

Lesson 7

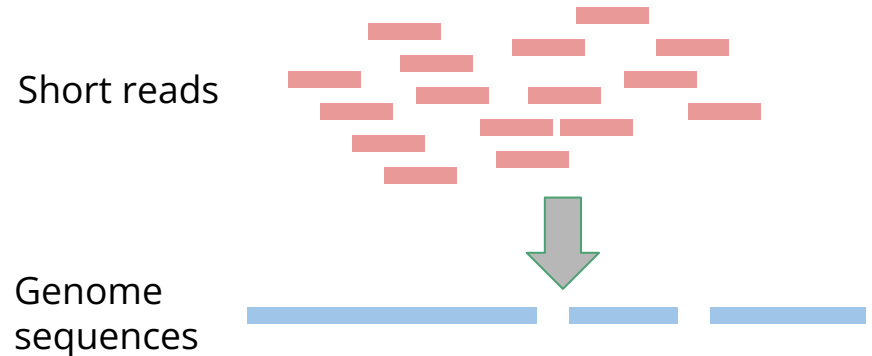
De novo genome assembly

By the end of this lesson you will...

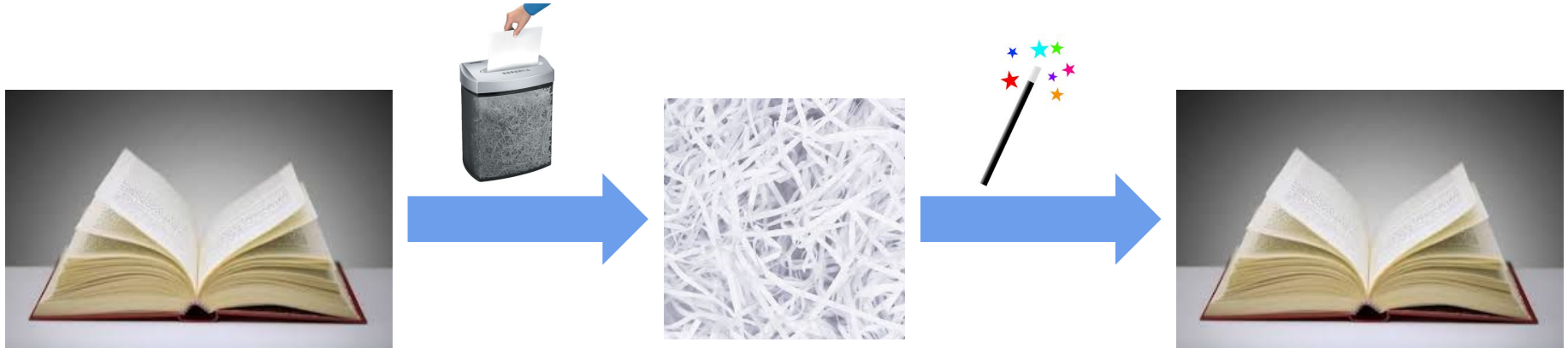
- Understand the basics of de novo genome assembly
- Be familiar with the DeBruijn graph method
- Know several methods and metrics for genome assembly QA
- Be able to perform assembly QA using QUAST and BUSCO

What is de novo genome assembly?

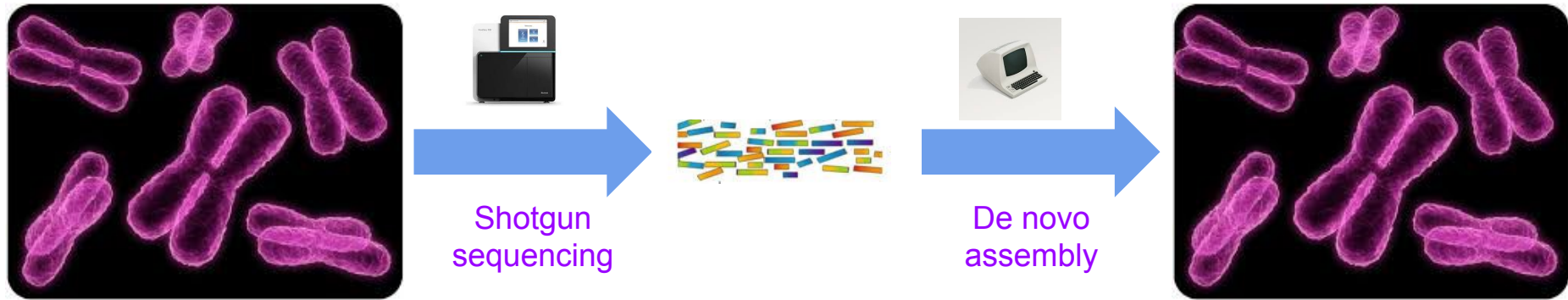
- Genome assembly - constructing long genomic sequences from shorter ones
- De novo = “from scratch”
- In NGS context - short reads → whole genome, without any external reference



The assembly problem

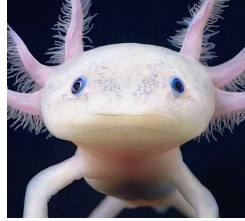


The assembly problem



Why do we need de novo assembly?

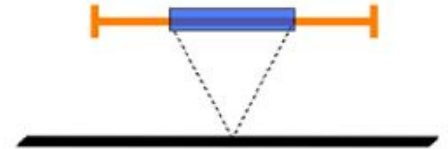
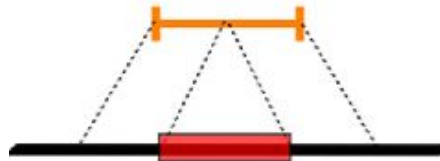
- Completely new organism



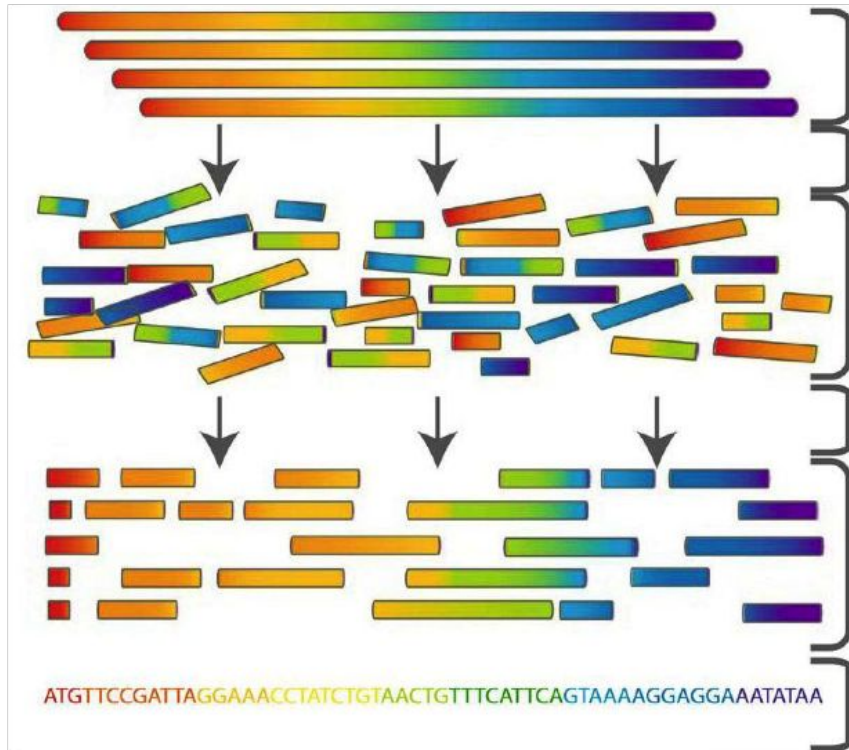
- Existing reference is too different from what we are interested in
- Identify large structural variation



- Detect novel sequences not present in the reference
- Cancer genomics



Genome assembly by reads overlap - challenges



Genomic DNA

Fragmentation + Sequencing

Sequence reads

Assembly

Connection between reads
found

Consensus sequence

Suffix Prefix Matching

TCTATATCTCGGCTCTAGG

TATCTCGACTCTAGGCC

Suffix Prefix Matching

TCTATATCTCGGCTCTAGG

TATCTCGACTCTAGGCC

Suffix Prefix Matching

TCTATATCTCGGCTCTAGG
GGCGTCTATATCTCGGCTCTAGGCCCTCATT
TATCTCGACTCTAGGCC

If a suffix of read A is similar to a prefix of read
then A and B might overlap in the genome

Suffix Prefix Differences

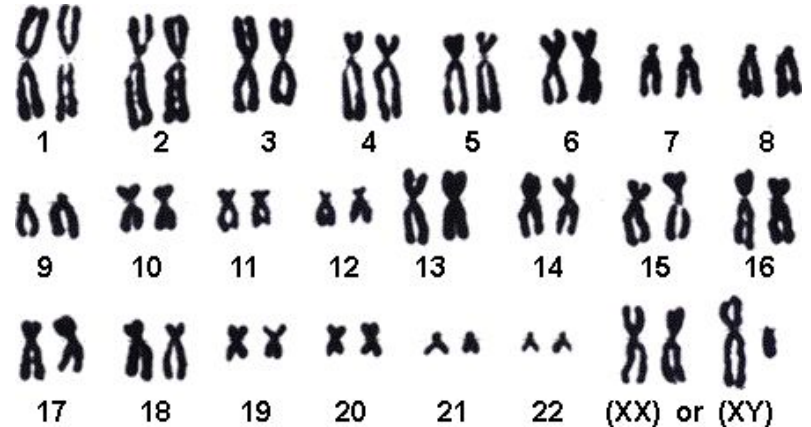
TCTATATCTCGGCTCTAGG

TATCTCGACTCTAGGCC



Why the differences?

1. Sequencing errors
2. Polyploidy



High and Low Coverage

CTAGGCCCTCAATTTT
GGCTCTAGGCCCTCATTTTT
CTCGGCTCTAGGCCCTCATTT
TATCTCGACTCTAGGCC
TCTATATCTCGGCTCTAGG
GGCGTCTATATCTCG
GGCGTCTATATCT
GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
CTAGGCCCTCAATTTT
TATCTCGACTCTAGGCCCTCA
GGCGTCTATATCT

More coverage

Less coverage

More coverage leads to more and longer overlaps

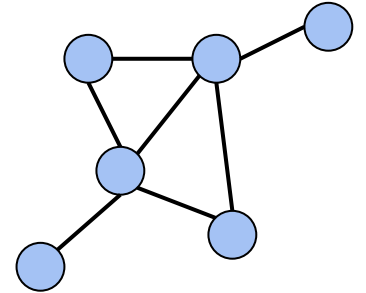
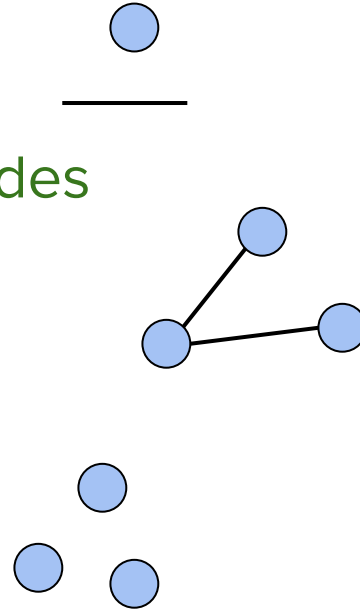
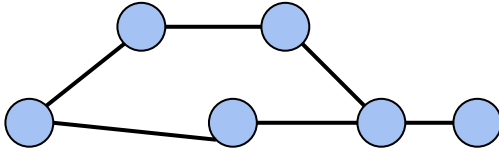
Overlap Consensus

What is the best representation to the set of sequences?

Graph

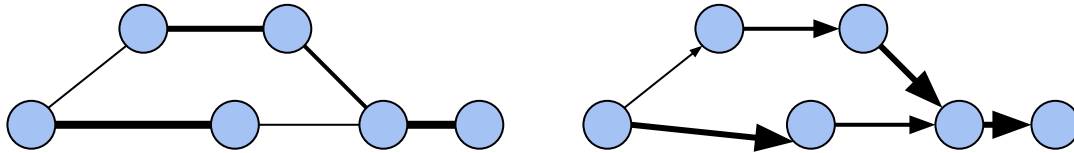
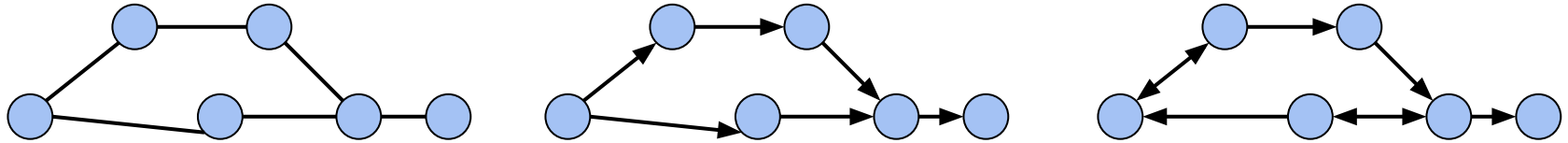
A graph is a set of:

- Nodes (vertices)
- Edges - connecting two nodes



Directed and Weighted Graphs

Graphs can be **directed** or **non-directed**



We can assign **weights** to edges

Overlap Consensus Graph

Nodes: all 6-mers in **GTACGTACGAT**

Edges: overlaps of length > 3

GTACGT

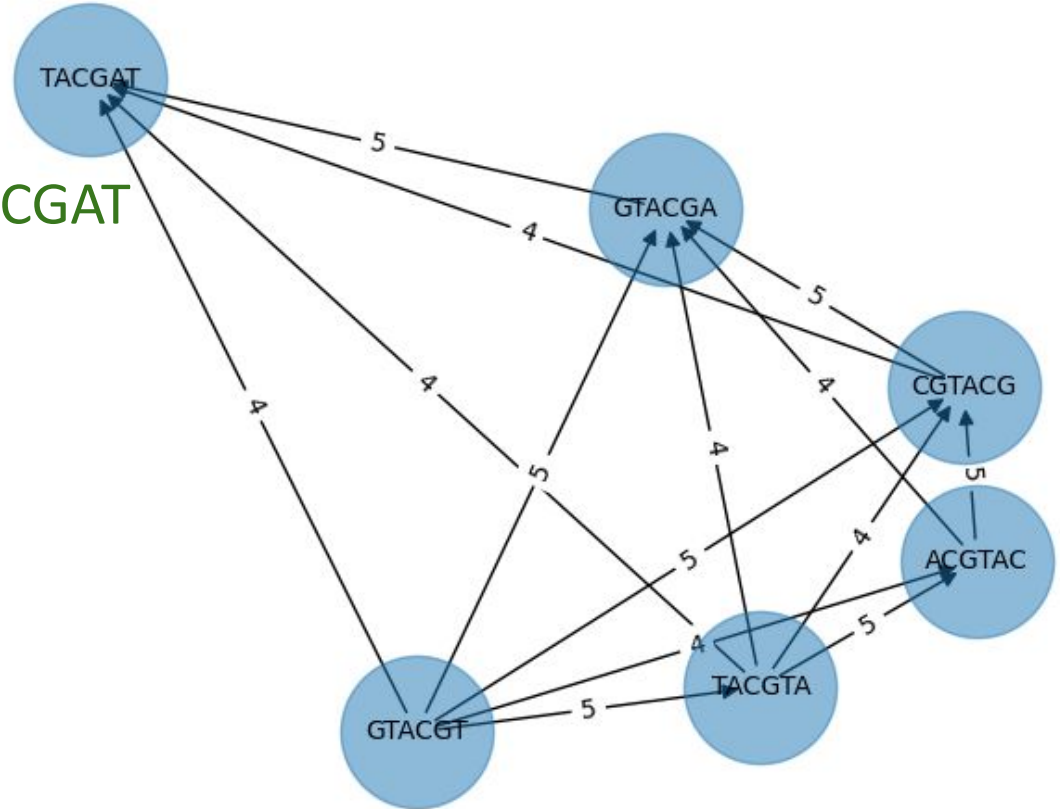
TACGTA

ACGTAC

CGTACG

GTACGA

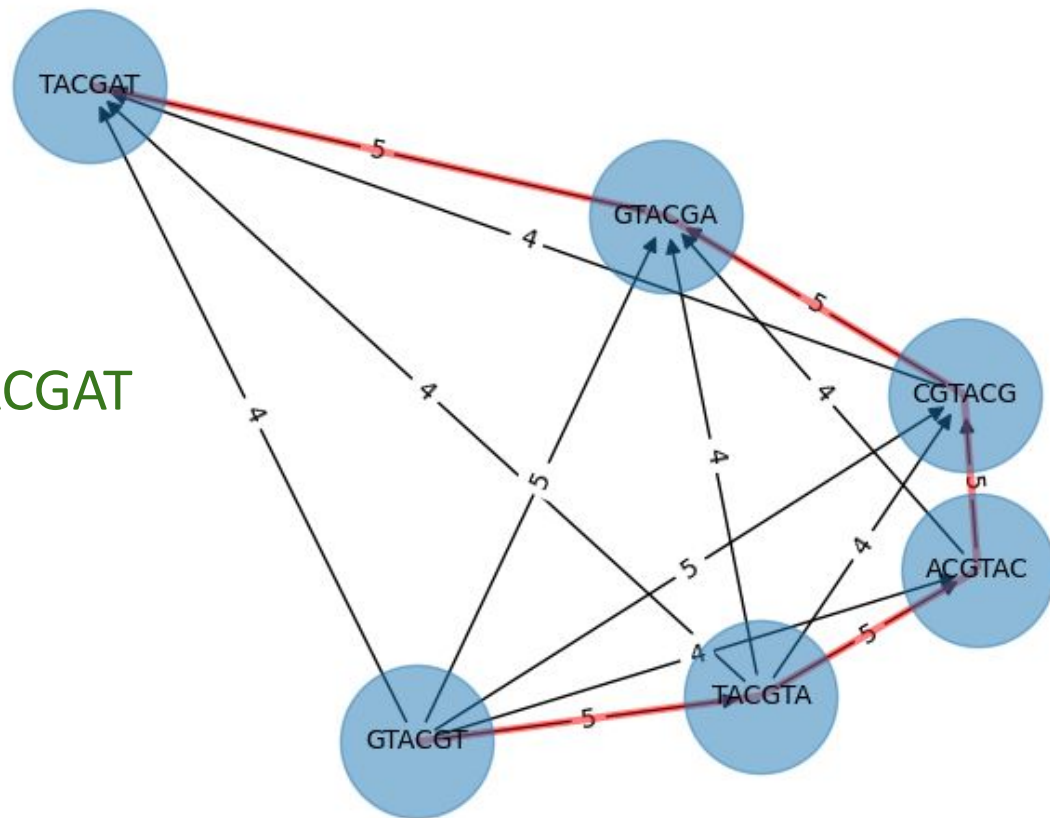
TACGAT



Overlap Consensus Graph

Nodes: all 6-mers in **GTACGTACGAT**

Edges: overlaps of length > 3



Shortest Common Superstring

The Shortest Common Superstring problem (SCS) aim to find the shortest possible string that contains every string in a given set as substrings

Example: BAA AAB BBA ABA ABB BBB AAA BAB

Concatenation: BAA AAB BBA ABA ABB BBB AAA BAB

AAA

AAB

ABB

BBB

BBA

BAB

ABA

BAA

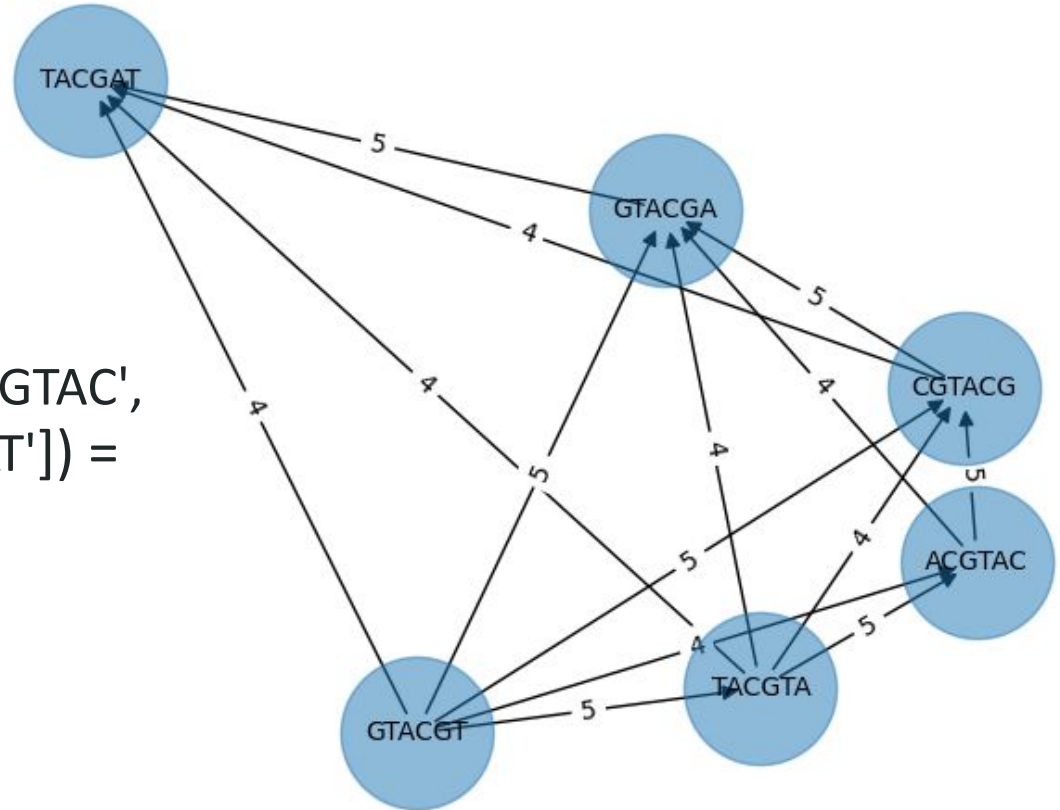
24

SCS: AAABBBBABAA

10

Shortest Common Superstring

SCS(['GTACGT', 'TACGTA', 'ACGTAC',
'CGTACG', 'GTACGA', 'TACGAT']) =
GTACGTACGAT



Shortest Common Superstring

Brute Force

Order 1: AAA AAB ABA ABB BAA BAB BBA BBB
AAABABBAABABBABBB Superstring 1

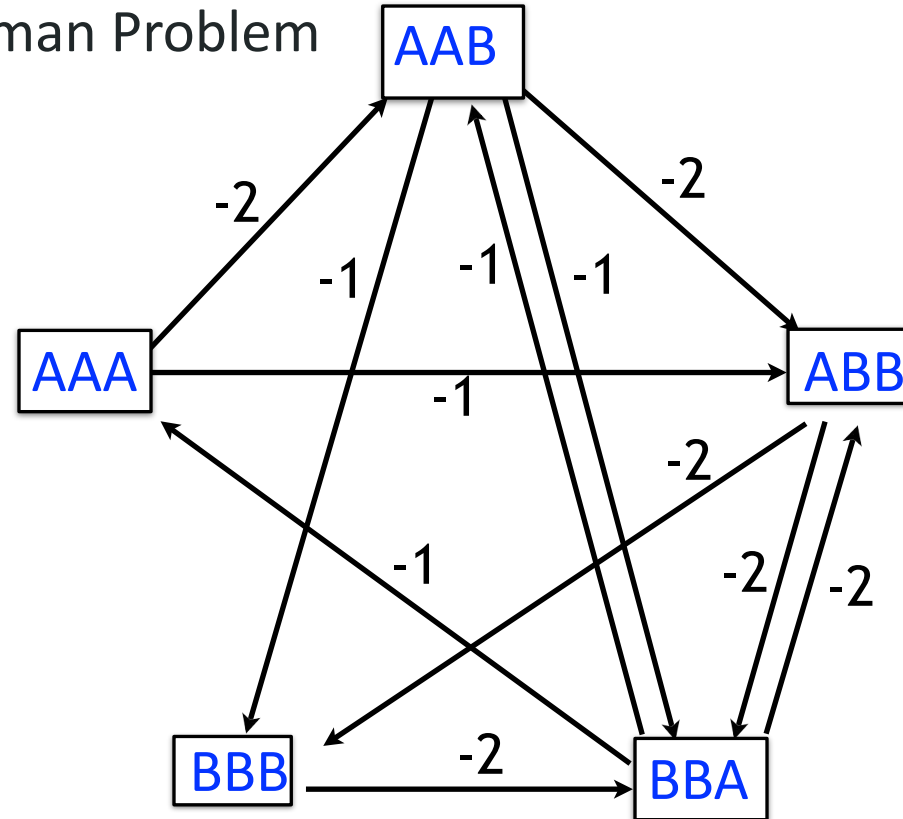
Order 2: AAA AAB ABA BAB ABB BBB BAA BBA
AAABABBBAABBA Superstring 2

$O(n!)$

Shortest Common Superstring

Traveling Salesman Problem

Modified overlap graph where each edge has cost = - (length of overlap)
SCS corresponds to a path that visits every node once, minimizing total cost along path



Shortest Common Superstring

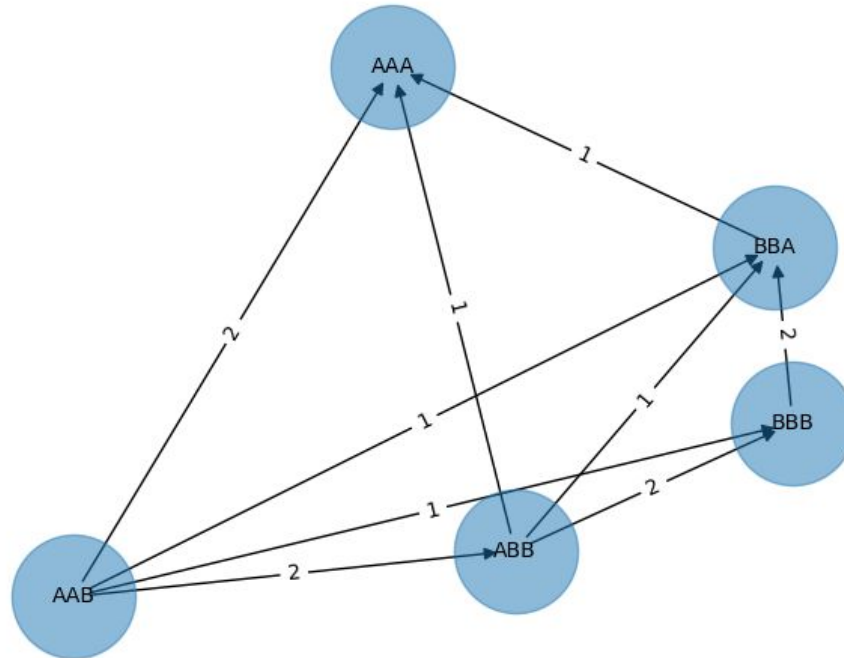
NP-complete

No efficient solution algorithm has been found

Shortest Common Superstring

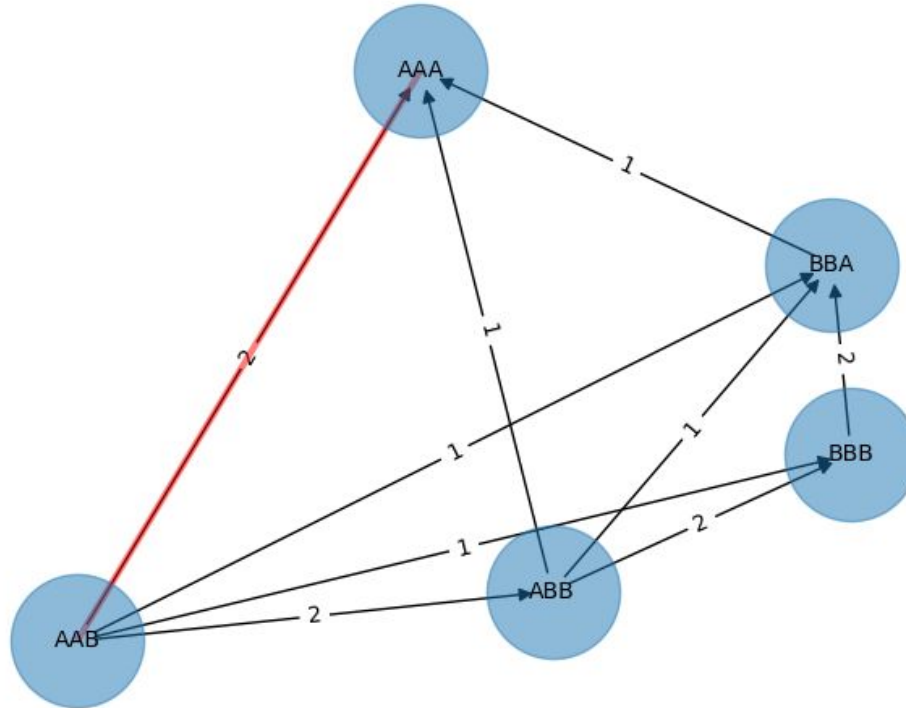
Greedy

Example: BAA AAB BBA ABA ABB BBB AAA BAB



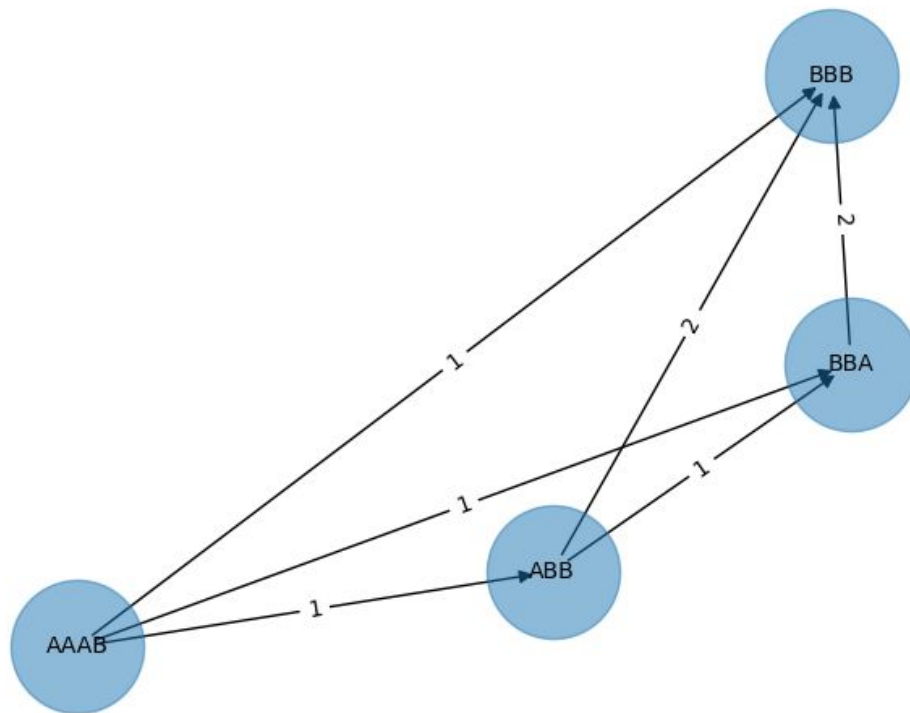
Shortest Common Superstring

Greedy



Shortest Common Superstring

Greedy



Shortest Common Superstring

Greedy



Superstring, length 7

Alternative Shuffling



Superstring, length 9

Greedy answer isn't necessarily optimal

Shortest Common Superstring

Greedy

Greedy algorithm is not guaranteed to choose overlaps yielding SCS

But greedy algorithm is a good approximation; i.e. the superstring yielded by the greedy algorithm won't be more than ~ 2.5 times longer than true SCS

Gusfield, Dan. "Algorithms on Strings, Trees, and Sequences - Computer Science and Computational Biology." (1997).

Shortest Common Superstring

Greedy

a_long_long_long_time !=6

ng_lon _long_ a_long long_l ong_ti ong_lo long_t g_long g_time ng_tim

5 ng_time ng_lon long_ a_long long_l ong_ti ong_lo long_t g_long

5 ng_time g_long_ ng_lon a_long long_l ong_ti ong_lo long_t

5 ng_time long_ti g_long_ ng_lon a_long long_l ong_lo

5 ng_time ong_lon long_ti g_long_ a_long long_l

5 ong_lon long_time g_long_ a_long long_l

5 long_lon long_time g_long_ a_long

5 long_lon g_long_time a_long

5 long_long_time a_long

4 a_long_long_time

a_long_long_time

Missing a _long

Shortest Common Superstring

Repeats often foil assembly. They certainly foil SCS, with its “shortest” criterion!

Reads might be too short to “resolve” repetitive sequences. This is why sequencing vendors try to increase read length.

Algorithms that don’t pay attention to repeats (like our greedy SCS algorithm) might *collapse* them

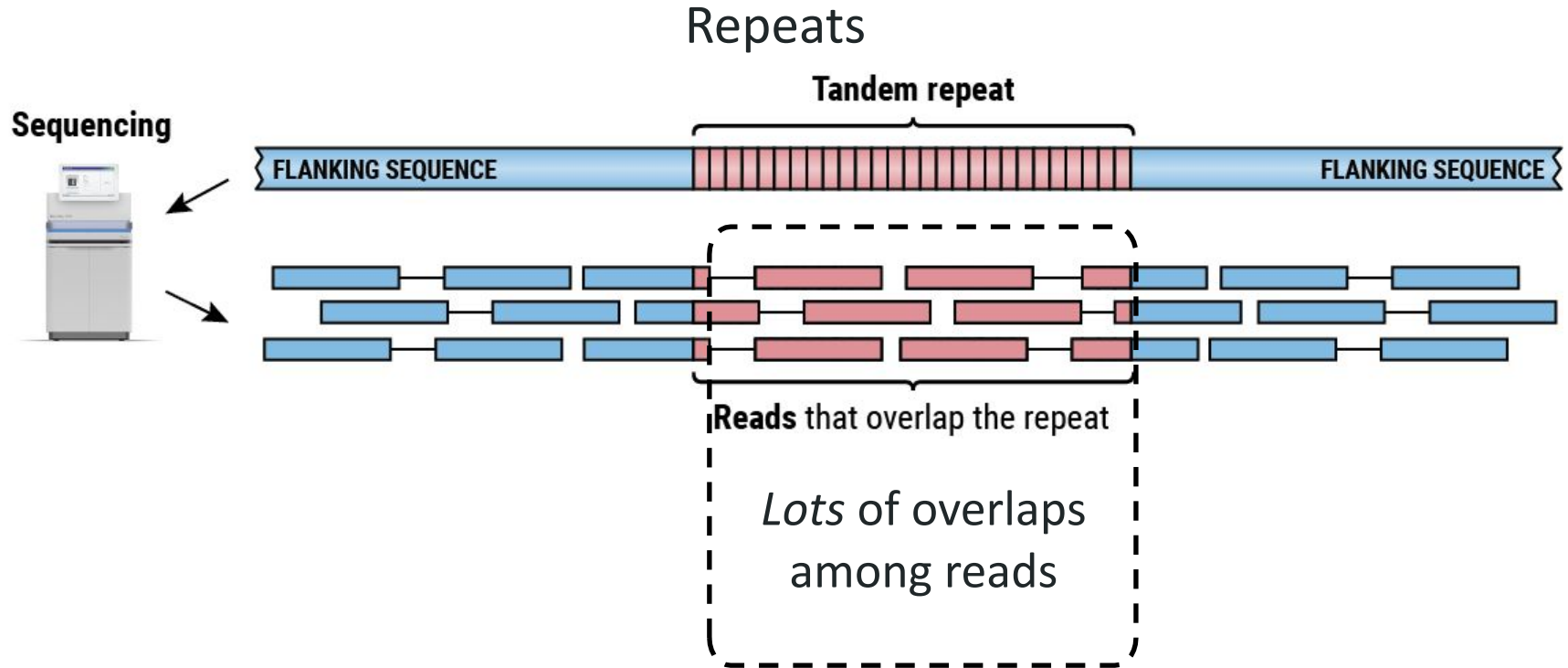
a_long_long_long_time

collapse

a_long_long_time

The human genome is ~ 50% repetitive!

Shortest Common Superstring

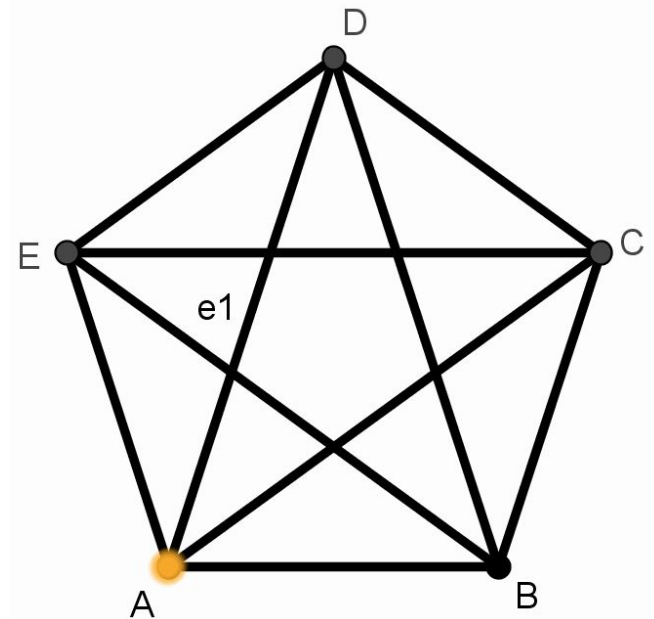


De Bruijn Graph assembly

- Currently the most popular assembly method
- Many software tools use variants of the method
 - SPAdes
 - SOAPdenovo
 - AbySS
 - MEGAHIT
- Based on graph theory
- Specifically **k-mer graphs**

The Eulerian path

- A path in a graph that visits every **edge** exactly once
 - Must visit **all** edges
 - Can't visit an edge twice
 - Can visit a node more than once
- A common problem in graph theory



K-mers

AGATCCAGCGAGGTCGCTATCCGTTAATTG

5-mers

AGATC

GATCC

ATCCA

...

AATTG

K-mers

AGATCCAGCGAGGTCGCTATCCGTTAATTG

5-mers

AGATC
GATCC
ATCCA
...
AATTG

7-mers

AGATCCA
GATCCAG
ATCCAGC
...
TTAATTG

How many 21-mers
are in a 100 bp
read?

De Bruijn Graph assembly - the basic algorithm

Genome (G=30): **A**GA**T**CC**A**GC**G**AG**G**TC**G**CT**A**TC**C**GT**T****A**ATT**G**



Reads (L=10): **A**GA**T**CC**A**GC**G**

AGC**G**AG**G**TC**G**

GCT**A**TC**C**GT**T**

CCG**T**T**A**ATT**G**

Break into k-mers (k=4):

AGA**T**CC**A**GC**G** → **A**GA**T** **G**AT**C** **A**TC**C** **T**CCA**C**CAG **C**AG**C** **A**GC**G**

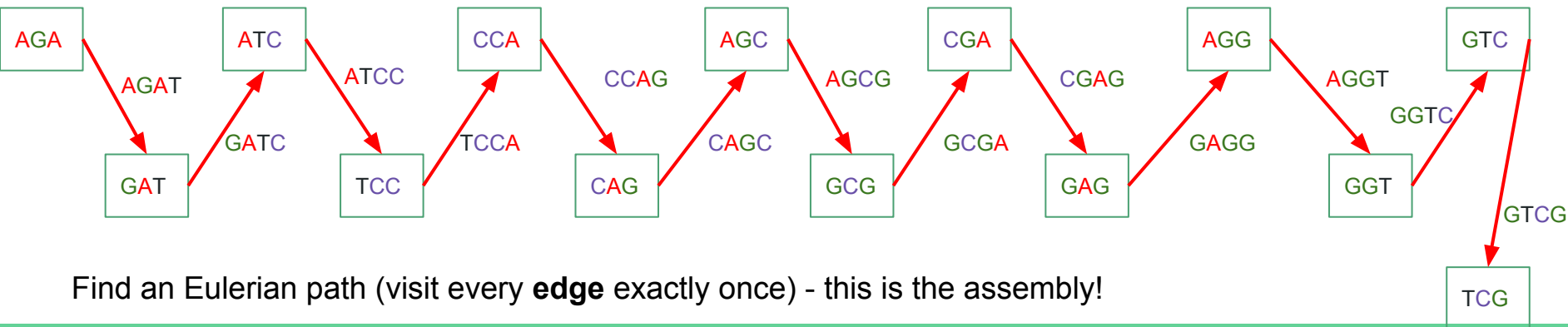
AGC**G**AG**G**TC**G** → **A**GC**G** **G**CG**A** **C**GAG **G**AG**G** **A**GG**T** **G**GT**C** **G**TC**G**

...

Create graph nodes - each **unique prefix and suffix** of length $k-1$ of k -mers



Add directed edge between node x and node y if k -mer exists with prefix x and suffix y



Find an Eulerian path (visit every **edge** exactly once) - this is the assembly!

De Bruijn Graph

A procedure for making a De Bruijn graph for a genome

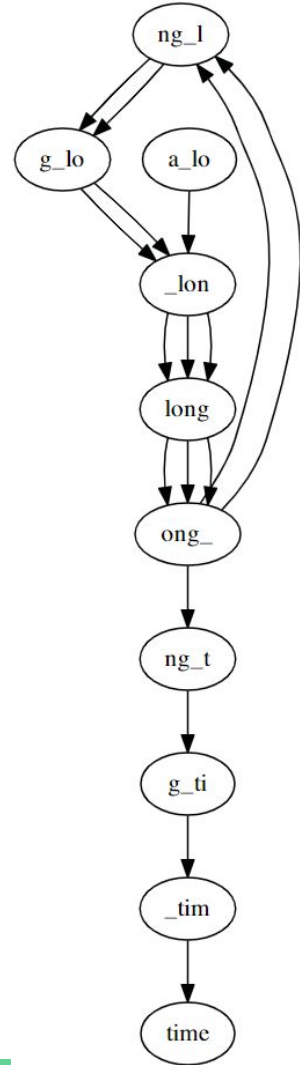
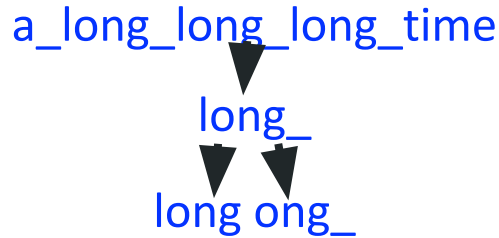
Assume perfect sequencing where each length-k substring is sequenced exactly once with no errors

Pick a substring length k: 5

Start with each read:

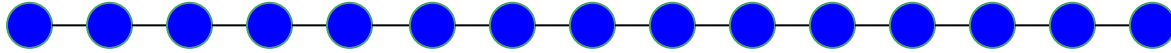
Take each k mer and split into left and right k-1 mers

Add k-1 mers as nodes to De Bruijn graph (if not already there), add edge from left k-1 mer to right k-1 mer

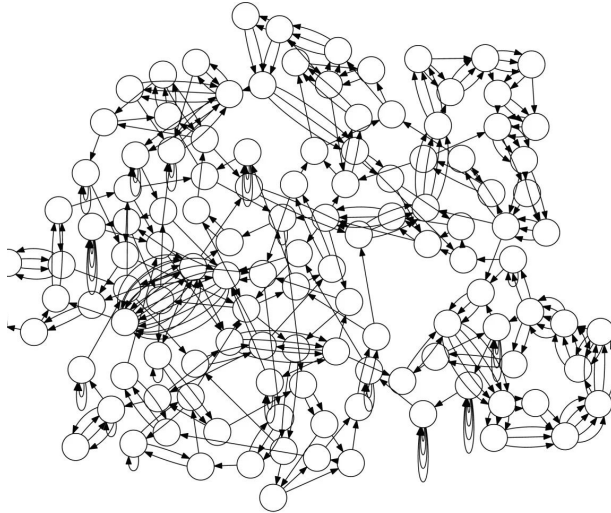


If only life was that simple...

Ideally, we want our graph to look like this:

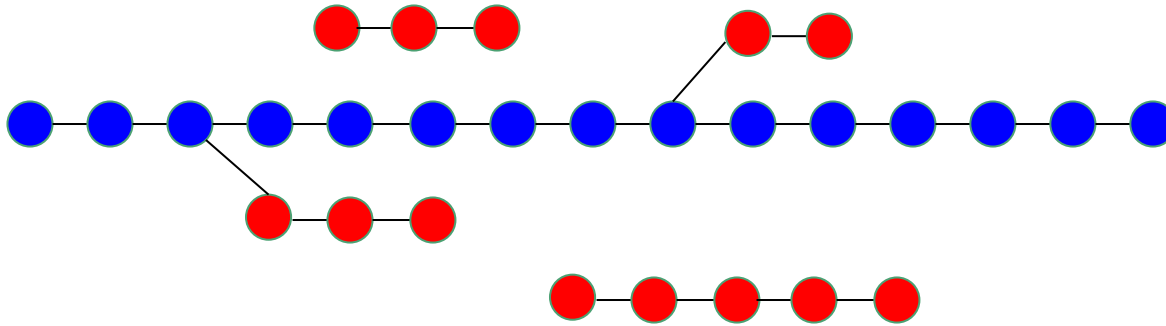


But in practice:



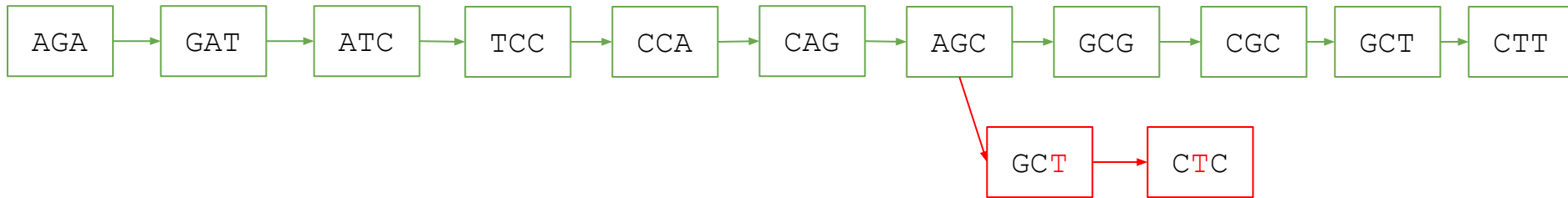
Sequencing errors

- Side branches
- Disconnected bits



Sequencing errors

AGATCCAGCG → AGAT GATC ATCC TCCA CCAG CAGC AGCG
GATCCAGC**T**C → GATC ATCC TCCA CCAG CAGC AGC**T** GC**T**C
TCCAGCGCTT → TCCA CCAG CAGC AGCG GCGC CGCT GCTT



Genomic repeats

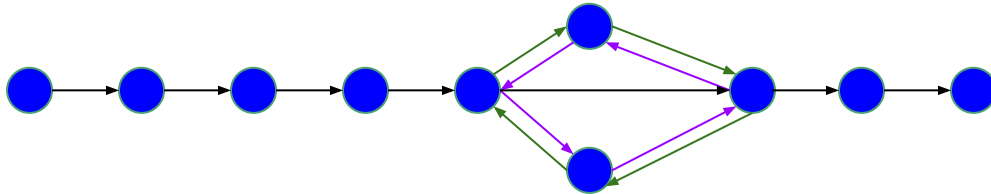
- Tandem duplication



- Interspersed duplications



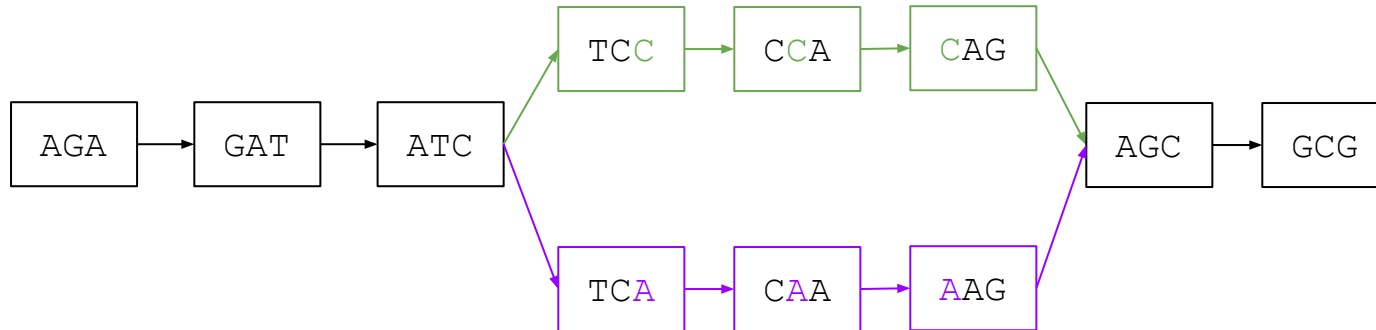
- Might create ambiguity - multiple possible eulerian paths



Heterozygosity

- Creates “bubbles” in the graph

Maternal AGATC**C**AGCG → AGAT GATC AT**C** TC**C**A CCAG **C**AGC AGCG
Paternal AGATC**A**AGCG → AGAT GATC AT**C**A TC**A**A CA**A**G **A**AGC AGCG

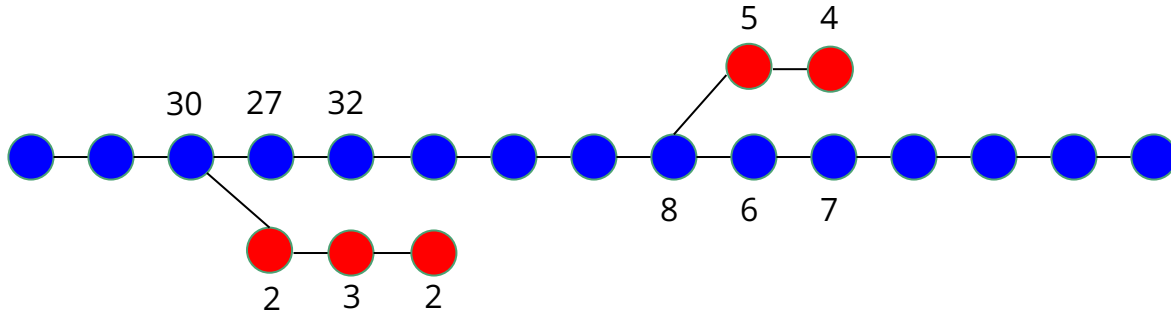


Uneven sequencing depth

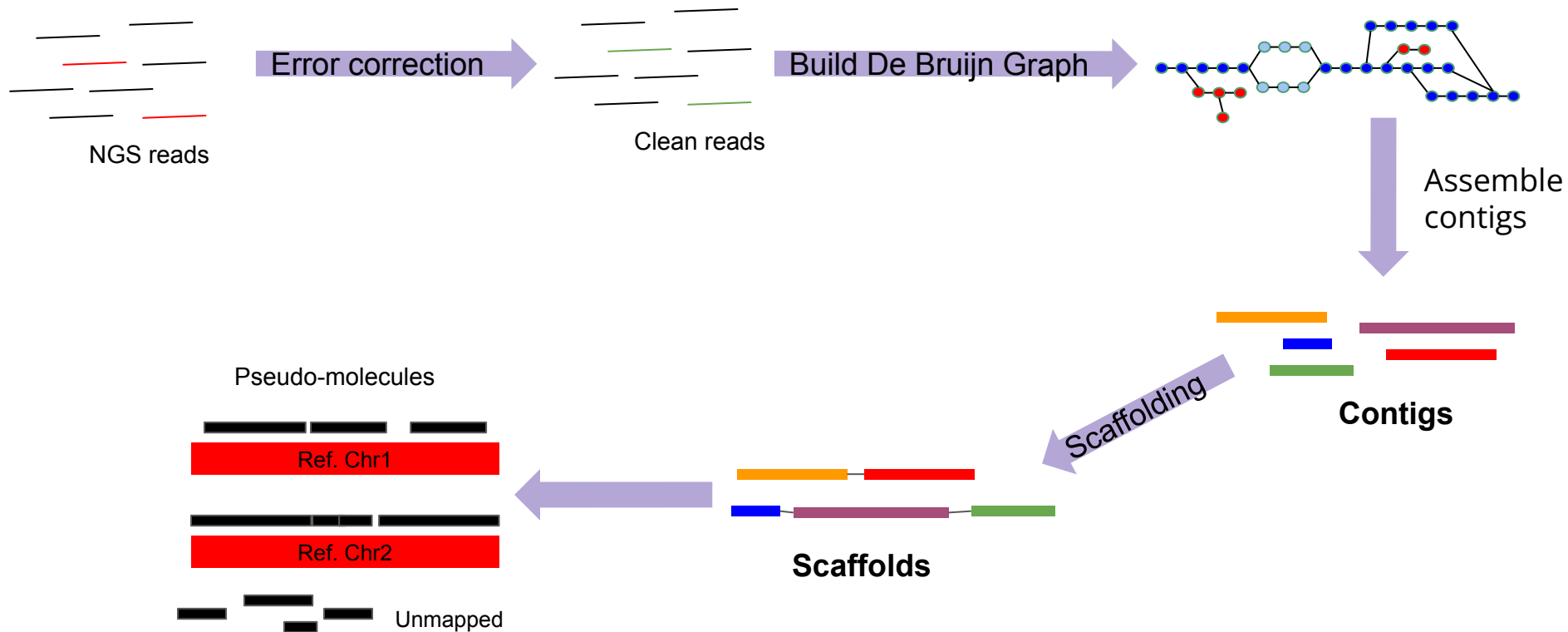
Very low depth (e.g. low complexity regions) can fragment the graph:



Uneven depth can make it hard to determine which branches are true and which are noise



General assembly workflow

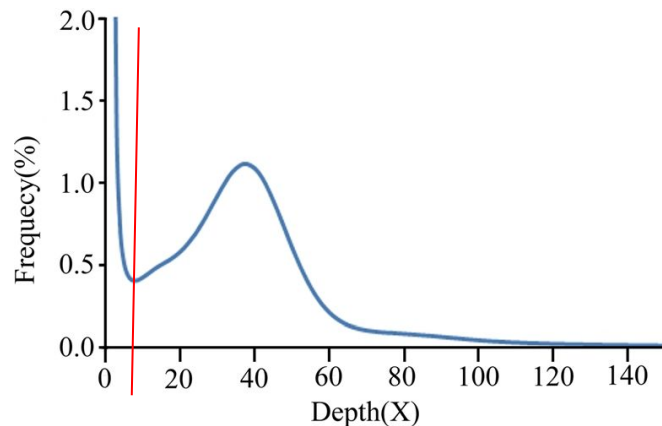


Error correction/filtration

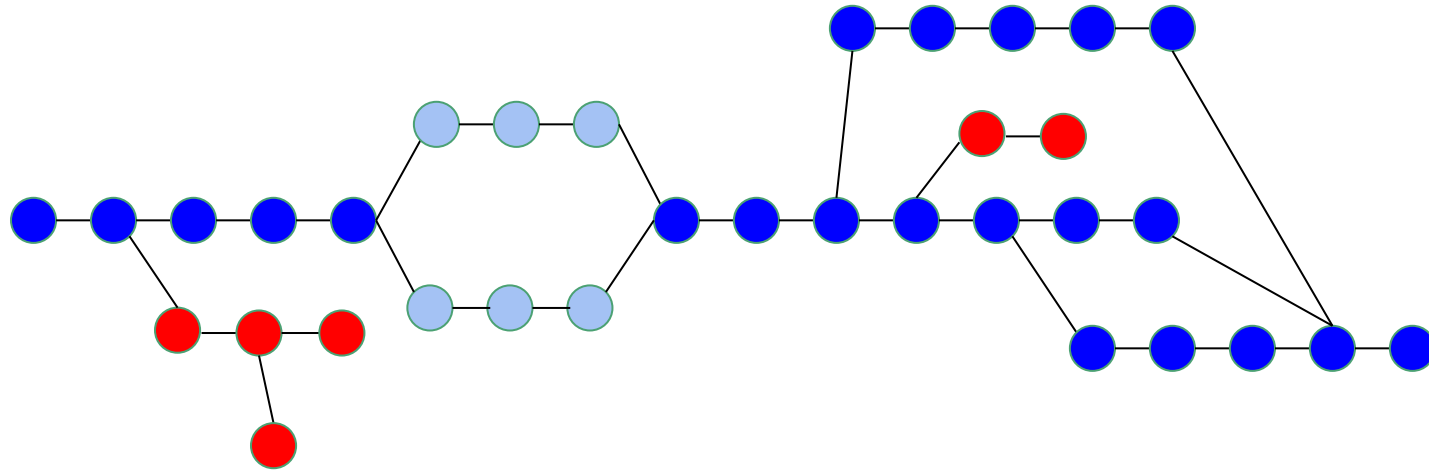
1. Extract all k-mers from all reads
2. Count how many times each k-mer was observed
3. Label rare k-mers as error k-mers
4. Find reads from which error k-mers came
5. Discard or correct the reads with error k-mers

Count(GGATAGGCACCAGTTAT) = 30

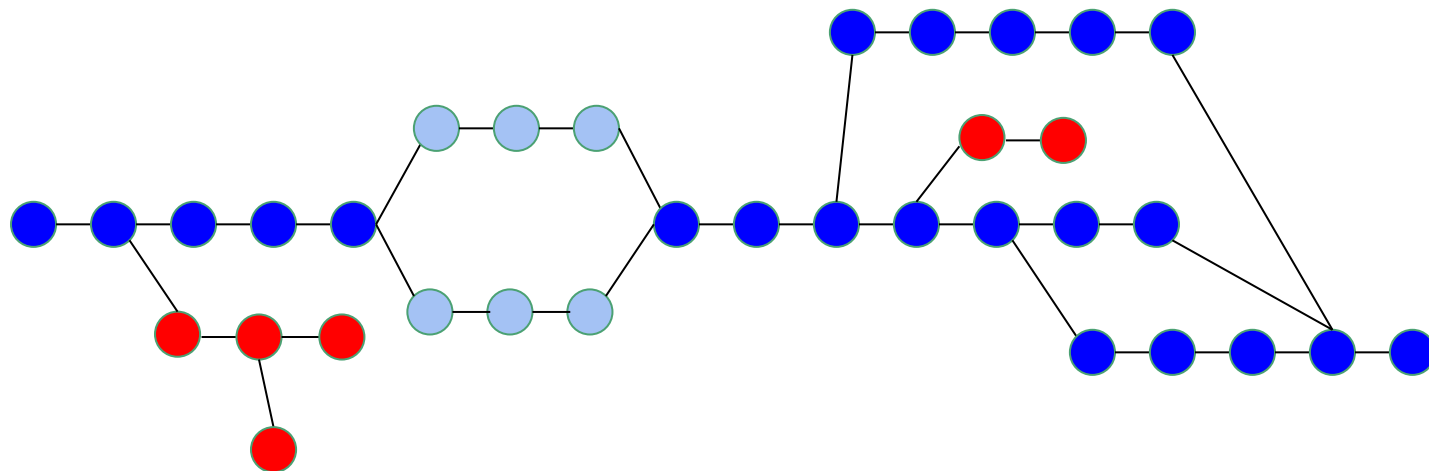
Count(GGATAGG**T**ACCAGTTAT) = 1



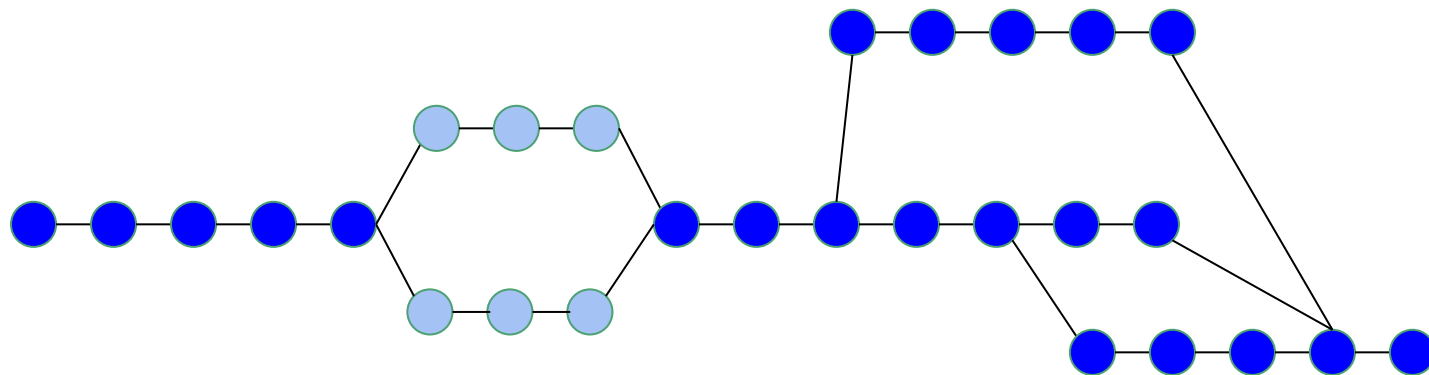
Build De Bruijn graph



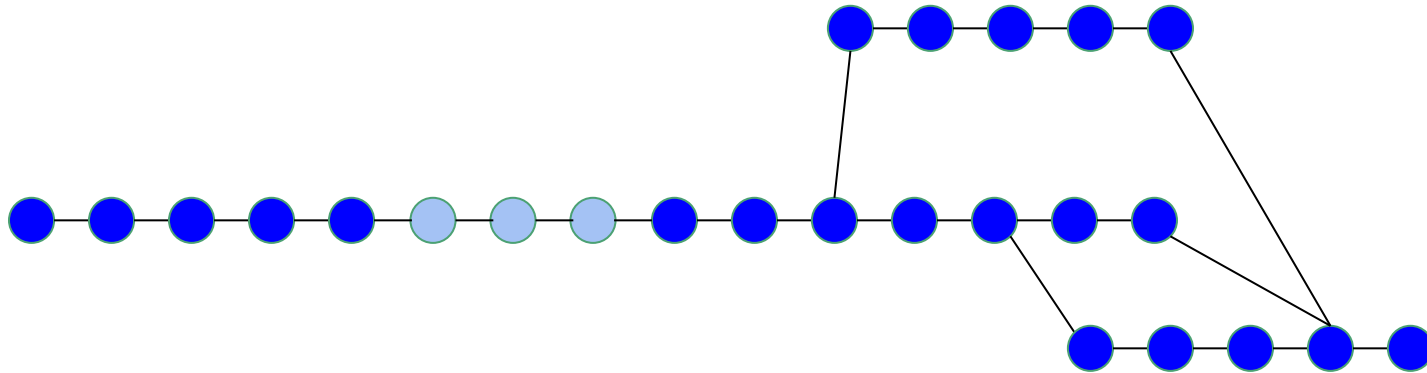
Prune error branches



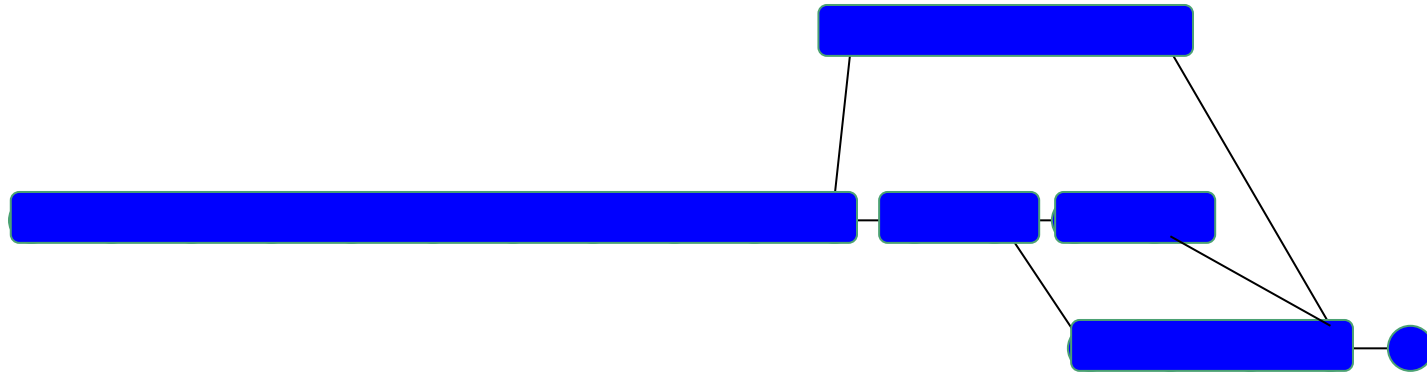
Resolve bubbles



Resolve bubbles



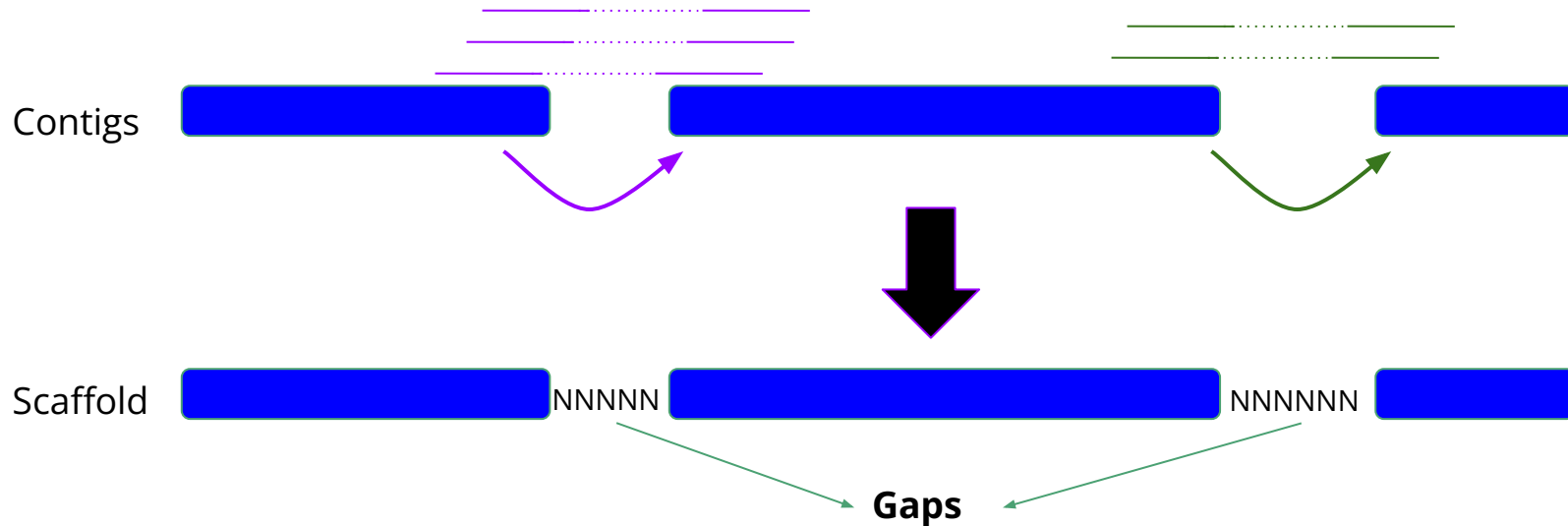
Create contigs - break on branching points



Scaffolding

Mate Pair Sequencing

- Use paired end information
- Look for evidence of two contigs linked together



The SPAdes assembler

- A popular assembler for small-medium size genomes
- Can use multiple types of data
 - Short reads - paired/single end
 - Mate pairs
 - Long reads
- Includes several modules:
 - Error correction
 - Contigs assembly
 - Scaffolding
- Does not require k-mer choice

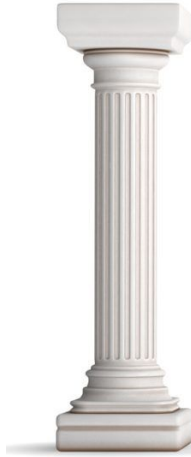


Beyond scaffolds?

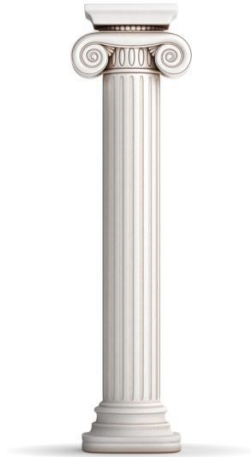
- Scaffolds are the best we can do *de novo*
- We can always produce more data - won't always help
- To get to pseudomolecules we need external data
 - A similar reference genome that we can map to
 - Hi-C data
 - Long reads



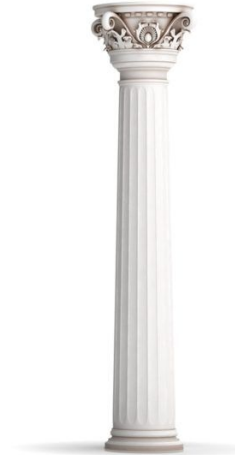
The “three pillars” of assembly quality



Completeness



Contiguity



Accuracy

A good assembly is...

- **Complete** - contains all or most of the genome
- **Contiguous** - built of large scaffolds
- **Accurate** - includes few mis-assemblies

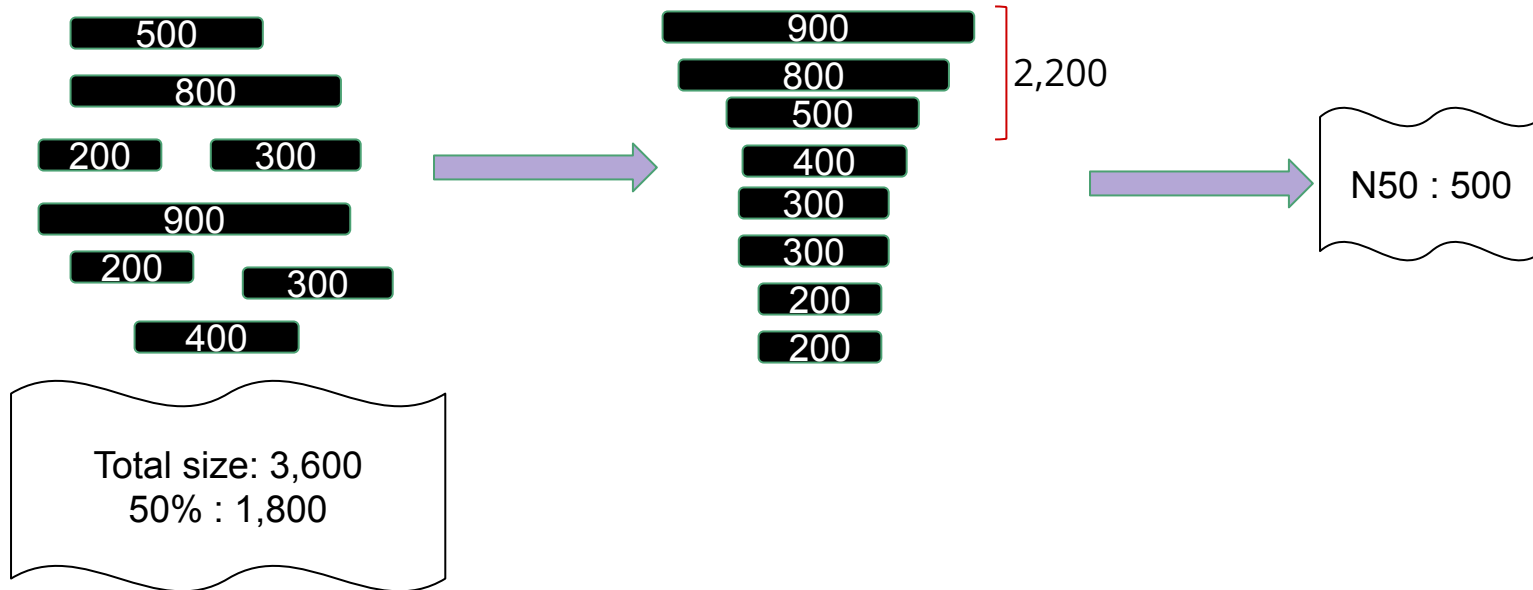
How can we assess these for a given assembly?

Assembly statistics

- **Completeness**
 - Assembly size - compared to expected genome size
 - % of gaps in assembly ('N' bases)
- **Contiguity**
 - Number of contigs/scaffolds
 - Mean/median scaffold/contig size
- **Always consider all stats together**

Completeness - N50

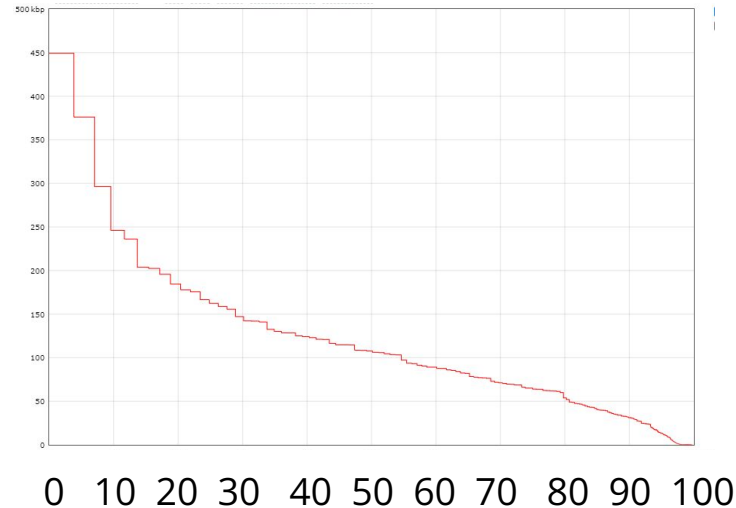
Length of the shortest contig for which the total size of all longer contigs is 50% of the assembly size.



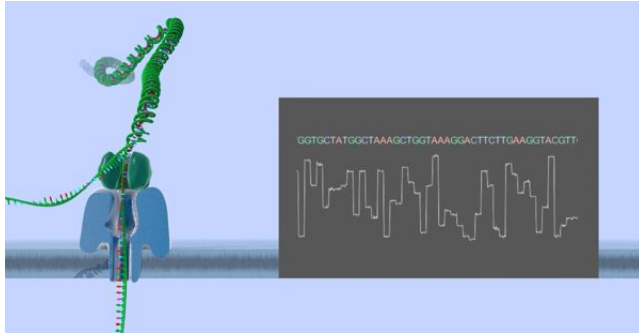
N50 et al.

- N90 - contig length required to cover 90% of assembly size
- L50 - Number of contigs longer than N50
- NG50 - use expected genome size instead of assembly size

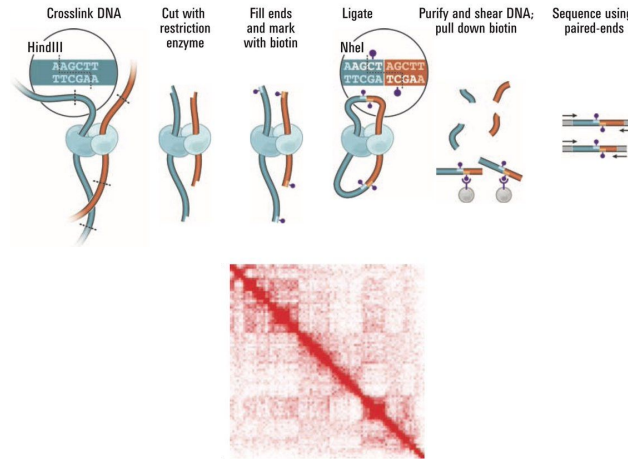
Nx plot



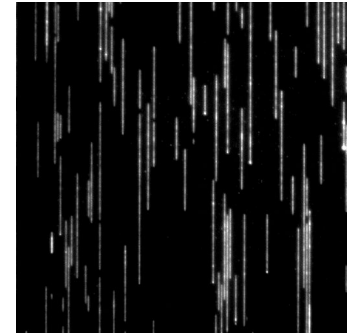
Emerging technologies for genome assemblies



Long reads



Hi-C



Optical mapping