

A little tour of assembly methods

Camille Marchet & Antoine Limasset

Université de Lille, CNRS, France

EVOMICS Workshop on Genomics 2020 - Český Krumlov



Content of this course

- How to reconstruct a genome with sequencing data?
- What are the main challenges?
- Which solutions have been proposed?



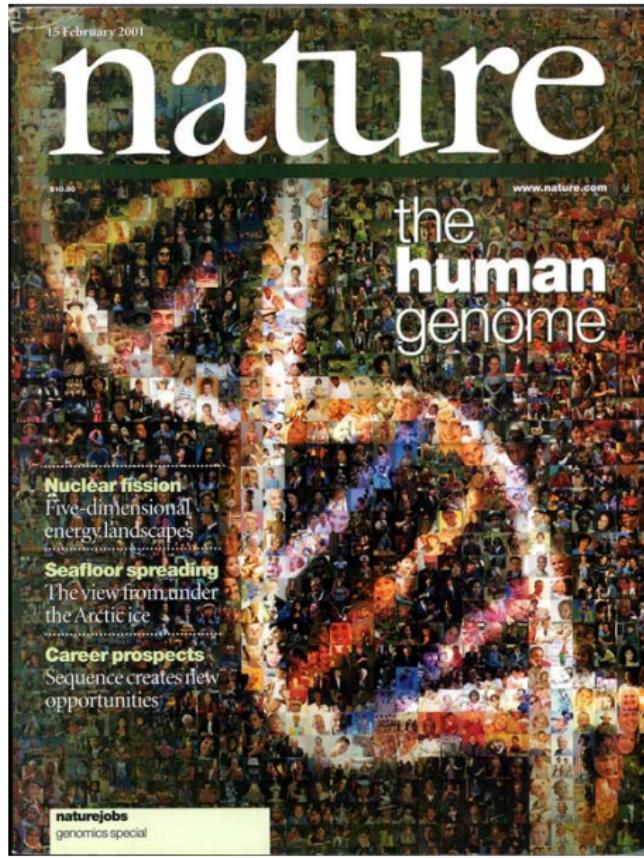
genome size: ~ 32 gigabases

Content of this course

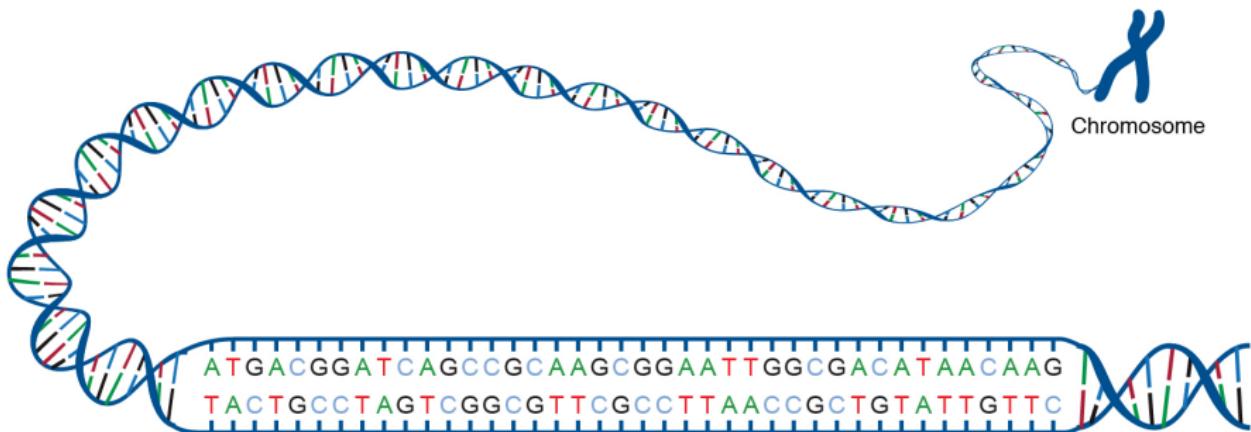
- What is assembly?
- What are the main problematics to bear in mind when conducting an assembly?
- What are the current limits of assembly?



Assembly

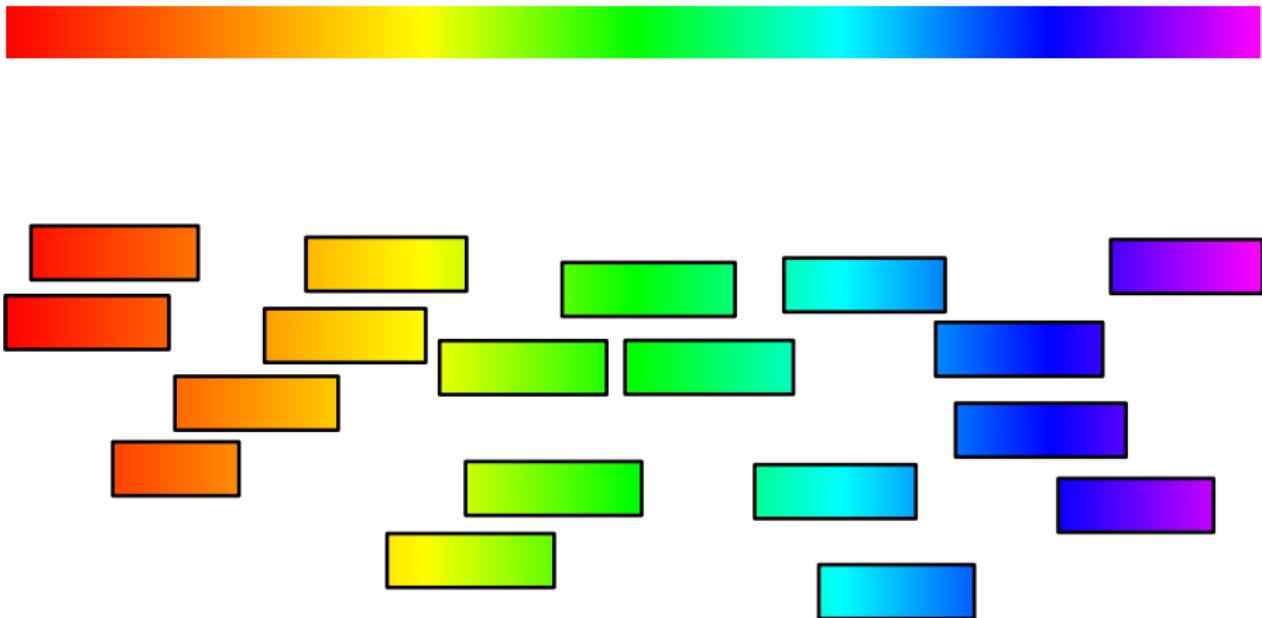


Accessing a genome

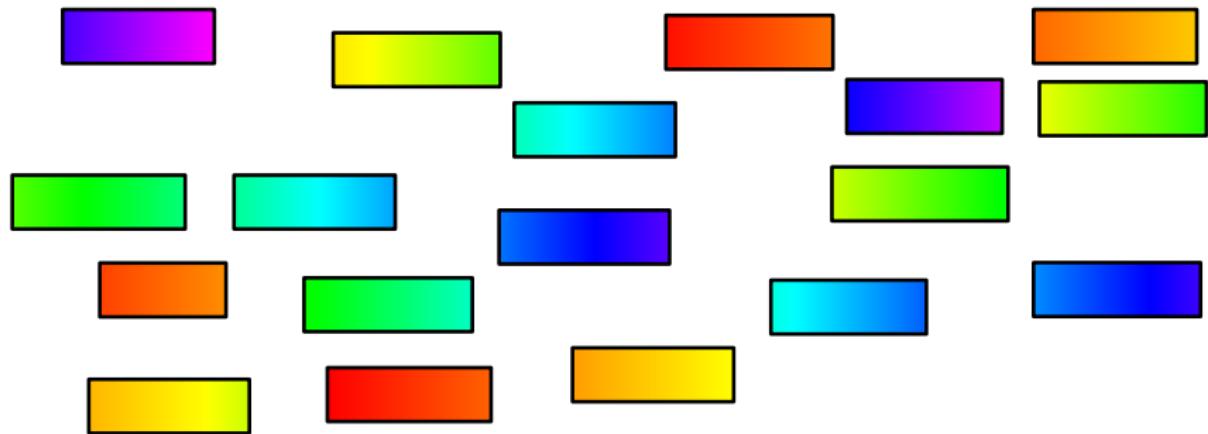


From www.genome.gov/genetics-glossary/acgt

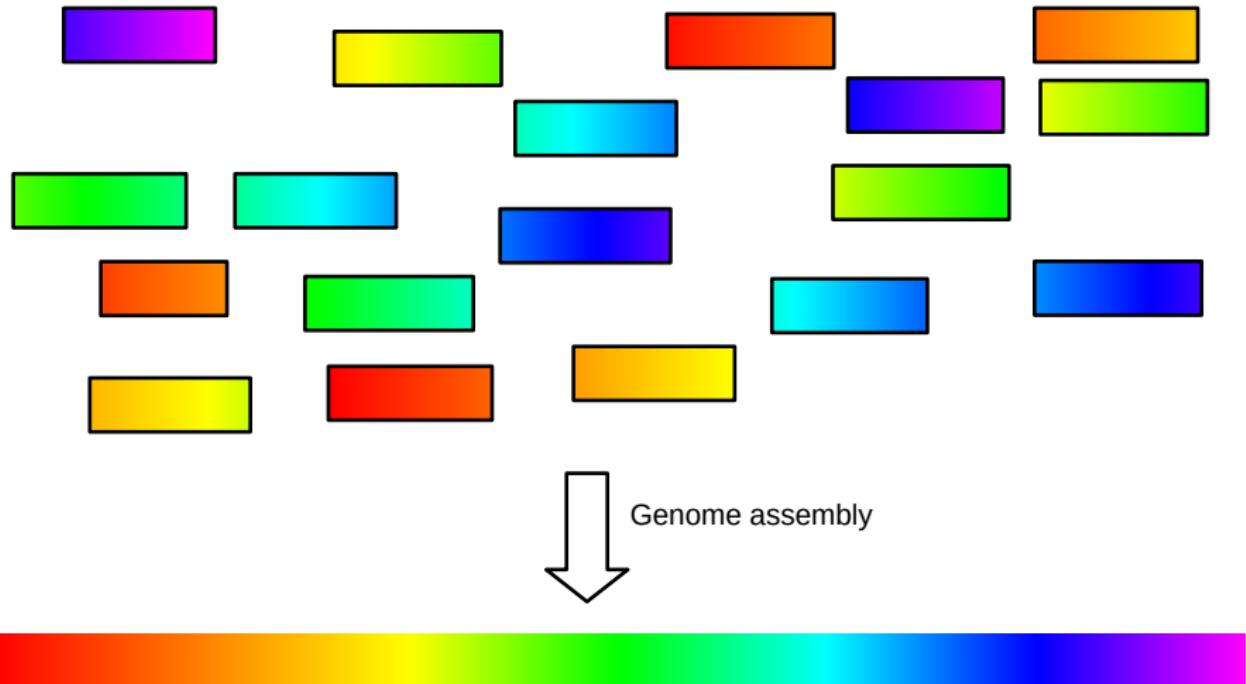
Reads are words from the genome



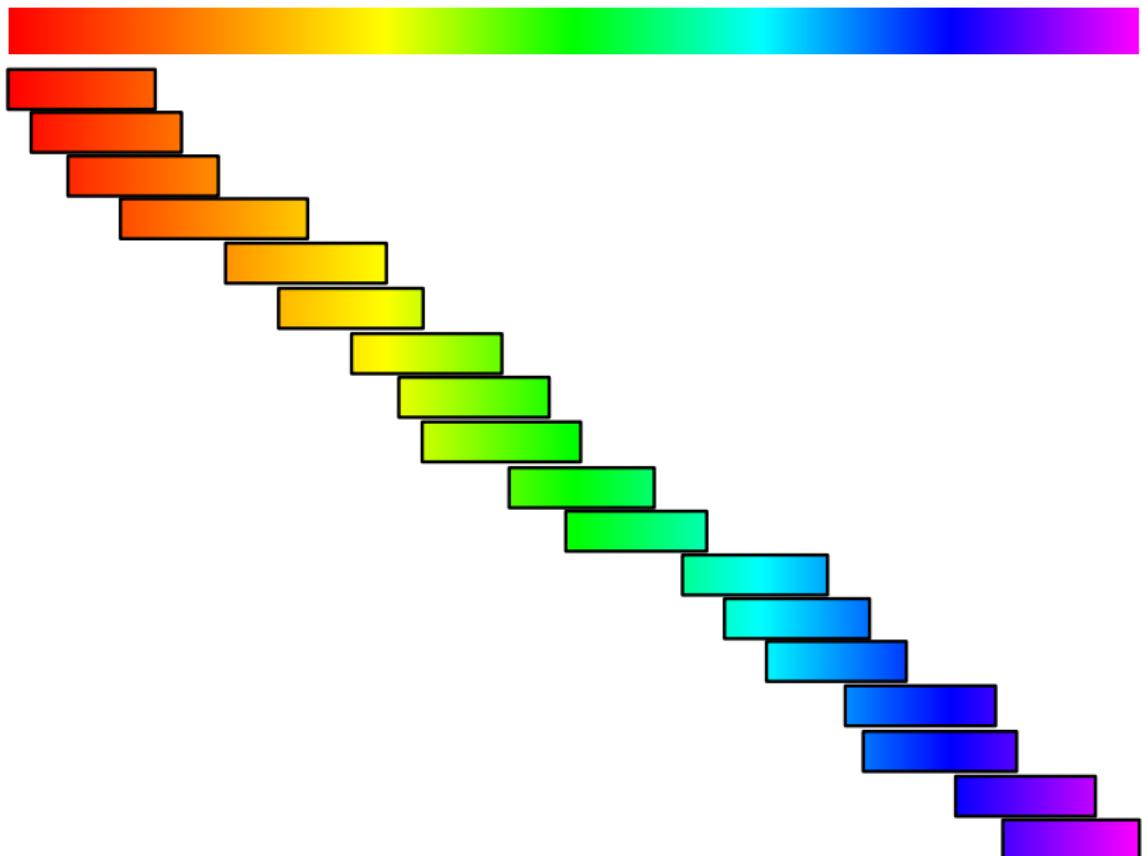
Reads are **shuffled** words from the genome



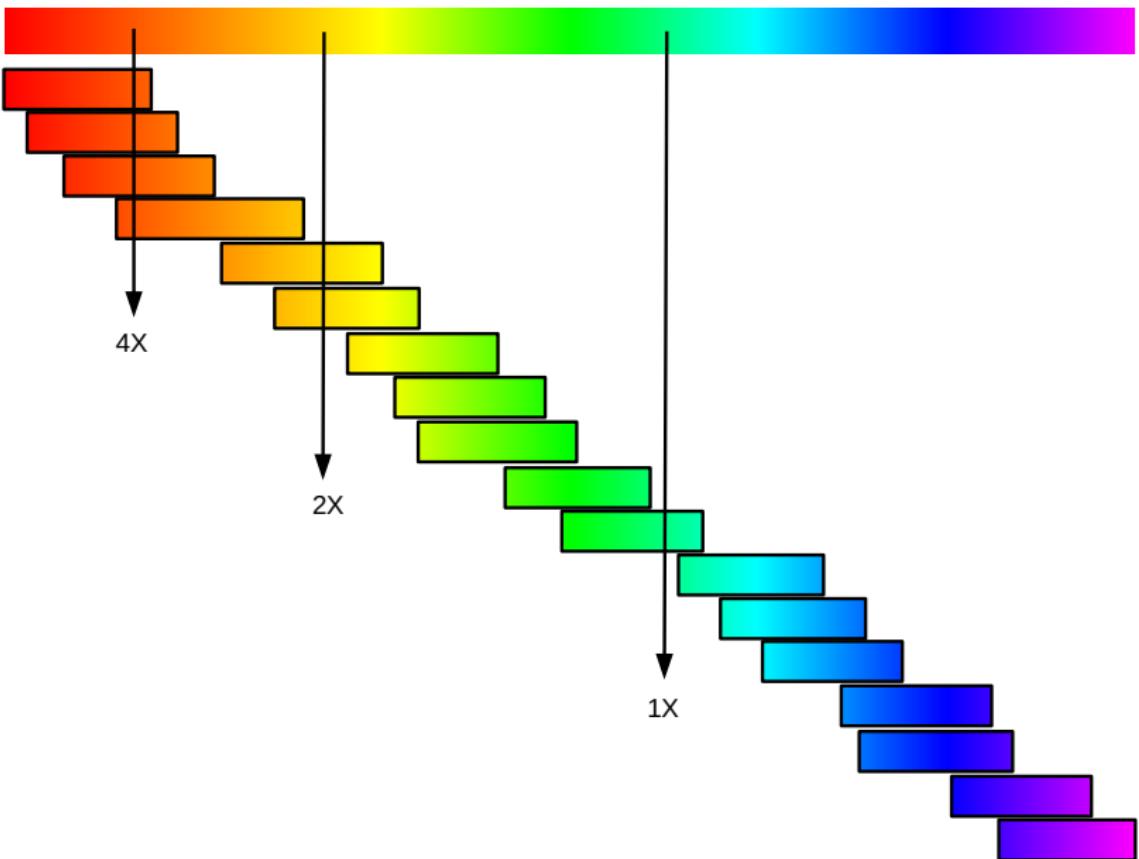
Genome assembly task



Using read overlaps



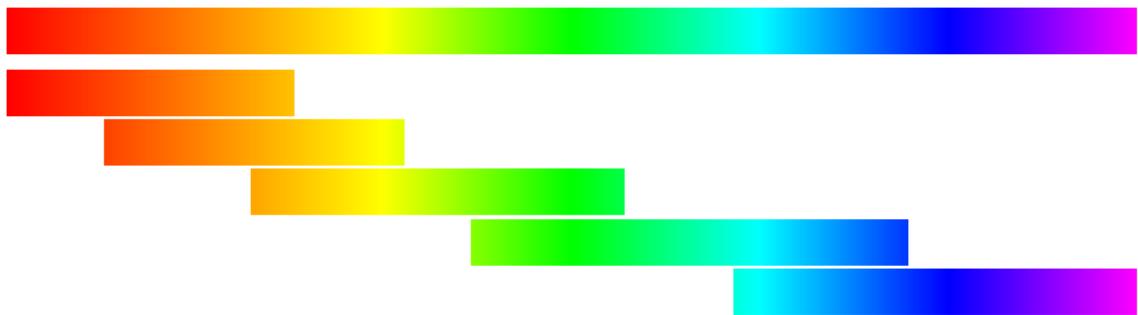
Genome sequencing: coverage



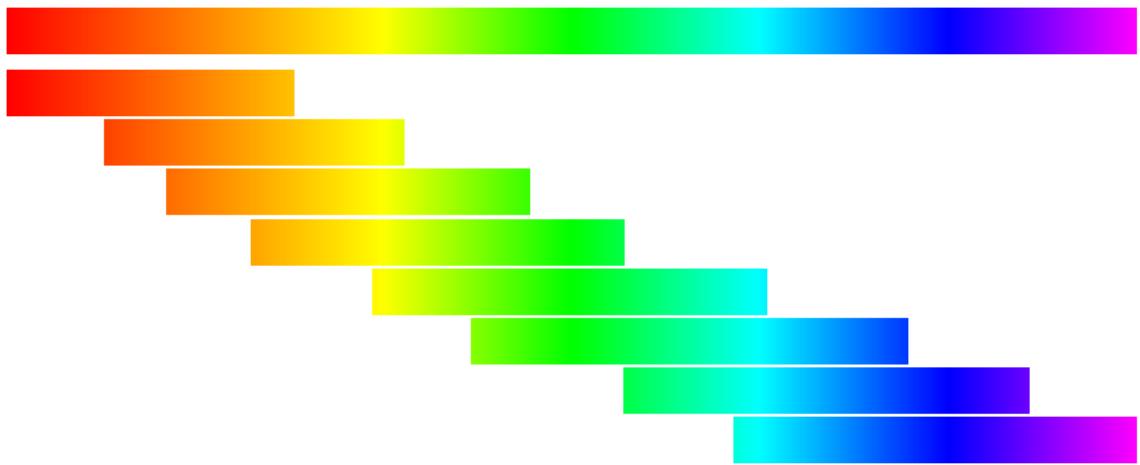
Genome sequencing: coverage



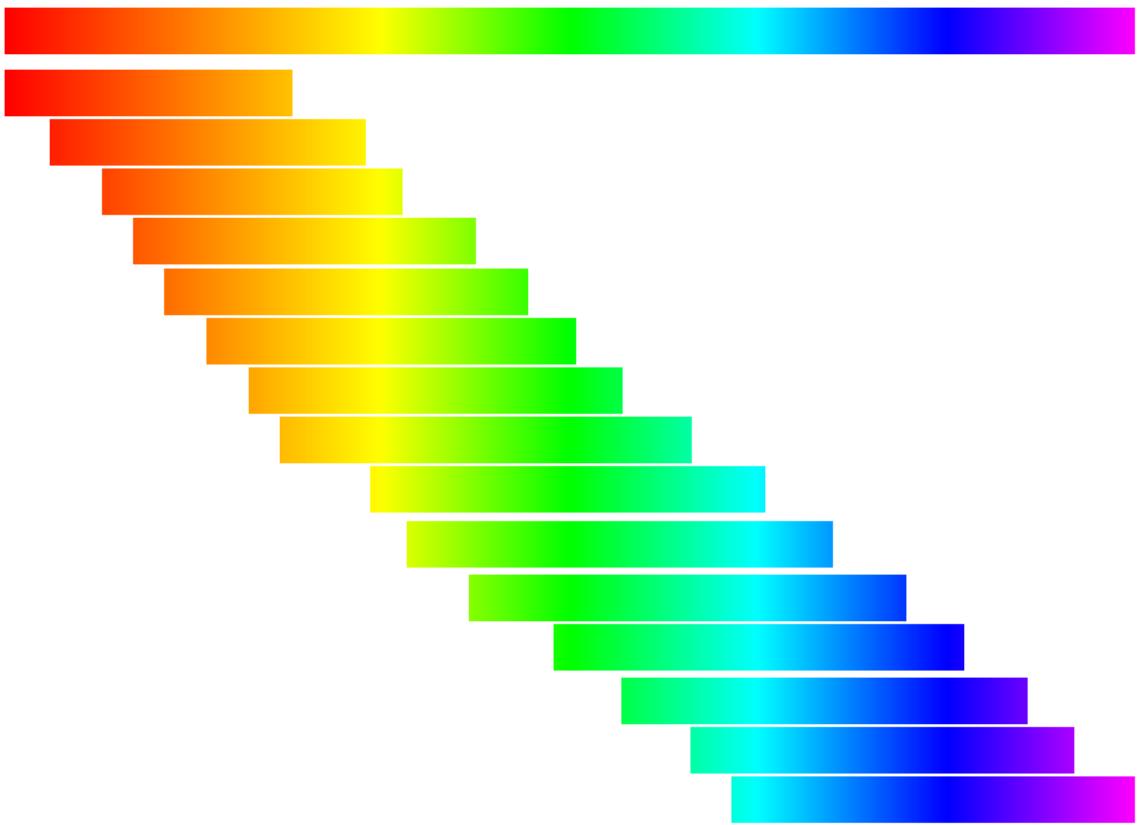
Genome sequencing: coverage



Genome sequencing: coverage



Genome sequencing: coverage



Assembly idea number 1: Select the longest overlaps

Reads:

1: ATCGGTATCG
2: GGTATCGTTA
3: ATCGTTACGG
4: GTTACGGTAT
5: ACGGTATAACC

1: ATCGGTATCG
2: GGTATCGTTA

Overlap length: 7

1: ATCGGTATCG
3: ATCGTTACGG

Overlap length: 4

1: ATCGGTATCG
4: GTTACGGTAT

Overlap length: 1

1: ATCGGTATCG
5: ACGGTATAACC

No Overlap

Assembly idea number 1: assemble the longest overlaps

Reads:

1: ATCGGTATCG
2: GGTATCGTTA
3: ATCGTTACGG
4: GTTACGGTAT
5: ACGGTATACC

Best overlaps:

1: ATCGGTATCG
2: GGTATCGTTA
3: ATCGTTACGG
4: GTTACGGTAT
5: ACGGTATACC

Output “genome”:

ATCGGTATCG + TTA + CGG + TAT + ACC
ATCGGTATCGTTACGGTATACC

Your time to shine!

Let assemble this genome!

Your read set:

1: A**T**TTA**C**G**GG**T
2: T**T**A**C**GG**GG**T**GG**
3: **A**CG**GG**T**C**CTT
4: G**T**C**C**TTT**C**TT
5: TTT**C**TT**A**CG**G**

For each read:

Find the best overlap (length>5)
Merge the two reads

The Greedy solution

The best overlaps:

ATTTACGGGT
TTACGGGTGG

ACGGGTCCTT
GTCCTTTCTT
TTTCTTACGG

Output “genome”

ATTTACGGGTGG
ACGGGTCCTTTCTTACGG

The actual solution

The actual genome:

ATTTC_ACGGGT_CC_TTTT_CTT_ACGGGT_GG

How the reads should be ordered:

ATTTC_ACGGGT_CC_TTT_ACGGGT_GG
ACGGGT_CC_TTT_ACGGGT_GG
GTC_CCTT_TT_CTT_ACGGGT_GG
TTT_CTT_ACGGGT_GG
TT_ACGGGT_GG

What happened?

The actual genome:

ATTTACGGGTCCCTTCTTACGGGTGG

How the reads should be ordered:

ATTTACGGGT
ACGGGTCCCTT
GTCCCTTCTT
TTTCTTACGG
8 → TTACGGGTGG

ATTTACGGGT
TTACGGGTGG

ATTTACGGGTGG
Not in the genome

ACGGGTCCCTT
GTCCCTTCTT
TTTCTTACGG

ACGGGTCCCTTCTTACGG
Not in the genome

Do we expect many repeats?

Probability to have NO repeated word of size 31 in a 5 megabases genome

Input interpretation:

$$\left(\frac{4^{31} - 1}{4^{31}}\right)^{1/2 \times (5 \times 10^6 \times 5 \times 10^6 - 1)}$$

Decimal approximation:

0.999997289498784302383172055421363836712023171938932024106...

From en.wikipedia.org/wiki/Birthday_problem

The burden of assembly: genomic repeats

Amount of repeats larger than a given size in *E. coli* genome

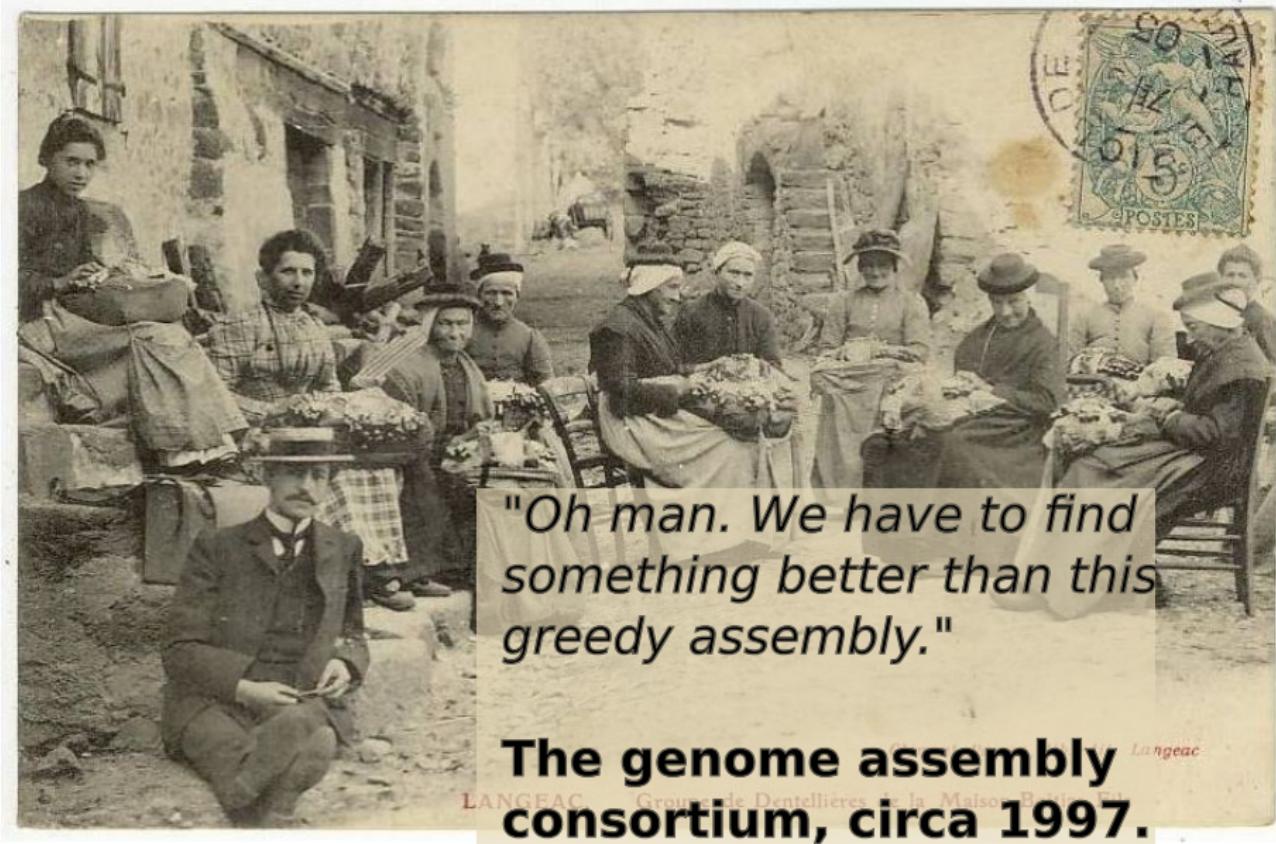
- 15: 44,994
- 21: 1,169
- 31: 559
- 41: 323
- 51: 225
- 61: 192

Genomic repeats are NOT random events

Greedy assemblers

- Simple and efficient scheme
- Rely on **local** best choice (greedy)
- May create errors because of local choices when there are repeats

History from the last century

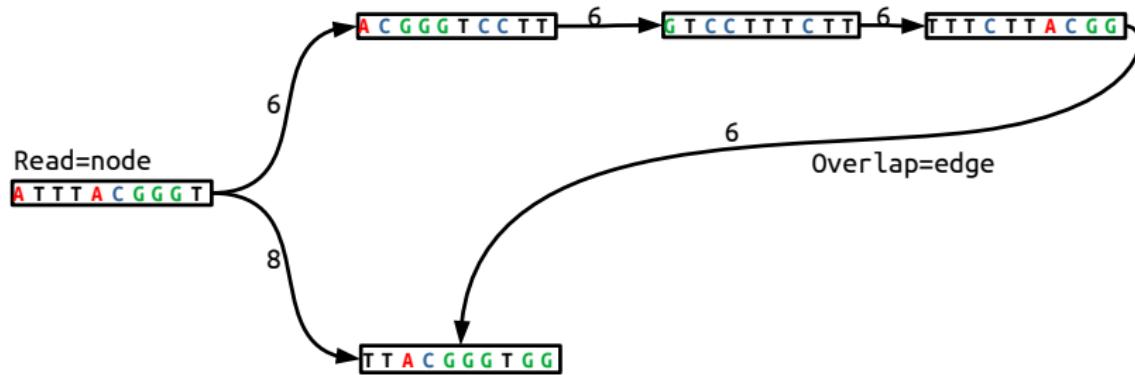


Assembly idea number 2: consider all overlaps

Genome:

ATTTACGGGT|CCTTTCTTACGGGTGG

Overlap graph:



Greedy solution

Genome:

ATTTACGGGTCCCTTCTTACGGGTGG

Overlap graph:



Read=node

ATTTACGGGT

Overlap=edge

8

TTACGGGTGG

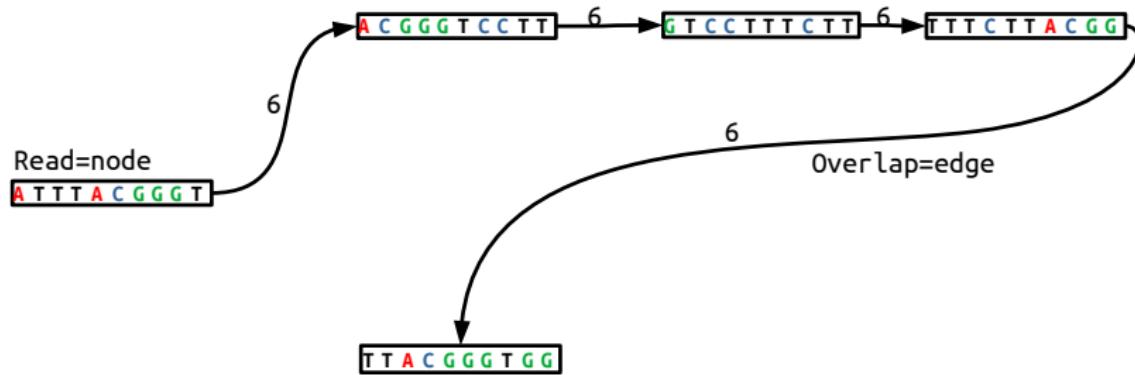
Greedy assembly output:
ATTTACGGGTGG
ACGGGTCCCTTCTTACGG

One piece solution

Genome:

ATTTACGGGT|CCTTTCTTACGGGTGG

Overlap graph:



Overlap graph output:

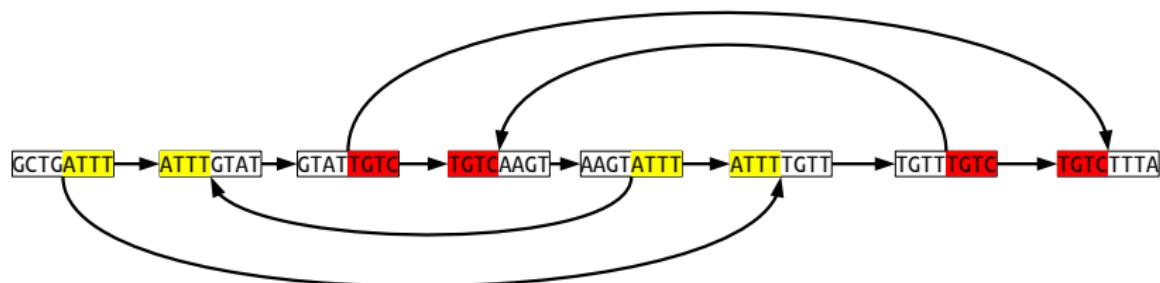
ATTTACGGGT|CCTTTCTTACGGGTGG

Multiple repeats

Reads:

GCTGATT
ATTTGTAT
GTATTGTC
TGTCAAGT
AAGTATT
ATTTTGT
TGTTTGTC
TGTCTTTA

Overlap graph:



First solution

Reads:

GCTGATT
ATTTGTAT
GTATTGTC
TGTCAAGT
AAGTATT
ATTTTGTT
TGTGGTC
TGTCTTTA

Overlap graph:



Possible assemblies:

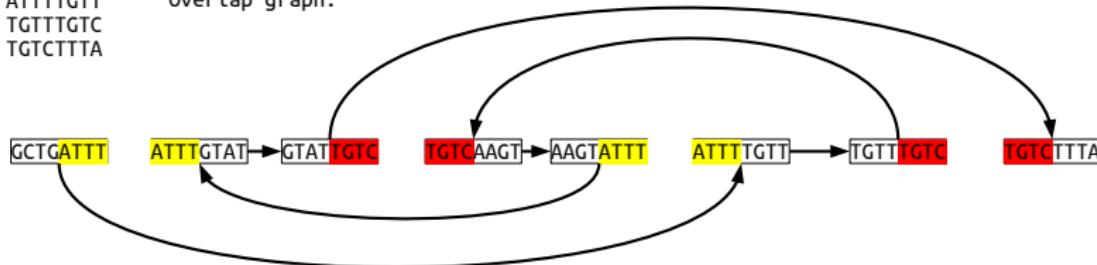
GCTGATT GTAT TGTCAAGT ATTTTGTT TGTCTTTA

Second solution

Reads:

GCTGATT
ATTTGTAT
GTATTGTC
TGTCAAGT
AAGTATT
ATTTTGTT
TGTTTGTC
TGTCTTTA

Overlap graph:



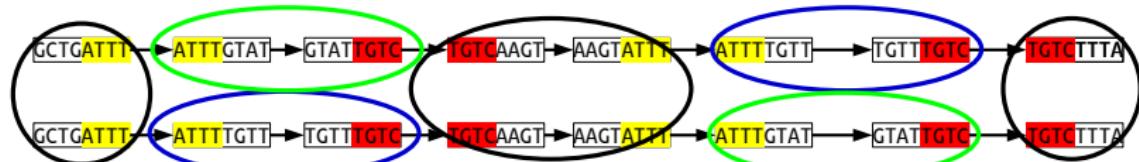
Possible assemblies:

GCTGATTGTATTGTCAAGTATTTTGTTTGTCCTTA
GCTGATTGTGTGTCAAGTATTTGTATTGTCTTTA

Those two solutions are indistinguishable

Parsimonious solution: do not assemble

Possible assemblies:

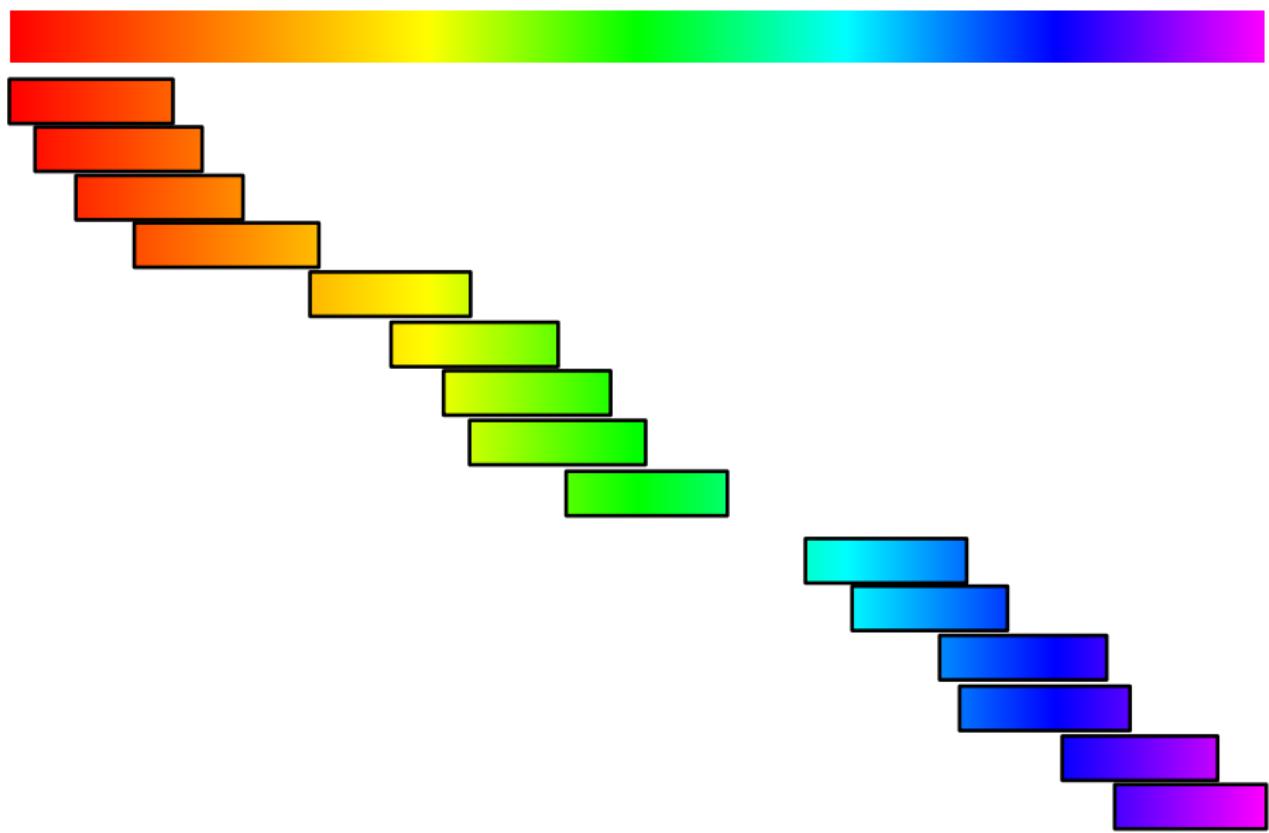


Genome pieces:

GCTGATTI ATTTGTATTGTC TGTC AAGTATTI ATTTGTATTGTC

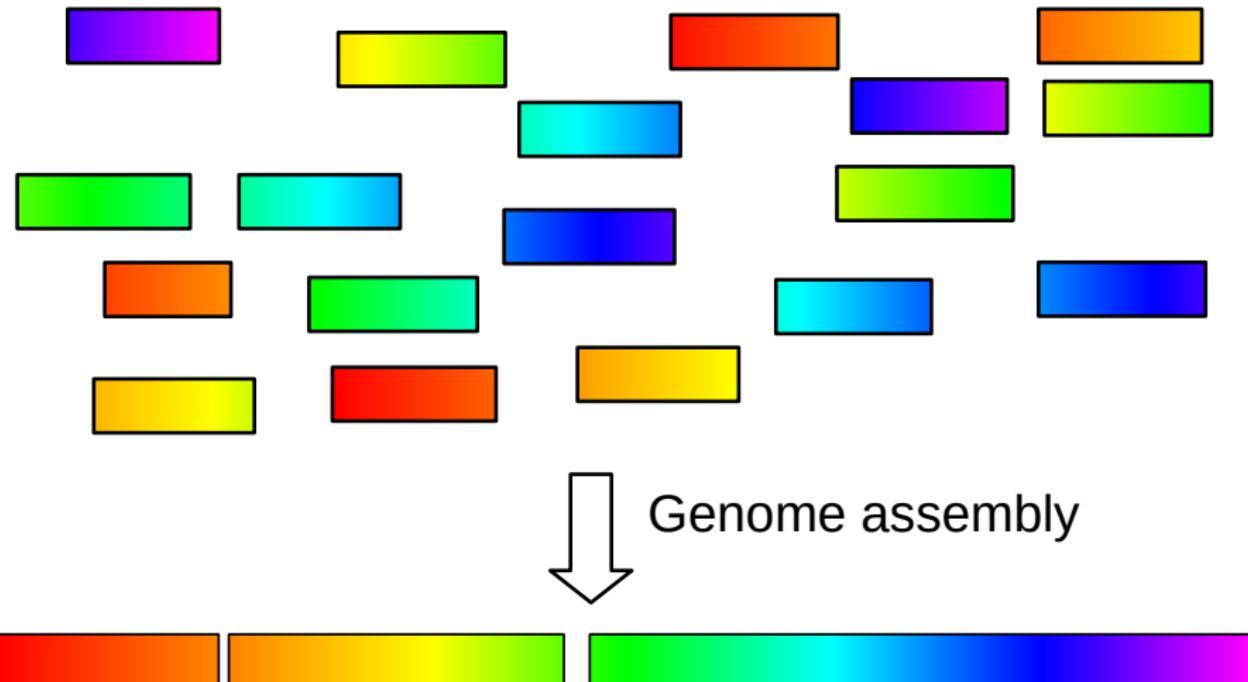
Repeats lead to the fragmentation of the assembly

Missing information also fragments the assembly

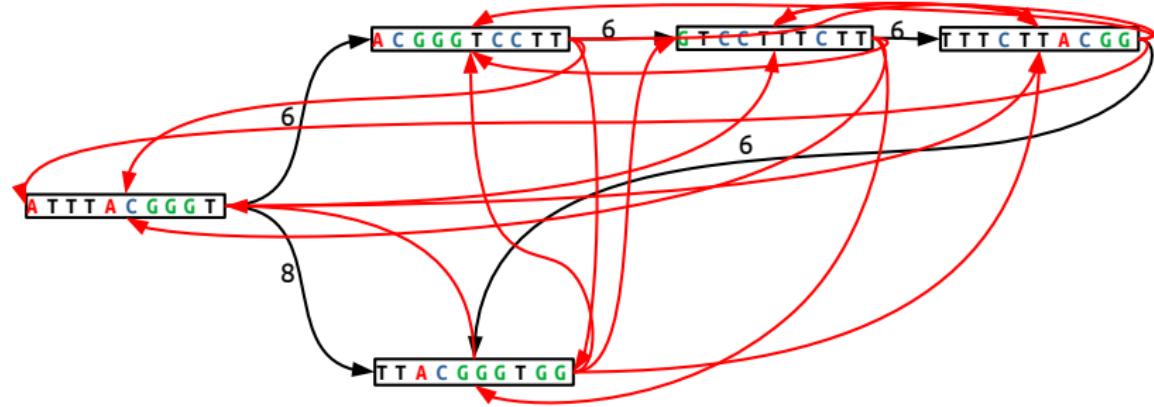


Assembly concession number 1: output fragments

In the real world, assemblers often provide pieces of genomes rather than complete ones

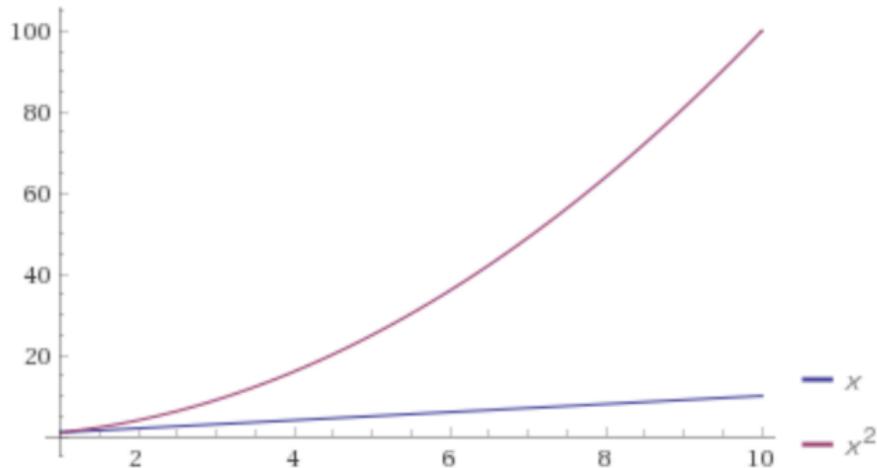


Overlap graph prerequisite : all overlaps



Overlap graph burden: number of reads

$n(n - 1)/2 = \mathcal{O}(n^2)$ possible overlaps for n reads



Linear: 2X data 2X time

Quadratic: 2X data 4X time

Overlap graph burden: number of reads

$n(n - 1)/2 = \mathcal{O}(n^2)$ possible overlaps for n reads

# Reads	# Overlaps
1000	499,500
10,000	50 million
100,000	5 billion
1 million	500 billion
10 million	50 trillion...

Most overlaps are too small to be considered...

The overlap computation is not linear

Overlap graphs in a nutshell

- Graphs of overlaps between the reads
- Can provide a global solution for assembly
- Can be difficult in real cases because it requires a lot of computation (overlaps)



S. cerevisiae, *D. melanogaster*, human could be assembled using overlap graphs approaches (Celera (Myers et al. 2000), SGA (Simpson & Durbin 2011), ...)

Fast forward



Assembly idea number 3: Focus on genome words

AGATACA  GTAGCATGCTAACTGTGACGGCATTAC
GCCATGACC read

a read is a word from the genome

words of size 7 from the genome:

AGATACA
GATACAG
ATACAGC
TACAGCC
...

k -mer definition

A k -mer is a word of size k

Exercise

List all **distinct** 4-mers and 5-mers of this set of reads:

CATAATGCA

TGCACATAA

CATAATGCA

Solution

CATAATGCA, TGCACATAA, CATAATGCA

5-mers:

CATAA, ATAAT, TAATC, AATGC, ATGCA, TGCAC, GCACA, CACAT, ACATA

CATAATGCA, TGCACATAA, CATAATGCA

4-mers:

CATA, ATAA, TAAT, AATG, ATGC, TGCA, GCAC, CACA, ACAT

Genome words / read words

Let's select a word from the genome:

AGATAACAGCCATGACCGTAGCATGCTAACTGTGACGGCATTAC

in the genome, after TAGCAT

only AGCATG appears

Genome words / read words

In real cases we don't have the genome

AGATAACAGCCATGACCGTAGCATGCTAACTGTGACGGCATTAC

in the genome, after TAGCAT

only AGCATG appears

but we have the reads

ATGACCGTAGCATGCT
ATGACCGTAGCATGCT
GACCGTAGCATGCTAA

in the reads, after TAGCAT

only AGCATG appears

Reconstitute larger genomic words

AGATAACAGCCATGACC GTAGC ATGCTAACTGTGACGGCATTAC

AGATACA

GATACAG

ATACAGC

TACAGCC

ACAGCCA

CAGCCAT

AGCCATG

AGATACA + G + C + C + A + T + G

AGATAACAGCCATG

a sequence from the genome

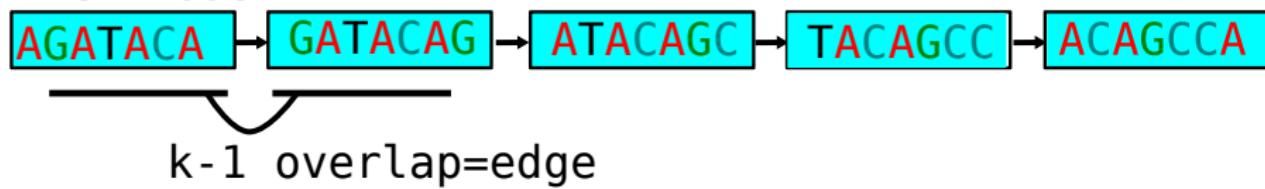
The De Bruijn graph

Read

AGATACAGCCA

De Bruijn graph

Kmer=node



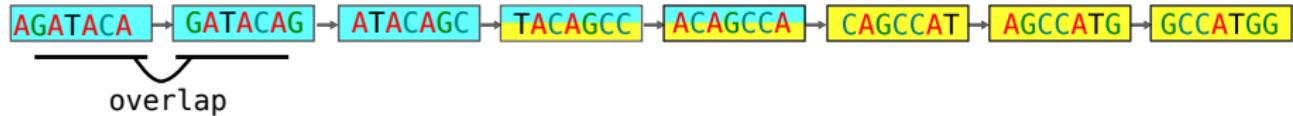
AGATACA + G + C + C + A
= AGATACAGCCA

De Bruijn graph assembly

Overlapping reads

AGATACAGCCA
TACAGCCATGG

De Bruijn graph



Resulting sequence

AGATACAGCCAATGG

De Bruijn graph time!

Reads

GCCATGGGTT
TACAGCCATGG
AGCCATGGGTT
GCCATGGGTT
AGATACAGCCA
ACAGCCATGGG
GATACAGCCAATG
CATGGGTTTAA
ACAGCCATGGG
GATACAGCCAATG
CATGGGTTTAA
CAGCCAATGGGT

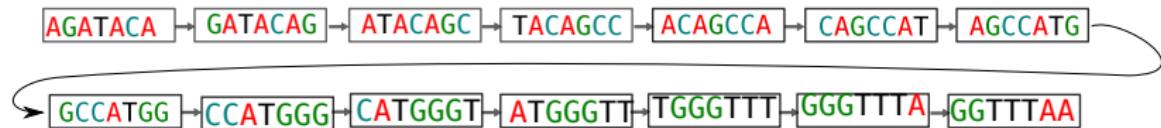
Hint: Use 7-mers

Solution

Overlapping reads

AGATAACAGCCA
GATACAGCCATG
GATACAGCCATG
TACAGCCATGG
ACAGCCATGGG
ACAGCCATGGG
CAGCCATGGGT
AGCCATGGGTT
GCCATGGGTTT
GCCATGGGTTT
CATGGGTTTAA
CATGGGTTTAA

De Bruijn graph



Resulting sequence

AGATAACAGCCATGGGTTTAA

De Bruijn graph versus overlap graph



AGATACA
GATA₁CAG TACAGCC AGCCATG words from the reads
ATACAGC ACAGGCCA CAGCCAT
...

word graph (De Bruijn graph)



Overlap graph from the reads



De Bruijn graphs abstract redundancy

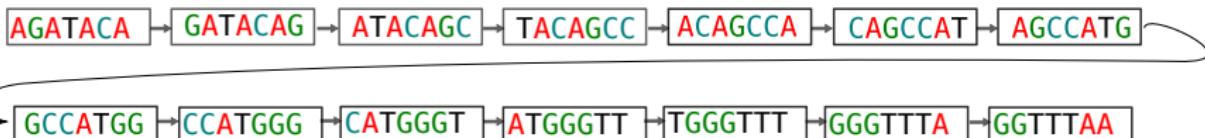
Overlapping reads

AGATACAGCCA
GATACAGCCATG
GATACAGCCATG
TACAGCCATGG
ACAGCCATGGG
ACAGCCATGGG
CAGCCATGGGT
AGCCATGGGTT
GCCATGGGTTT
GCCATGGGTTT
CATGGGTTTAA
CATGGGTTTAA

62 (**non distinct**) 7-mers in the reads

De Bruijn graph

14 **distinct** 7-mers in the De Bruijn graph



De Bruijn graphs only rely on $k - 1$ overlaps

Overlapping reads

AGATACAGCCA
GATAACAGCCATG
GATAACAGCCATG
TACAGCCATGG
ACAGCCATGGG
CAGCCATGGGT
GCCATGGGTT
CATGGGTTT
TGGGTTTA
GGTTTAA

Overlap length: 9

Overlap length: 10

Overlap length: 7

Overlap length: 8

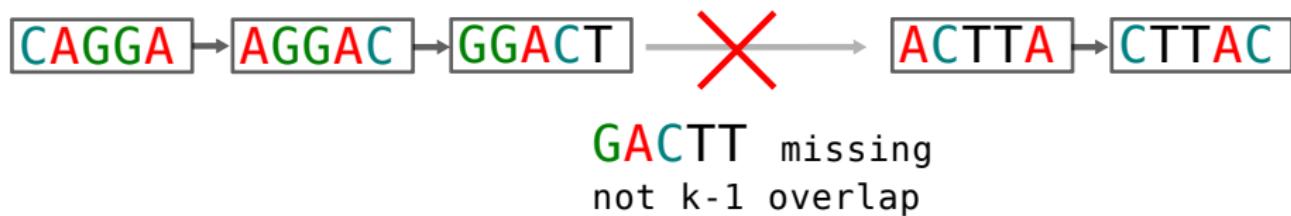
Overlap length: 6

De Bruijn graph overlap length: 6



De Bruijn graphs limitation

Fixed overlaps



De Bruijn graphs limitation

Repeats...

...TACAGGACTTA... ...TATAGGACTGA...



De Bruijn graph limitation

...TACAGGACTTA... ...TATAGGACTGA...



genome pieces

...TATAGGA

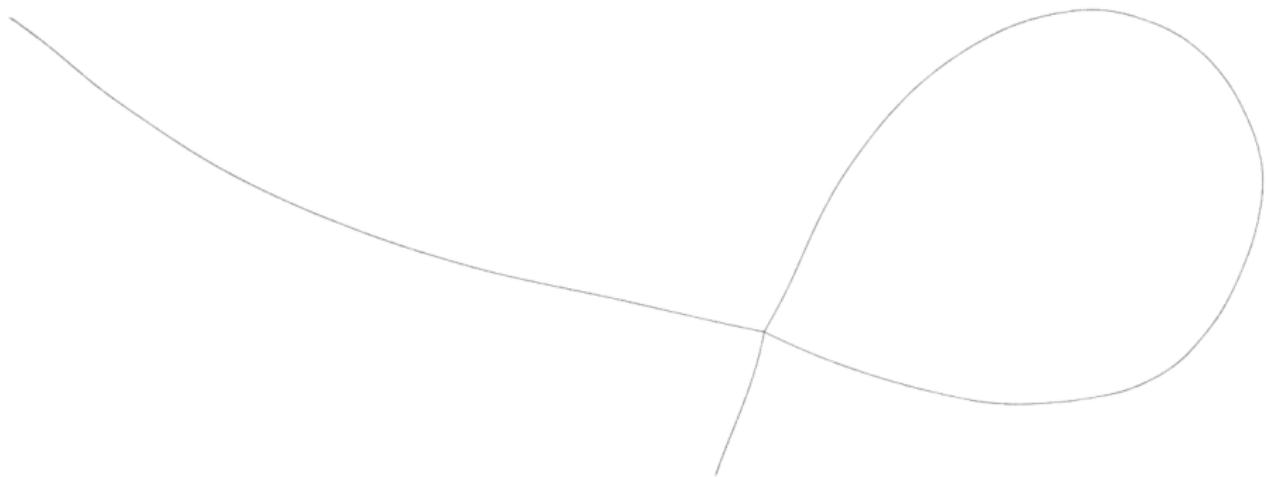
GACTGA...

AGGACT

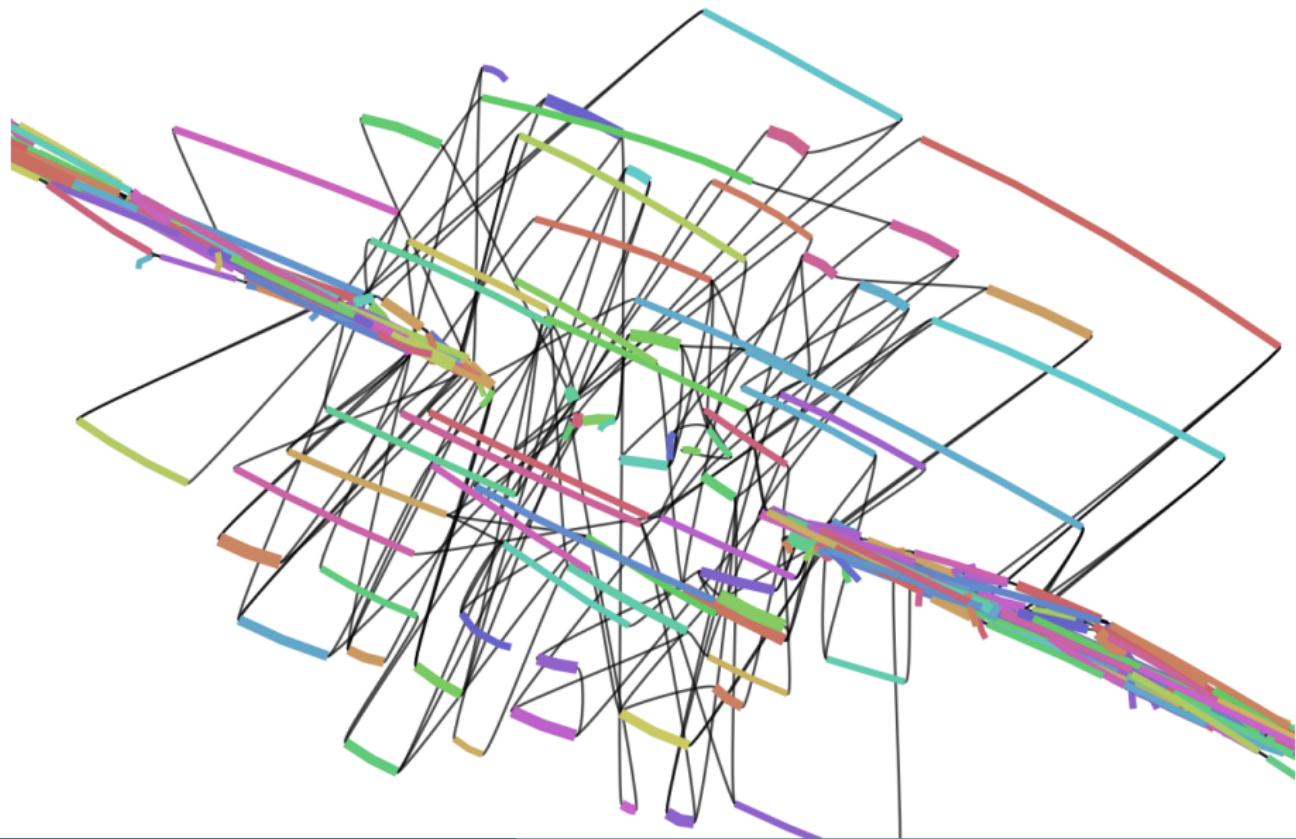
...TACAGGA

GACTTA...

De Bruijn graph on a real dataset



De Bruijn graph on real dataset ZOOMED IN

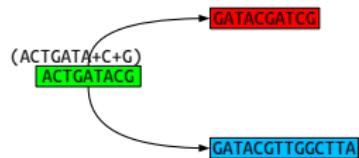


On the representation of De Bruijn graphs

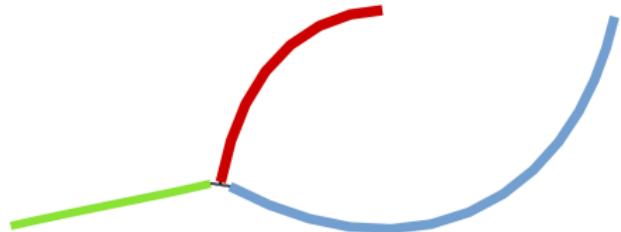
De Bruijn graph:



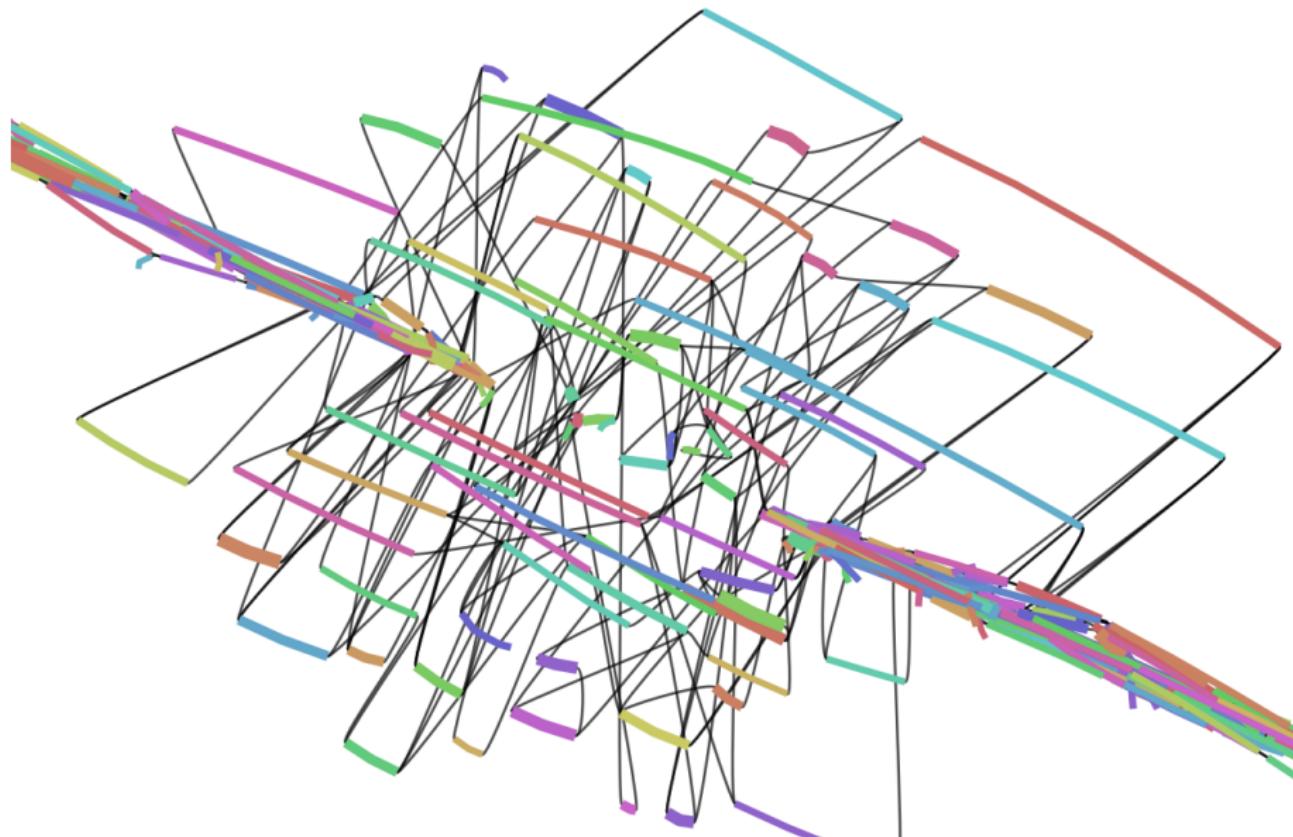
Compacted De Bruijn graph:



Graphical representation (.gfa plot using Bandage):



De Bruijn graph on a real dataset ZOOMED IN



Dealing with sequencing errors

Genome:

AT**C**GGT**A**T**C**G**T**T**A**CG**G**T**A**T**A**CC

Reads:

AT**C**G**C**T**A**TCG
GGTT**T**CG**T**TA
AT**C**GA**T**ACGG

TCG**C**TA
GGTT**T**C
AT**C**GA**T**

...

Are not genomic kmers...

Erroneous k -mers vs genomic k -mers

Genome:

TAAGAAAGCTCTGAATCAACGGACTGCGACA

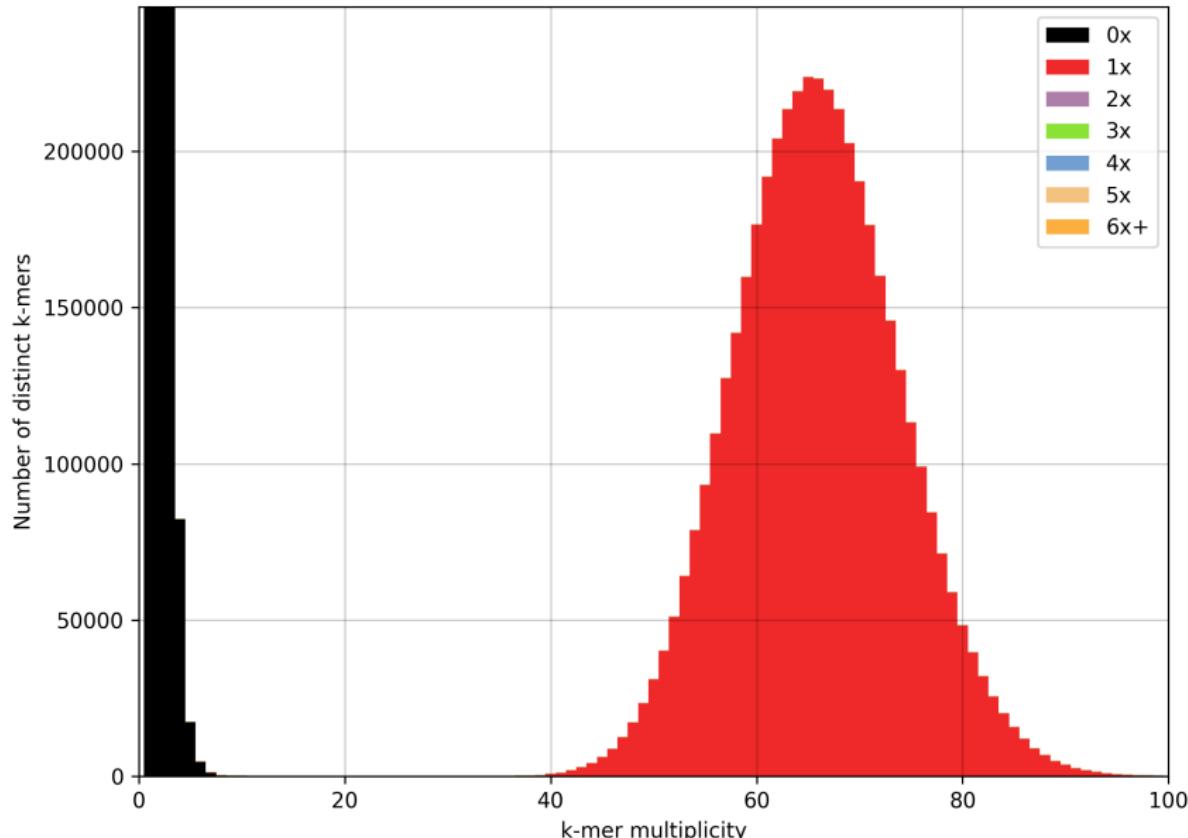
Reads:

TAAGAAAGCTCTGAATCA		
AAGAAAGCTCTAAATCAAC		
AGAAAGCTCTGAATCAACG	9 times	TCTGAAT
GAAAGCTCTGAATCAACGGA	1 time	TCTAAAT
AAAGCTCTGAATCAACGGAC		
AAGCTCTGAATCAACGGACT	6 times	CAACGGA
AGCTCTGAATCAACGGACTG	1 time	CAACGGT
GCTCTGAATCAACGGTCTGC		
CTCTGAATCAACGGACTGCG		
TCTGAATCAACGGACTGCGA		

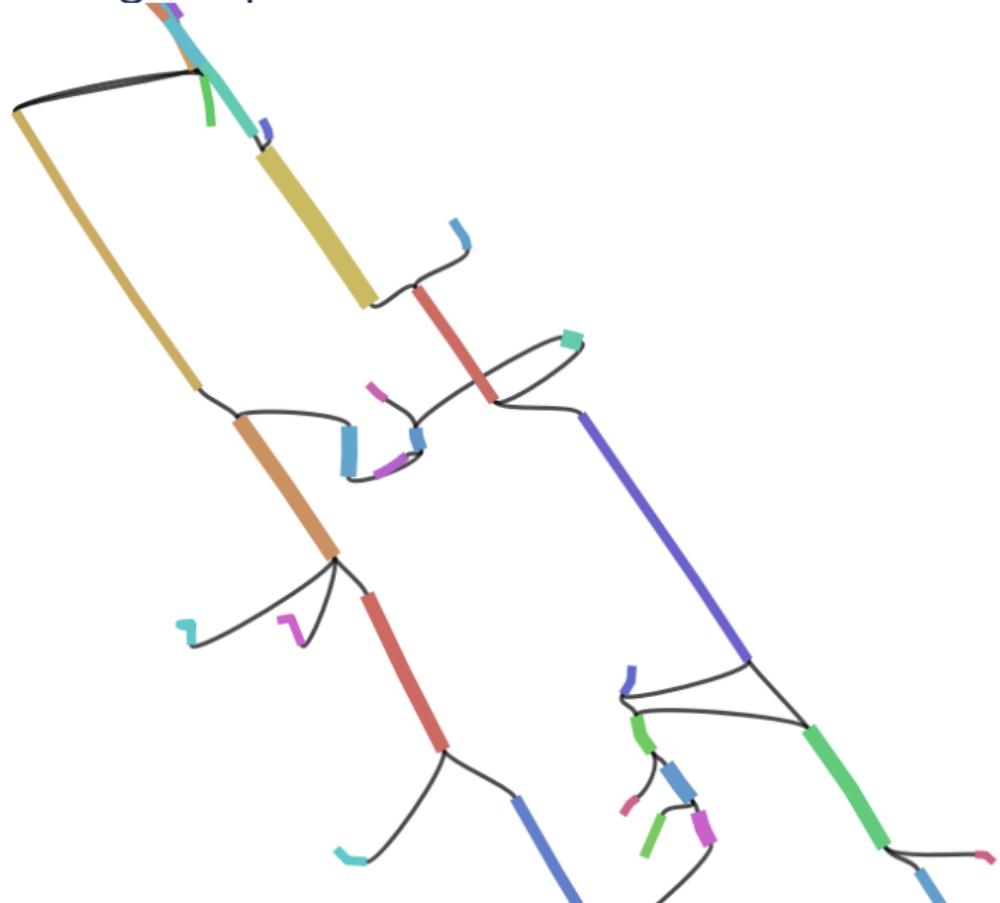
Erroneous k -mers are seen less than genomic ones

K-mer histogram

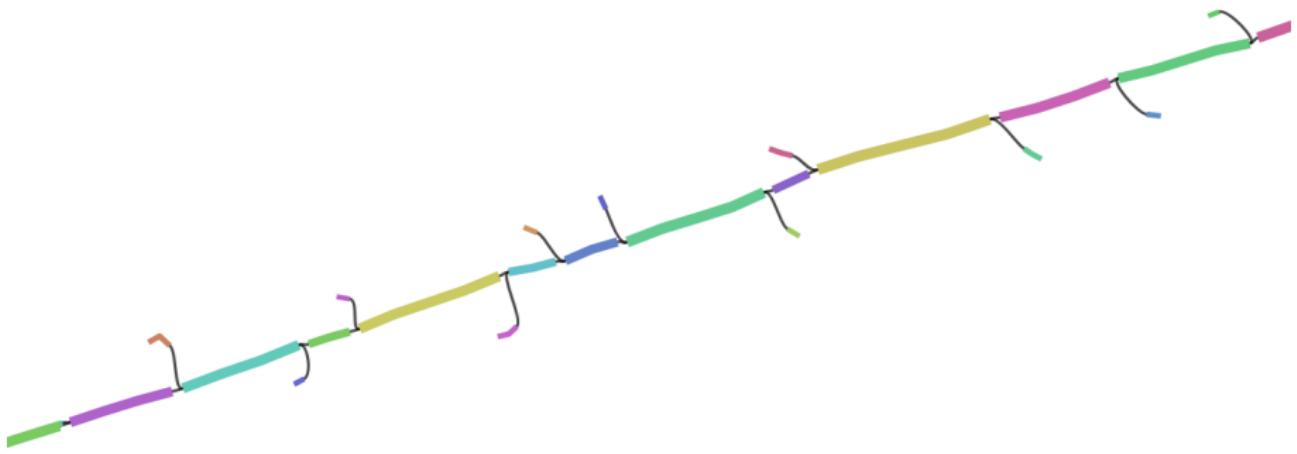
K-mer comparison plot



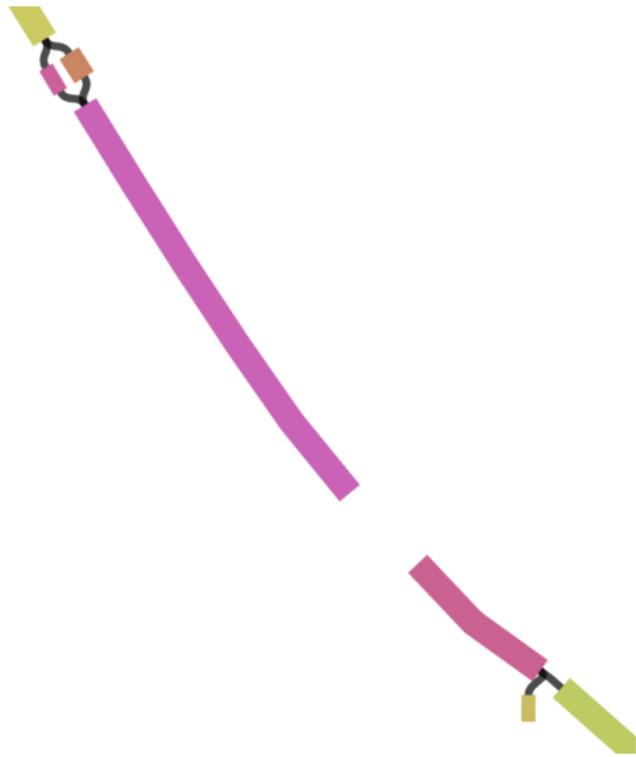
Removing unique k -mers



Removing k -mers seen less than 3 times



Removing k -mers seen less than 4 times



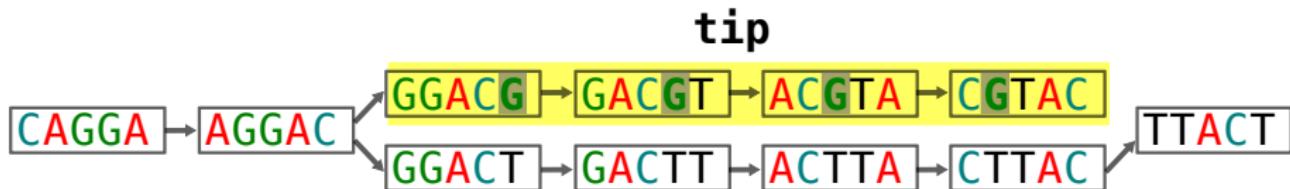
Errors in De Bruijn graphs



Errors in De Bruijn graphs

...TACAGGACTTACTGA... genome

reads CAGGACTTA
AGGACGTAC ← **sequencing error**
AGGACTTAC
GGACTTACT

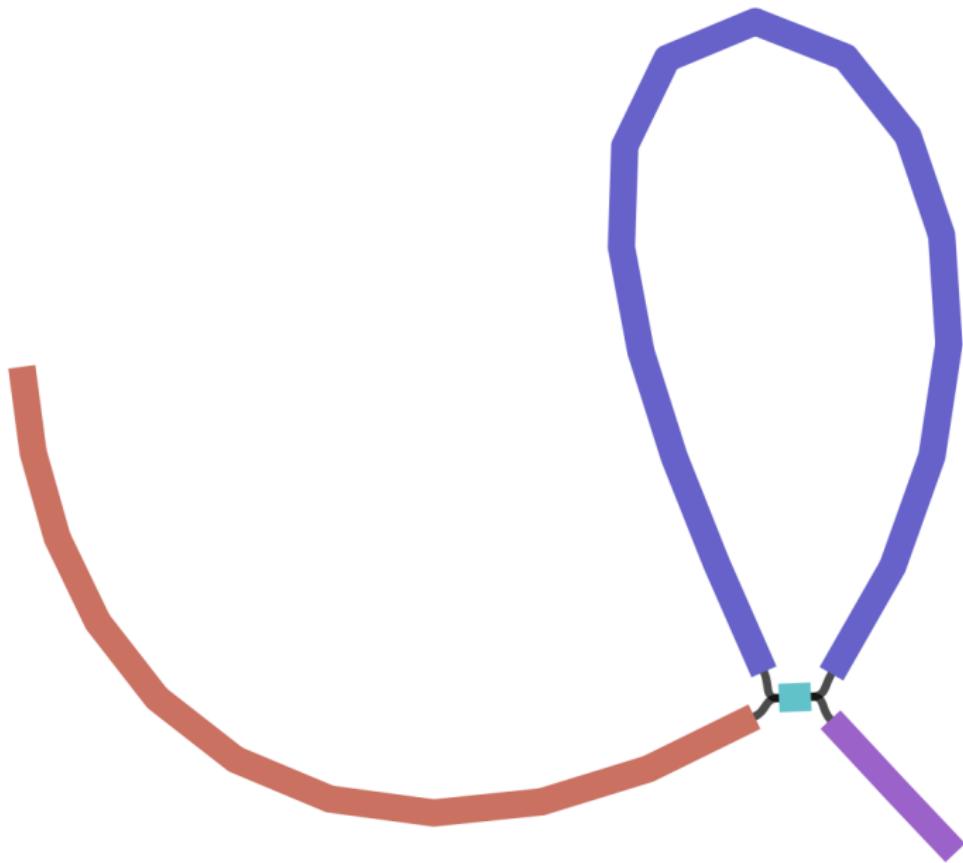


Errors in De Bruijn graphs

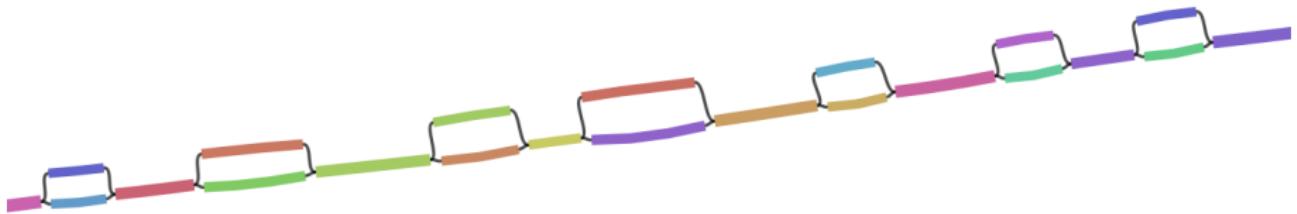
... TACAGGACTTACTGA... genome
reads ACAGGACTTA
 CAGGAATTAC ← **sequencing error**
 CAGGACTTAC
 AGGACTTACT



(Almost assembled phage !)



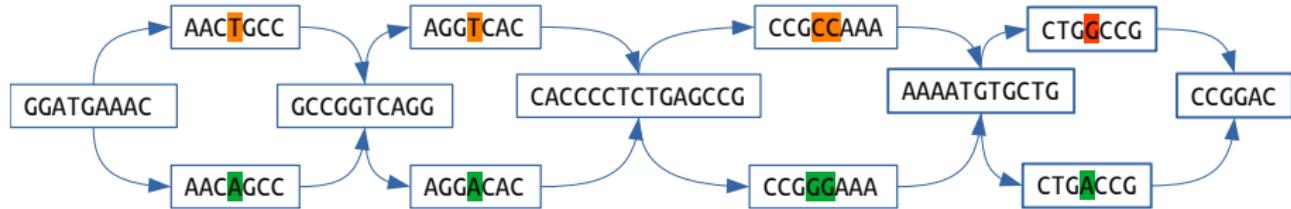
De Bruijn graph on an eukaryota



Two or more genomes per individual

♀ GGATGAAACTGCCGGTCAGGTACCCCCCTTGAGCCGCCAAAATGTGCTGGCCGGAC

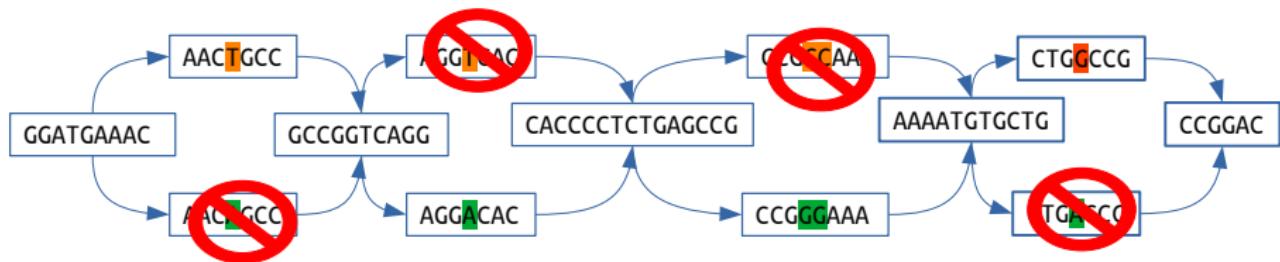
♂ GGATGAAACAGCCGGTCAGGAACCCCCCTTGAGCCGGAAAATGTGCTGACCGGAC



Two or more genomes per individual

♀ GGATGAAACTGCCGGTCAGGTACCCCTCTGAGCCGCCAAAATGTGCTGCCGGAC

♂ GGATGAAACAGCCGGTCAGGACACCCCTCTGAGCCGGAAAATGTGCTGACCGGAC



Assembly:

GGATGAAACTGCCGGTCAGGACACCCCTCTGAGCCGGAAAATGTGCTGCCGGAC

Assembly concession number 2: collapse variability

♀ GGATGAAACTGCCGGTCAGGTACCCCTCTGAGCCGCCAAAATGTGCTGCCGGAC

♂ GGATGAAACAGCCGGTCAGGAACCCCTCTGAGCCGCCAAAATGTGCTGACCGGAC

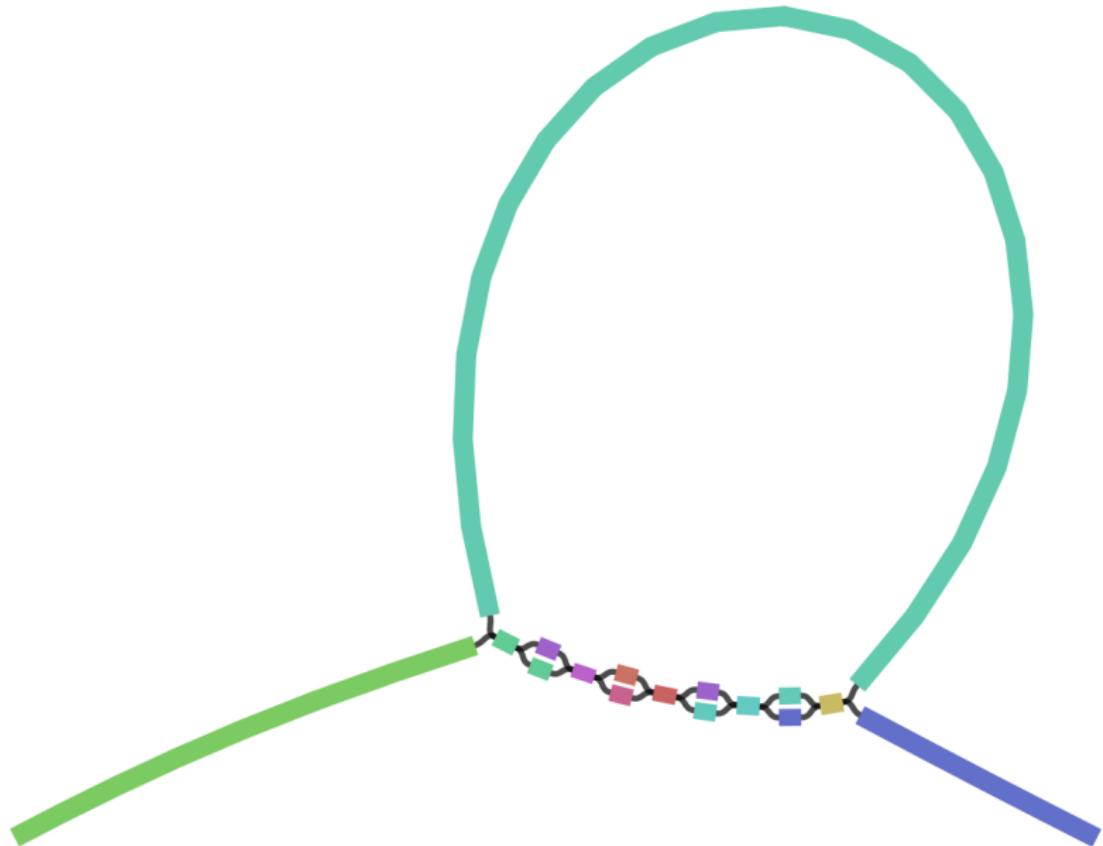
Assembly:

GGATGAAACTGCCGGTCAGGAACCCCTCTGAGCCGCCAAAATGTGCTGCCGGAC

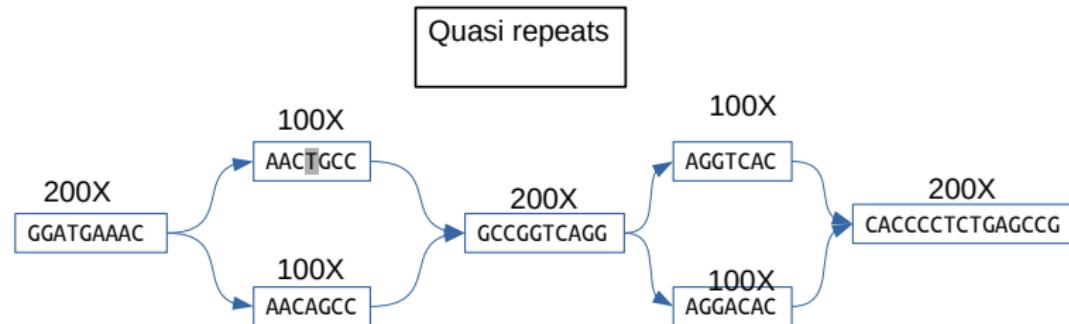
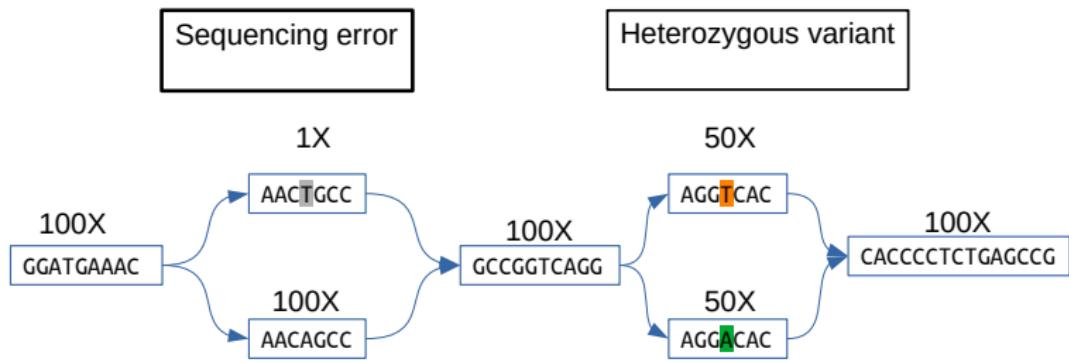
Reads:

GATGAAACTG
ATGAAACAGC
TGAAACAGCCG
GAAACTGCCGG
AAACTGCCGGT
AACAGCCGGTC
ACAGCCGGTCA
CTGCCGGTCAG

Paralog genes/repeats



Paralog genes/repeats in graph



An assembler is a set of heuristics

Graph cleaning heuristics

- Nodes coverage
- Graph local/global topology
- Reads that can be mapped on nodes
- Estimated coverage/genome size
- ...

An assembly is a model

Different tools can produce very similar assemblies

A single tool can produce very different assemblies with small changes of parameters(!)

De Bruijn graphs in a nutshell

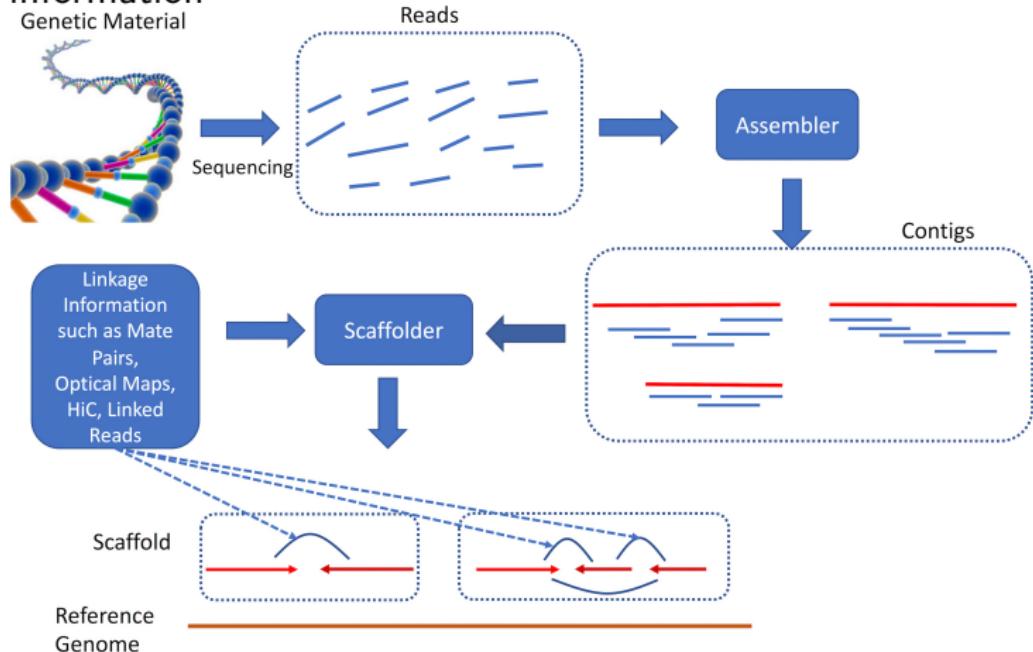
- Graph of words of size k , $k-1$ overlaps
- Collapses identical k -mers
- Very successful, have replaced the overlap graphs with high throughput sequencing data
- Still outputs fragments of the genome



white spruce, 20 gigabases

Scaffolding

Softwares can improve the assembly continuity by using other kinds of information



From "Modern technologies and algorithms for scaffolding assembled genomes" Plos Computational Biology

Second generation sequencing



NextSeq Series



HiSeq 4000 System



HiSeq X Series[‡]



NovaSeq 6000
System

- Short reads $\approx 150\text{bp}$
- Low error rate $\approx 1\%$
- High throughput (up to billion of reads per run)
- GC bias

Mainly assembled using De Bruijn graphs

State-of-the-art

Your toolkit for the practical session

- SPAdes
- Megahit
- Minia



Other notable assemblers

- SGA
- Discovar denovo
- Abyss

Third generation sequencing



- Long reads $\approx 10 - 100\text{ kbp}$
- High error rate $\approx 10 - 12\%$
- High throughput (up to millions of reads per run)

Nanopore VS Pacbio

Nanopore

- Portable
- Ultra long reads (100kbases, some reads reach the megabase level)
- Mostly deletions

Pacbio

- More mature
- HiFi reads (99% identity)
- Mostly insertions

Long reads killed the assembly star



Laura Landweber @LandweberLab · Jan 2

Our newest version of *Oxytricha*'s somatic genome is out ([rdcu.be/bZNfC](https://doi.org/10.1371/journal.pgen.1009371)) and has 18,617 distinct chromosomes. That's 2000 more than we previously published in [doi.org/10.1371/journal....](https://doi.org/10.1371/journal.pgen.1009371) PacBio captured most chromosomes in single reads: Genome sequence, No assembly required

Repeats spanning

Genome:

GGTAATGGTTTTTGGTGCTAATGCGTTTTTTCATGGATGTCGTAATTTTTATCTG

Reads:

GGTAATG TTTTTT GTGCTAAT GTTTTTT ATGGATG TTTTTTA
ATGGTTT AATGCGTT ATGTCGT CGTAATT TTTATCTG
TTTGGTG TTTCATG

Contexts of the repeat:

...ATGG

ATCT...

???TTTTT???

GGTG...

...TGC

CATG...

...GTAA

Repeats spanning

Genome:

GGTAATGGTTTTTGGTGCTAA_{TGCGTTTTT}CATGATGTCTGAATTTTTATCTG

Reads:

GGTAATGTTTTTGTGCTAA_{AATGCCTT}ATGGATGATGTCTGTTTTTA
ATGGCTT_{TTTGGTG}TTTTCATGCTGAATT_{TTTATCTG}

Long reads:

TGGTTTTTGGT TGCC_{TTTTTT}CAT TGAATTTTTATCT

Contexts of the repeats:

...ATGG → GGTG... ...TGCG → CATG... ...GTAA → ATCT...

Repeats spanning

Contexts of the repeats:

...ATGG → GGTG...
...TGC... → CATG...
...GTAA → ATCT...

Overlapping Reads:

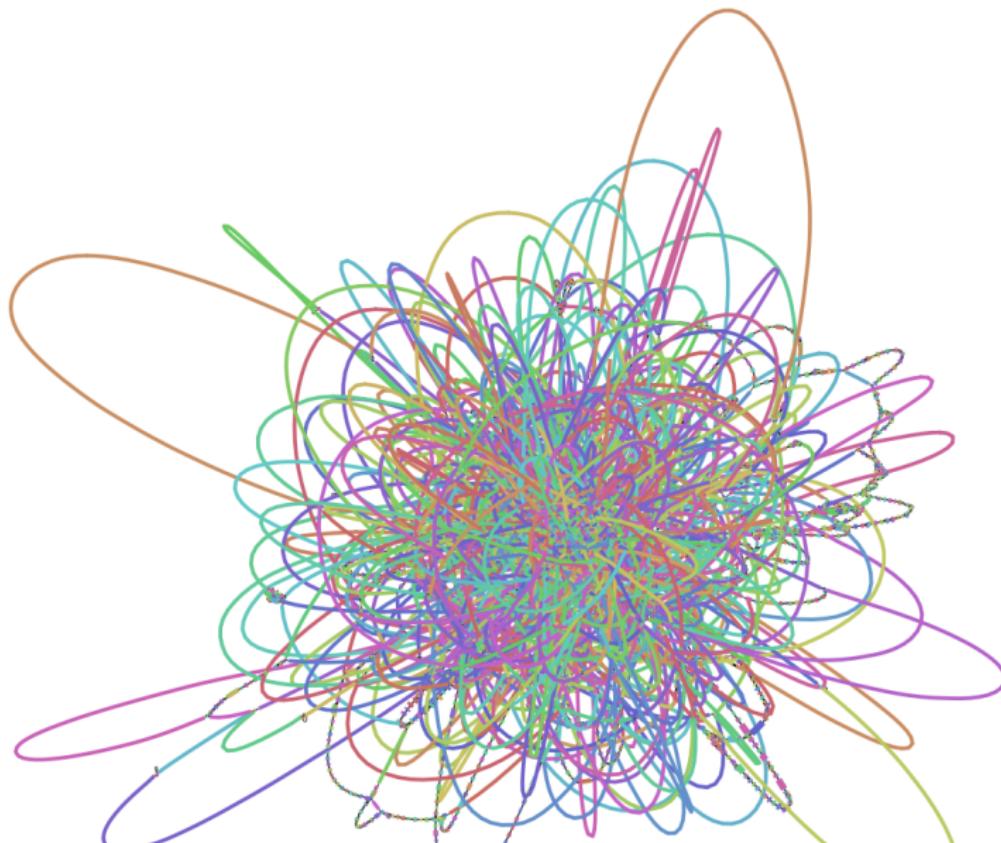
GGTAATG	GTGCTAAT		
ATGGTTT	AATGCGTT	ATGTCTG	
TGGTTTTTGGT	TGCCCTTTTCAT	CTGAATT	
TTGGTG	TTTCATG	TGAATT	TTTATCTG
	ATGGATG	TATCTG	

Output assembly:

GGTAATGGTTTTGGTGTAAATGCGTTTCATGGATGTCTGAATTATCTG

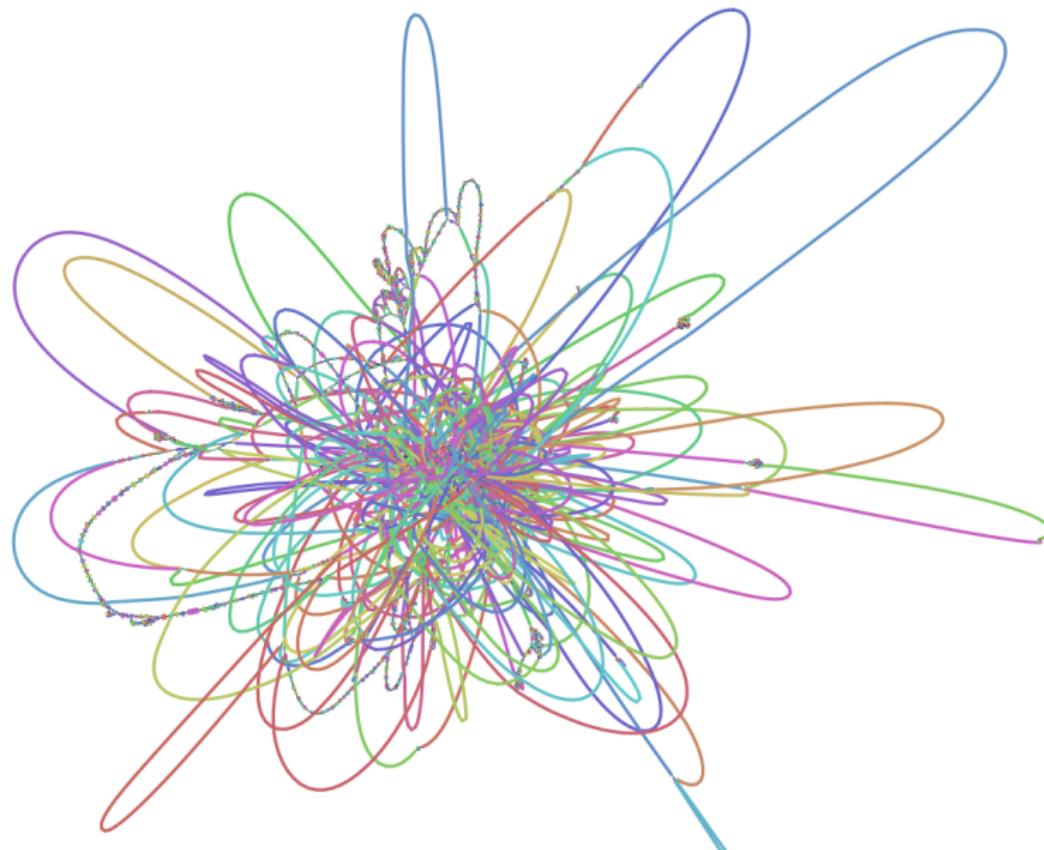
Read length matter

Read size=21



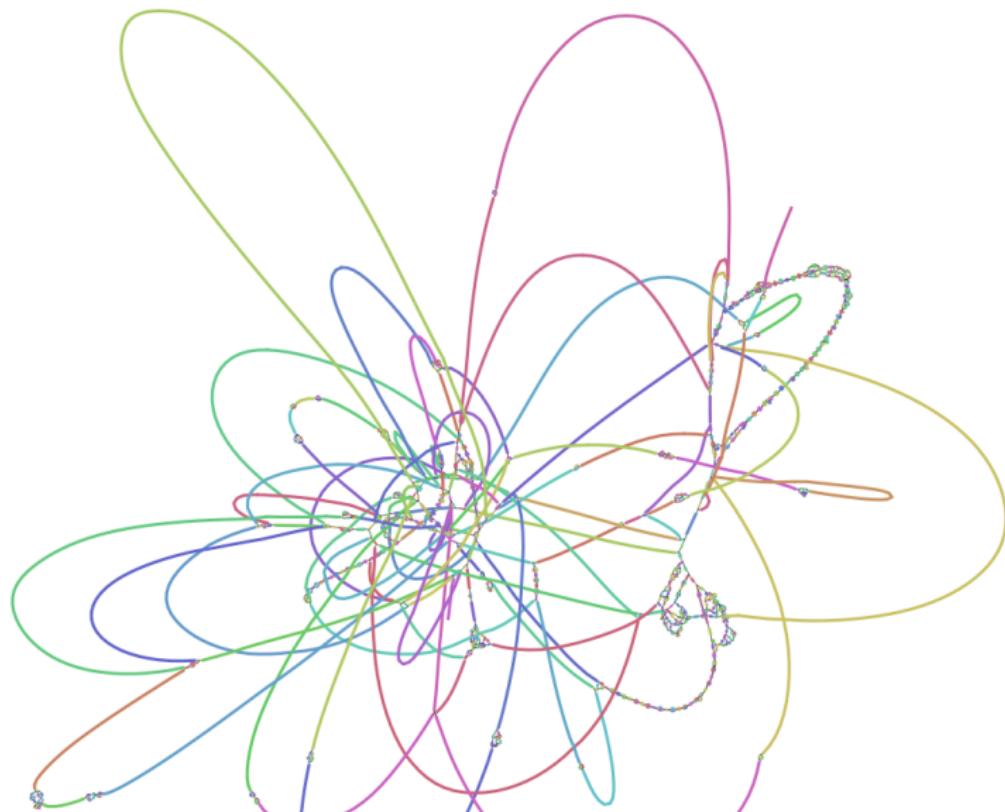
Read length matter

Read size=31



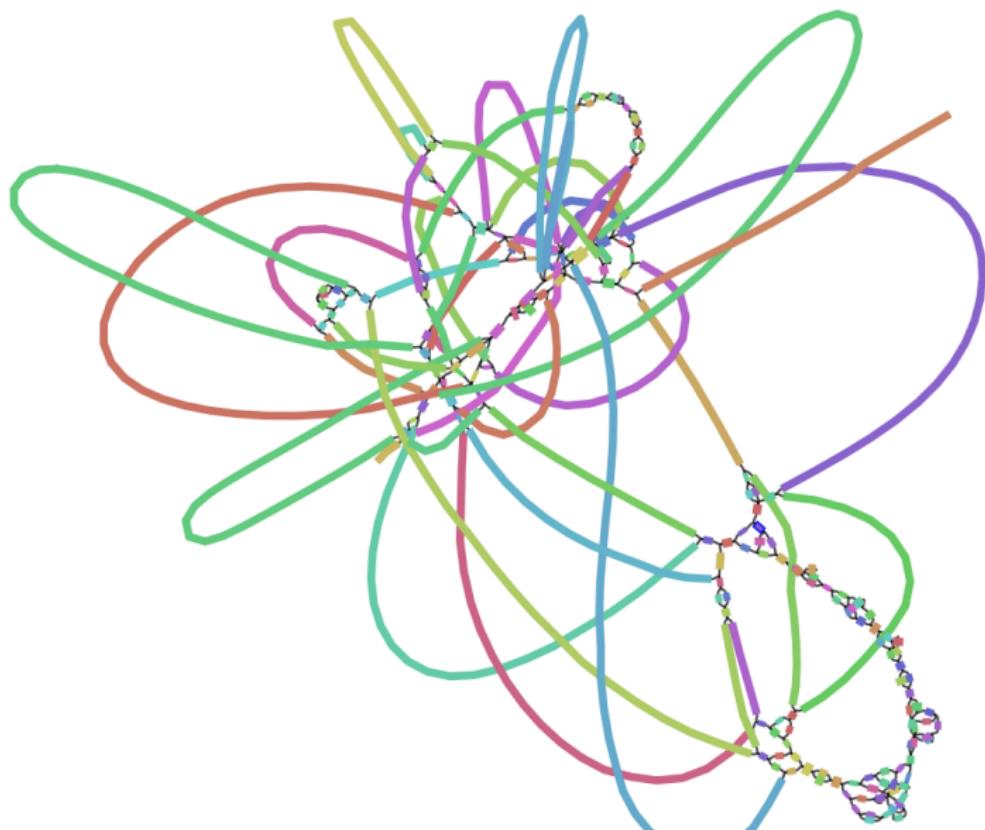
Read length matter

Read size=63



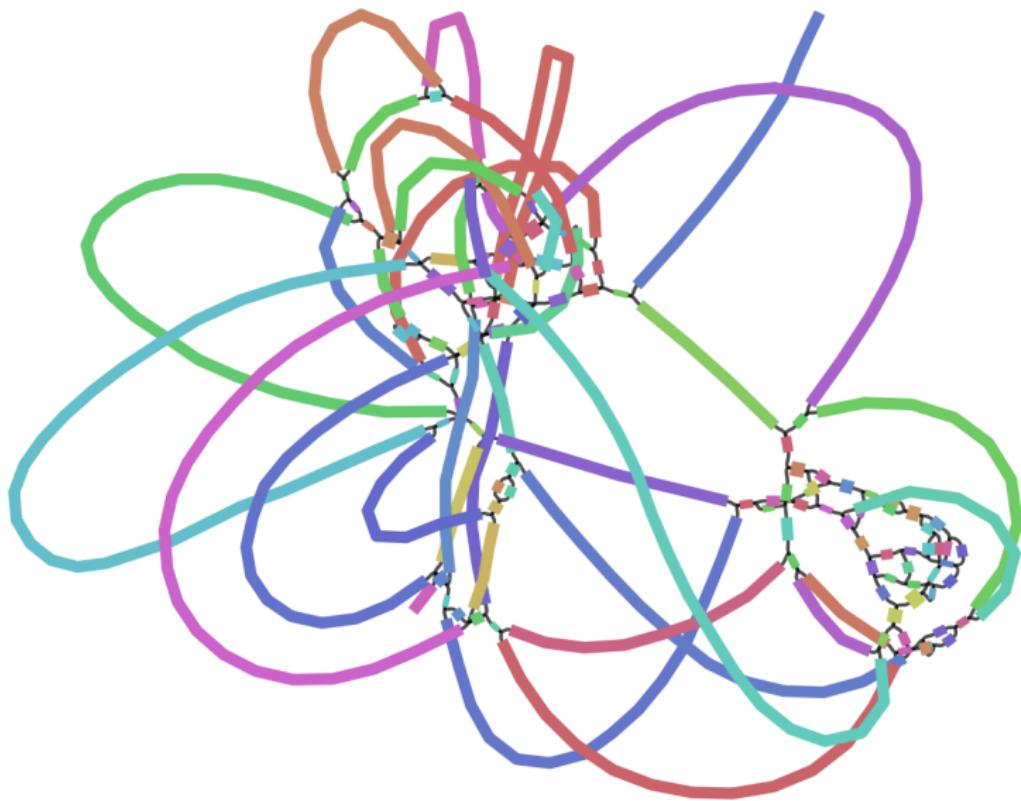
Read length matter

Read size=255



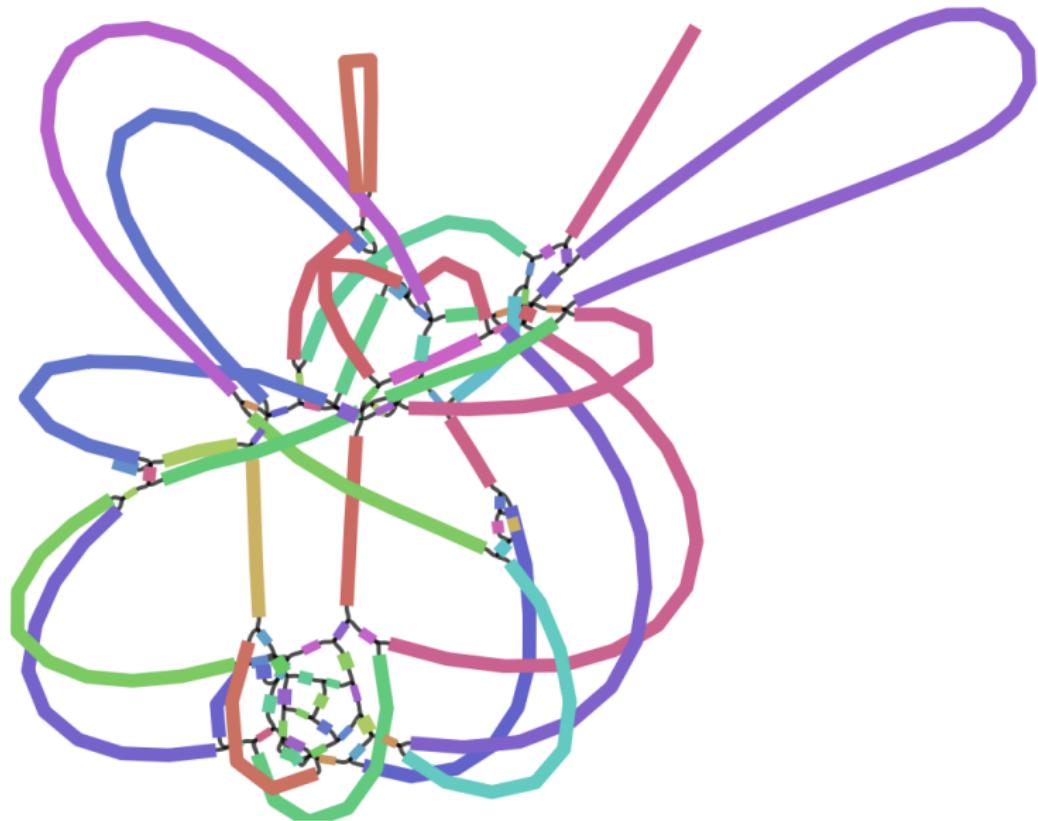
Read length matter

Read size=500



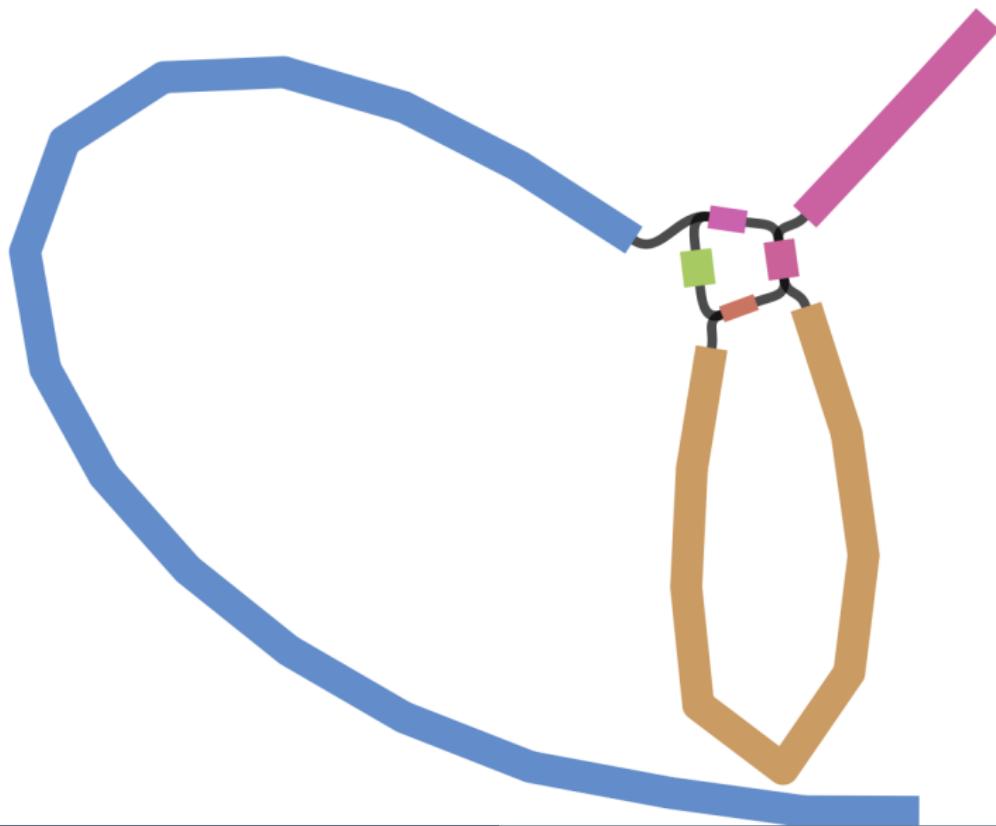
Read length matter

Read size=1000

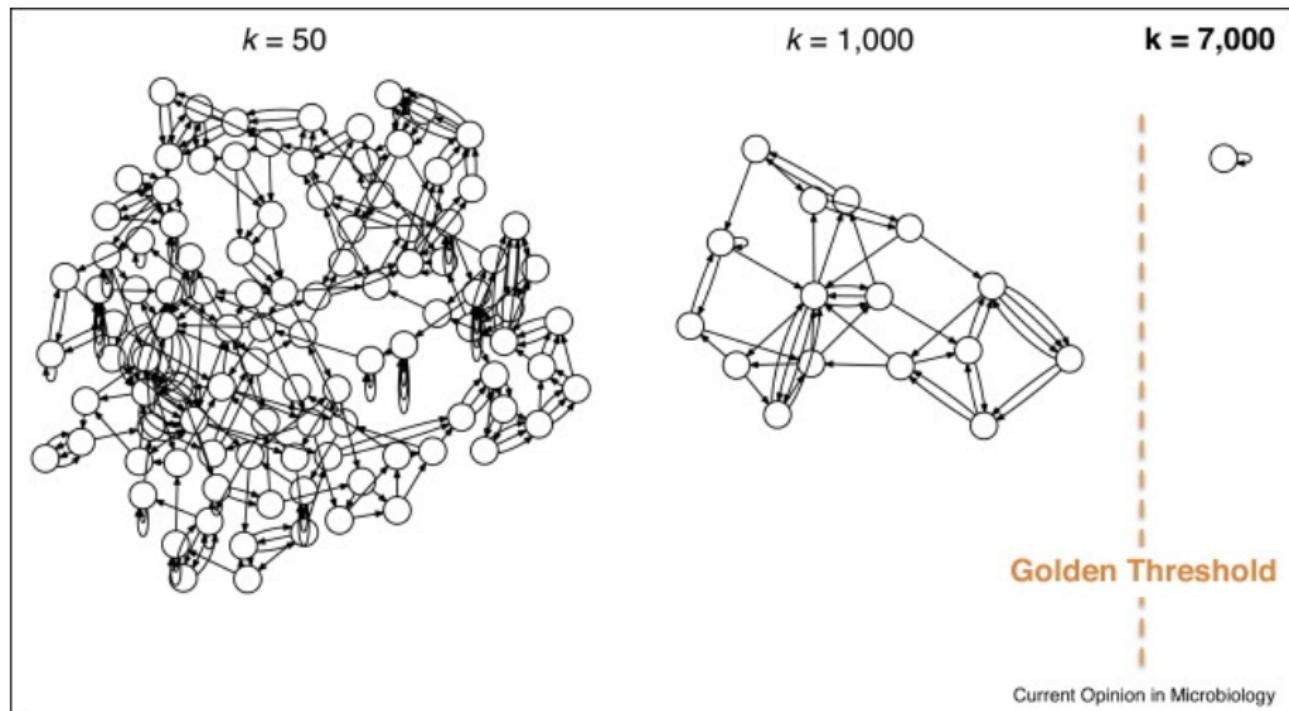


Read length matter

Read size=2000

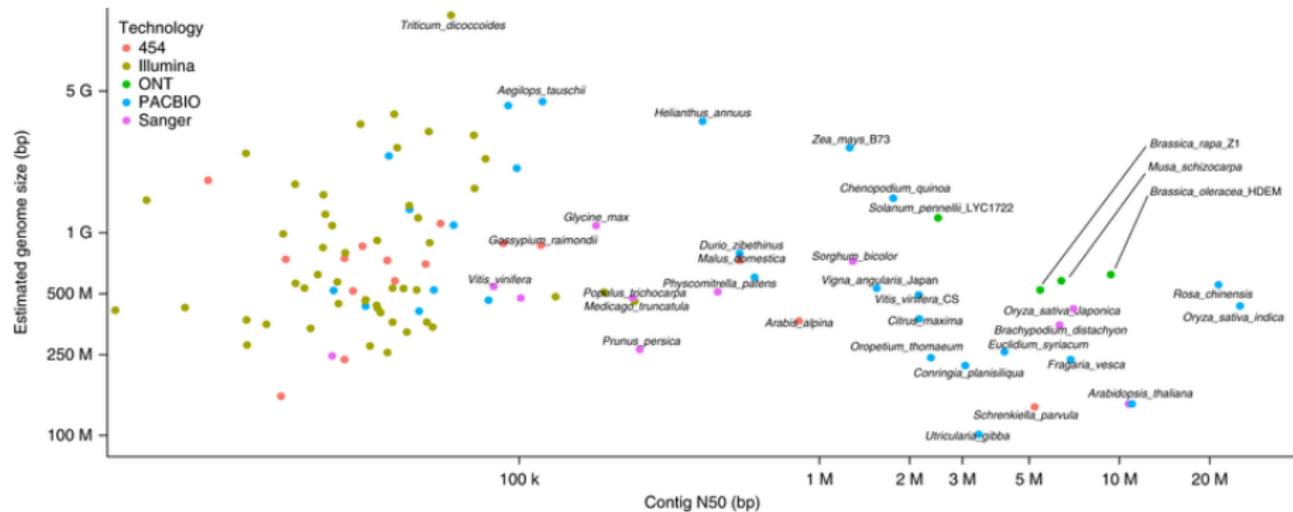


Great hope for assembly



From "One chromosome, one contig: complete microbial genomes from long-read sequencing and assembly" Current Opinion in Microbiology 2015

Great hope for assembly



From "Chromosome-scale assemblies of plant genomes using nanopore long reads and optical maps" Nature Plants 2018

Which assembly strategy is best suited?

- Long reads $\approx 10\text{ kbp}$
- High error rate $\approx 12\%$
- High throughput (up to millions of reads per run)

Based on long reads properties, which assembly solution would you choose and why?

Vote!

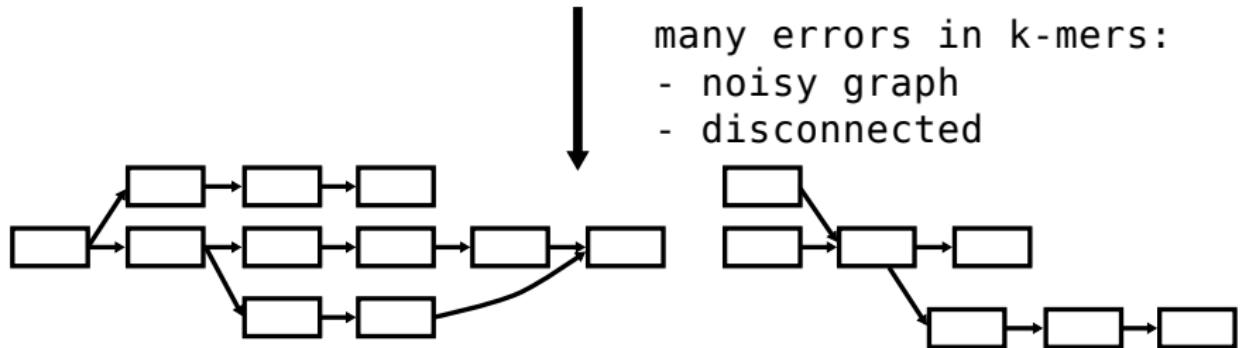
- Fish
- Greedy
- Overlap graph
- Birbs
- De Bruijn graph

Long reads for assembly: De Bruijn graph?

genome



reads



Long reads for assembly: overlap graph?

Supposed to be super expensive!

TAAGAAAGCTCTGAATCACGGACTCGCACAATAAGTGGTGTATCCAGAATTGTCACTTCAAGTAAAAACACCTCACGAGTTAACACCTAAGTTC
TAAGAAAGCT
AAGAAAGCTC
AGAAAGCTCT
GAAAGCTCTG
AAAGCTCTGA
AAGCTCTGAA
AGCTCTGAAT
GCTCTGAATC
CTCTGAATCA
TCTGAATCAA
CTGAATCAAC
TGAATCAACG
GAATCAACGG
AATCAACGGA
ATCAACGGAC
TCAACGGACT
CAACGGACTG
AACGGACTGC
ACGGACTGG
CGGACTCGGA
GGACTCGGAC
GACTGGGACA
ACTGGGACAA
CTGGGACAAT
TGGGACAATA
...

Average coverage: 10
Read length: 10
Average overlap: 9
Read number: 100

TAAGAAAGCTCTGAATCACGGACTCGCACAATAAGTGGTGTATCCAGAATTGTCACTTCAAGTAAAAACACCTCACGAGTTAACACCTAAGTTC
TAAGAAAGCTCTGAATCACGGACTCGCACA
GAAAGCTCTGAATCACGGACTCGCACAAT
AGCTCTGAATCACGGACTCGCACAATAAG
TCTGAATCACGGACTCGCACAATAAGTGG
GAATCAACGGACTCGCACAATAAGTGGTGT
TCAACGGACTCGCACAATAAGTGGTGTATCCA
ACGGACTCGCACAATAAGTGGTGTATCCAG
GACTCGCACAATAAGTGGTGTATCCAG
TGGGACAATAAGTGGTGTATCCAGAATTG
GACAATAAGTGGTGTATCCAGAATTG
ATAAAGTGGTGTATCCAGAATTGTC

Average coverage: 10
Read length: 30
Average overlap: 27
Read number: 33

Longer reads, better overlaps

- Less reads for the same coverage
- Larger overlaps

5Mb bacteria example with 100X coverage

Short reads

- 5 million 100bp reads
- 99 bp average overlap

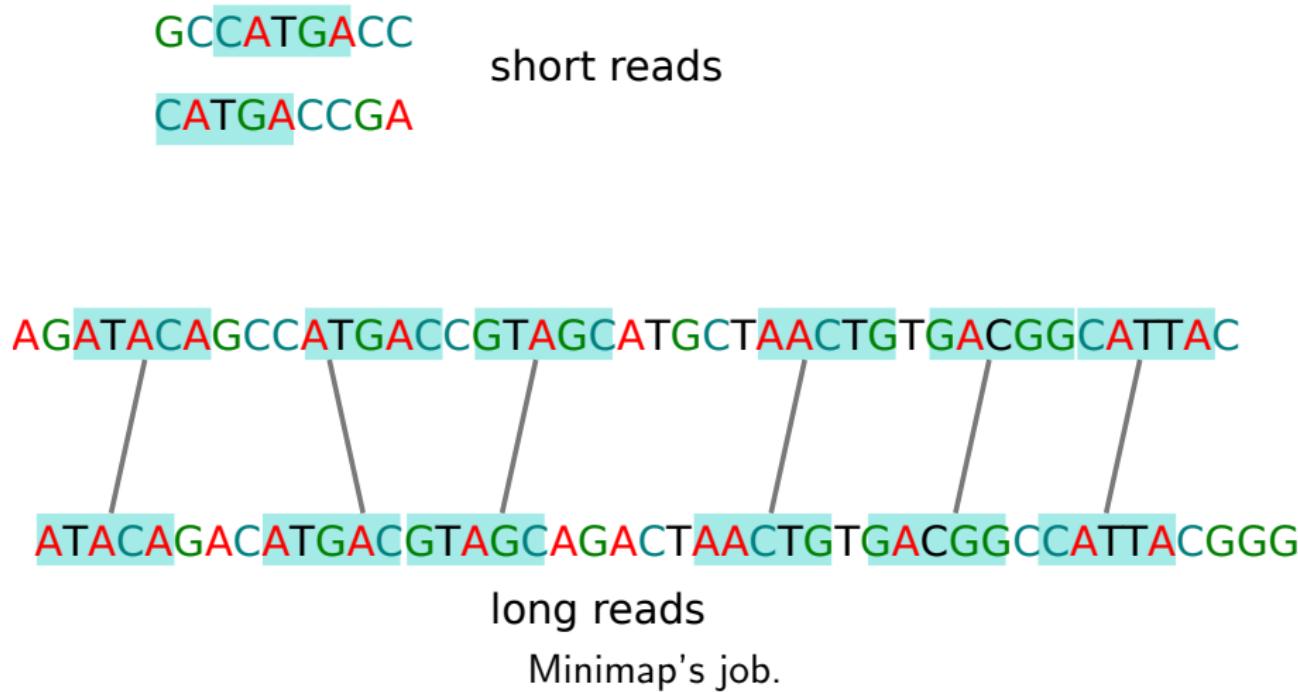
Long reads

- 50,000 10kbp reads
- 9,900 bp average overlap

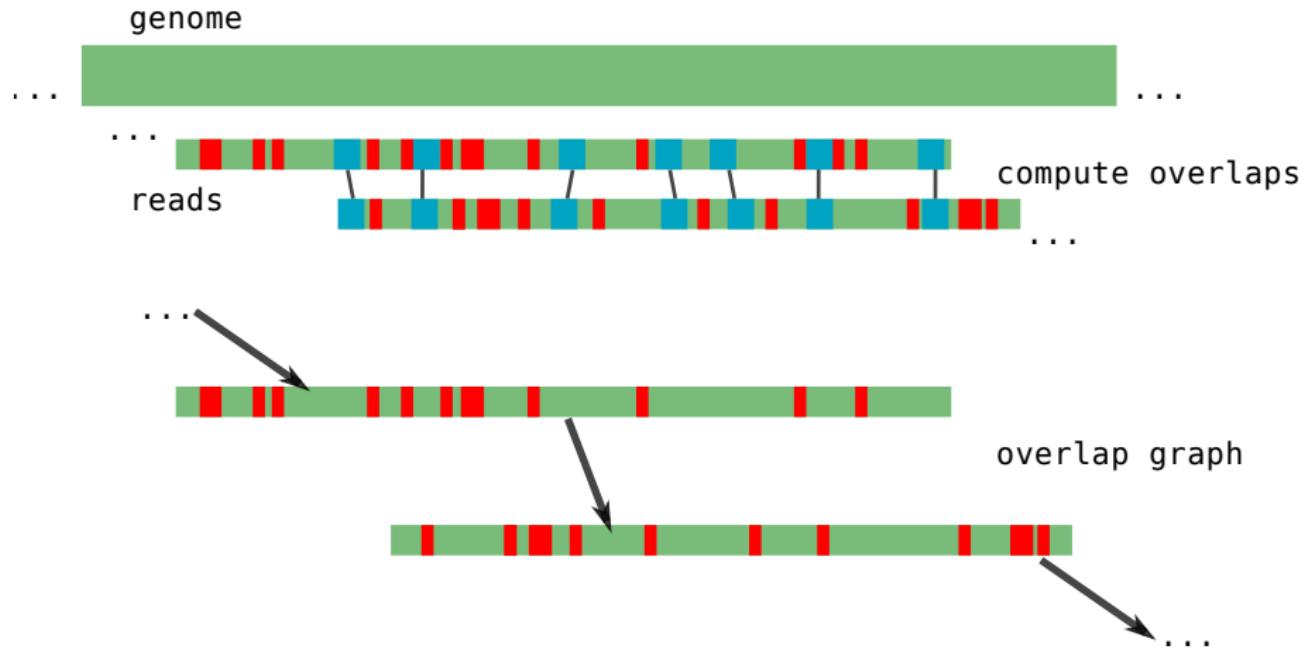
Very long reads

- 5,000 100kbp reads
- 99,000 bp average overlap

Anchors chaining in overlap graph



Long reads for assembly: overlap graphs



Sequencing errors

Genome:

ATCGGTATCGTTACGGTATAACC

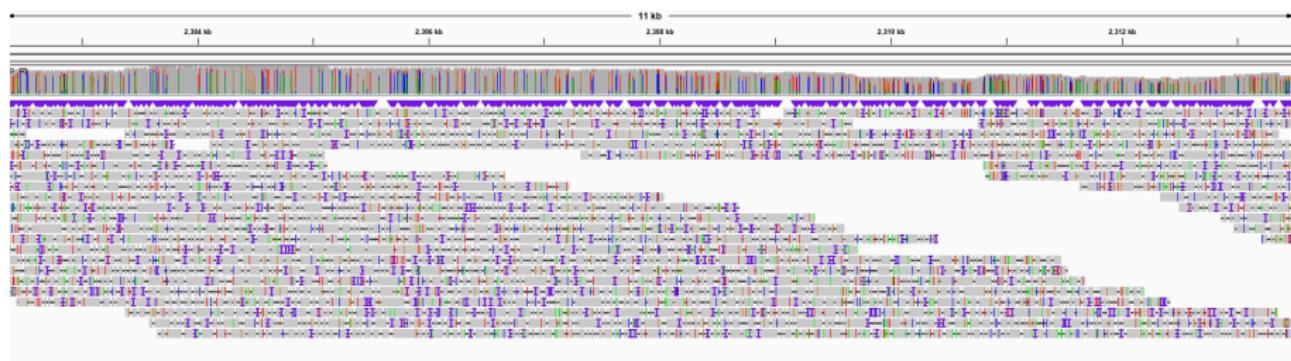
Reads:

ATCG**G**TATCG
GGTAT**C**GT**T**A
AT**T**GTTACGG

(Substitution)
(Insertion)
(Deletion)

Insertion and deletion made calling almost impossible

Sequencing errors



Miniasm

Long reads can be assembled without taking care of the sequencing errors ("Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences", Bioinformatics '16)

Genome characteristics and structure can be quickly estimated

Using coverage to remove noise: Consensus

Genome:

TAAGAAAAGCTCTGAATCAACGGACTGCGACAATAAGTGGTGGTATCCAGAATTTCACCTT

Reads:

AAAGAAAAGCACTGAATCATGGGACTTCGAG
GAAAGCTCTCAACCAACGGACTGCGACTTT
ACCTCTCAAGCAACGGACTGCGACAAAAAAG
TCTGAATCACCGGACTGCGTCAAAAAGTGC
GAATCACCGGACTGCGACAGTTGTGGTGG
TCAACCGCACTGCGACAATAAGTCCTGGTAT
ACGGACTGCGACAAAAAGTGTGGTATCCA
GAECTGCCACAAAAAGTGGTGGTATCCAG
TCCGACAAAAAGTGGGGTATCCAGAAT
GACAATAAGGGGGTGGTATCCAATTTG
AAAAAGGGGTGGTATCCAGAATTTC
TAAGTGGGGTATCCAATTTTCAGTT

Consensus:

AAAGATAGCTCTGAATCAACGGACTGCGACAAAAAGTGGTGGTATCCAGAATTTCAGTT
1/1 4/7 9/10 6/11 3/4

Exercise: Perform a consensus

Erroneous reads:

TAAGAAAGCTCTGAATCAAACGGTACTGCGA
GAAAGCTTGAATCAAACGGACTGCGACAA
AGCTCTGAATCAACGGACTGCGACAATAA

Contig to polish:

TAAGAAAGCTTGAATCAACGGAATGGCGACAATAA

Exercise: Perform a consensus - solution

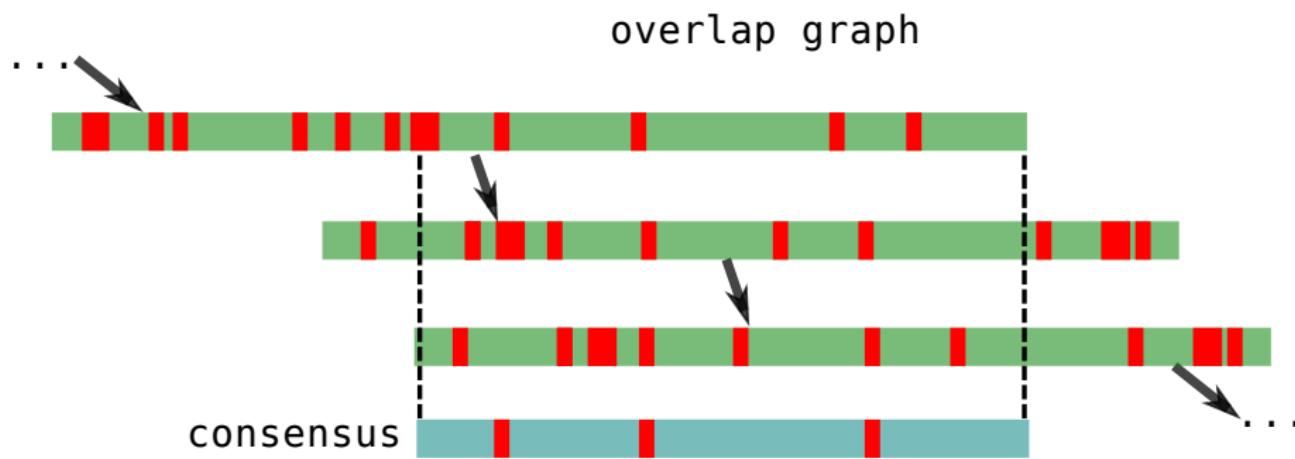
Correct contig:

TAAGAAAGCTCTGAATCAA-CGGACTG-CGACAATAA

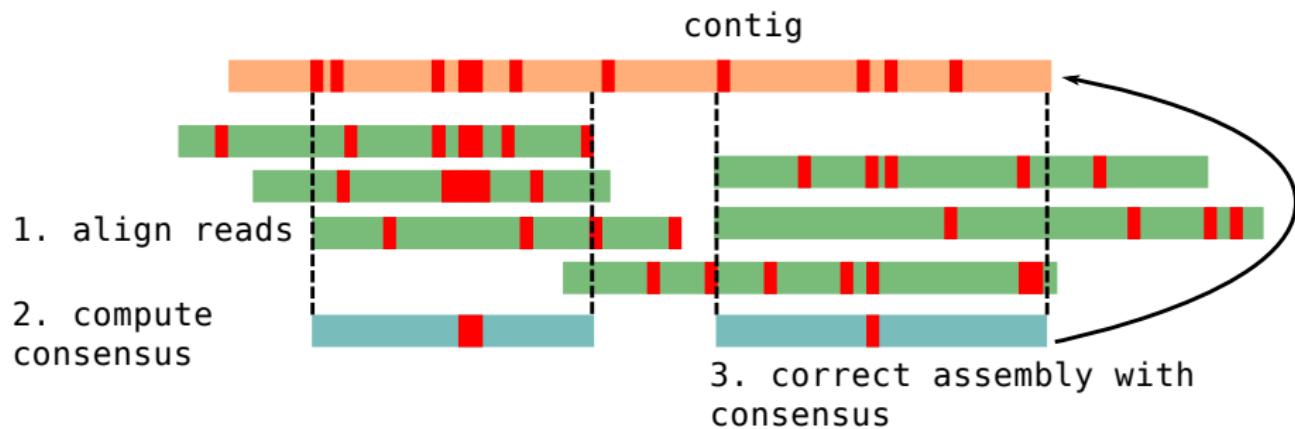
Aligned reads:

TAAGAAAGCTCTGAATCAAACGGT**A**CT**G**CGA
GAAAG**C**T-TGAATCAAACGGACTG-CGACAA
AGCTCTGAATCAA-CGGACTG-CGACAATAA

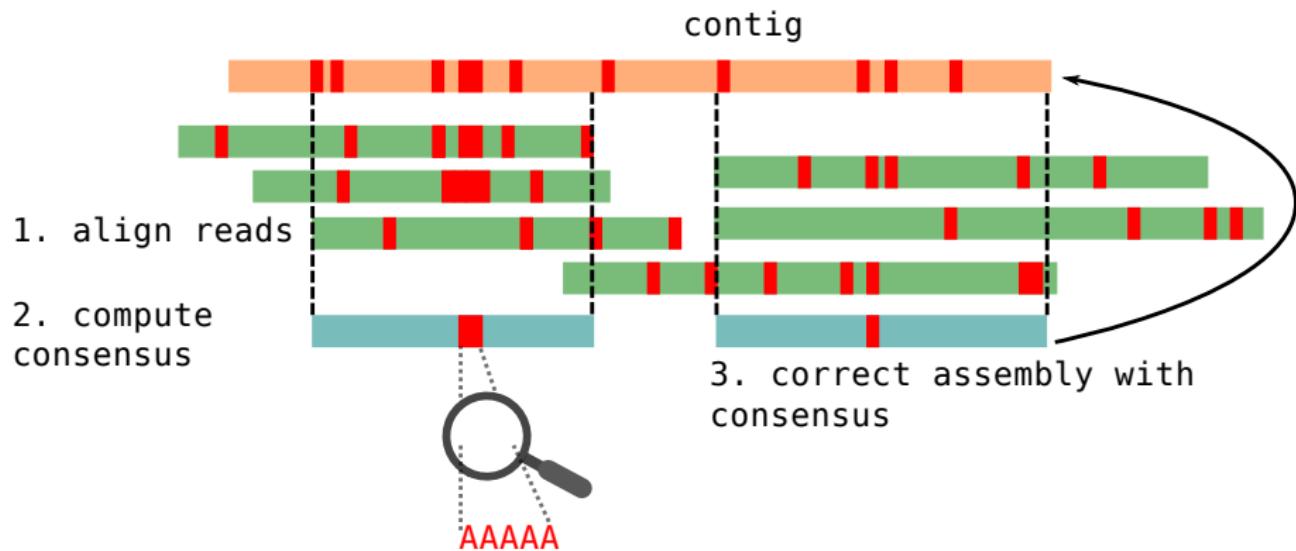
Consensus during assembly



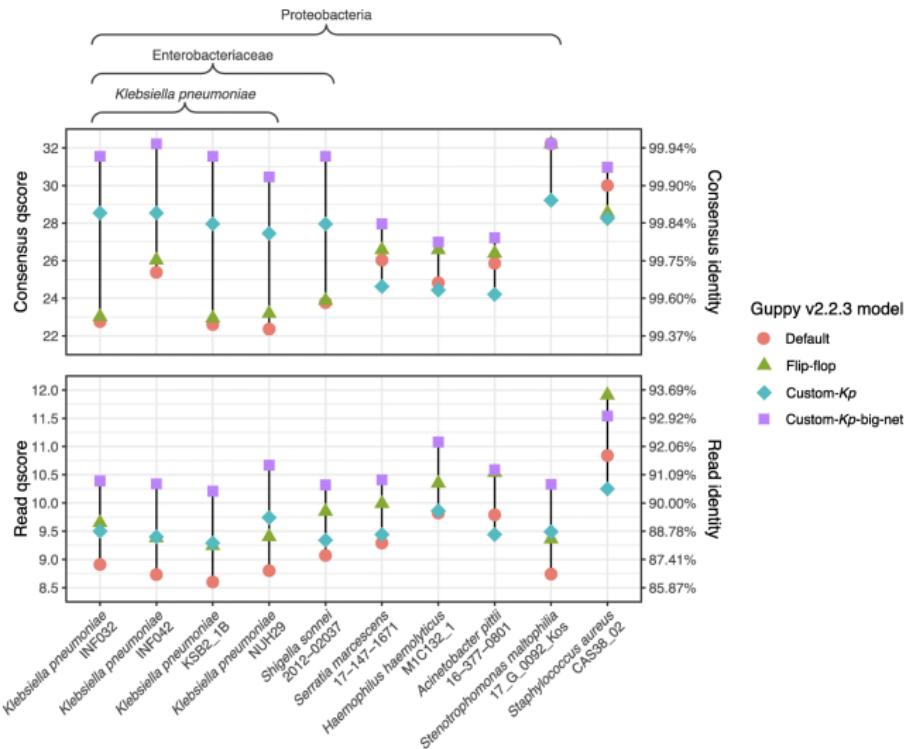
Consensus after assembly: polishing



Consensus after assembly: polishing

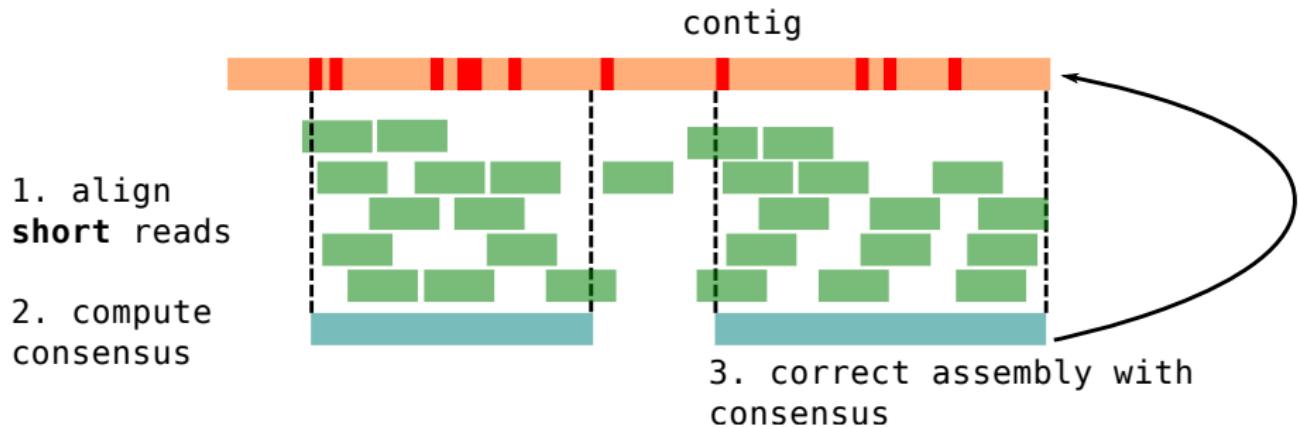


Systematic errors



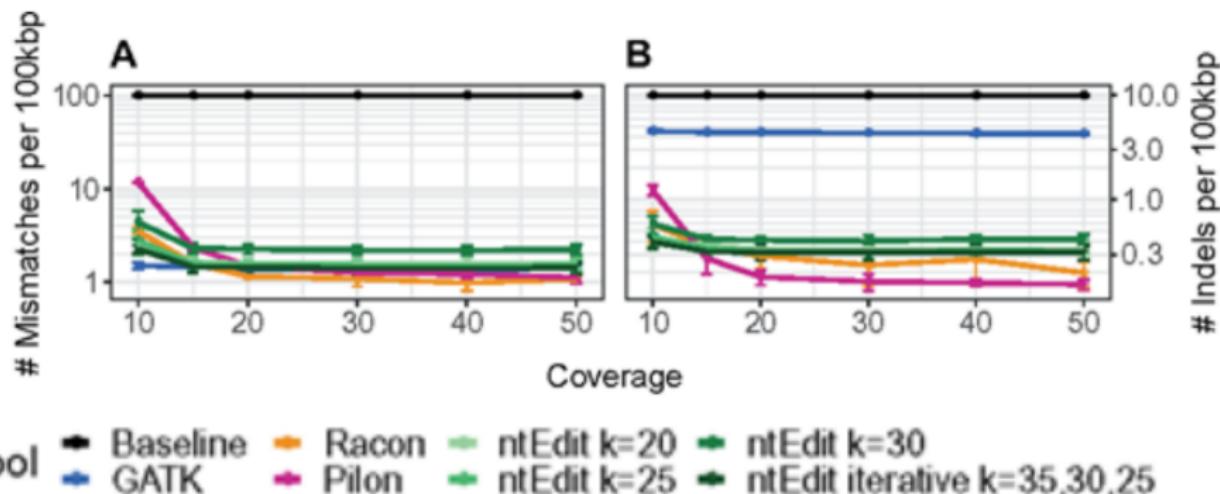
From "Performance of neural network basecalling tools for Oxford Nanopore sequencing" Genome biology 2019

Polishing using accurate reads



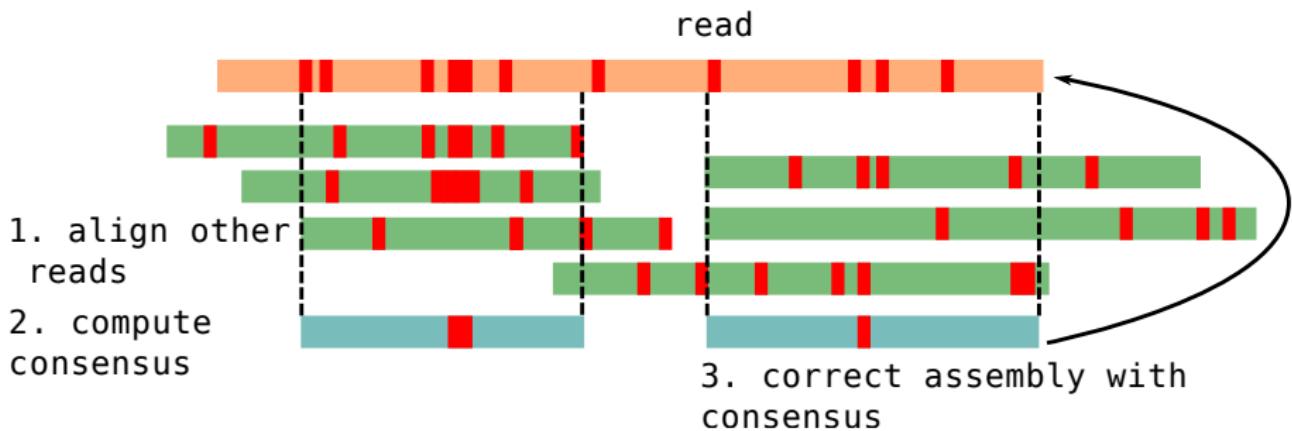
Systematics errors

Polishing with Illumina data can improve the final error rate



From ntEdit: scalable genome sequence polishing Bioinformatics 2019

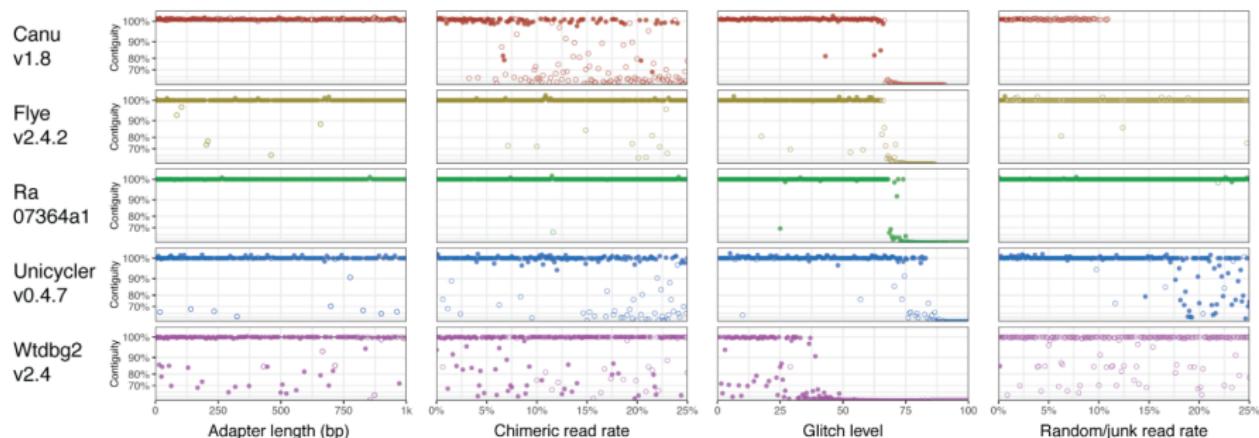
Correction before assembly



Long reads for assembly: assembly solved?

Assembly is not solved yet

Sometime the software fail



From github.com/rrwick/Long-read-assembler-comparison

Long reads for assembly: assembly solved?

Assembly is not solved yet

Sometimes the data cannot solve the problem

- Very large repeated region
- Low local coverage
- Chimeric/noisy reads

Long reads assemblers

Your toolkit for the practical

- Flye
- Miniasm
- Raven

Other notable assemblers

- Canu
- Mecat
- Redbean
- ...

Long read assembly summary

- Not resolved: correction before or after assembly (polishing)
- Overlap graphs with quick overlap computation
- Long reads can span repeats and improve assemblies
- Methods to polish contigs

Conclusions

Take home messages

- Short reads: De Bruijn graphs / Long reads: Overlap graph
- Repeats are the core issue
- Output fragments of genomes (**contigs**)
- Several steps and heuristics in practice

Challenges

- Difficult to reconstruct haplotypes
- Scaling on large genomes
- Robustness to noisy data

Topics we did not review today

- Read correction
- Multi k assembly
- Use of paired-end, mate-pair, 10X, Hi-C ...)

The end



PopGenGoogling
@popgengoogling

i trust you to figure out your own genome

[Traduire le Tweet](#)

3:01 AM · 14 déc. 2019 · [Twitter for iPhone](#)

37 Retweets **267** J'aime

Practical session

Datasets

Short Illumina reads 150bp

Long PacBio reads

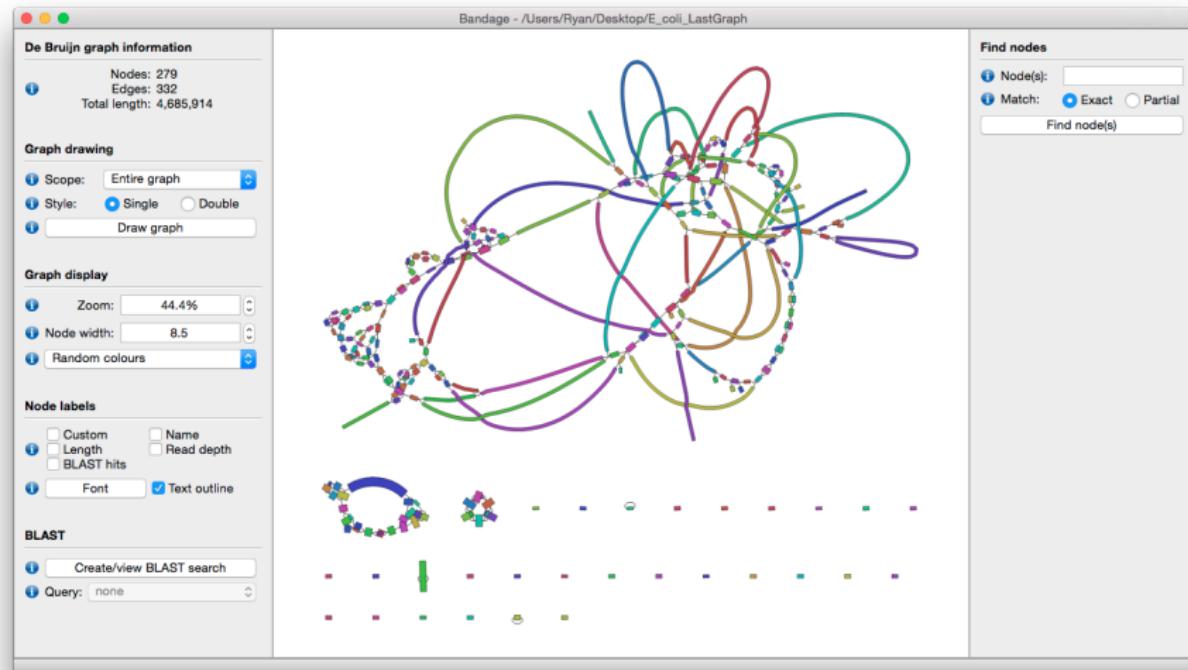
Steps

- Perform a draft assembly
- Evaluate it
- (Try to) Perform a better assembly

You will read the assembler website and papers to decide which one you want to use

Visualize assembly

Bandage tool can visualize assembly graphs (GFA)



From rwick.github.io/Bandage

Evaluate assembly

Contigs can be mapped and compared to a reference/closely related genome



From quast.bioinf.spbau.ru/manual.html

Assembly continuity

N50

N50 can be described as a weighted median statistic such that 50% of the entire assembly is contained in contigs or scaffolds equal to or larger than this value.

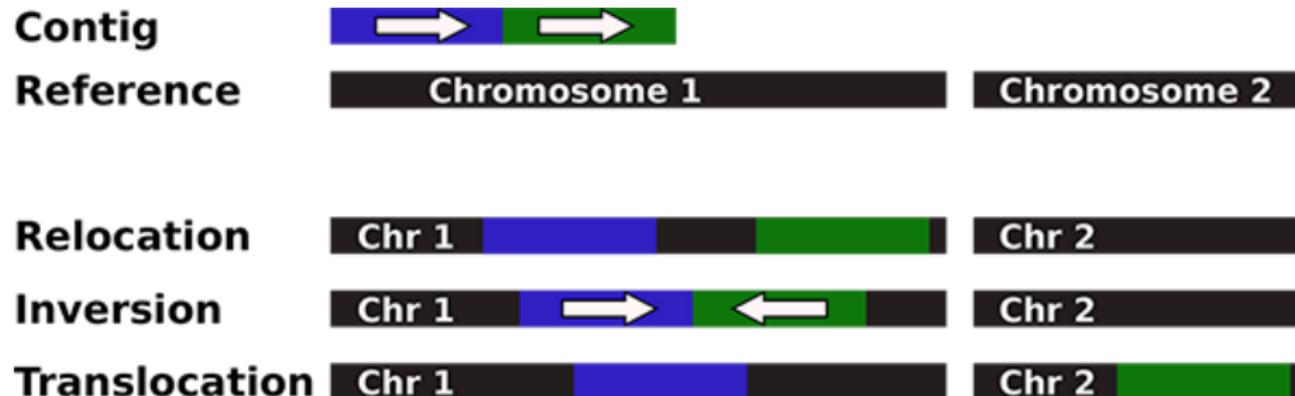
N75

N75 is the same statistic for 75% of the assembly

NGA50

Similar to the N50 but only takes into account contigs/scaffolds that can be aligned on the reference genome and consider 50% of the **genome size** instead of the assembly size

Misassemblies



Multiple k assembly

Most De Bruijn graph assemblers can now perform several assemblies with different k -mer sizes to produce an improved "super" assembly
(will be discussed in metagenomic session)

Advice

Using a single size of k -mer will allow the assembly to go way faster during the practical

Fig. 1. The workflow of MEGAHIT

