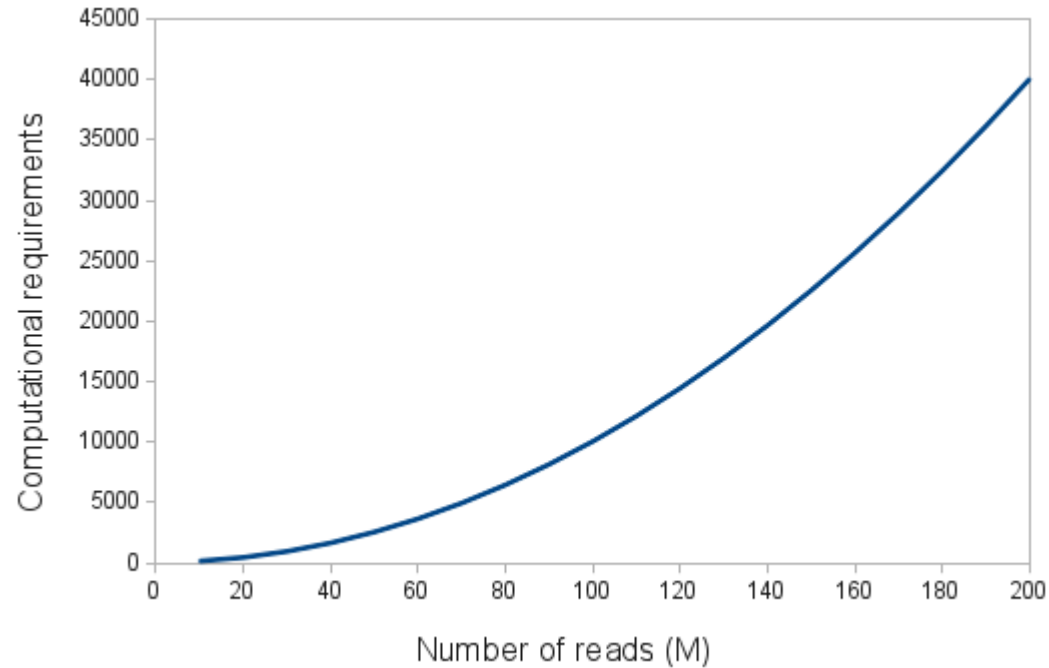
- Mate pairs help resolve ambiguity in overlap patterns

## Scaffolds

- Every contig is a scaffold
- Every scaffold contains one or more contigs
- Scaffolds **can contain** sequence **gaps** of any size (including negative)

# Overlap-Consensus-Layout limitation

---



30 Million 5Kb (long) reads

Number of alignments: 30 M reads x 30 M reads = **900 trillion alignments (It does not scale!!)**

Main requirements:

Memory.

Time.

The kmer-based assemblers requirements depend on the genome complexity and not so much on the reads.

# Algorithm II

kmer-based

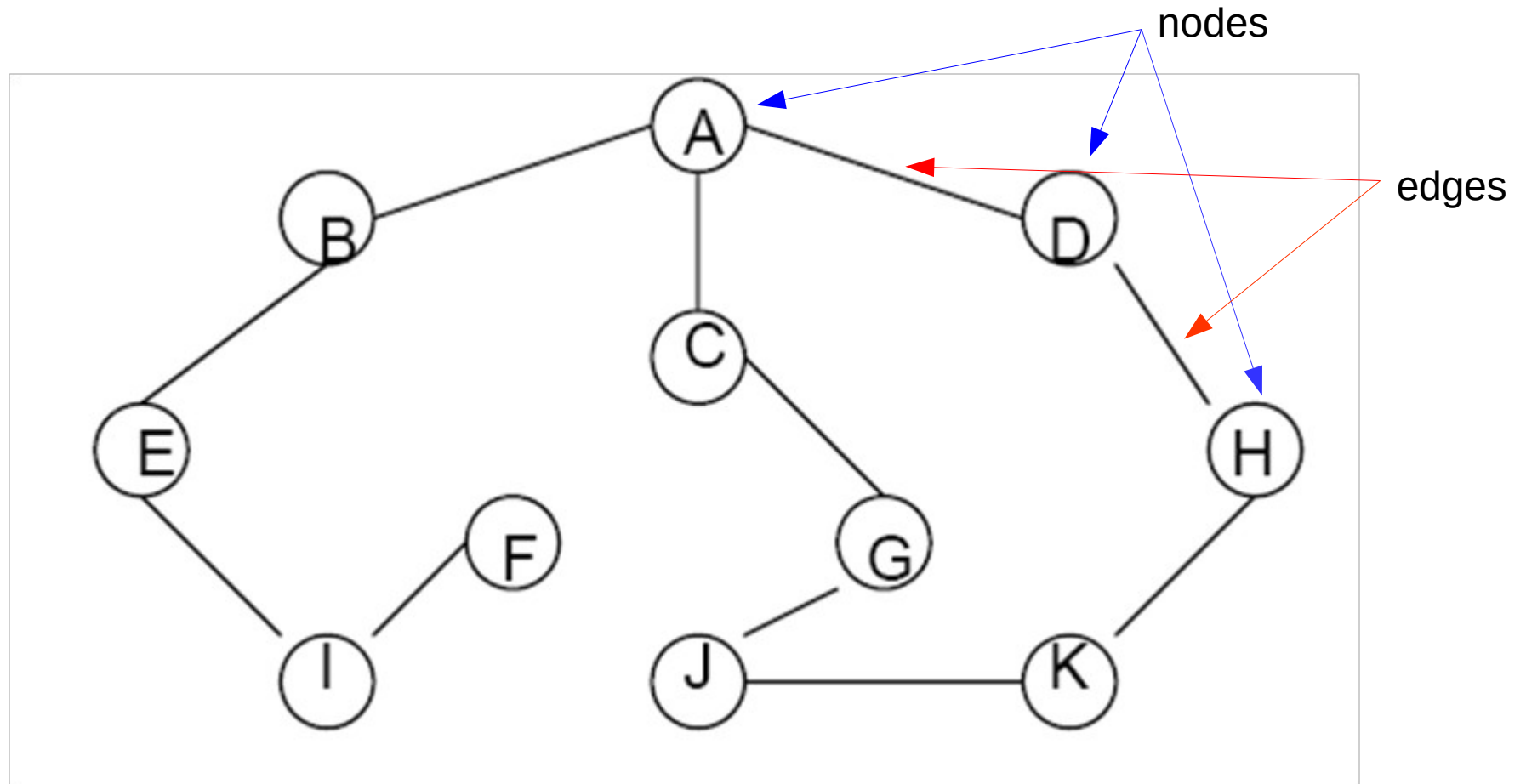
# K-mer based assembly

---

Seq		ACTGGTCAT
K-mers	5pb	ACTGG
		CTGGT
		TGGTC
		GGTCA
		GTCAT

# K-mer based assembly

K-mer graphs are de Bruijn **graphs**.



# K-mer based assembly

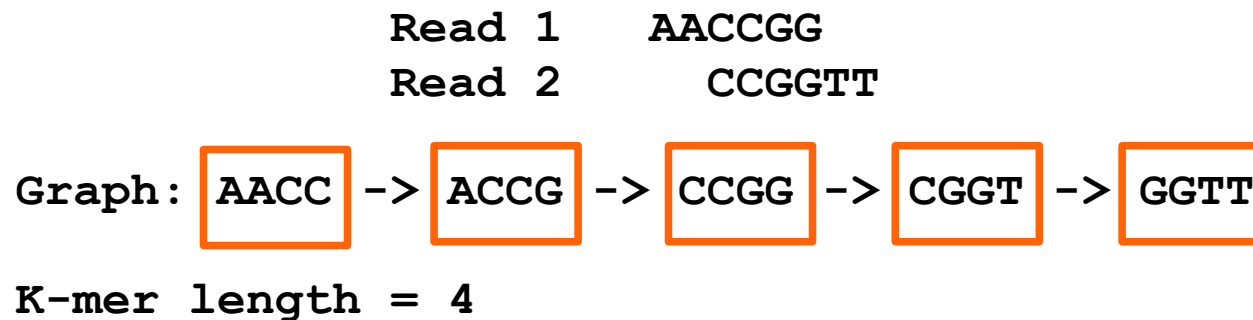
---

K-mer graphs are de Bruijn graph.

Nodes represent all K-mers.

Edges join consecutive K-mers.

Assembly is reduced to a graph reduction problem.



# K-mer based assembly

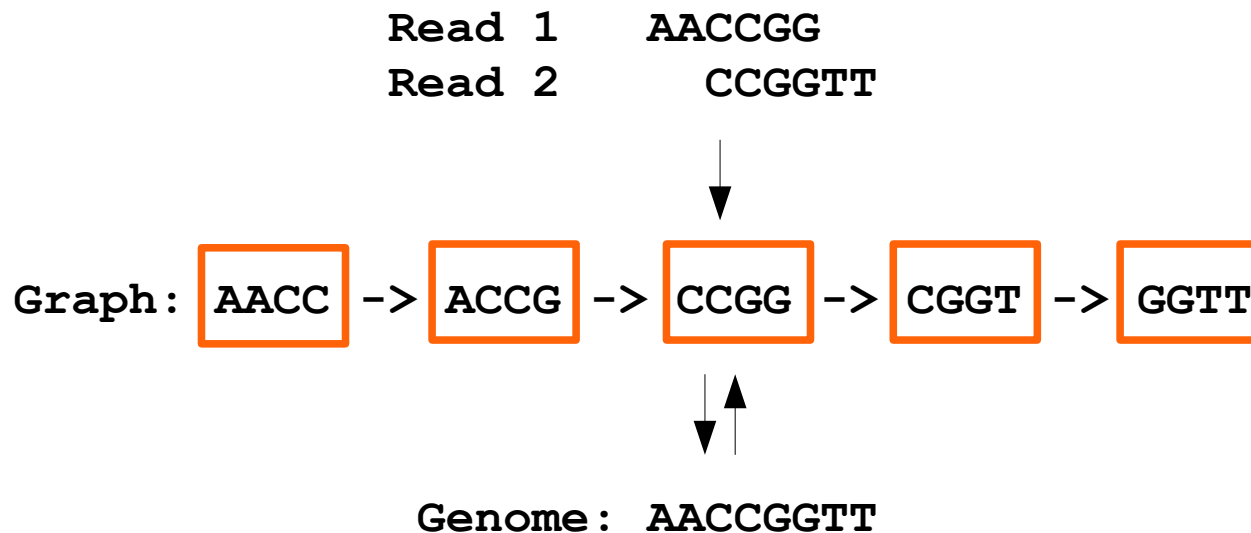
---

K-mer graphs derived from reads and genome is very similar

From the graph we can recreate the genome sequence

K-mer graph size depends, mostly on genome size, not on number of reads

Reads are only read once to prepare the kmer graph



# K-mer based assembly

---

reads

A C T G A T A  
A T A C G T T C  
G T T C C A G G



Kmer graph ACTG→CTGA→TGAT→GATA→ATAC→TACG→ACGT→CGTT→GTTC→TTCC→TCCA→CCAG→CAGG



Genomic A C T G A T A C G T T C C A G G



# K-mer bubbles

Genomic A C T G A T A **C/G** G T T C C A G G



reads

A C T G A T A  
 G A T A **C** G T T C  
 G A T A **G** G T T C  
 G T T C C A G G



Kmer graph

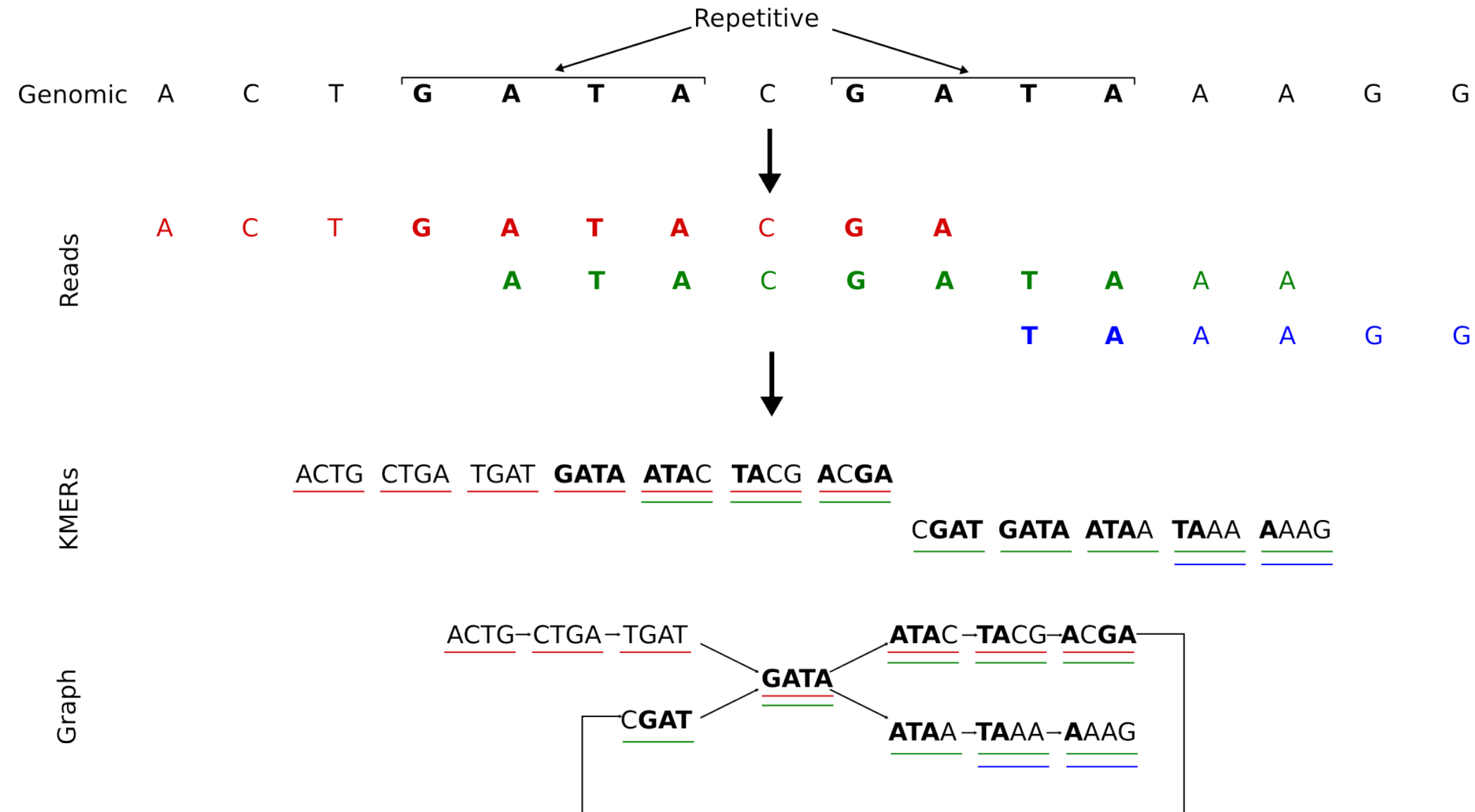
ACTG→CTGA→TGAT→GATA { ATAC→TACG→ACGT→CGTT } GTTC→TTCC→TCCA→CCAG→CAGG  
 { ATAG→TAGG→AGGT→GGTT }

Het. bubble

ACTG→CTGA→TGAT→GATA { ATAC→TACG→ACGT→CGTT } GTTC→TTCC→TCCA→CCAG→CAGG  
 { ATAG→TAGG→AGGT→GGTT }

Seq. error bubble

# K-mer repetitive



# Graph problems

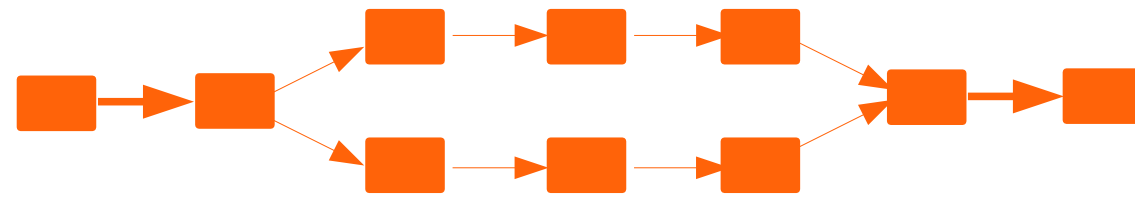
---

Repeats, sequencing errors and polymorphisms increase graph complexity, leading to tangles difficult to resolve

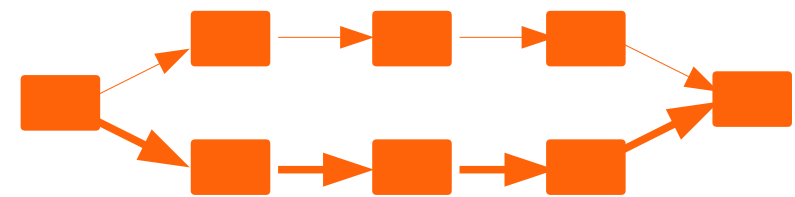
K-mer graphs are more sensitive to repeats and sequencing errors than overlap based methods.

Optimal graph reductions algorithms are NP-hard, so assemblers use heuristic algorithms

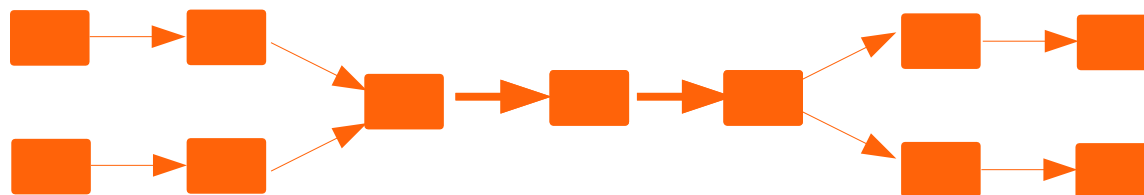
Polymorphism



Sequencing error



Repeat



# Quality Assessment

# N50

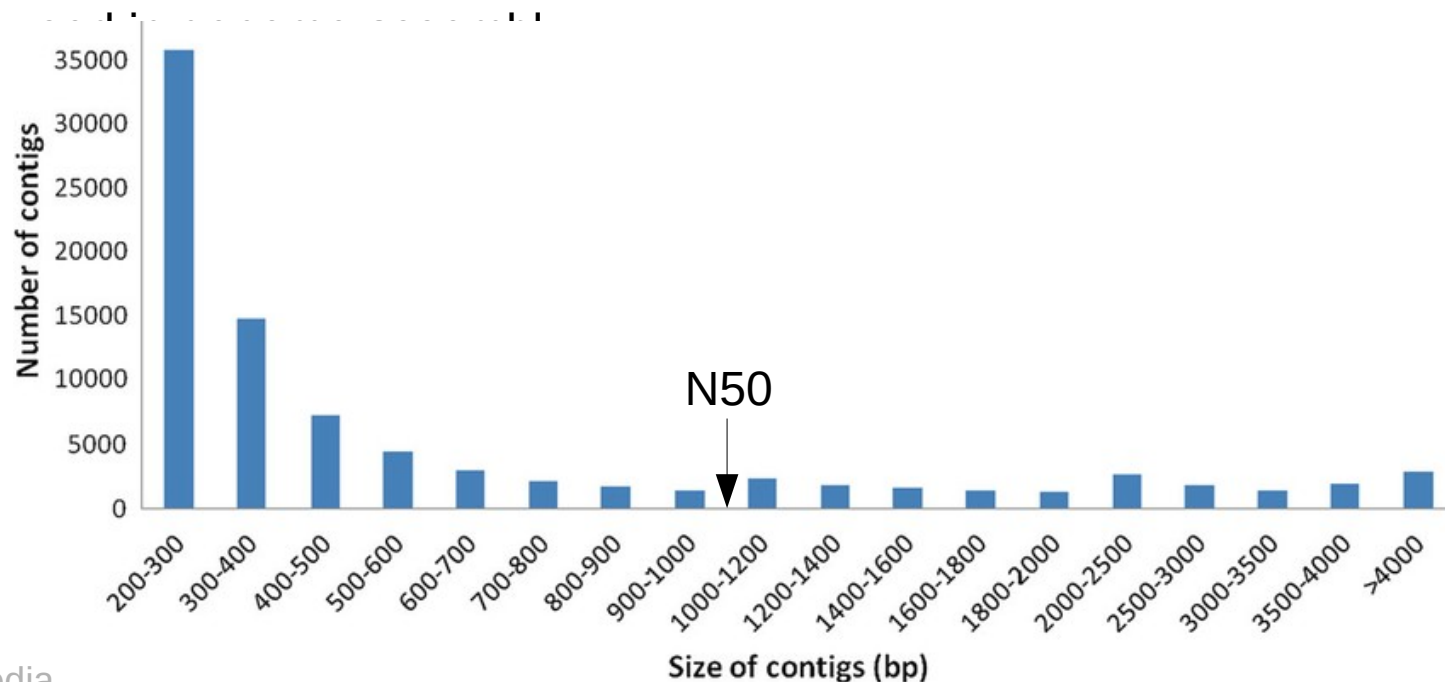
---

N50 is defined as the contig length such that using equal or longer contigs produces half the bases of the assembly.

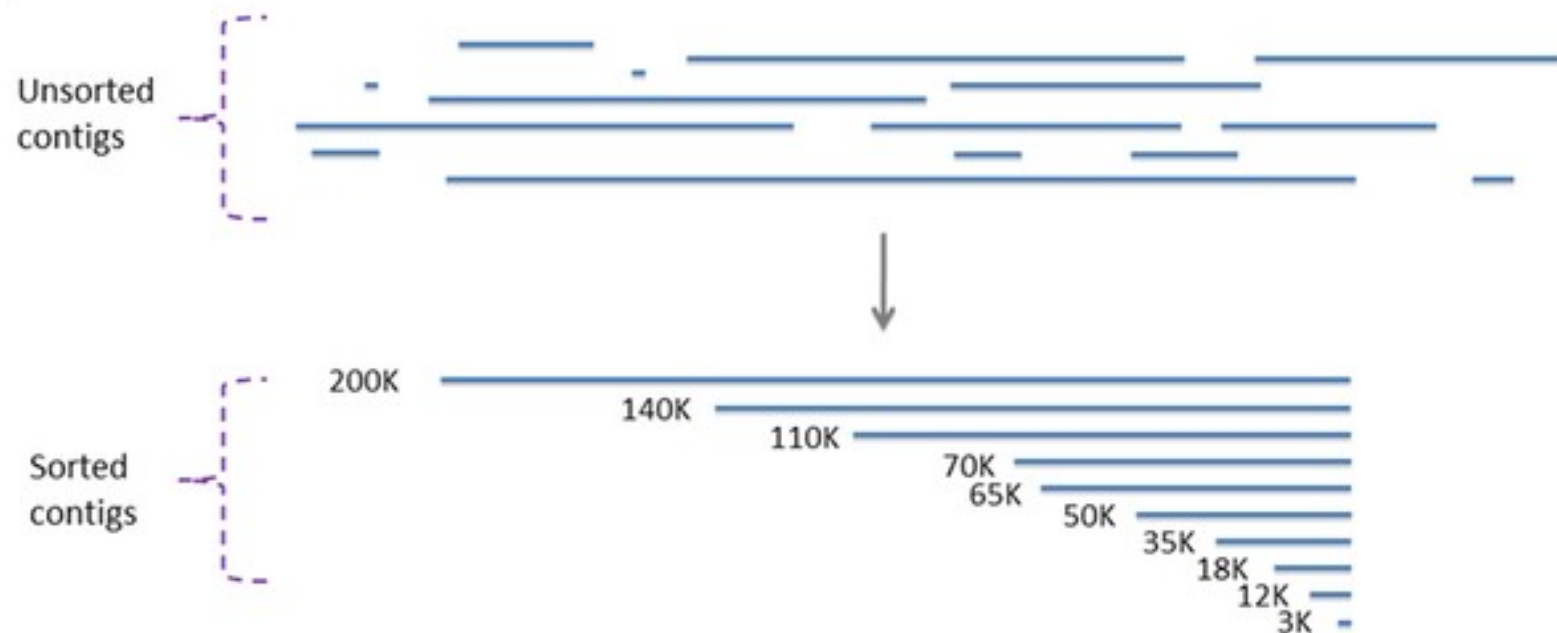
- NG50: 50% of the genome (if available)

The N50 statistic is a measure of the average length of a set of sequences, with greater weight given to longer sequences.

It is widely



# N50



Total contig length =  $200\text{K} + 140\text{K} + 110\text{K} + 70\text{K} + 65\text{K} + 50\text{K} + 35\text{K} + 18\text{K} + 12\text{K} + 3\text{K} = 703\text{K}$

50% total contig length =  $703\text{K} \times 50\% = 351.5\text{K}$

$\therefore 200\text{K} + 140\text{K} + 110\text{K} > 351.5\text{K}$ ,  $\therefore \text{N50} = 110\text{K}$

# N50

---

short N50



Long N50



# N50 for scaffolds and contigs

---

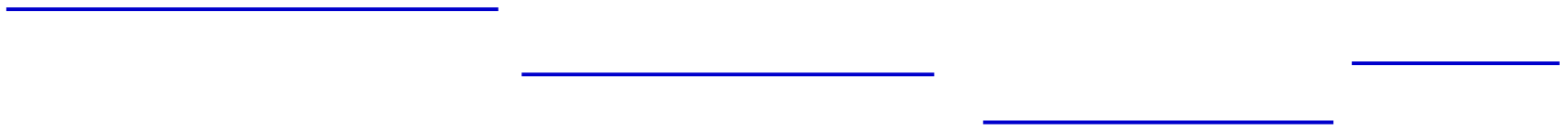
Long Scaffold N50



Short Contig N50



short Scaffold N50



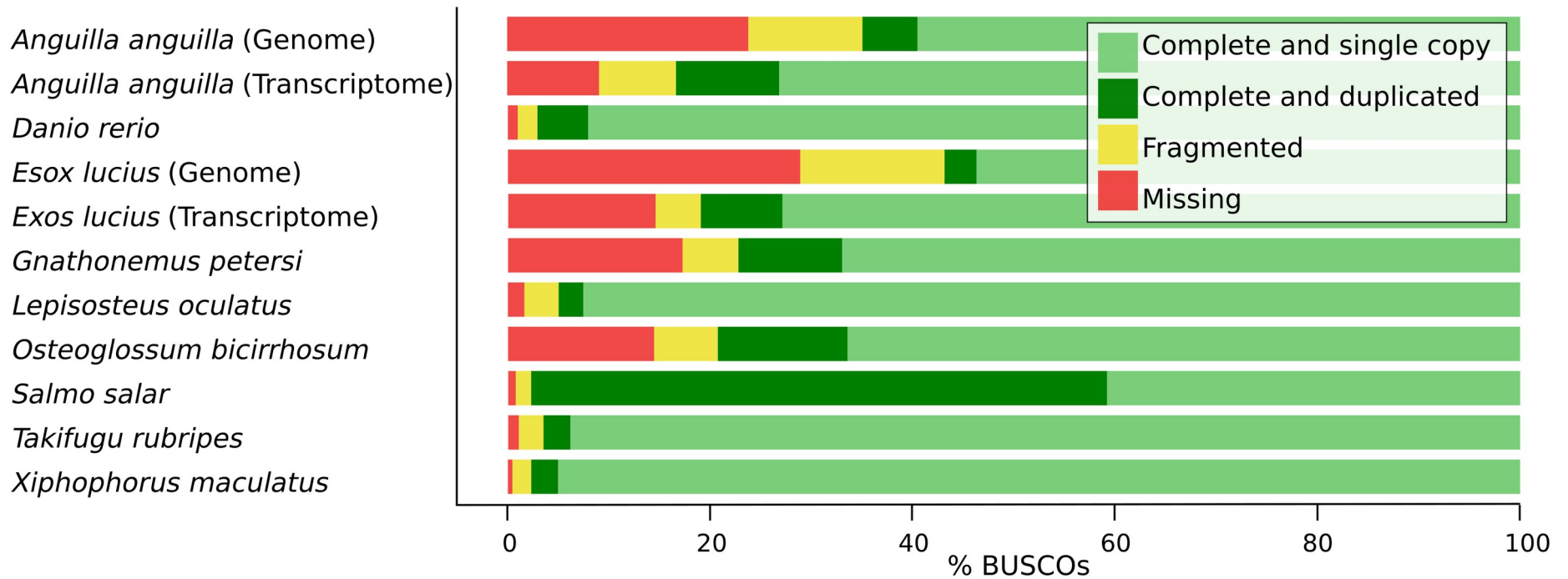
Long Contig N50





# BUSCO orthologs

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



# Mapping reads against the assembly

---



Reads used to create the assembly

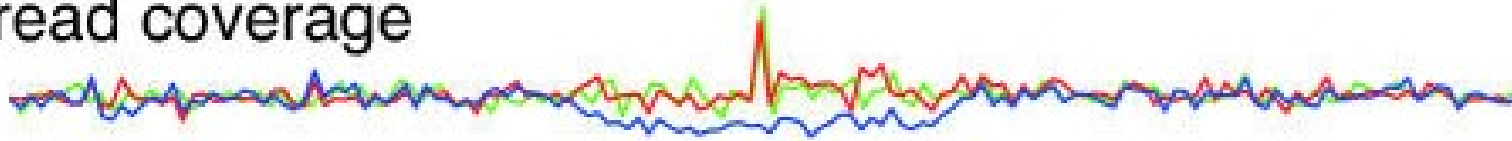
Other reads from the same species

Percentage of reads mapped

- Most reads should properly map

# Mapping reads against the assembly

## i read coverage



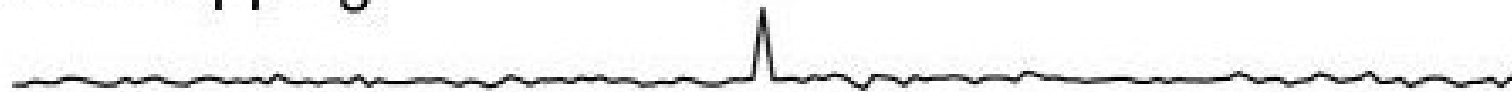
## ii type of read coverage, on each strand



Properly paired

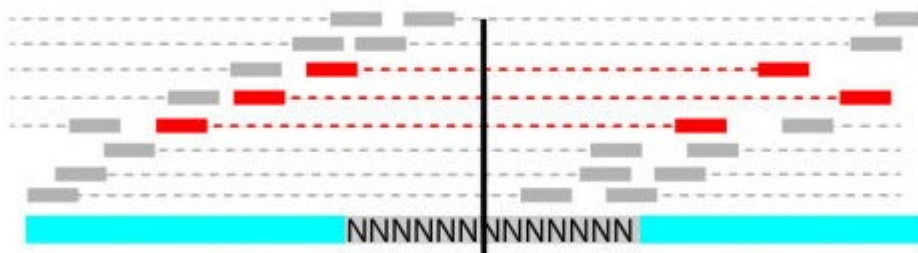
Orphaned

## iii read clipping



Reads only  
partially aligned

If the base of interest lies in a gap



REAPR corrected N50

# Main assembly hurdles

---

## Main Problems:

- Short, inaccurate sequence reads
- genomic repeats
  - Difficulty of the assembly depends a lot on the genome: easy for bacteria, very difficult for long highly repetitive polyploids

## Main solutions:

- Long reads
- Low sequencing error rates:
  - Simplify graph problem
  - Complex graphs typically yield worse assemblies
  - Allow to differentiate between repeats

# Common assembly errors

---

## Collapsed repeats

- Align reads from distinct (polymorphic) repeat copies
- Fewer repeats in assembly than in genome
- Fewer tandem repeat copies in assembly than genome

## Missed joins

- Missed overlaps due to sequencing error
- Contradictory evidence from overlaps
- Contradictory or insufficient evidence from pairs

## Chimera

- Enter repeat at copy 1 but exit repeat at copy 2
- Assembly joins unrelated sequences

# Ingredients for Good Assembly

---

## Coverage and read length

- 20X for (corrected) PacBio, 150X for Illumina

## Paired reads

- Read lengths long enough to place uniquely
- Inserts longer than long repeats
- Pair density sufficient to traverse repeat clusters
- Tight insert size variance
- Diversity of insert sizes

## Read Quality:

- Sequencing errors
- No vector contamination
- Contamination: mitochondrial or chloroplastic

The requirements depend greatly on the quality: it is not the same a draft that high quality genome.

# Common assemblers

---

## Genomic

- SOAPdenovo, Illumina
- Canu: Pacbio and Illumina
- 

Staden for Sanger reads.

Transcriptomic assemblers are specialized software

# Assembly comparison

---

Empirical evaluation of methods for de novo genome assembly

- <https://peerj.com/articles/cs-636/>

The best assemblers and parameters depend a lot on the genome and on the information available

*Arabidopsis thaliana*:

- SOAPdenovo2 produces much larger contigs than any other assembler
- it has many assembly errors
- SPAdes appeared to be preferable

*Bacillus cereus*

- SOAPdenovo2 has the fewest contigs, as well as a low N50 value
- its error rate was among the best

Human genome

- HiFiasm outperformed other assemblers
- problems with mismatches and misassemblies.

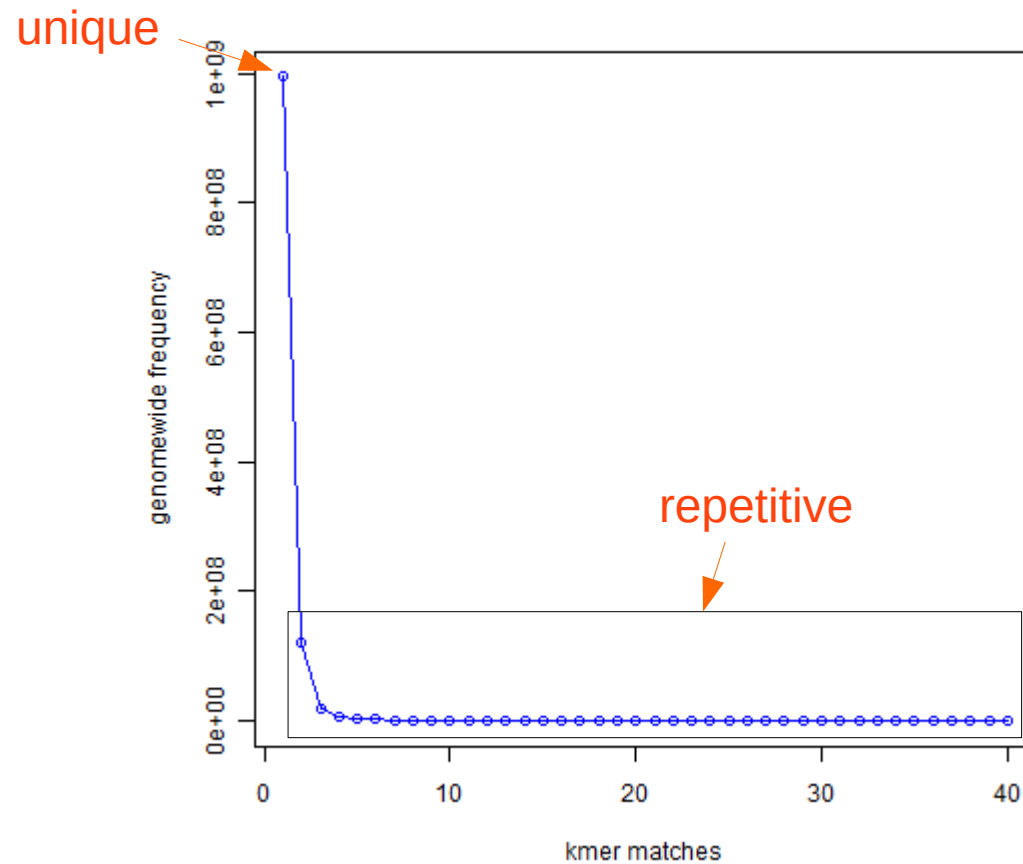
Flye has consistently demonstrated superior output based on contig size, with the trade-off between scale and error rate. Hinge and Canu, though they demonstrate more errors than Flye, both performed reasonably well too.



# **K-mer analysis**

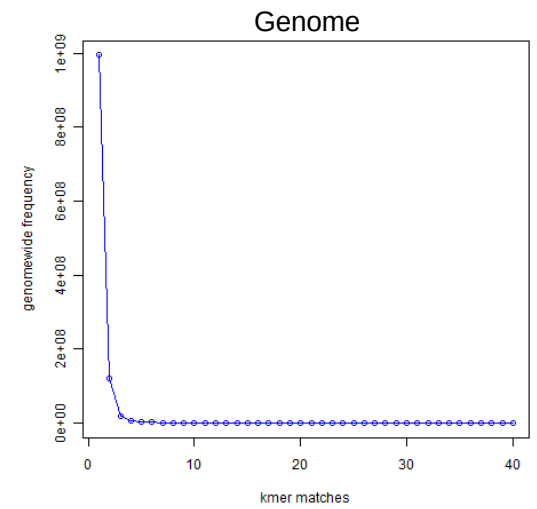
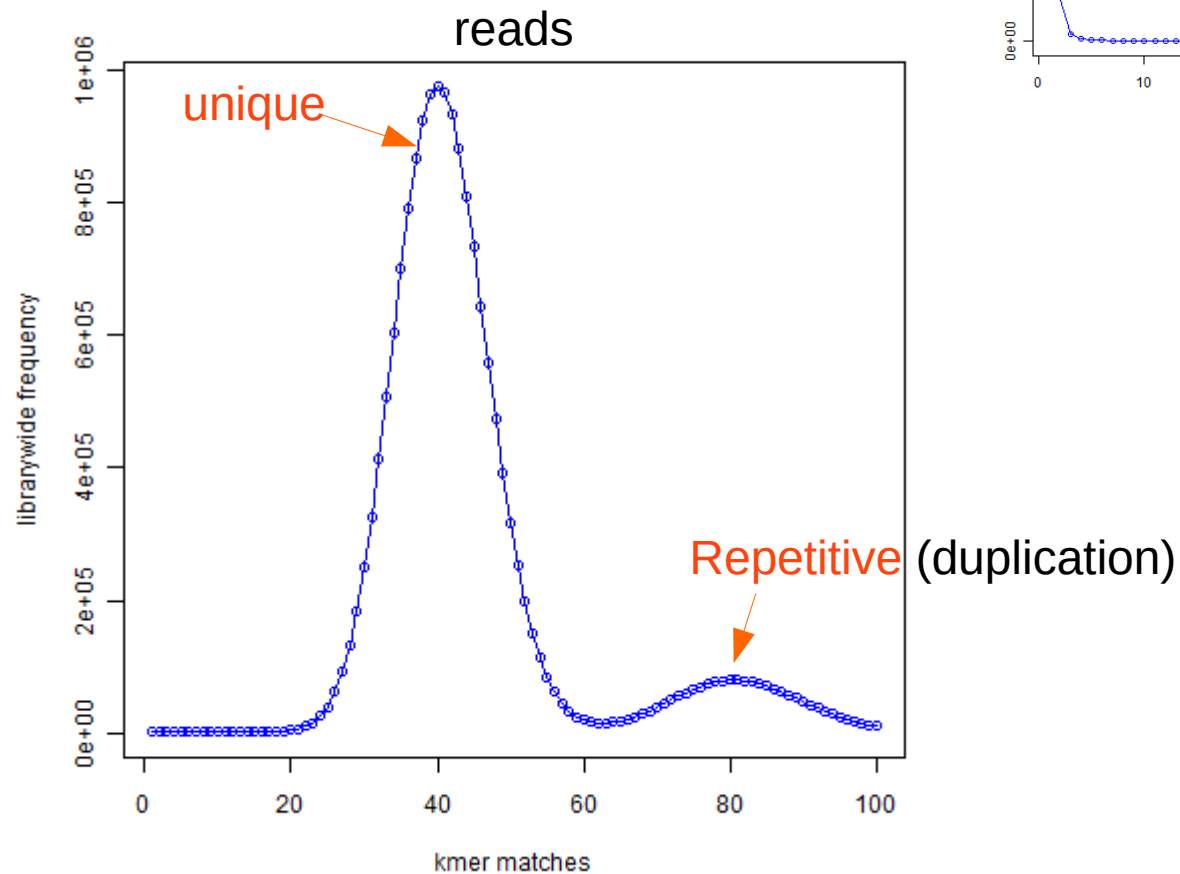
# K-mer analysis

21-mer distribution of sea urchin  
(*Strongylocentrotus purpuratus*) genome (900  
Mbase long)



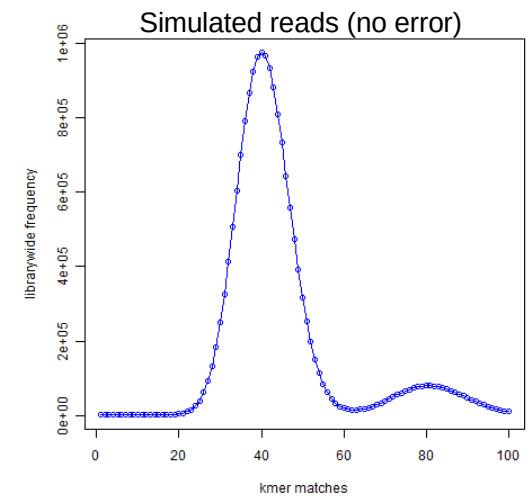
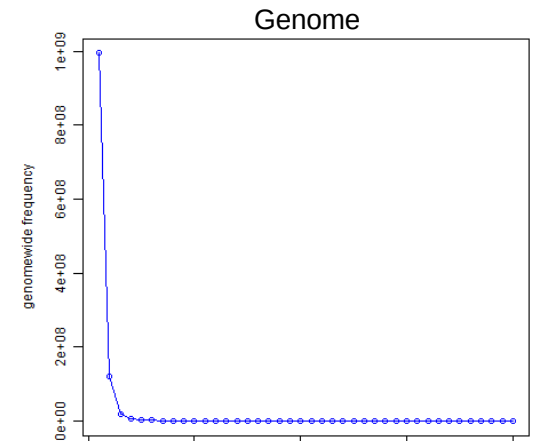
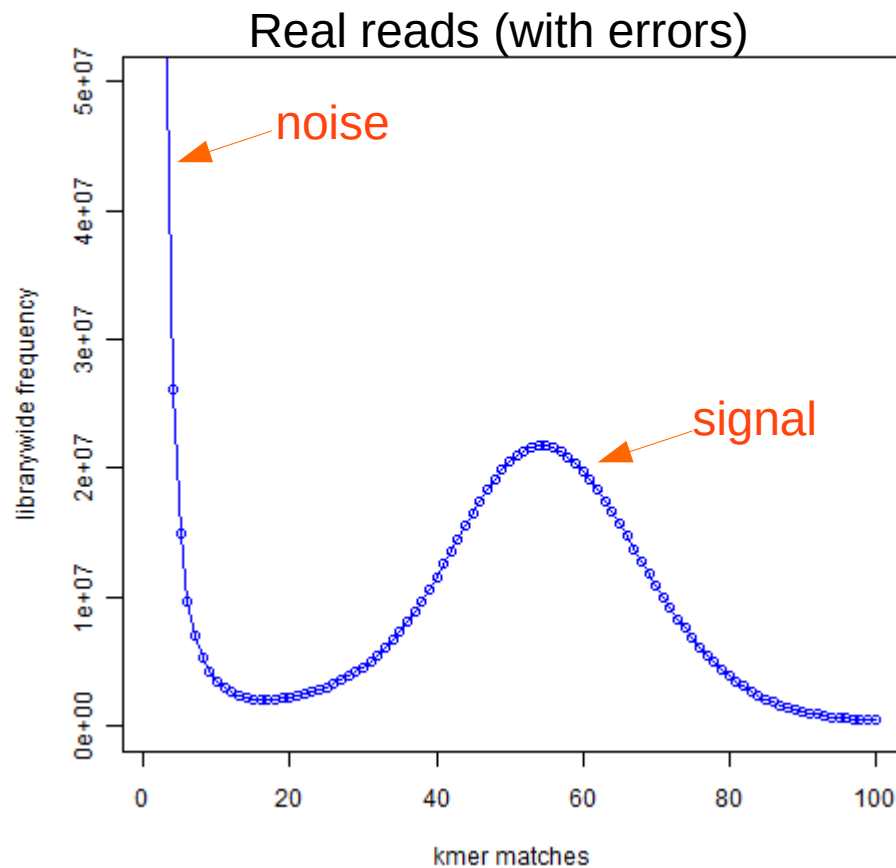
# K-mer analysis

K-mer distribution  
Simulated reads.  
40X coverage  
No errors



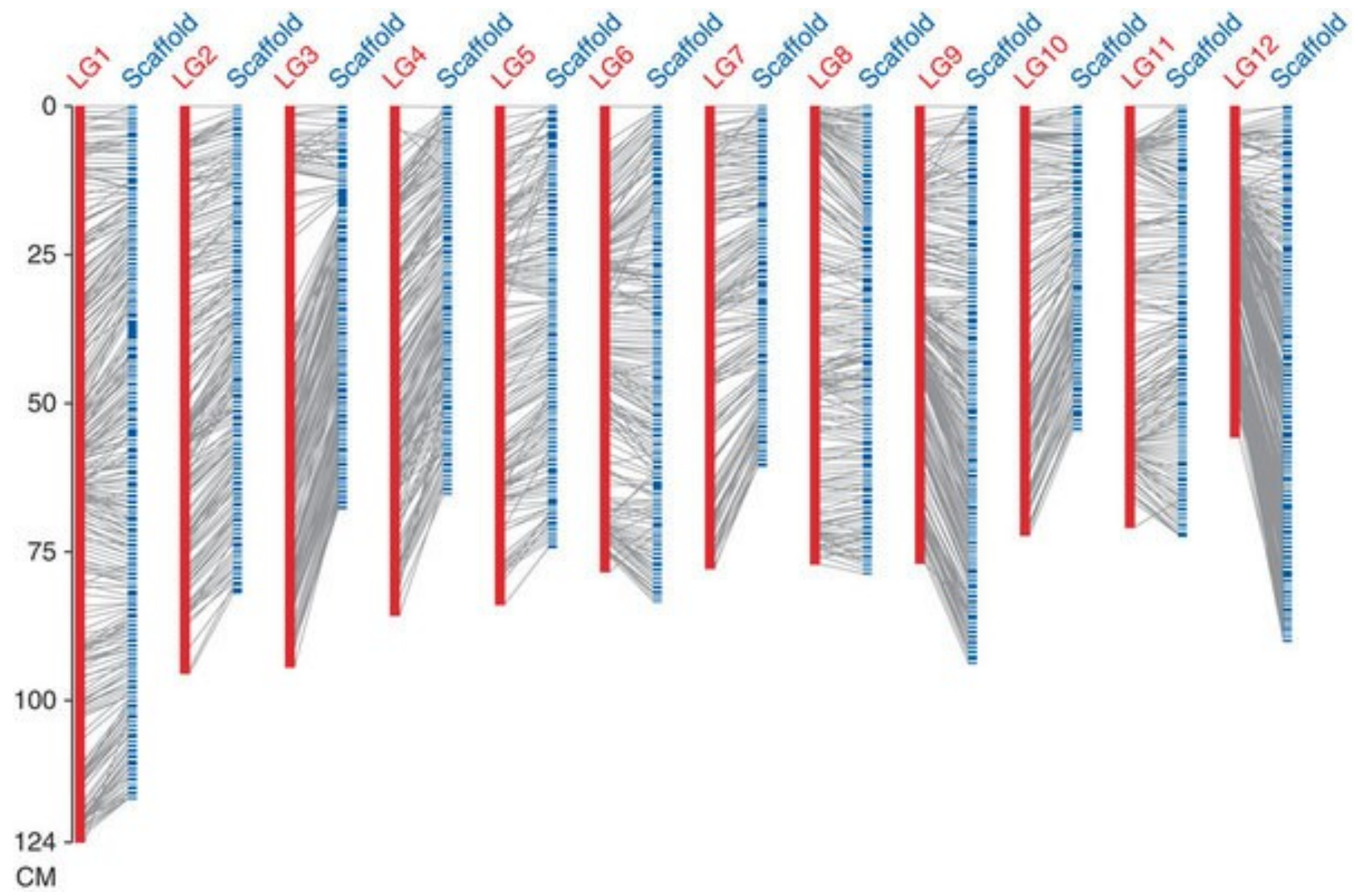
# K-mer analysis

K-mer distribution  
real reads.  
50X coverage  
With errors



# **From scaffolds to chromosomes**

# Genetic maps



# Bionano

## 1 Sequence-specific labeling

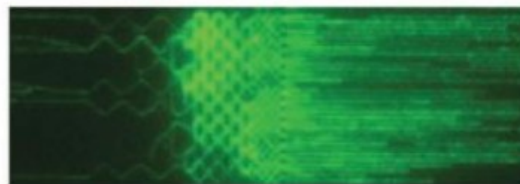
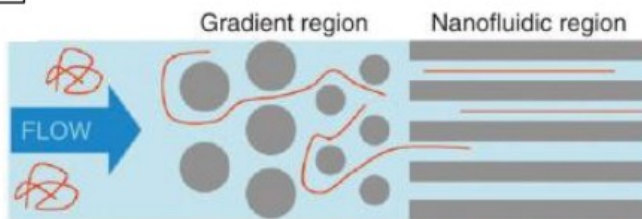
Nickase (Nt.BspQI)

5'-ATGC**GCTCTTC**CATGAATGCGAGC-3'  
3'-TACG**CGAGAAG**GTACTTACGCTCG-5'

Nick  
labeling

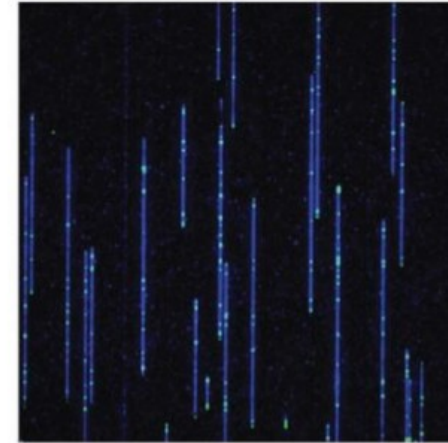
5'-ATGC**GCTCTTC**CATGAATGCGAGC-3'  
3'-TACG**CGAGAAG**GTGCTTACGCTCG-5'

## 2 DNA linearization

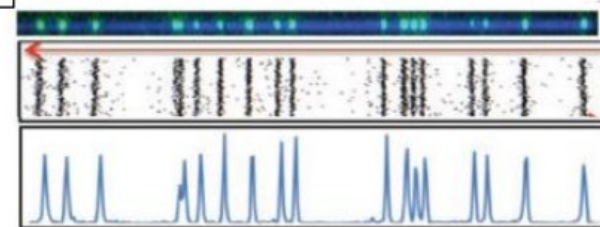


Lam et al., Nat. Biotechnol. 30(8) 2012

## 3 Fluorescence imaging



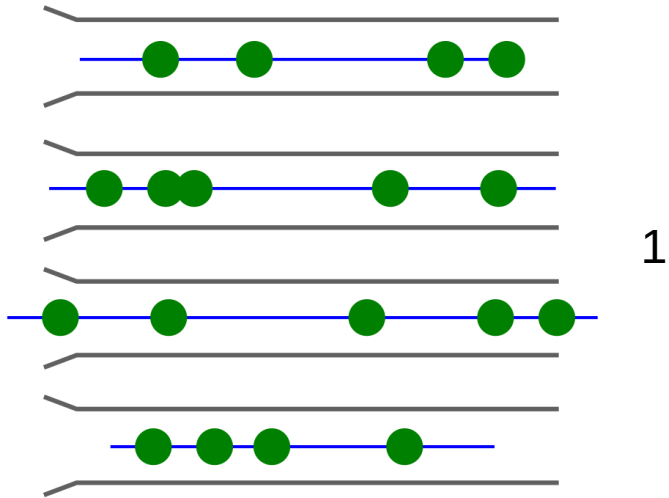
## 4 Map construction



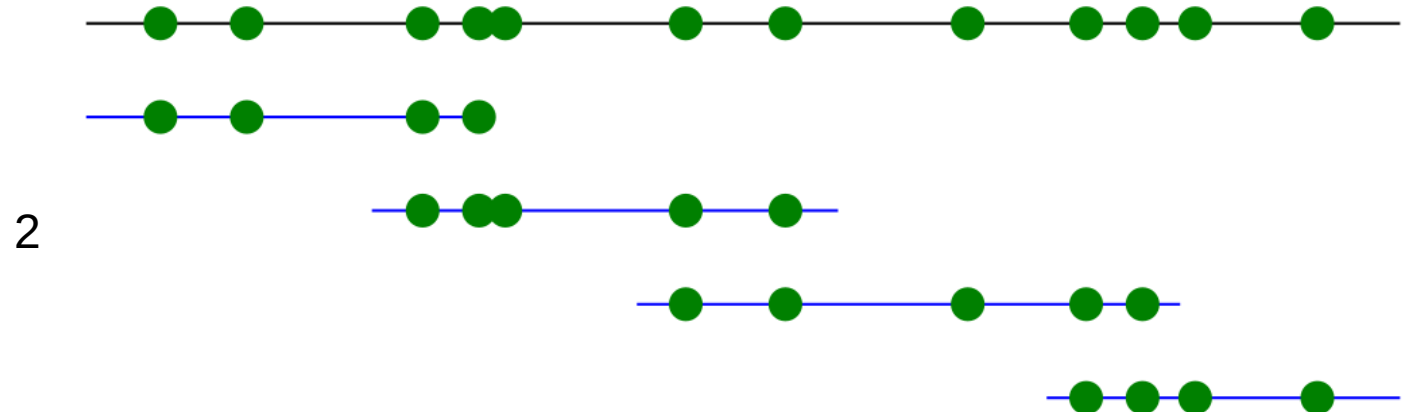
**Output:**  
sequence motif (GCTCTTC) map  
along hundreds-kb to megabase  
DNA stretches

# Bionano

---

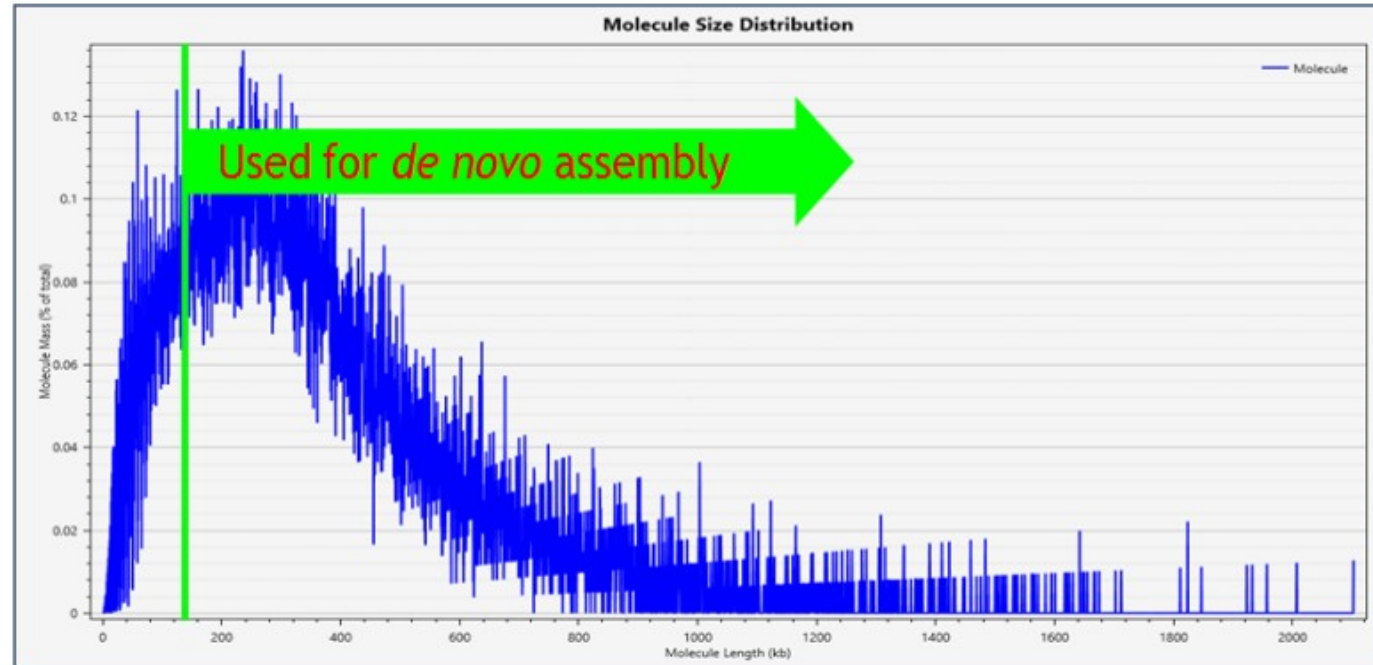
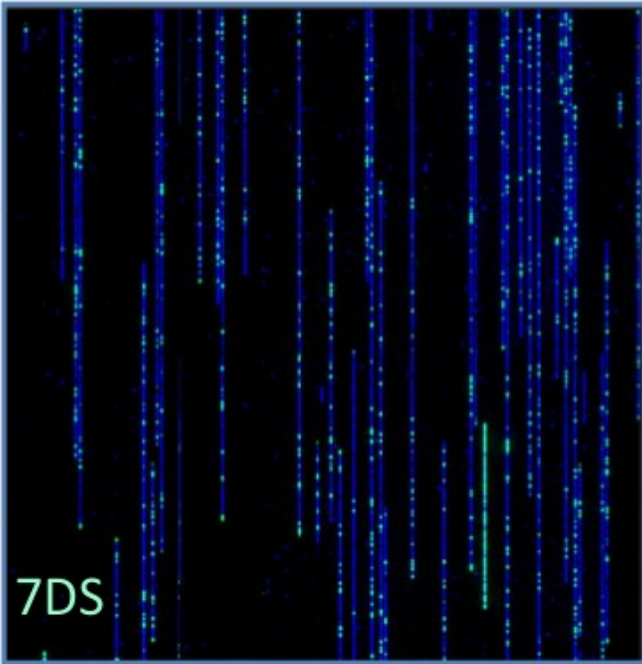


Bionano helps you solve **long range** structures  
(although not chromosome wide assemblies)





# Bionano

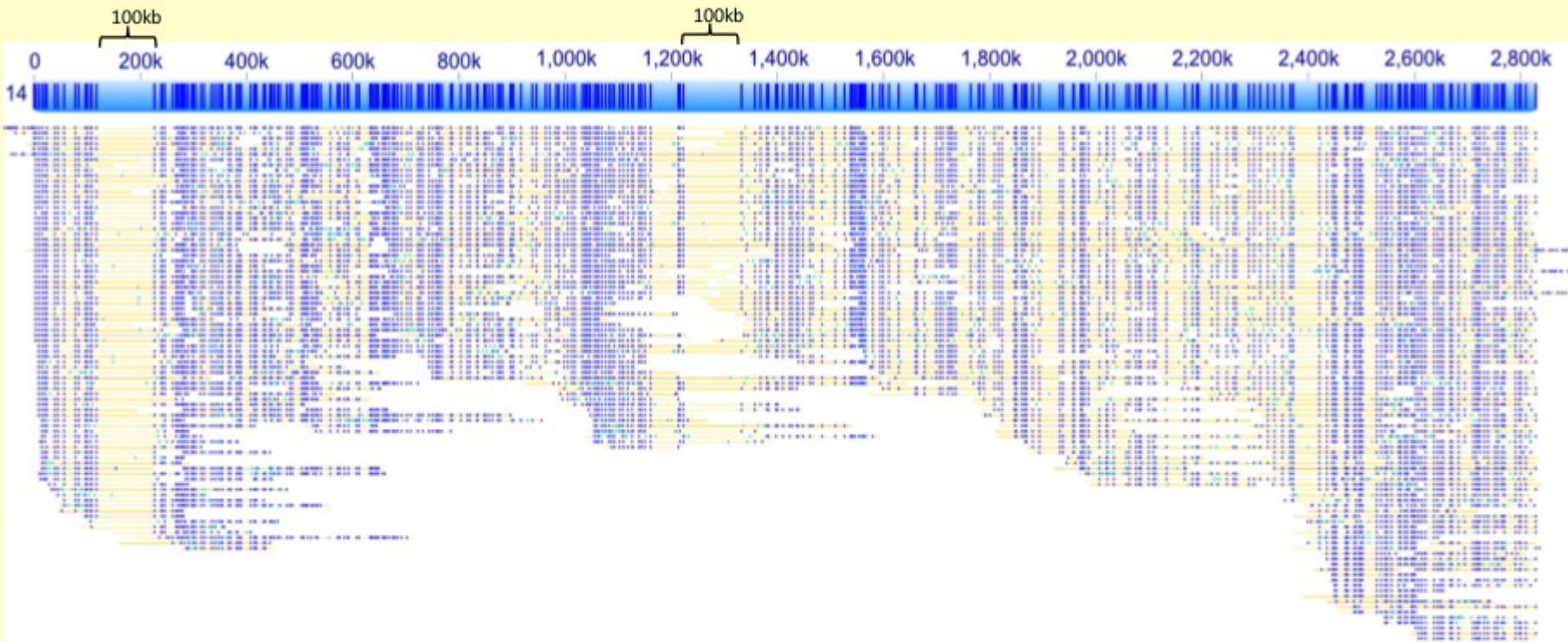


<http://olomouc.ueb.cas.cz/>

# Bionano

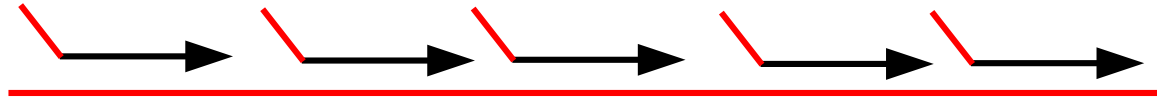
## Genome map 14

Single molecule maps



# 10X linked reads

---

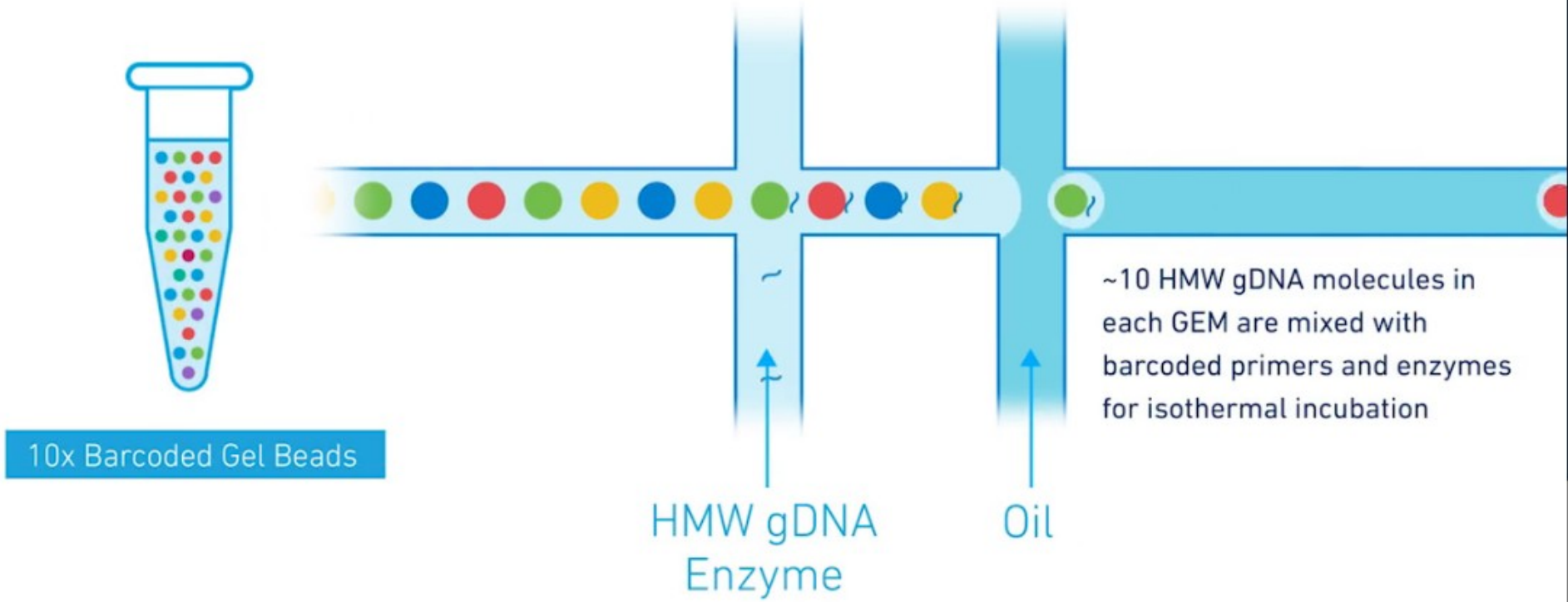


Reads generated from the same long DNA sequence are tagged with the same index and sequenced using Illumina

Reads with the same index convey **medium range** information (Up to 100Kb)

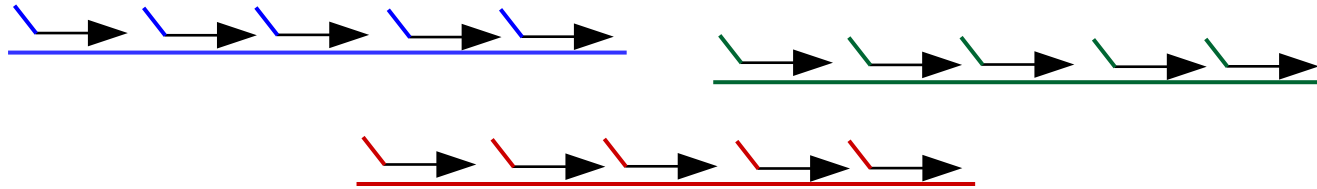
# 10X GemCode

## 10x GemCode™ Technology for Partitioning High Molecular Weight gDNA



# 10X linked reads assembly

---



# 10x possible uses

---

*De novo* genome assembly

Haplotype phasing

Structural Variants:

- Linked read sequencing resolves complex genomic rearrangements in gastric cancer metastases (<https://doi.org/10.1186/s13073-017-0447-8>)
- Integrative analysis of genomic alterations in triple-negative breast cancer in association with homologous recombination deficiency. (DOI: 10.1371/journal.pgen.1006853)

# Tomato assembly example

---

	PacBio	10X	PacBio-Bionano	10X-Bionano
N75	1.789.934	813.655	16.105.568	6.758.801
L75	141	288	18	40
N's per 100 Kb	0	6071	2325	14885
Cost	\$\$\$	\$	\$\$\$\$	\$\$

# Pacbio HiFi vs Nanopore ultralong

---

Comparison of the two up-to-date sequencing technologies for genome assembly: HiFi reads of Pacific Biosciences Sequel II system and ultralong reads of Oxford Nanopore

- <https://doi.org/10.1093/gigascience/giaa123>

Comparison with rice genome

ONT ultralong:

- 92 Gb data (230×) with an N50 of 41,473 bp
- higher contiguity
  - 18 contigs of which 10 were assembled into a single chromosome compared to 394 contigs and 3 chromosome-level contigs for the PacBio assembly
  - prevented assembly errors caused by long repetitive regions. PacBio assembly: over- or underestimation of the gene families

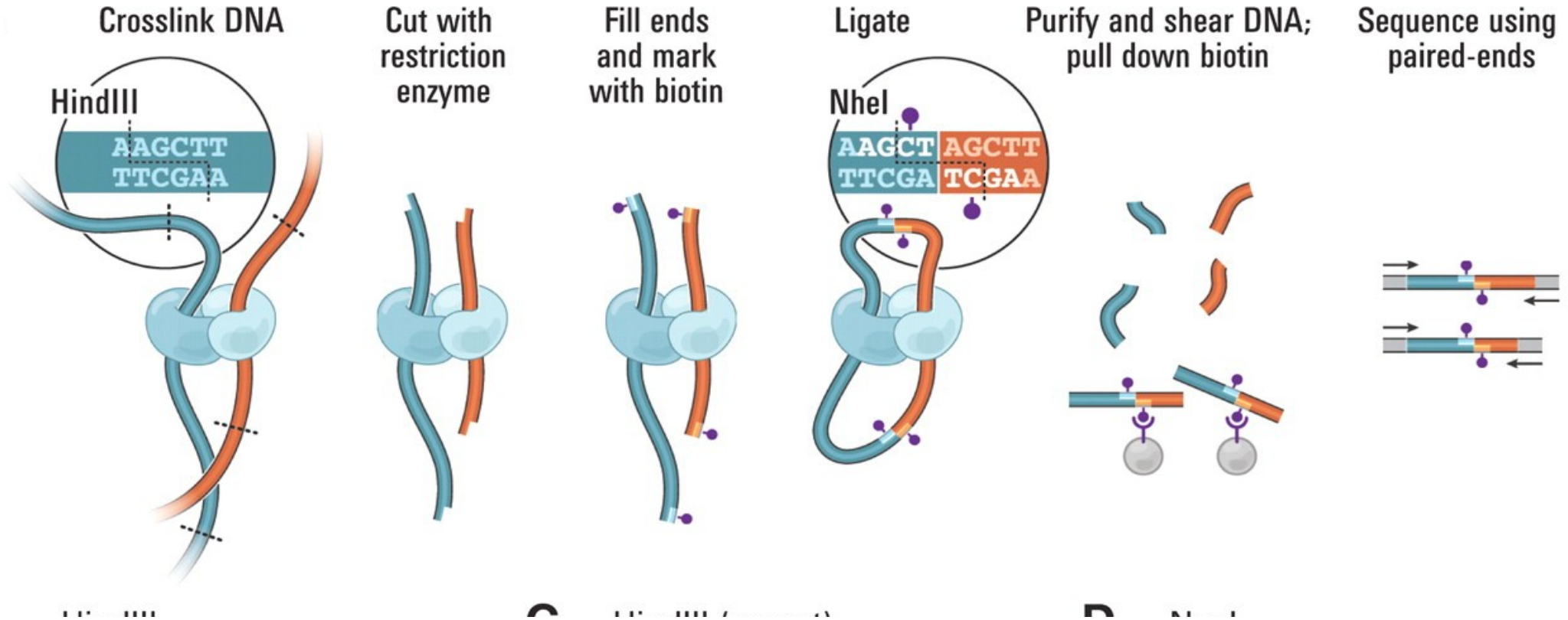
Pacbio HiFi:

- 20 Gb HiFi reads (50×) average length 13,363 bp
- fewer errors at the level of single nucleotides and small insertions and deletions

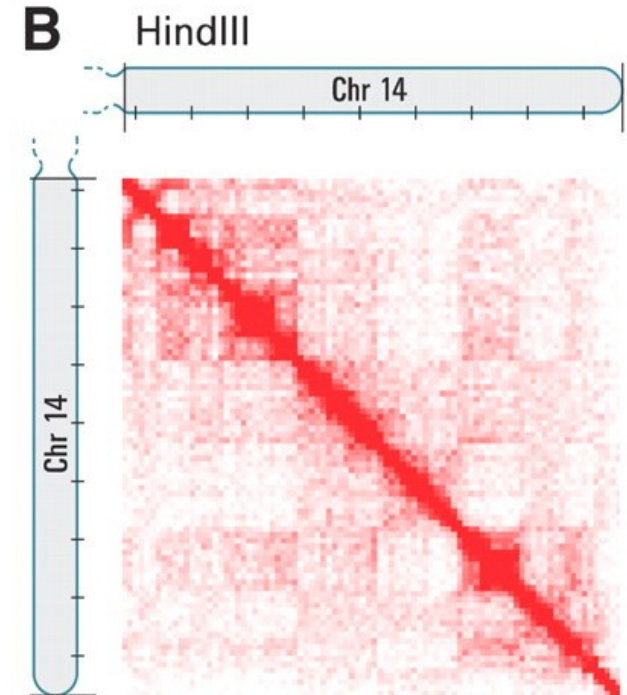
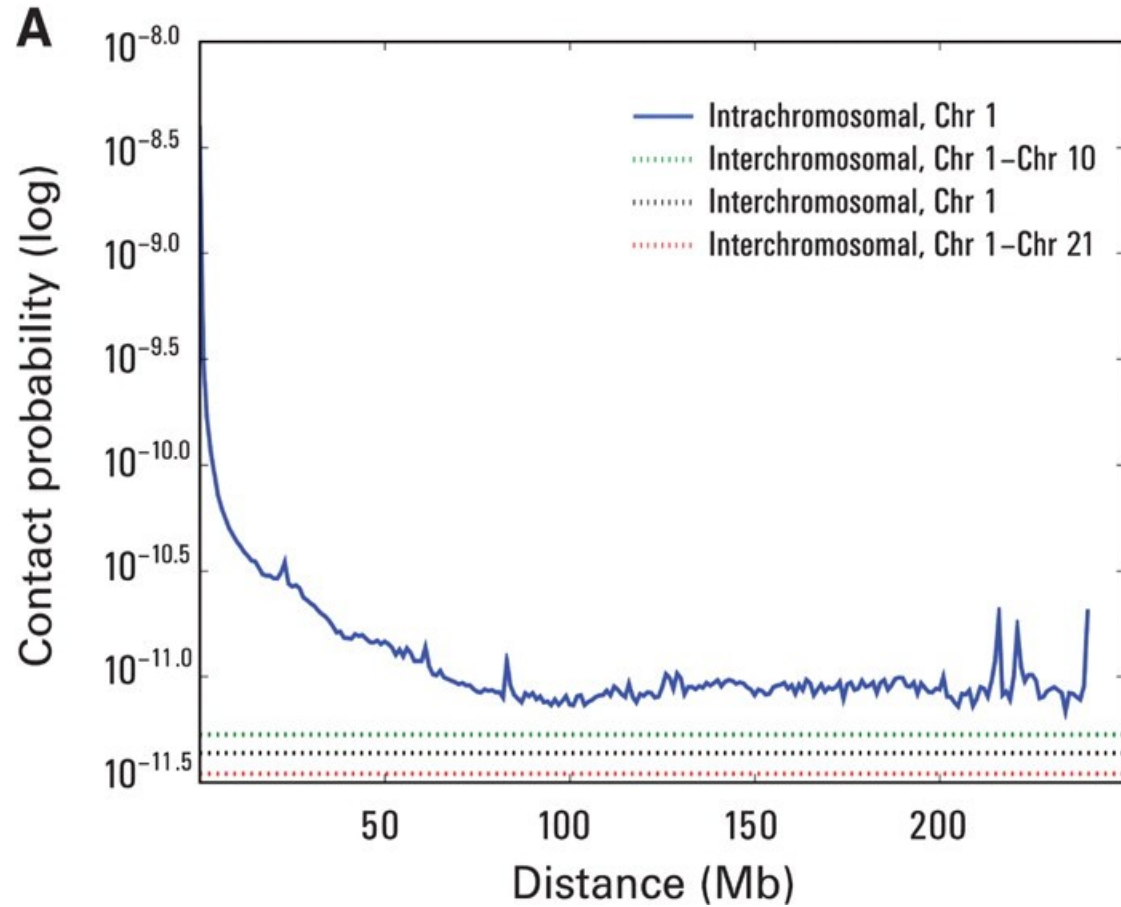


# Hi-C

Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome. (DOI: 10.1126/science.1181369)



# Hi-C



# Hi-C Assembly

---

De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds. (DOI: [10.1126/science.aal3327](https://doi.org/10.1126/science.aal3327))

Hi-C data provide links across a variety of length scales

Unlike mate-pair and pair-end reads Hi-C contact spans an unknown length and may connect loci on different chromosomes

Human genome assembly:

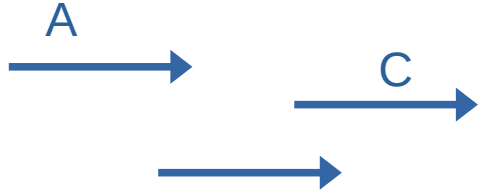
- Pair-end Illumina reads (67X coverage)
- Hi-C data (6.7 coverage)
- 23 scaffold that span the 99.5% of the 23 human chromosomes
- Errors remain in the ordering of short distances than could be fixed by mate-pairs, long reads or bionano

# Haplotypes

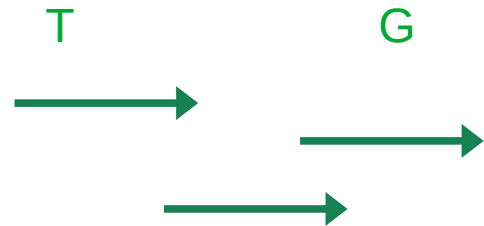
# The haplotype problem

---

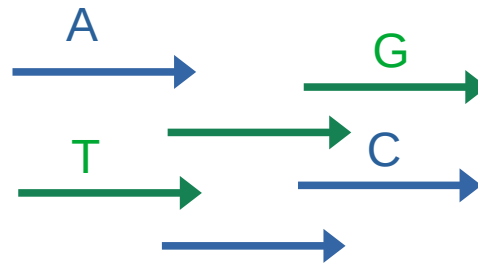
Haplotype 1



Haplotype 2



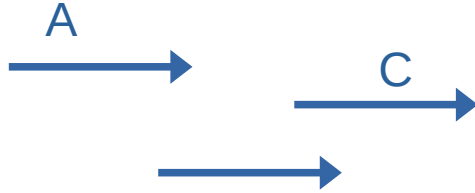
reads



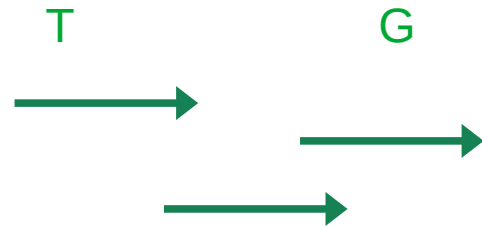
# The haplotype problem

---

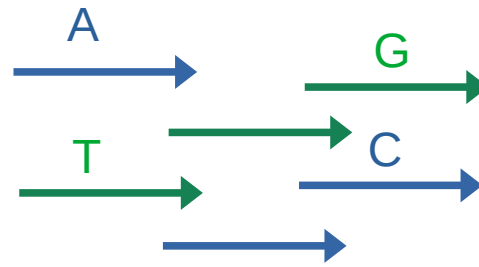
Haplotype 1



Haplotype 2



reads



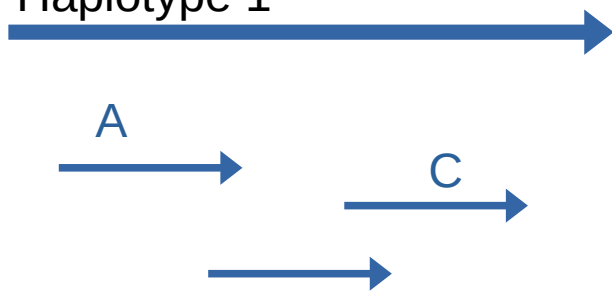
Standard assembly



# The haplotype problem

---

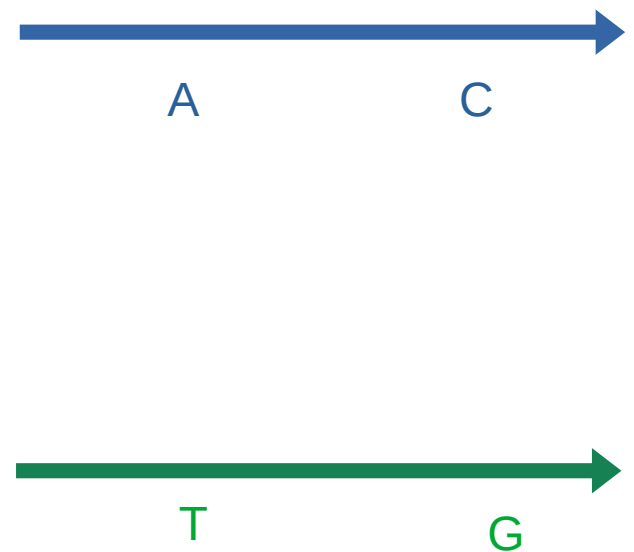
Haplotype 1



Haplotype 2



Ideal assembly



# Source of info 1: Good quality long reads

---

## Haplotype-resolved de novo assembly using phased assembly graphs with **hifiasm**

Haoyu Cheng, Gregory T. Concepcion, Xiaowen Feng, Haowen Zhang & Heng Li [✉](#)

*Nature Methods* **18**, 170–175 (2021) | [Cite this article](#)

**8349** Accesses | **29** Citations | **117** Altmetric | [Metrics](#)

### Abstract

Haplotype-resolved de novo assembly is the ultimate solution to the study of sequence variations in a genome. However, existing algorithms either collapse heterozygous alleles into one consensus copy or fail to cleanly separate the haplotypes to produce high-quality phased assemblies. Here we describe hifiasm, a de novo assembler that takes advantage of **long high-fidelity** sequence reads to faithfully represent the haplotype information in a phased assembly graph. Unlike other graph-based assemblers that only aim to maintain the contiguity of one haplotype, hifiasm strives to preserve the contiguity of all haplotypes. This feature enables the development of a graph trio binning algorithm that greatly advances over standard trio binning. On three human and five nonhuman datasets, including California redwood with a ~30-Gb hexaploid genome, we show that hifiasm frequently delivers better assemblies than existing tools and consistently outperforms others on haplotype-resolved assembly.





# Long distance haplotype info

---

10X

Hi-C

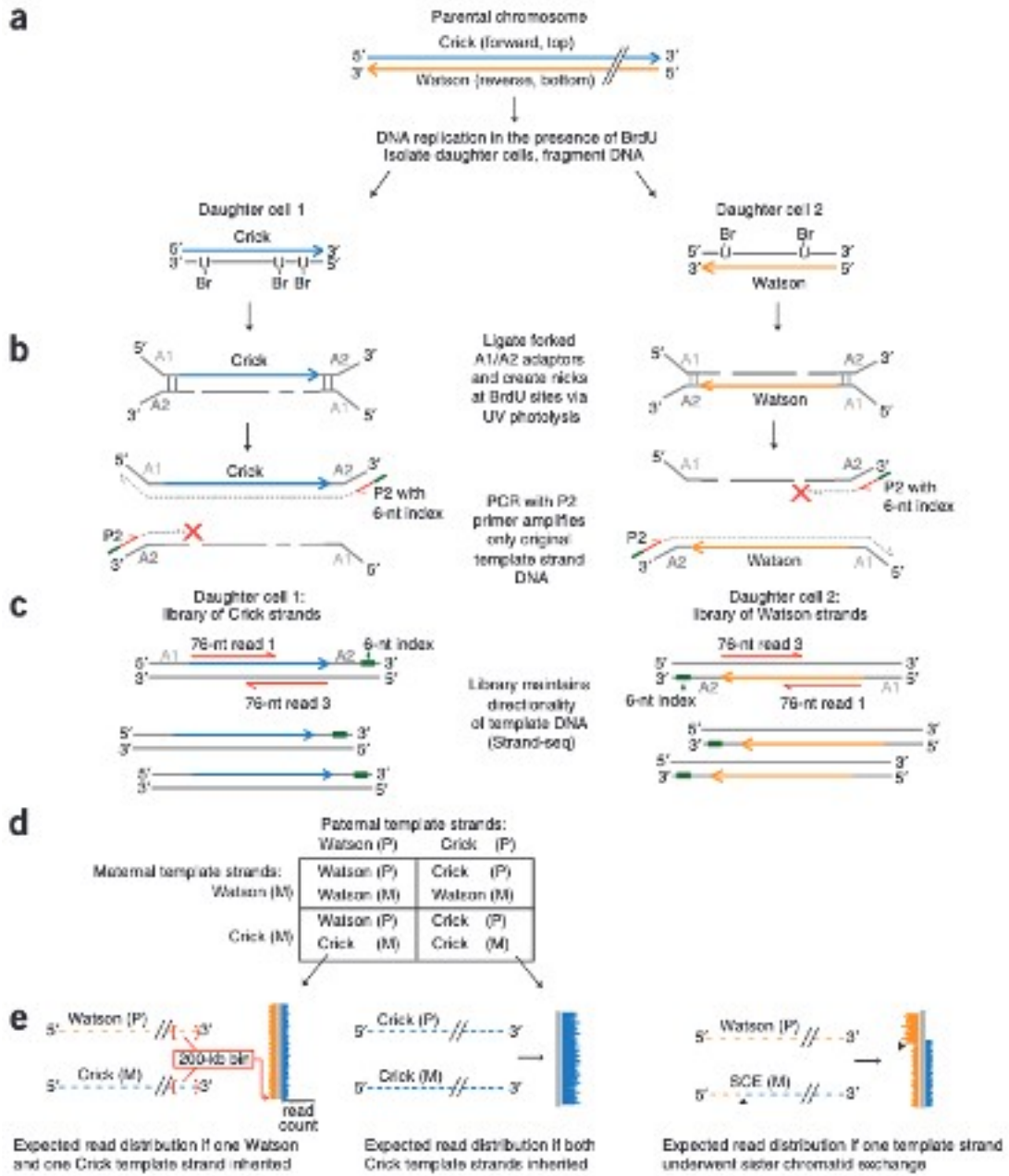
Single strand seq

# Single strand seq

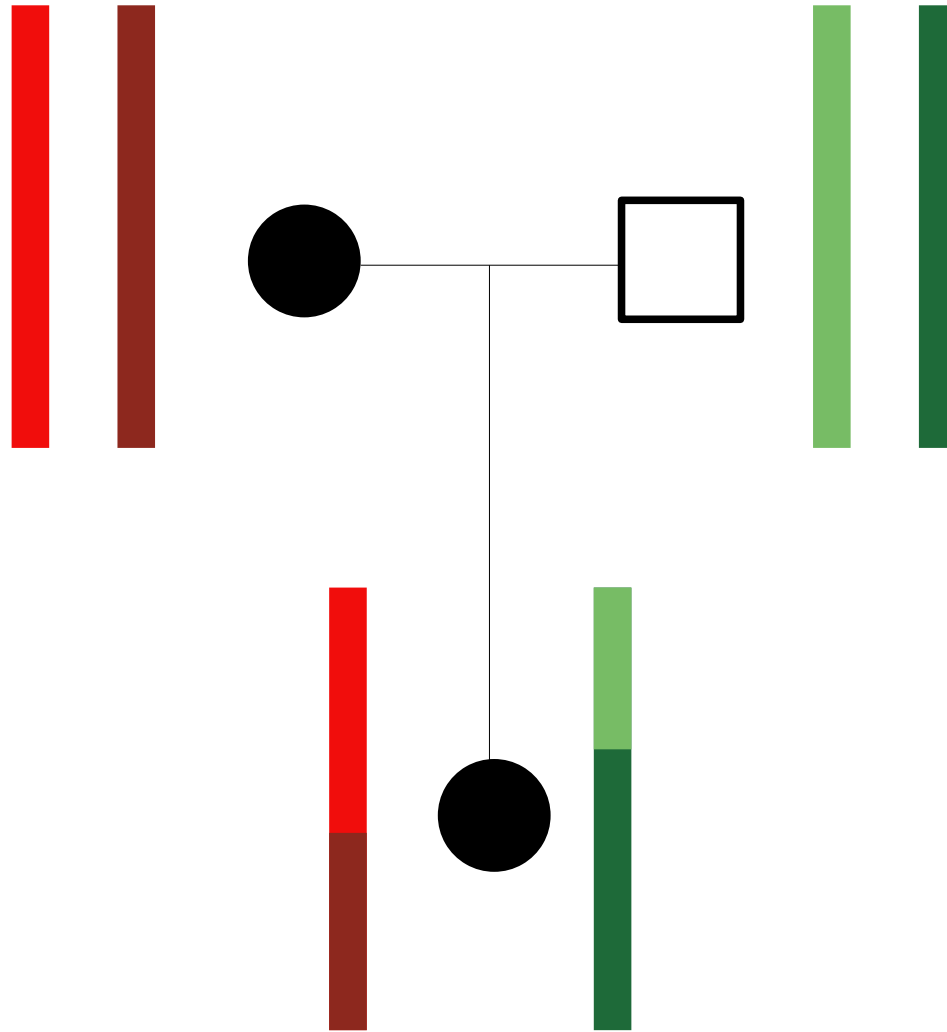
short-read, single-cell sequencing method  
 thymidine analog (BrdU) to selectively label  
 and remove one of the DNA strands

References:

- Fully phased human genome assembly without parental data using single-cell strand sequencing and long reads
- Single-cell template strand sequencing by Strand-seq enables the characterization of individual homologs



# Long distance haplotype info: parental trio



# Haplotype assembly with hifiasm

---

Additional short read data:

- Hi-C pair ends or Strand-seq
  - improves contiguity and phasing accuracy
- Parental trios
  - Improves phasing accuracy
  - Advisable specially for high heterozygosity

Hifiasm does not perform scaffolding for now

Checked with human and California redwood with a ~30-Gb hexaploid genome

Similar to HiCanu

Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. Nat Methods, 18:170-175. <https://doi.org/10.1038/s41592-020-01056-5>

---

This work is licensed under the Creative Commons Attribution 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by/3.0/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

Jose Blanca  
COMAV institute  
[bioinf.comav.upv.es](http://bioinf.comav.upv.es)